



CZECH TECHNICAL UNIVERSITY IN PRAGUE

Faculty of Civil Engineering
Department of Mechanics

FFT-based method for homogenization of periodic media: Theory and applications

DOCTORAL THESIS

Mgr. Ing. Jaroslav Vondřejc

Ph.D. Programme: Civil Engineering
Branch of study: Mathematics in Civil Engineering

Supervisor: doc. Ing. Jan Zeman, Ph.D.

Prague, 2013



CZECH TECHNICAL UNIVERSITY IN PRAGUE
Faculty of Civil Engineering
Thákurova 7, 166 29 Prague 6

DECLARATION BY CANDIDATE

Author: Jaroslav Vondřejc

Title: FFT-based method for homogenization of periodic media: Theory and applications

I hereby declare that this thesis is my own work and effort and that it has not been submitted anywhere for any award. Where other sources of information have been used, they have been acknowledged.

In Prague, 2nd January 2013

Signature:

Abstract

This dissertation is devoted to an FFT-based homogenization scheme, a numerical method for the evaluation of the effective (homogenized) matrix of periodic linear heterogeneous materials. A problem is explained and demonstrated on a scalar problem modeling electric conduction, heat conduction, or diffusion.

Originally, FFT-based homogenization, that was proposed by Moulinec and Suquet in [32], is a numerical algorithm derived from Lippmann-Schwinger equation. Its equivalence to a corresponding weak formulation is shown; it eliminates a reference homogeneous material, a parameter of Lippmann-Schwinger equation. Next, Galerkin approximation with numerical integration is introduced to produce Moulinec-Suquet algorithm; trigonometric polynomials are taken as the trial space [47]. Convergence of approximate solutions to the solution of weak formulation is provided using a standard finite element approach together with approximation properties of trigonometric polynomials stated in [43].

Then, the solution of assembled non-symmetric linear system by Conjugate gradients, proposed by Zeman et al. in [57], is clarified.

Next, we study arbitrary accurate guaranteed bounds of homogenized matrix introduced by Dvořák in [12, 13] for a scalar problem and later independently by Wieǳowski in [55] for linear elasticity. This approach is also applicable for FFT-based homogenization. A general technique is proposed to allow for efficient calculation by FFT algorithm and to maintain the upper-lower bound structure. Dual formulation is employed to obtain lower bounds — for odd number of discretization points, the solution of dual formulation can be avoided. A general number of discretization points leads to a more complicated theory in both discretization and numerical treatment.

Finally, applications of FFT-based homogenization to real-world problems are demonstrated. The method is used to calculate homogenized matrix for cement paste, gypsum and aluminum alloy with local data obtained from nanoindentation. Next, it is employed as a part of two-step homogenization for a highly porous aluminium foam.

Keywords: homogenization, Fourier transform, FFT, discretization, finite element method, convergence, guaranteed bounds

Abstrakt

Tato práce je věnována homogenizační metodě založené na FFT, metodě počítající efektivní (homogenizovanou) matici periodického heterogenního materiálu. Problematika je vysvětlena a demonstrována na skalárním problému modelujícím elektrickou vodivost, tepelnou vodivost nebo difúzi.

Původně byla FFT homogenizace založena na řešení Lippmannovy-Schwingerovy rovnice a z něho navrženého numerického algoritmu podle Moulineca a Suqueta v [32].

Ukazujeme, že Lippmannova-Schwingerova rovnice je ekvivalentní příslušné slabé formulaci; to umožňuje eliminovat referenční homogenní materiálovou hodnotu, parametr rovnice. Dále jsme navrhli Galerkinovskou aproximaci s numerickou integrací tak, aby produkovala Moulinecův-Suquetův algoritmus; trigonometrické polynomy jsou použity jako konečně dimenzionální prostor [47]. Dokazujeme konvergenci přibližných řešení k řešení slabé formulace pomocí klasického přístupu z konečných prvků s využitím odhadů podle [43].

Dále je vysvětleno řešení vzniklé nesymetrické soustavy lineárních rovnic pomocí sdružených gradientů, které bylo navrženo Zemanem et al. v [57].

Věnujeme se libovolně přesným zaručeným mezím efektivní matice, jak byly navrženy Dvořákem v [12, 13] pro skalární problém a nezávisle Wieçkowskim v [55] pro lineární elasticitu. Tato metoda je využitelná i pro homogenizaci založenou na FFT. Navrhujeme obecnou metodu pro výpočet mezí tak, aby bylo možno využít efektivního algoritmu FFT a aby byla zachována struktura dolních a horních mezí. Duální formulace problému je využita pro dolní mez — ukazujeme, že pro lichý počet diskretizačních bodů se lze vyhnout řešení duální úlohy.

Nakonec je ukázána aplikace FFT homogenizace na lineární elasticitu. Metoda je využita k počítání efektivní matice pro cementovou pastu, sádku a hliníkovou slitinu s lokálními daty získanými z nanoindentace. Metoda je dále využita v dvoustupňové homogenizaci pro výpočet efektivních materiálových parametrů vysoce porézní hliníkové pěny.

Klíčová slova: homogenizace, Fourierova transformace, FFT, diskretizace, metoda konečných prvků, konvergence, zaručené meze

Acknowledgments

First of all, I would like to express my deepest gratitude and thanks to my supervisor doc. Ing. Jan Zeman, Ph.D for his ingenuous support and endless patience during the whole study and for the freedom that I have obtained in research directions. Special gratitude is to prof. RNDr. Ivo Marek, DrSc. for his support and critical comments in the fields of mathematical and numerical analysis.

I thank to co-authors of papers attached to the work, namely to already mentioned supervisors, Ing. Vlastimil Králík, doc. Ing. Jiří Němeček, Ph.D., and Ing. Jan Novák, Ph.D.

I also thank to prof. Ing. Milan Jirásek, DrSc. besides other things for his course Modeling of Localized Inelastic Deformation that I have held at the beginning of my Ph.D. studies. I also appreciate the courses, that I have taken at the Faculty of Mathematics and Physics at Charles University and I thank to their organizers: prof. RNDr. Vít Dolejší, Ph.D., DSc.; prof. RNDr. Miloslav Feistauer, DrSc., dr.h.c.; doc. RNDr. Jiří Felcman, CSc.; prof. RNDr. Jaroslav Haslinger, DrSc.; doc. RNDr. Petr Holický, CSc.; doc. Dr. Mgr. Petr Knobloch; doc. RNDr. Martin Kružík, Ph.D.; doc. Mgr. Milan Pokorný, Ph.D.; prof. Ing. Tomáš Roubíček, DrSc.; and doc. RNDr. Jan Zítko, CSc. The courses have provided me the opportunity to obtain the fundamental knowledge of mathematics that I have fully utilized during my works.

Most importantly, I would like to thank my wife Katka and daughter Justýnka for their patience and support during my studies and writing the works.

Last but not least, the financial support of this work provided by the Czech Science Foundation through project No. GAČR P105/12/0331, GAČR 103/09/1748, GAČR 103/09/P490 and by the Grant Agency of the Czech Technical University in Prague through project No. SGS 12/027/OHK1/1T/11, SGS10/124/OHK1/2T/11 is gratefully acknowledged.

“The mathematical facts worthy of being studied are those which, by their analogy with other facts, are capable of leading us to the knowledge of a physical law.” (Henri Poincare)

“Empirical evidence can never establish mathematical existence — nor can the mathematician’s demand for existence be dismissed by the physicist as useless rigor. Only a mathematical existence proof can ensure that the mathematical description of a physical phenomenon is meaningful.” (Richard Courant)

Content

I	Integrating text	1
1	Motivation	1
2	State of the art	2
2.1	Variational formulation	3
2.2	Formulation based on the Lippmann-Schwinger equation	5
2.3	Guaranteed bounds on homogenized matrix	7
3	Methodology	8
4	Results	8
4.1	Weak formulation and Lippmann-Schwinger equation	9
4.2	Discretization via trigonometric polynomials	10
4.2.1	Trigonometric polynomials	10
4.2.2	Galerkin approximation with numerical integration	13
4.3	Guaranteed bounds by FFT-based homogenization	14
4.3.1	Connection of primal and dual formulation	15
4.3.2	Calculation of bounds	15
4.4	Numerical experiments	17
4.4.1	Acceleration by Conjugate gradients	17
4.4.2	Guaranteed bounds	19
4.5	Applications	19
5	Conclusions	20
6	References	21
7	List of thesis papers	26
II	Paper 1	27
III	Paper 2	35
IV	Paper 3	44
V	Paper 4	53
VI	Paper 5	91
VII	Paper 6	120

Part I

Integrating text

1 Motivation

In the recent decades, the computational mechanics — together with its tremendous development, increase in computer performance and hardware availability — has delivered extensive societal benefits. Simulations, replaced with computer models, have dramatically reduced the cost of the engineering design process, reduced design cycle times, and have made it possible to address scientific and engineering questions beyond the capabilities of experiments, [39, 2].

Nevertheless, the extensive engineering problems and their modeling, treated all in one with its complexity, promptly reach the limit of accessible computer capacity, particularly for heterogeneous materials with fully resolved microstructure. For example, concrete — the basic construction material in civil engineering structures — looks and behaves as a homogeneous substance at a macroscale, the scale of design object. However, a closer investigation into microscale discovers aggregates, sand, and binder; under another closer investigation, they are still heterogeneous. Concurrently, the heterogeneities are crucial, for example, in crack propagation problems; thus the design difficulty lies in a detail that is required. This justifies the development of new effective and reliable computational models, methods, and algorithms dealing with heterogeneities.

This work is dedicated to numerical homogenization of linear periodic materials or to the problem of finding effective (homogenized) material properties — electric conductivity, heat conductivity, or stiffness — by taking into account the heterogeneities at microscale. It consists of the calculation of microscale fields satisfying linear elliptic partial differential equation. Contrary to well-established h , p , hp -versions of Finite Element Method (FEM) based on weak formulation, see e.g. [1, 12, 13, 55, 58], they can be found by a method based on the Fourier Transform — Fast Fourier Transform homogenization (FFTH) — introduced in the form of numerical algorithm in [32]. The latter method is based on the solution of the Lippmann-Schwinger type of integral equation incorporating the Green function of an auxiliary homogeneous problem with a reference material property, the parameter of the method.

The method has become, next to analytical homogenizations [31, 8, 16] and already mentioned FEM, broadly used by engineers — especially for its simplicity in implementation, efficiency resulting from the application of the FFT algorithm, and direct usage of material geometry defined as pixel or voxel images. Although it is a rather special topic of interest, it has a wide area of applications; among others, recent relevant examples are listed: evaluation of effective properties for real structures [45, 54, 48] with data from nanoindentation [37, 38], microstructure modeling with visco-elastic [17, 53] and visco-plastic [21, 23, 44] material, conductivity with imperfect interfaces [20], micromechanical behavior of polycrystals [22, 21, 44], permeability of porous medium [30], non-local fracture (damage) models [24, 25], and microstructure reconstruction with Wang tilings [36].

In addition, for the last two decades, FFTH method has gone through extensive

investigation. The convergence of numerical scheme, based on the Neumann expansion, has been proposed in [33]. The accelerated schemes have emerged in [15, 49]. The algorithm has been extended for voided materials by augmented Lagrangians in [27, 28] and by application of the dual formulation in [29]; some of the methods were compared in [34]. Recently, the discretization of Lippmann-Schwinger equation, based on piecewise constant basis functions, with convergence results was suggested in [6, 7]. Another approach is based on approximation of the Green function with a FEM solution in [56].

Nevertheless, no rigorous theory — providing the resolvability of Lippmann-Schwinger equation for various parameters, comparing the Lippmann-Schwinger equation to variational formulation, establishing discretization and convergence of approximate solutions to the continuous one — has been published. It is the role of this work allowing the additional research within the standard mathematical instruments and contributing to reliability and confidence in engineering design, especially with guaranteed bounds of homogenized material properties.

Note that this dissertation is based on the compilation of six papers [57, 50, 37, 38, 51, 52] that are attached as Parts II–VII in the chronological order, see also Section 7 List of thesis paper. Paper 1 has been published in ISI journal, then Paper 2 was published as a peer-reviewed conference paper. Paper 3 is accepted for publication, while Paper 4 is in peer review stage. Finally, Papers 5 and 6, that are fundamental for the thesis, are to be sent for publication. Since each paper needs to be self-contained, the dissertation may contain repetitions.

2 State of the art

This section briefly overviews the state of the art of numerical homogenization of linear elliptic partial differential equation, namely the evaluation of homogenized matrix by FFT-based homogenization method. A scalar problem of electric conductivity, equally to diffusion or heat transfer, is chosen as a model case.

The variational formulation of homogenization problem is described in Section 2.1, followed by an alternative formulation with Lippmann-Schwinger equation including Moulinec-Suquet numerical solution [32] in Section 2.2 and the theory for guaranteed bounds [12, 13] in Section 2.3 are summarized.

For the later use in this dissertation, the following notation is introduced. Note that some articles can differ in some details, nevertheless, the most recent two [51, 52] should be consistent.

The letter d denotes the dimension of the problem, assuming $d = 2, 3$; the Greek letters α, β are reserved to indices relating dimension, thus ranging $1, \dots, d$ (the range is for simplicity often omitted).

The sets \mathbb{C}^d and \mathbb{R}^d are spaces of complex and real vectors with canonical basis $\{\epsilon_\alpha\}$ and are equipped with the Lebesgue measure $d\mathbf{x}$. We denote by $|\Omega|_d$ the d -dimensional Lebesgue measure of a measurable set $\Omega \subset \mathbb{R}^d$. The norm $\|\cdot\|_2$ on \mathbb{C}^d is induced by scalar product $(\mathbf{u}, \mathbf{v})_{\mathbb{C}^d} = \sum_\alpha u_\alpha \bar{v}_\alpha$ for $\mathbf{u}, \mathbf{v} \in \mathbb{C}^d$.

The set $\mathbb{R}_{\text{spd}}^{d \times d}$ denotes the space of symmetric positive definite matrices of size $d \times d$ with norm $\|\mathbf{C}\|_2 = \max_{\mathbf{x} \in \mathbb{R}^d, \|\mathbf{x}\|=1} \|\mathbf{C}\mathbf{x}\|_2$ that equals to a largest eigenvalue.

A function $f : \mathbb{R}^d \mapsto \mathbb{R}$ is \mathbf{Y} -periodic (with period $\mathbf{Y} \in \mathbb{R}^d$) if $f(\mathbf{x} + \mathbf{Y} \odot \mathbf{k}) = f(\mathbf{x})$ for arbitrary $\mathbf{x} \in \mathbb{R}^d, \mathbf{k} \in \mathbb{Z}^d$, where operator \odot denotes element-wise

multiplication. The \mathbf{Y} -periodic functions are sufficient to define only on a periodic unit cell (PUC), set to $\mathcal{Y} := (-Y_\alpha, Y_\alpha)_{\alpha=1}^d \subset \mathbb{R}^d$. Two integrable functions which are almost everywhere equal are identified. The mean value of function $\mathbf{v} \in L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)$ over periodic unit cell \mathcal{Y} is denoted as $\langle \mathbf{v} \rangle := \frac{1}{|\mathcal{Y}|^d} \int_{\mathcal{Y}} \mathbf{v}(\mathbf{x}) \, d\mathbf{x} \in \mathbb{R}^d$.

We define space $C_{\text{per}}(\mathcal{Y}; \mathbb{X})$ of continuous \mathbf{Y} -periodic functions $\mathbb{R}^d \mapsto \mathbb{X}$, where \mathbb{X} is some finite dimensional vector space, e.g. \mathbb{C} , \mathbb{R} , \mathbb{C}^d , or \mathbb{R}^d . Vector valued functions, for $\mathbb{X} = \mathbb{C}^d$ or $\mathbb{X} = \mathbb{R}^d$, are denoted with small bold letters, e.g. \mathbf{v} with components v_α .

The spaces $L^2_{\text{per}}(\mathcal{Y}; \mathbb{X})$ or $L^\infty_{\text{per}}(\mathcal{Y}; \mathbb{R}^{d \times d}_{\text{spd}})$ are composed of functions $\mathbf{v} : \mathbb{R}^d \mapsto \mathbb{X}$ or $\mathbf{A} : \mathbb{R}^d \mapsto \mathbb{R}^{d \times d}_{\text{spd}}$ having \mathbf{Y} -periodic, measurable components v_α or $A_{\alpha\beta}$ and having finite norm, i.e. $\|\mathbf{v}\|_{L^2_{\text{per}}(\mathcal{Y}; \mathbb{X})} < \infty$ or $\|\mathbf{A}\|_{L^\infty_{\text{per}}(\mathcal{Y}; \mathbb{R}^{d \times d}_{\text{spd}})} < \infty$; the first norm is generated by scalar product $(\mathbf{u}, \mathbf{v})_{L^2_{\text{per}}(\mathcal{Y}; \mathbb{X})} = \frac{1}{|\mathcal{Y}|^d} \int_{\mathcal{Y}} (\mathbf{u}(\mathbf{x}), \mathbf{v}(\mathbf{x}))_{\mathbb{X}} \, d\mathbf{x}$ while the second norm is defined as $\|\mathbf{A}\|_{L^\infty_{\text{per}}(\mathcal{Y}; \mathbb{R}^{d \times d}_{\text{spd}})} = \text{esssup}_{\mathbf{x} \in \mathcal{Y}} \|\mathbf{A}(\mathbf{x})\|_2$. If there is no ambiguity, both the norms and the scalar products are denoted with subscript L^2_{per} or L^∞_{per} rather than $L^2_{\text{per}}(\mathcal{Y}; \mathbb{X})$ or $L^\infty_{\text{per}}(\mathcal{Y}; \mathbb{R}^{d \times d}_{\text{spd}})$.

Next, we introduce the Helmholtz decomposition $L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d) = \mathcal{U} \oplus^\perp \mathcal{E} \oplus^\perp \mathcal{J}$ to the spaces of constant, curl-free with zero mean, and divergence free with zero mean fields

$$\mathcal{U} = \{\mathbf{v} \in L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d) : \mathbf{v}(\mathbf{x}) = \text{const.}\}, \quad (2.1a)$$

$$\mathcal{E} = \{\mathbf{v} \in L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d) : \nabla \times \mathbf{v} = \mathbf{0}, \langle \mathbf{v} \rangle = \mathbf{0}\}, \quad (2.1b)$$

$$\mathcal{J} = \{\mathbf{v} \in L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d) : \nabla \cdot \mathbf{v} = 0, \langle \mathbf{v} \rangle = \mathbf{0}\}, \quad (2.1c)$$

where differential operator $\nabla = (\frac{\partial}{\partial x_\alpha})_{\alpha=1}^d$ is meant in the distributional sense. For dimension $d \neq 3$, the curl-free condition in (2.1b) means $(\nabla \times \mathbf{v})_{\alpha\beta} := \frac{\partial v_\alpha}{\partial x_\beta} - \frac{\partial v_\beta}{\partial x_\alpha} = 0$. Since space \mathcal{U} consists of constant functions, we identify the space \mathcal{U} with \mathbb{R}^d ; this validates the operations such as $\mathbf{E} + \mathbf{v} \in L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)$ for $\mathbf{E} \in \mathbb{R}^d$, $\mathbf{v} \in L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)$ and $\mathbf{C}\mathbf{J} \in \mathbb{R}^d$ for $\mathbf{C} \in \mathbb{R}^{d \times d}$ and $\mathbf{J} \in \mathcal{U}$.

2.1 Variational formulation

This section begins with a homogenization problem defining a homogenized or effective matrix for a scalar linear elliptic problem, particularly the problem of electric conductivity.

Here and in the sequel, $\mathbf{A} \in L^\infty_{\text{per}}(\mathcal{Y}, \mathbb{R}^{d \times d}_{\text{spd}})$ denotes symmetric¹ and uniformly elliptic² material coefficients of electric conductivity, $\mathbf{e} \in \mathcal{E}$ and $\mathbf{j} \in \mathcal{J}$ perturbation of electric field and electric current, and $\mathbf{E}, \mathbf{J} \in \mathcal{U}$ are their macroscopic counterparts.

Definition 2.1 (Homogenization problem). *The primal and dual homogenization problem states: find homogenized matrix $\mathbf{A}_{\text{eff}} \in \mathbb{R}^{d \times d}$ satisfying for arbitrary fixed*

¹For almost all $\mathbf{x} \in \mathcal{Y}$, equality $\mathbf{A}(\mathbf{x}) = \mathbf{A}(\mathbf{x})^T$ holds.

²There exists positive constant $c_A > 0$ such that for almost all $\mathbf{x} \in \mathcal{Y}$ and all nonzero $\mathbf{u} \in \mathbb{R}^d$, inequality $c_A \|\mathbf{u}\|_2^2 \leq (\mathbf{A}(\mathbf{x})\mathbf{u}, \mathbf{u})_{\mathbb{R}^d}$ holds.

macroscopic quantities $\mathbf{E} \in \mathbb{R}^d$

$$(\mathbf{A}_{\text{eff}} \mathbf{E}, \mathbf{E})_{\mathbb{R}^d} = \inf_{\mathbf{e} \in \mathcal{E}} (\mathbf{A}(\mathbf{E} + \mathbf{e}), \mathbf{E} + \mathbf{e})_{L^2_{\text{per}}} \quad (2.2a)$$

$$(\mathbf{A}_{\text{eff}}^{-1} \mathbf{J}, \mathbf{J})_{\mathbb{R}^d} = \inf_{\mathbf{j} \in \mathcal{J}} (\mathbf{A}^{-1}(\mathbf{J} + \mathbf{j}), \mathbf{J} + \mathbf{j})_{L^2_{\text{per}}} \quad (2.2b)$$

Remark 2.2. For particular macroscopic load \mathbf{E} , a minimizer in previous definition exists and is unique according to the direct method in the calculus of variation since material coefficients \mathbf{A} are bounded and uniformly elliptic, cf Rem. 2.5. Next, homogenized matrix \mathbf{A}_{eff} coincide in both formulations, is symmetric and positive definite, see e.g. [3, 11].

Remark 2.3. Space \mathcal{E} contains curl-free functions — such function \mathbf{e} possess potential U such that $\nabla U = \mathbf{e}$. The reformulation of homogenization problem with potentials avoids a trouble with constructing finite element spaces with curl-free basis functions, however, FFT-based FEM naturally overcome this difficulty by employing a certain projection operators, see Lem. 4.2.

Remark 2.4. The previous homogenization formulas in Eq. (2.2) are obtained from an asymptotic behavior as $\varepsilon \rightarrow 0$ of a linear elliptic problem: for $\varepsilon > 0$, find $u^\varepsilon \in H_0^1(\Omega) := \{v \in L^2(\Omega), v|_{\partial\Omega} = 0, \frac{\partial v_\alpha}{\partial x_\alpha} \in L^2(\Omega)\}$ such that³

$$(\mathbf{A}^\varepsilon \nabla u^\varepsilon, \nabla v)_{L^2(\Omega)} = F(v), \quad \forall v \in H_0^1(\Omega) \quad (2.3)$$

where Ω is a bounded open set in \mathbb{R}^d with the Lipschitz boundary $\partial\Omega$. Oscillating material coefficients are defined as $\mathbf{A}^\varepsilon(\mathbf{x}) := \mathbf{A}\left(\frac{\mathbf{x}}{\varepsilon}\right)$ for prescribed material coefficients \mathbf{A} and linear functional F contains boundary conditions and source terms.

Since the material coefficients are uniformly elliptic and bounded, the unique solutions u^ε of weak formulation (2.3) exists, are uniformly bounded for ε , and thus weakly converges in $H_0^1(\Omega)$ to some function $u_{\text{eff}} \in H_0^1(\Omega)$ representing the averaged field of oscillating solutions u^ε .

Various methods, see e.g. the notion of H -convergence in [11, 46] or formal method of asymptotic expansion [3], reveal that weak limit u_{eff} satisfies variational formulation

$$(\mathbf{A}_{\text{eff}} \nabla u_{\text{eff}}, \nabla v)_{L^2(\Omega)} = F(v), \quad \forall v \in H_0^1(\Omega) \quad (2.4)$$

with homogeneous material coefficients \mathbf{A}_{eff} obtained by homogenization formula (2.2a).

Hence the homogenized matrix represents a limit state, however, in reality it is a reliable value if the periodic unit cell is sufficiently small compared to the macroscale of a designed object.

The necessity for homogenization theories is observed from Eq. (2.3); a direct discretization, for some small ε , is inappropriate because of highly oscillating solutions — it requires huge number of degrees of freedom. The homogenization thus splits elliptic problem (2.3) to the evaluation of homogenized matrix (2.2a) and to the solution of Eq. (2.4) without the oscillatory part.

³The values on boundary $v|_{\partial\Omega}$ are meant in the sense of trace operator, partial derivatives $\frac{\partial v_\alpha}{\partial x_\alpha}$ are meant in the sense of distributional derivatives.

In the sequel, an attention is focused on the primal homogenization problem in Def. 2.1; the dual homogenization problem is only utilized later in Sec. 4.3 for the guaranteed bounds of the homogenized matrix.

Remark 2.5. *Since material coefficients \mathbf{A} are symmetric, the homogenization problem in Eq. 2.2a is equivalent to a weak formulation⁴*

$$(\mathbf{A}\tilde{\mathbf{e}}^{(\mathbf{E})}, \mathbf{v})_{L^2_{\text{per}}} = -(\mathbf{A}\mathbf{E}, \mathbf{v})_{L^2_{\text{per}}}, \quad \forall \mathbf{v} \in \mathcal{E} \quad (2.5)$$

describing the minimizer point. Since it has a linear structure, minimizer $\tilde{\mathbf{e}}^{\mathbf{E}} \in \mathcal{E}$ for macroscopic field $\mathbf{E} \in \mathbb{R}^d$ is obtained from unitary minimizers $\tilde{\mathbf{e}}^{(\alpha)}$, see the next definition, as

$$\tilde{\mathbf{e}}^{\mathbf{E}} = \sum_{\alpha} E_{\alpha} \tilde{\mathbf{e}}^{(\alpha)}.$$

Definition 2.6 (auxiliary problems). *We say that $\tilde{\mathbf{e}}^{(\alpha)} \in \mathcal{E}$ are unitary minimizers if they satisfy the weak formulations with unitary macroscopic loads $\boldsymbol{\epsilon}_{\alpha}$, i.e.*

$$(\mathbf{A}\tilde{\mathbf{e}}^{(\alpha)}, \mathbf{v})_{L^2_{\text{per}}} = -(\mathbf{A}\boldsymbol{\epsilon}_{\alpha}, \mathbf{v})_{L^2_{\text{per}}}, \quad \forall \mathbf{v} \in \mathcal{E}, \quad (2.6)$$

Unitary microscopic fields are noted without tilde, i.e. $\mathbf{e}^{(\alpha)} := \boldsymbol{\epsilon}_{\alpha} + \tilde{\mathbf{e}}^{(\alpha)} \in \mathcal{U} \oplus^{\perp} \mathcal{E}$.

Remark 2.7. *The unitary minimizers also serve to evaluate the components of the homogenized matrix; formula*

$$\mathbf{A}_{\text{eff}, \alpha\beta} = (\mathbf{A}(\boldsymbol{\epsilon}_{\beta} + \tilde{\mathbf{e}}^{(\beta)}), (\boldsymbol{\epsilon}_{\alpha} + \tilde{\mathbf{e}}^{(\alpha)}))_{L^2_{\text{per}}}. \quad (2.7)$$

follows from the linear structure.

2.2 Formulation based on the Lippmann-Schwinger equation

This section is dedicated to an alternative formulation of the auxiliary problems in Def. 2.6, namely to the Lippmann-Schwinger equation incorporating the Green function for a reference homogeneous problem.

Definition 2.8 (Lippmann-Schwinger equation). *Let parameter $\mathbf{A}^0 \in \mathbb{R}_{\text{spd}}^{d \times d}$ be a symmetric positive definite matrix. We say that $\mathbf{e}_{\text{LS}}^{(\mathbf{E})} \in L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)$ is the solution of Lippmann-Schwinger equation if it satisfies*

$$\mathbf{e}_{\text{LS}}^{(\mathbf{E})}(\mathbf{x}) + \int_{\mathcal{Y}} \Gamma^0(\mathbf{x} - \mathbf{y})(\mathbf{A}(\mathbf{y}) - \mathbf{A}^0)\mathbf{e}_{\text{LS}}^{(\mathbf{E})}(\mathbf{y}) \, d\mathbf{y} = \mathbf{E}, \quad \text{for almost all } \mathbf{x} \in \mathcal{Y} \quad (2.8)$$

where the convolution integral is defined with the help of the Fourier series

$$\int_{\mathcal{Y}} \Gamma^0(\mathbf{x} - \mathbf{y})\mathbf{v}(\mathbf{y}) \, d\mathbf{y} := \sum_{\mathbf{k} \in \mathbb{Z}^d \setminus \{\mathbf{0}\}} \frac{\boldsymbol{\xi}(\mathbf{k}) \otimes \boldsymbol{\xi}(\mathbf{k})}{(\mathbf{A}^0 \boldsymbol{\xi}(\mathbf{k}), \boldsymbol{\xi}(\mathbf{k}))_{\mathbb{R}^d}} \hat{\mathbf{v}}(\mathbf{k}) \varphi_{\mathbf{k}}(\mathbf{x}) \quad (2.9)$$

⁴Left-hand side represents continuous symmetric uniformly elliptic bilinear form for $\tilde{\mathbf{e}}^{(\mathbf{E})}$ and \mathbf{v} . Right-hand side represents continuous linear functional for \mathbf{v} . The existence of an unique solution is provided by Lax-Milgram lemma or simply by Riesz representation theorem.

where $\boldsymbol{\xi}(\mathbf{k}) \in \mathbb{R}^d$ is a vector with components $\xi_\alpha(\mathbf{k}) = \frac{k_\alpha}{Y_\alpha}$, operator \otimes denotes the tensor product⁵, and $\hat{\mathbf{v}}(\mathbf{k})$ are the Fourier coefficients with components $\hat{v}_\alpha(\mathbf{k}) := (v_\alpha, \varphi_{\mathbf{k}})_{L^2_{\text{per}}}$ for trigonometric polynomial $\varphi_{\mathbf{k}} = \exp(i\pi \sum_\alpha \frac{x_\alpha k_\alpha}{Y_\alpha})$.

Lippmann-Schwinger equation is formulated for microscopic field $\mathbf{e}_{\text{LS}}^{(\mathbf{E})}$ contrary to weak formulation (2.5) written for the perturbation part $\tilde{\mathbf{e}}^{(\mathbf{E})} = \mathbf{e}^{(\mathbf{E})} - \mathbf{E}$; both of the solutions, if coincide, differ by constant \mathbf{E} as $\langle \mathbf{e}_{\text{LS}}^{(\mathbf{E})} \rangle = \mathbf{E}$ and $\langle \tilde{\mathbf{e}}^{(\mathbf{E})} \rangle = \mathbf{0}$.

Remark 2.9. *Lippmann-Schwinger equation is deduced from a strong formulation⁶: for prescribed macroscopic load $\mathbf{E} \in \mathbb{R}^d$, find $\tilde{\mathbf{e}}$ with continuous partial derivatives satisfying*

$$\nabla \cdot [\mathbf{A}(\mathbf{E} + \tilde{\mathbf{e}})] = 0, \quad \nabla \times \tilde{\mathbf{e}} = 0, \quad \langle \tilde{\mathbf{e}} \rangle = 0. \quad (2.10)$$

The problem is reformulated for a reference homogeneous material $\mathbf{A}^0 \in \mathbb{R}_{\text{spd}}^{d \times d}$, the parameter of Lippmann-Schwinger equation, to a homogeneous problem

$$\nabla \cdot [\mathbf{A}^0 \mathbf{e}(\mathbf{x})] = \mathbf{f}(\mathbf{x})$$

where $\mathbf{f}(\mathbf{x}) = -\nabla \cdot [(\mathbf{A}(\mathbf{x}) - \mathbf{A}^0) \mathbf{e}(\mathbf{x})]$ is called a divergence of polarization current and $\mathbf{e} = \mathbf{E} + \tilde{\mathbf{e}}$ is microscopic field with $\tilde{\mathbf{e}}$ satisfying the curl-free and zero mean condition. It is then transformed with the technique of the Fourier transform⁷ to the system of algebraic equations, whose solution yields the Lippmann-Schwinger equation (2.8).

The resulting equation is already defined in a way to have a good sense and to be easily analyzed, see Sec. 4.1 and particularly Theorem 4.1. It will become apparent that the detailed derivation of Lippmann-Schwinger equation can be omitted; the only important part is the convolution integral providing projection on curl-free fields with zero mean, see Lem. 4.2.

Remark 2.10 (Solution of Lippmann-Schwinger equation). *Lippmann-Schwinger equation (2.8) can be written in operator form $(I + \mathcal{B})\mathbf{e}_{\text{LS}}^{(\mathbf{E})} = \mathbf{E}$ with the obvious definition of operators I and \mathcal{B} . The inverse of $(I + \mathcal{B})$ is then expressed using Neumann series expansion $(I + \mathcal{B})^{-1} = \sum_{k=0}^{\infty} \mathcal{B}^k$ as $\|\mathcal{B}\| < 1$ for the special choice of parameter $\mathbf{A}^0 := \frac{1}{2}(\|\mathbf{A}^{-1}\|_{L^\infty_{\text{per}}} + \|\mathbf{A}\|_{L^\infty_{\text{per}}})\mathbf{I}$, see [28].*

Remark 2.11 (FFT-based method according to Moulinec and Suquet in [32]). *The solution of Lippmann-Schwinger equation obtained from Neumann series $\mathbf{e}_{\text{LS}}^{(\mathbf{E})} = \sum_{k=0}^{\infty} \mathcal{B}^k[\mathbf{E}]$ is the limit in L^2_{per} norm as $k \rightarrow \infty$ of iteration algorithm $\mathbf{e}_{(k+1)} = \mathcal{B}[\mathbf{e}_{(k)}] + \mathbf{E}$ with initial value $\mathbf{e}_{(0)} := \mathbf{E}$. The iteration algorithm serves as a foundation for the numerical algorithm proposed by Moulinec and Suquet in [32, 33]*

⁵The tensor product of two vectors $\mathbf{u}, \mathbf{v} \in \mathbb{R}^d$ defines matrix $\mathbf{u} \otimes \mathbf{v} \in \mathbb{R}^{d \times d}$ with components $(\mathbf{u} \otimes \mathbf{v})_{\alpha\beta} = u_\alpha v_\beta$.

⁶The solution of the strong formulation also satisfies the weak formulation in Eq. (2.5); it is shown by the multiplication of a test function, by integration over PUC \mathcal{Y} , and by application of Green's theorem (boundary term vanishes because of periodicity).

⁷The most important properties are: derivative is transformed to the multiplication with a Fourier variable and the convolution of two functions is transformed to their multiplication. The inverse Fourier transform is provided with an analogical formula.

who approximated solution \mathbf{e} with the function values at a regular grid corresponding to pixels or voxels in the images of periodic unit cell. The convolution integral of Lippmann-Schwinger equation is then replaced with the Discrete Fourier Transform⁸, its inverse, and the multiplication by the integral kernel at the Fourier space; in a detail, the algorithm is described in Paper 1.

2.3 Guaranteed bounds on homogenized matrix

The upper-lower bounds obtained from a posteriori error estimates were introduced by Dvořák [12, 13] for a scalar problem and independently by Wiećkowski [55] for linear elasticity. This section provides a summary of results in [12]. In this section, we work with some conforming approximations of unitary minimizers $\tilde{\mathbf{e}}^{(\alpha)}$ and $\tilde{\mathbf{j}}^{(\alpha)}$, namely $\tilde{\mathbf{e}}_{\mathbf{N}}^{(\alpha)} = (\mathbf{e}_{\mathbf{N}}^{(\alpha)} - \boldsymbol{\epsilon}_\alpha) \in \mathcal{E}$ and $\tilde{\mathbf{j}}_{\mathbf{N}}^{(\alpha)} = (\mathbf{j}_{\mathbf{N}}^{(\alpha)} - \boldsymbol{\epsilon}_\alpha) \in \mathcal{J}$; parameter \mathbf{N} represents the inverse of discretization size of FEM or the number of discretization points in the case of FFT-based method, for details see Sec. 4.2 or Papers 5 and 6. In what follows, relation $\mathbf{C} \preceq \mathbf{D}$ between symmetric and positive definite matrices $\mathbf{C}, \mathbf{D} \in \mathbb{R}_{\text{spd}}^{d \times d}$ stands for ordering in the sense of quadratic forms; it is equivalent to $\mathbf{E} \cdot \mathbf{C} \mathbf{E} \leq \mathbf{E} \cdot \mathbf{D} \mathbf{E}$ for all $\mathbf{E} \in \mathbb{R}^d$. Upper bound of homogenized matrix is obtained from the primal homogenization problem (2.2a). The replacement of minimizers with approximate minimizers leads to an increase in the value of the quadratic form with the homogenized matrix

$$(\mathbf{A}_{\text{eff}} \mathbf{E}, \mathbf{E})_{\mathbb{R}^d} = \inf_{\mathbf{e} \in \mathcal{E}} (\mathbf{A}(\mathbf{E} + \mathbf{e}), \mathbf{E} + \mathbf{e})_{L_{\text{per}}^2} \leq (\mathbf{A}(\mathbf{E} + \mathbf{e}_{\mathbf{N}}^{(\mathbf{E})}), \mathbf{E} + \mathbf{e}_{\mathbf{N}}^{(\mathbf{E})})_{L_{\text{per}}^2}.$$

The last term then defines the upper bound of the homogenized matrix.

Definition 2.12. We say that matrix $\overline{\mathbf{A}}_{\text{eff}, \mathbf{N}} \in \mathbb{R}^{d \times d}$ defined as

$$(\overline{\mathbf{A}}_{\text{eff}, \mathbf{N}})_{\alpha\beta} = (\mathbf{A}(\boldsymbol{\epsilon}_\beta + \tilde{\mathbf{e}}_{\mathbf{N}}^{(\beta)}), \boldsymbol{\epsilon}_\alpha + \tilde{\mathbf{e}}_{\mathbf{N}}^{(\alpha)})_{L_{\text{per}}^2}$$

is the upper bound on homogenized matrix \mathbf{A}_{eff} .

The upper bound of homogenized matrix $\mathbf{A}_{\text{eff}} \preceq \overline{\mathbf{A}}_{\text{eff}, \mathbf{N}}$ can be completed with the upper bound of inverse homogenized matrix $\mathbf{A}_{\text{eff}}^{-1} \preceq \underline{\mathbf{A}}_{\text{eff}, \mathbf{N}}^{-1}$ obtained from dual formulation (2.2b). Both relations lead to the upper-lower bound structure of the homogenized matrix

$$\underline{\mathbf{A}}_{\text{eff}, \mathbf{N}} \preceq \mathbf{A}_{\text{eff}} \preceq \overline{\mathbf{A}}_{\text{eff}, \mathbf{N}}. \quad (2.11)$$

Remark 2.13 (Element-wise upper-lower bounds). *Upper-lower bounds (2.11) imply element-wise bounds, the bounds on components of the homogenized matrix.*

The mean of upper-lower bounds $\mathbf{A}_{\text{eff}, \mathbf{N}} = \frac{1}{2}(\overline{\mathbf{A}}_{\text{eff}, \mathbf{N}} + \underline{\mathbf{A}}_{\text{eff}, \mathbf{N}})$ is a more reliable approximation of homogenized matrix with guaranteed error $\mathbf{D}_{\mathbf{N}} = \frac{1}{2}(\overline{\mathbf{A}}_{\text{eff}, \mathbf{N}} - \underline{\mathbf{A}}_{\text{eff}, \mathbf{N}})$. Following lemma provides the convergence⁹ of the error to zero if the approximates minimizers converge to minimizers as $\min_\alpha N_\alpha \rightarrow \infty$.

⁸Numerically, it is realized with FFT algorithm, for details see [9].

⁹The trace of a matrix is a norm on the set of positive definite matrices, particularly on $\mathbf{D}_{\mathbf{N}}$ for $\mathbf{N} \in \mathbb{N}^d$.

Lemma 2.14 (Rate of convergence of homogenized properties). *Guaranteed error \mathbf{D}_N satisfies following inequality*

$$\mathrm{tr} \mathbf{D}_N \leq C_1 \sum_{\alpha} \|\tilde{\mathbf{e}}^{(\alpha)} - \tilde{\mathbf{e}}_N^{(\alpha)}\|_{L^2_{\mathrm{per}}}^2 + C_2 \sum_{\alpha} \|\tilde{\mathbf{j}}^{(\alpha)} - \tilde{\mathbf{j}}_N^{(\alpha)}\|_{L^2_{\mathrm{per}}}^2$$

where constants C_1, C_2 are independent of N .

3 Methodology

The goal of this dissertation consists of analysis of Lippmann-Schwinger equation (2.8) and particularly to Moulinec-Suquet algorithm [32], see Rem. 2.11. The relevant questions arising in the analysis can be summarized as:

1. Is there the unique solution of Lippmann-Schwinger equation for various parameters \mathbf{A}^0 ; the existence has been shown only for particular choice of parameter \mathbf{A}^0 , see Rem. 2.10. Does the solution of Lippmann-Schwinger equation equal to the solution of weak formulation (2.5)? Do the solutions coincide for various parameters \mathbf{A}^0 ?
2. (a) Is there some consistent approximation of weak formulation (2.5) or of Lippmann-Schwinger equation (2.8) that is equivalent to Moulinec-Suquet numerical algorithm? Do approximate solutions, if exist, converge to the solution of weak formulation?
- (b) It has been observed by Zeman et al. in [57] that Moulinec-Suquet algorithm, Rem. 2.11, is equivalent to the solution of a system of linear equations. It appears that this non-symmetric system can be successfully solved using Conjugate gradient algorithm. What is the reason for that?
3. In [12, 13] and later independently in [55], the guaranteed bounds of the homogenized matrix were calculated using the standard Finite element methods. Is the approach applicable to the FFT-based homogenization?

The analysis of the method is provided with standard mathematical instruments and techniques. Namely, it is based on the following subjects with the list of literature: homogenization theory [11, 40, 18, 3], mathematical analysis [41, 42, 26], modern theory of partial differential equations [19, 5, 35], finite element method [10, 4], and trigonometric collocation method [43].

Numerical calculations stated in Sec. 4.4, 4.5 and attached papers were provided with the own software written in MATLAB[®] and PYTHON programming language; it is available at <http://mech.fsv.cvut.cz/~vondrejic/publications.php#SW>.

4 Results

This section provides the summary of the results obtained in six papers [57, 50, 37, 38, 51, 52], referenced as Papers 1-6, that are attached in the chronological order in Parts II–VII. The section is split into five subsections; The first three correspond to the tasks in Methodology Sec. 3 while the last two contain numerical examples and applications of FFT-based method to linear elasticity.

4.1 Weak formulation and Lippmann-Schwinger equation

This section start with theorem describing solvability of Lippmann-Schwinger equation, the main theorem about continuous formulation of homogenization problem.

Theorem 4.1 (Equivalence of weak formulation and Lippmann-Schwinger equation, Theorem 2.29 in Paper 5). *Let $\mathbf{A}^0 \in \mathbb{R}_{\text{spd}}^{d \times d}$ be symmetric and positive definite, then weak formulation (2.5) and Lippmann-Schwinger equation (2.8) are equivalent in the sense that the solution coincide in both formulations.*

The theorem shows that the unique solution of weak formulation (2.5) is the solution of Lippmann-Schwinger equation and contrary. It reveals the existence of the unique solution of Lippmann-Schwinger equation for various parameters \mathbf{A}^0 .

The proof is based on a next lemma providing a projection on \mathcal{E} , the space of curl-free with zero mean fields that is both a trial space and the space of test functions in weak formulation (2.5).

Lemma 4.2 (Lem. 2.28 in Paper 5). *Operator $\mathcal{G}[\cdot] : L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d) \rightarrow L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)$ defined as*

$$\mathcal{G}[\mathbf{f}](\mathbf{x}) = \int_{\mathcal{Y}} \Gamma^0(\mathbf{x} - \mathbf{y}) \mathbf{A}^0 \mathbf{f}(\mathbf{y}) \, d\mathbf{y} \quad (4.1)$$

is a projection on \mathcal{E} and orthogonal for $\mathbf{A}^0 = \lambda \mathbf{I}$ with $\lambda > 0$.

The proof is based on the expression of convolution integral in the Fourier space, cf. Def. 2.8. It is shown that the operator is continuous and real valued. Since matrix $\frac{\boldsymbol{\xi}(\mathbf{k}) \otimes \boldsymbol{\xi}(\mathbf{k})}{\mathbf{A}^0 \boldsymbol{\xi}(\mathbf{k}) \cdot \boldsymbol{\xi}(\mathbf{k})} \mathbf{A}^0$ is simply a projection, it shows that operator \mathcal{G} is. Besides, we show that it is the projection on \mathcal{E} with the help of a potential that exists in this case as the periodic unit cell is simply connected.

Proof outline of Theorem 4.1. The proof is outlined for a special choice of parameter \mathbf{A}^0 equal to the identity matrix $\mathbf{A}^0 = \mathbf{I}$. First, we show that the solution of the weak formulation is the solution of Lippmann-Schwinger equation. We enlarge space \mathcal{E} in the weak formulation by adding the projection operator

$$(\mathbf{A}\tilde{\mathbf{e}}, \mathcal{G}[\mathbf{v}])_{L^2_{\text{per}}} = -(\mathbf{A}\mathbf{E}, \mathcal{G}[\mathbf{v}])_{L^2_{\text{per}}}, \quad \forall \mathbf{v} \in L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d).$$

In this special case, operator \mathcal{G} is self-adjoint leading to

$$(\mathcal{G}[\mathbf{A}\tilde{\mathbf{e}}], \mathbf{v})_{L^2_{\text{per}}} = -(\mathcal{G}[\mathbf{A}\mathbf{E}], \mathbf{v})_{L^2_{\text{per}}}, \quad \forall \mathbf{v} \in L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d).$$

Finally, a following equation is deduced

$$\mathcal{G}[\mathbf{A}\tilde{\mathbf{e}}] = -\mathcal{G}[\mathbf{A}\mathbf{E}], \quad \text{for almost all } \mathbf{x} \in \mathcal{Y}$$

by removing the testing functions and the rest is the consequence of elementary algebra and properties of projection \mathcal{G} .

The opposite implication is done analogically by splitting the solution $\mathbf{e}_{\text{LS}}^{(E)}$ with projection operator \mathcal{G} . The general case for $\mathbf{A}^0 \in \mathbb{R}_{\text{spd}}^{d \times d}$ follows similar ideas. \square

4.2 Discretization via trigonometric polynomials

This section is dedicated to the discretization of the homogenization problem (2.2) in a way to produce the identical linear system like in Moulinec-Suquet algorithm [32], cf. Rem. 2.11.

The space of trigonometric polynomials is used as a finite dimensional space for discretization. Relating definitions and properties are summarized in Sec. 4.2.1 according to [47, 43], however, a slight modification is performed for the non-odd number of discretization points to assure conforming approximations that are required for guaranteed bounds of the homogenized matrix, see Sec. 4.3 or Papers 5-6 for more details.

Then, approximate solutions are defined through Galerkin approximation (GA) and Galerkin approximation with numerical integration (GAwNI). Together with their convergence to the minimizers and the numerical solution of the linear system, it is described in Sec. 4.2.2 and in a detail in Papers 5 and 6.

Notation 4.3. *In the sequel, let $\mathbf{N} \in \mathbb{N}^d$ be reserved for a number of discretization points with the number of degrees of freedom $|\mathbf{N}|_{\Pi} := \prod_{\alpha} N_{\alpha}$; if N_{α} is odd (even) for all α we talk about odd (even) number of discretization points, otherwise about non-odd ones. The reduced and full index set state for*

$$\mathbb{Z}_{\mathbf{N}}^d = \left\{ \mathbf{k} \in \mathbb{Z}^d : |k_{\alpha}| < \frac{N_{\alpha}}{2} \right\}, \quad \underline{\mathbb{Z}}_{\mathbf{N}}^d = \left\{ \mathbf{k} \in \mathbb{Z}^d : -\frac{N_{\alpha}}{2} \leq k_{\alpha} < \frac{N_{\alpha}}{2} \right\}. \quad (4.2)$$

A multi-index notation is employed, in which $\mathbb{R}^{\mathbf{N}}$ represents $\mathbb{R}^{N_1 \times \dots \times N_d}$. Set $\mathbb{X}_{\mathbf{N}}$ represents the space of vectors \mathbf{v} with components $v_{\alpha}^{\mathbf{n}}$ and $\mathbb{X}_{\mathbf{N}}^2$ the space of matrices \mathbf{A} with components $A_{\alpha\beta}^{\mathbf{nm}}$ for α, β and $\mathbf{n}, \mathbf{m} \in \underline{\mathbb{Z}}_{\mathbf{N}}^d$. Next, $\mathbf{v}^{\mathbf{n}} \in \mathbb{R}^d$ for $\mathbf{n} \in \underline{\mathbb{Z}}_{\mathbf{N}}^d$ and $\mathbf{v}_{\alpha} \in \mathbb{R}^{\mathbf{N}}$ for α represent subvectors of \mathbf{v} with components $v_{\alpha}^{\mathbf{n}}$; analogically the submatrices $\mathbf{A}^{\mathbf{nm}} \in \mathbb{R}^{d \times d}$ and $\mathbf{A}_{\alpha\beta} \in \mathbb{R}^{\mathbf{N} \times \mathbf{N}}$ can be defined. A scalar product on $\mathbb{X}_{\mathbf{N}}$ is defined as $(\mathbf{u}, \mathbf{v})_{\mathbb{X}_{\mathbf{N}}} := \sum_{\alpha} \sum_{\mathbf{n} \in \underline{\mathbb{Z}}_{\mathbf{N}}^d} u_{\alpha}^{\mathbf{n}} v_{\alpha}^{\mathbf{n}}$ and matrix \mathbf{A} by vector \mathbf{v} multiplication as $(\mathbf{A}\mathbf{v})_{\alpha}^{\mathbf{n}} := \sum_{\beta} \sum_{\mathbf{m} \in \underline{\mathbb{Z}}_{\mathbf{N}}^d} A_{\alpha\beta}^{\mathbf{nm}} v_{\beta}^{\mathbf{m}}$. Matrix \mathbf{A} is symmetric positive definite if $A_{\alpha\beta}^{\mathbf{mn}} = A_{\beta\alpha}^{\mathbf{nm}}$ holds for all components and $(\mathbf{A}\mathbf{v}, \mathbf{v})_{\mathbb{X}_{\mathbf{N}}} > 0$ applies for arbitrary $\mathbf{v} \in \mathbb{X}_{\mathbf{N}}$. We use a serif font for vectors \mathbf{v} and matrices \mathbf{A} to distinguish from vectors $\mathbf{E} \in \mathbb{R}^d$ and matrices $\mathbf{A}_{\text{eff}} \in \mathbb{R}^{d \times d}$ and from vector valued functions $\mathbf{v} \in L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)$. In order to distinguish vectors and matrices for different number of discretization points \mathbf{N} , we write them with subscript, i.e. $\mathbf{v}_{\mathbf{N}}$ and $\mathbf{A}_{\mathbf{N}}$.

Operator \oplus^{\perp} denotes the direct sum of mutually orthogonal subspaces, e.g. $\mathbb{R}^d = \boldsymbol{\epsilon}_1 \oplus^{\perp} \boldsymbol{\epsilon}_2 \oplus^{\perp} \dots \oplus^{\perp} \boldsymbol{\epsilon}_d$.

4.2.1 Trigonometric polynomials

This section presents the trigonometric polynomials with their properties; it is also well described in Paper 5 for the odd number of discretization points and in Paper 6 for the general number of discretization points, see also [43].

First, the trigonometric polynomials can be expressed as the linear combination of the Fourier coefficients and the Fourier basis functions, see Eq. (4.3) and (4.5). Second, they can be expressed as the linear combination of function values at nodal

points and shape basis functions, see Def. 4.5 and Eq. (4.6); Fig. 1 shows the examples of the nodal points and the shape basis function. Both formulations are related through the Discrete Fourier Transform (DFT) see Rem. 4.12.

Dirac delta property of shape basis functions $\varphi_{N,m}(\mathbf{x}_N^n) = \delta_{mn}$ ensures the uniqueness of trigonometric polynomials representation and allows us to define an interpolation operator through the nodal points, see Def. 4.8 and 4.5.

Both definitions of trigonometric polynomials \mathcal{T}_N and $\tilde{\mathcal{T}}_N$ in Def. 4.7 coincide if the number of discretization points N is odd for all components, cf. Rem. 4.13.

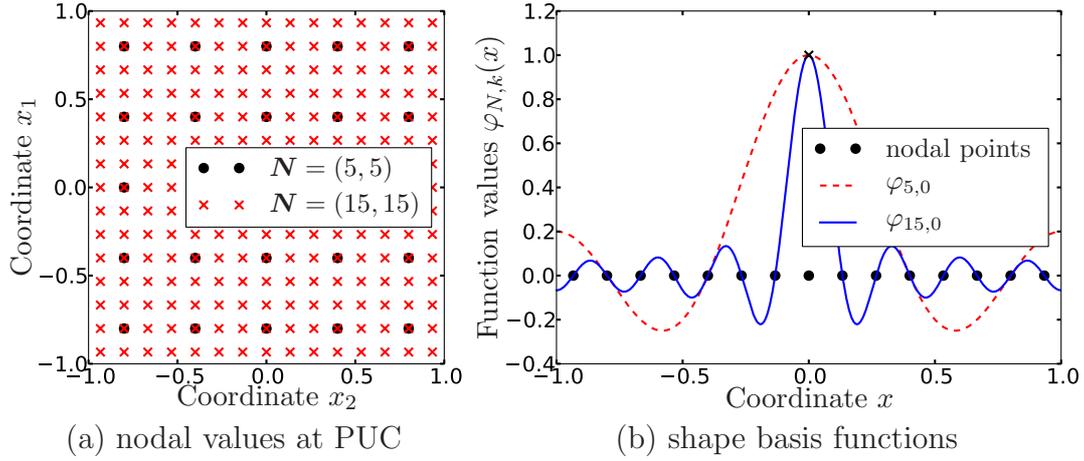


Figure 1: Nodal values for $d = 2$ and shape basis function for $d = 1$

Notation 4.4 (DFT). For $N \in \mathbb{N}^d$ we define, up to constant, unitary matrices $\mathbf{F}_N, \mathbf{F}_N^{-1} \in \mathbb{C}^{d \times d \times N \times N}$ of the Discrete Fourier transform (DFT) and its inverse (iDFT) as

$$\mathbf{F}_N = \frac{1}{|N|_{\Pi}} (\delta_{\alpha\beta} \omega_N^{-mn})_{\alpha,\beta=1,\dots,d}^{m,n \in \mathbb{Z}_N^d} \quad \mathbf{F}_N^{-1} = (\delta_{\alpha\beta} \omega_N^{mn})_{\alpha,\beta=1,\dots,d}^{m,n \in \mathbb{Z}_N^d}$$

where $\delta_{\alpha\beta}$ is Kronecker delta and $\omega_N^{mn} = \exp\left(2\pi i \sum_{\alpha=1}^d \frac{m_{\alpha} n_{\alpha}}{N_{\alpha}}\right)$ with $\mathbf{m}, \mathbf{n} \in \mathbb{Z}^d$.

Definition 4.5 (nodal points, basis functions). Let $N \in \mathbb{N}^d$. For $\mathbf{n} \in \mathbb{Z}_N^d$, we define nodal points of the periodic unit cell $\mathbf{x}_N^n = \sum_{\alpha} \frac{2Y_{\alpha} n_{\alpha}}{N_{\alpha}} \boldsymbol{\epsilon}_{\alpha}$ and the Fourier and shape basis functions

$$\varphi_{\mathbf{n}}(\mathbf{x}) = \exp\left(\pi i \sum_{\alpha} \frac{n_{\alpha} x_{\alpha}}{Y_{\alpha}}\right), \quad \varphi_{N,\mathbf{n}}(\mathbf{x}) = \frac{1}{|N|_{\Pi}} \sum_{\mathbf{m} \in \mathbb{Z}_N^d} \omega_N^{-m\mathbf{n}} \varphi_{\mathbf{m}}(\mathbf{x}). \quad (4.3)$$

Lemma 4.6 (Properties of $\varphi_{\mathbf{m}}$ and $\varphi_{N,\mathbf{m}}$). Let $\mathbf{m}, \mathbf{n} \in \mathbb{Z}_N^d$, then

$$(\varphi_{\mathbf{m}}, \varphi_{\mathbf{n}})_{L^2_{\text{per}}} = \delta_{\mathbf{m}\mathbf{n}} \quad \varphi_{\mathbf{n}}(\mathbf{x}_N^{\mathbf{m}}) = \omega_N^{m\mathbf{n}} \quad (4.4a)$$

$$\varphi_{N,\mathbf{m}}(\mathbf{x}_N^{\mathbf{n}}) = \delta_{\mathbf{m}\mathbf{n}} \quad (\varphi_{N,\mathbf{m}}, \varphi_{N,\mathbf{n}})_{L^2_{\text{per}}} = \frac{\delta_{\mathbf{m}\mathbf{n}}}{|N|_{\Pi}} \quad (4.4b)$$

Definition 4.7 (Trigonometric polynomials). For $\mathbf{N} \in \mathbb{N}^d$, we define the spaces of trigonometric polynomials $\mathcal{T}_{\mathbf{N}}$, $\tilde{\mathcal{T}}_{\mathbf{N}}$ and their vector valued versions $\mathcal{T}_{\mathbf{N}}^d, \tilde{\mathcal{T}}_{\mathbf{N}}^d$ as

$$\mathcal{T}_{\mathbf{N}} = \left\{ \sum_{\mathbf{n} \in \mathbb{Z}_{\mathbf{N}}^d} \hat{\mathbf{v}}^{\mathbf{n}} \varphi_{\mathbf{n}} : \hat{\mathbf{v}}^{\mathbf{n}} \in \mathbb{C}, \hat{\mathbf{v}}^{\mathbf{n}} = \overline{\hat{\mathbf{v}}^{-\mathbf{n}}} \right\}, \quad \mathcal{T}_{\mathbf{N}}^d = \{ \mathbf{v} : v_{\alpha} \in \mathcal{T}_{\mathbf{N}} \}. \quad (4.5)$$

$$\tilde{\mathcal{T}}_{\mathbf{N}} = \left\{ \sum_{\mathbf{n} \in \mathbb{Z}_{\mathbf{N}}^d} \mathbf{v}^{\mathbf{n}} \varphi_{\mathbf{N}, \mathbf{n}} : \mathbf{v}^{\mathbf{n}} \in \mathbb{R} \right\}, \quad \tilde{\mathcal{T}}_{\mathbf{N}}^d = \{ \mathbf{v} : v_{\alpha} \in \tilde{\mathcal{T}}_{\mathbf{N}} \}. \quad (4.6)$$

where the index sets are given by

$$\mathbb{Z}_{\mathbf{N}}^d = \left\{ \mathbf{k} \in \mathbb{Z}^d : |k_{\alpha}| < \frac{N_{\alpha}}{2} \right\}, \quad \underline{\mathbb{Z}}_{\mathbf{N}}^d = \left\{ \mathbf{k} \in \mathbb{Z}^d : -\frac{N_{\alpha}}{2} \leq k_{\alpha} < \frac{N_{\alpha}}{2} \right\}. \quad (4.7)$$

Definition 4.8 (Interpolation projection). We define interpolation operator $Q_{\mathbf{N}} : C_{\text{per}}(\mathcal{Y}; \mathbb{R}^d) \rightarrow L_{\text{per}}^2(\mathcal{Y}; \mathbb{C}^d)$ as

$$Q_{\mathbf{N}}[f] = \sum_{\mathbf{m} \in \underline{\mathbb{Z}}_{\mathbf{N}}^d} f(\mathbf{x}_{\mathbf{N}}^{\mathbf{m}}) \varphi_{\mathbf{N}, \mathbf{m}}.$$

Lemma 4.9. Interpolation operator $Q_{\mathbf{N}}$ is projection and its image is $\tilde{\mathcal{T}}_{\mathbf{N}}^d$.

Definition 4.10. The operator $\mathcal{I}_{\mathbf{N}} : \tilde{\mathcal{T}}_{\mathbf{N}}^d \rightarrow \mathbb{R}^{d \times N}$ stocks the values of trigonometric polynomials at nodal points to a vector $\mathcal{I}_{\mathbf{N}}[\mathbf{v}_{\mathbf{N}}] = (\mathbf{v}_{\mathbf{N}, \alpha}(\mathbf{x}_{\mathbf{N}}^{\mathbf{n}}))_{\alpha=1, \dots, d}^{\mathbf{n} \in \underline{\mathbb{Z}}_{\mathbf{N}}^d}$.

Lemma 4.11. The operator $\mathcal{I}_{\mathbf{N}}$ from previous definition is isomorphism.

Remark 4.12 (Connection of representations). The trigonometric polynomial $\mathbf{v}_{\mathbf{N}} \in \tilde{\mathcal{T}}_{\mathbf{N}}^d$ can be uniquely expressed using both the Fourier coefficients and the function values at nodal points

$$\mathbf{v}_{\mathbf{N}} = \sum_{\mathbf{m} \in \underline{\mathbb{Z}}_{\mathbf{N}}^d} \mathbf{v}_{\mathbf{N}}(\mathbf{x}_{\mathbf{N}}^{\mathbf{m}}) \varphi_{\mathbf{N}, \mathbf{m}} = \sum_{\mathbf{n} \in \underline{\mathbb{Z}}_{\mathbf{N}}^d} \hat{\mathbf{v}}_{\mathbf{N}}(\mathbf{n}) \varphi_{\mathbf{n}} \quad (4.8)$$

with connection through the DFT $\hat{\mathbf{v}}_{\mathbf{N}} = \mathbf{F}_{\mathbf{N}} \mathbf{v}_{\mathbf{N}}$, where $\mathbf{v}_{\mathbf{N}} = \mathcal{I}_{\mathbf{N}}[\mathbf{v}_{\mathbf{N}}]$ and $\hat{\mathbf{v}}_{\mathbf{N}} = (\hat{\mathbf{v}}_{\mathbf{N}, \alpha}(\mathbf{m}))_{\alpha=1, \dots, d}^{\mathbf{m} \in \underline{\mathbb{Z}}_{\mathbf{N}}^d}$. Thus, space $\tilde{\mathcal{T}}_{\mathbf{N}}^d$ can be possibly characterized with the Fourier coefficients as $\tilde{\mathcal{T}}_{\mathbf{N}}^d = \{ \sum_{\mathbf{m} \in \underline{\mathbb{Z}}_{\mathbf{N}}^d} \hat{\mathbf{v}}_{\mathbf{N}}^{\mathbf{m}} \varphi_{\mathbf{m}} : \hat{\mathbf{v}}_{\mathbf{N}} \in \mathbf{F}_{\mathbf{N}}(\mathbb{R}^{d \times N}) \}$.

Remark 4.13. The trigonometric polynomials are real valued if the Fourier coefficients obey conjugate symmetry $\hat{\mathbf{v}}(\mathbf{n}) = \overline{\hat{\mathbf{v}}(-\mathbf{n})}$, $\mathbf{n} \in \mathbb{Z}^d$; from definition, it is valid for the trigonometric polynomials $\mathcal{T}_{\mathbf{N}}^d \subset L_{\text{per}}^2(\mathcal{Y}; \mathbb{R}^d)$.

The peculiar situation occurs for the space $\tilde{\mathcal{T}}_{\mathbf{N}}^d$. If \mathbf{N} is odd, both spaces coincide $\mathcal{T}_{\mathbf{N}}^d = \tilde{\mathcal{T}}_{\mathbf{N}}^d$ as the index sets do $\mathbb{Z}_{\mathbf{N}}^d = \underline{\mathbb{Z}}_{\mathbf{N}}^d$; generally, the inclusion $\mathcal{T}_{\mathbf{N}}^d \subseteq \tilde{\mathcal{T}}_{\mathbf{N}}^d$ holds. Unfortunately, the space $\tilde{\mathcal{T}}_{\mathbf{N}}^d$ fails to be real valued $\tilde{\mathcal{T}}_{\mathbf{N}}^d \not\subseteq L_{\text{per}}^2(\mathcal{Y}; \mathbb{R}^d)$ because the Fourier coefficients with frequencies $\mathbf{n} \in \underline{\mathbb{Z}}_{\mathbf{N}}^d \setminus \mathbb{Z}_{\mathbf{N}}^d$ miss opposite counterpart with frequencies $-\mathbf{n}$.

4.2.2 Galerkin approximation with numerical integration

Approximate solutions are gained from Galerkin approximation (GA) and Galerkin approximation with numerical integration (GAwNI); both methods can be defined either for weak formulation (2.5) or minimization formulation (2.2). In all cases, formulations are based on trial space \mathcal{E} or \mathcal{J} in the primal or dual formulation resp. The following remark clarifies their finite dimensional relatives.

Remark 4.14. *We confine the description to the problem for the odd number of discretization points \mathbf{N} for which trigonometric polynomials $\tilde{\mathcal{T}}_{\mathbf{N}}^d$ and $\mathcal{T}_{\mathbf{N}}^d$ coincide; thus it is possible to directly follow the results in Paper 5. The case of general number of discretization points can be found in Section 3.2 in Paper 6.*

Hence, we define a conforming finite dimensional spaces as $\mathcal{E}_{\mathbf{N}} := \mathcal{T}_{\mathbf{N}}^d \cap \mathcal{E}$ and $\mathcal{J}_{\mathbf{N}} := \mathcal{T}_{\mathbf{N}}^d \cap \mathcal{J}$ that satisfy the finite dimensional Helmholtz decomposition

$$\tilde{\mathcal{T}}_{\mathbf{N}}^d = \mathcal{U} \oplus^{\perp} \mathcal{E}_{\mathbf{N}} \oplus^{\perp} \mathcal{J}_{\mathbf{N}}.$$

Fully discrete spaces defined as $\mathbb{U}_{\mathbf{N}} = \mathcal{I}_{\mathbf{N}}[\mathcal{U}]$, $\mathbb{E}_{\mathbf{N}} = \mathcal{I}_{\mathbf{N}}[\mathcal{E}_{\mathbf{N}}]$, $\mathbb{J}_{\mathbf{N}} = \mathcal{I}_{\mathbf{N}}[\mathcal{J}_{\mathbf{N}}]$, with isomorphism from Def. 4.10, are their associates and thus satisfy an analogue

$$\mathcal{I}_{\mathbf{N}}[\tilde{\mathcal{T}}_{\mathbf{N}}^d] = \mathbb{R}^{d \times \mathbf{N}} = \mathbb{U}_{\mathbf{N}} \oplus^{\perp} \mathbb{E}_{\mathbf{N}} \oplus^{\perp} \mathbb{J}_{\mathbf{N}}. \quad (4.9)$$

Definition 4.15 (Galerkin approximation, Def. 3.20 in Paper 5). *Galerkin approximation of the auxiliary problems in Def. 2.6 states: for unitary macroscopic load $\boldsymbol{\epsilon}_{\alpha}$, find minimizer $\tilde{\mathbf{e}}_{\mathbf{N}}^{(\alpha)}$ satisfying*

$$(\mathbf{A}\tilde{\mathbf{e}}_{\mathbf{N}}^{(\alpha)}, \mathbf{v}_{\mathbf{N}})_{L_{\text{per}}^2} = -(\mathbf{A}\boldsymbol{\epsilon}_{\alpha}, \mathbf{v}_{\mathbf{N}})_{L_{\text{per}}^2}, \quad \forall \mathbf{v}_{\mathbf{N}} \in \mathcal{E}_{\mathbf{N}}. \quad (4.10)$$

Remark 4.16 (Solution and convergence). *The unique approximate solutions of Galerkin approximation are provided by Lax-Milgram lemma as material coefficients \mathbf{A} are uniformly elliptic and bounded. The convergence of the approximate solutions to the minimizers is provided by Cea's lemma*

$$\|\mathbf{e}^{(\alpha)} - \mathbf{e}_{\mathbf{N}}^{(\alpha)}\|_{L_{\text{per}}^2} \leq C \inf_{\mathbf{v}_{\mathbf{N}} \in \mathcal{E}_{\mathbf{N}}} \|\mathbf{e}^{(\alpha)} - \mathbf{v}_{\mathbf{N}}\|_{L_{\text{per}}^2} \quad (4.11)$$

together with density of trigonometric polynomials in space L_{per}^2 . Although, the convergence can be arbitrarily slow, more regular minimizers permit the order of convergence, see Sec. 3.4 in Paper 5.

Approximate solutions from Galerkin approximation behaves meaningfully, unfortunately, the linear systems determined by GA cannot be obtained in the closed form for general coefficients \mathbf{A} . Thus, the need for numerical integration arises.

Definition 4.17 (GAwNI). *Let the material coefficients \mathbf{A} be continuous. Galerkin approximation with numerical integration of the primal homogenization problem, Def. 2.1, states as: for arbitrary fixed macroscopic load $\mathbf{E} \in \mathbb{R}^d$, find discrete homogenized matrix $\mathbf{A}_{\text{eff}, \mathbf{N}}^{\text{FFTH}} \in \mathbb{R}^{d \times d}$ satisfying*

$$(\mathbf{A}_{\text{eff}, \mathbf{N}}^{\text{FFTH}} \mathbf{E}, \mathbf{E})_{\mathbb{R}^d} = \inf_{\mathbf{e}_{\mathbf{N}} \in \mathcal{E}_{\mathbf{N}}} (Q_{\mathbf{N}}[\mathbf{A}(\mathbf{E} + \mathbf{e}_{\mathbf{N}})], \mathbf{E} + \mathbf{e}_{\mathbf{N}})_{L_{\text{per}}^2}. \quad (4.12)$$

It is shown in Lem. 3.28 in Paper 5 and in Lem. 3.22 in Paper 6, that GAwNI is equivalent to the fully discrete formulation.

Definition 4.18 (Fully discrete formulation of GAwNI). *Find $\mathbf{A}_{\text{eff},N}^{\text{FFTH}} \in \mathbb{R}^{d \times d}$ such that for arbitrary fixed $\mathbf{E} \in \mathbb{R}^d$ holds*

$$(\mathbf{A}_{\text{eff},N}^{\text{FFTH}} \mathbf{E}, \mathbf{E})_{\mathbb{R}^d} = \frac{1}{|\mathbf{N}|_{\Pi}} \inf_{\mathbf{e}_N \in \mathbb{E}_N} (\mathbf{A}_N(\mathbf{E}_N + \mathbf{e}_N), \mathbf{E}_N + \mathbf{e}_N)_{\mathbb{R}^{d \times N}} \quad (4.13)$$

where $\mathbf{E}_N = \mathcal{I}_N[\mathbf{E}] \in \cup_N$ and $\mathbf{A}_N \in \mathbb{R}^{d \times d \times N \times N}$ with components assembled as $\mathbf{A}_{N,\alpha\beta}^{mn} = \mathbf{A}_{\alpha\beta}(\mathbf{x}_N^m) \delta_{mn}$ for α, β and $\mathbf{m}, \mathbf{n} \in \underline{\mathbb{Z}}_N^d$.

Remark 4.19 (Solution of GAwNI by Conjugate gradients). *In Papers 2 and 5, it is explained that the minimizers of fully discrete formulation can be obtained by Conjugate gradients that minimize a quadratic functional on a subspace $\mathbb{E}_N = \mathcal{I}_N[\mathcal{E}_N]$. It is equivalent to the solution of linear system $\mathbf{C}\mathbf{x} = \mathbf{b}$ defined for particular α as*

$$\underbrace{\mathbf{G}_0^1 \mathbf{A}_N}_{\mathbf{C}} \underbrace{\mathbf{e}_N^{(\alpha)}}_{\mathbf{x}} = - \underbrace{\mathbf{G}_0^1 \mathbf{A}_N \mathbf{E}_N^{(\alpha)}}_{\mathbf{b}}$$

for the initial approximation $\mathbf{x}_{(0)} \in \mathbb{E}_N$ from the appropriate subspace. Matrix \mathbf{G}_0^1 is an orthogonal projection on \mathbb{E}_N . It follows from continuous projection (4.2) on \mathcal{E} , see Sec. 3.2 in Paper 6 providing the scheme of subspaces for both of the primal and the dual formulations.

Remark 4.20 (Convergence of discrete minimizers). *In Paper 5, the convergence of the approximate solutions of GAwNI to the minimizers of homogenization problem (2.2) is provided for sufficiently regular material coefficients \mathbf{A} . It incorporates first Strang's lemma — the standard approach in the finite element method — and the estimates of interpolation operator, Def. (4.8), provided in [43] for one and two-dimensional setting, see also Sec. 3.2 in Paper 5 for d -dimensional setting. Numerical examples confirming the rates of convergence can be found in Paper 6 in Sec. 4.2.*

Remark 4.21. *Regularity of material coefficients \mathbf{A} principally influences the rates of convergence. For rough coefficients, arbitrary slow convergence can be observed for GA. However, GAwNI is even difficult to define as interpolation operator \mathbf{Q}_N requires function values — it requires some type of continuity. In Paper 5, we remedy this trouble by smoothing of the material coefficients.*

4.3 Guaranteed bounds by FFT-based homogenization

This section summarizes results obtained in Paper 6 for guaranteed bounds on the homogenized matrix, see summary in Sec. 4.3, that were introduced by Dvořák in [12, 13] for a scalar problem and later independently by Więkowski [55] for linear elasticity.

The upper-lower bounds are based on the primal and dual formulations of homogenization problem in Def. 2.1. In Sec. 4.3.1, the connection of the primal and dual formulations is investigated in fully discrete setting. Then in Sec. 4.3.2, the method for an effective evaluation of the upper-lower bounds is provided.

Galerkin approximation with numerical integration of the homogenization problem in Def. (2.1) leads to a fully discrete homogenization problem in both primal and dual setting

$$(\tilde{\mathbf{A}}_{\text{eff},\mathbf{N}}^{\text{FFTH}} \mathbf{E}, \mathbf{E})_{\mathbb{R}^d} = \frac{1}{|\mathbf{N}|_{\Pi}} \inf_{\mathbf{e}_N \in \mathbb{E}_N} (\mathbf{A}_N(\mathbf{E}_N + \mathbf{e}_N), \mathbf{E}_N + \mathbf{e}_N)_{\mathbb{R}^d \times \mathbf{N}}, \quad (4.14a)$$

$$((\tilde{\mathbf{A}}_{\text{eff},\mathbf{D},\mathbf{N}}^{\text{FFTH}})^{-1} \mathbf{J}, \mathbf{J})_{\mathbb{R}^d} = \frac{1}{|\mathbf{N}|_{\Pi}} \inf_{\mathbf{j}_N \in \mathbb{J}_N} (\mathbf{A}_N^{-1}(\mathbf{J}_N + \mathbf{j}_N), \mathbf{J}_N + \mathbf{j}_N)_{\mathbb{R}^d \times \mathbf{N}}, \quad (4.14b)$$

compare to Def. 4.18. Minimizers $\mathbf{e}_N^{(\alpha)} := \mathcal{I}_N^{-1}[\mathbf{e}_N^{(\alpha)}]$ and $\mathbf{j}_N^{(\alpha)} := \mathcal{I}_N^{-1}[\mathbf{j}_N^{(\alpha)}]$ are then used to evaluate the upper-lower bounds of the homogenized matrix.

4.3.1 Connection of primal and dual formulation

The connection between formulations in (4.14) is summarized for odd number of discretization points.

Theorem 4.22 (Primal-dual formulation for odd number of discretization points, Corollary 3.27 in Paper 6). *Let the number of discretization points $\mathbf{N} \in \mathbb{N}^d$ be odd. Then, the fully discrete formulations satisfy:*

1. Both primal and dual homogenized matrices coincide $\mathbf{A}_{\text{eff},\mathbf{N}}^{\text{FFTH}} = \mathbf{A}_{\text{eff},\mathbf{D},\mathbf{N}}^{\text{FFTH}}$.
2. Primal and dual discrete minimizers $\tilde{\mathbf{e}}_N^{(\alpha)}$ and $\tilde{\mathbf{j}}_N^{(\alpha)}$ are related as

$$\epsilon_\beta + \tilde{\mathbf{j}}_N^{(\beta)} = \mathbf{A}_N \sum_{\alpha} E_{\alpha}(\epsilon_{\alpha} + \tilde{\mathbf{e}}_N^{(\alpha)}) \quad (4.15)$$

where $\mathbf{E} = (\mathbf{A}_{\text{eff},\mathbf{N}}^{\text{FFTH}})^{-1} \epsilon_{\beta}$.

This enables to avoid the solution of the dual formulation in order to obtain the dual minimizers. The proof is the consequence of perturbation duality theorem [14], that is mentioned for a discrete setting in Lem. 3.26 in Paper 6, and discrete version of the Helmholtz decomposition (4.9).

Unfortunately, the previous theorem fails to hold for the general number of discretization points as the discrete version of the Helmholtz decomposition (4.9) is no longer valid. A general theorem is provided in Cor. 3.27 in Paper 6.

4.3.2 Calculation of bounds

The calculation of the upper-lower bounds of the homogenized matrix consists of the integral evaluation of type $(\mathbf{A}\mathbf{e}_N^{(\alpha)}, \mathbf{e}_N^{(\beta)})_{L_{\text{per}}^2}$ occurring in Def. 2.12. Generally, the integral cannot be evaluated in a closed form because of non-specific material coefficients. The idea is to adjust the material coefficients to calculate the integrals accurately and efficiently and simultaneously keep the upper-lower bounds structure.

For an easier orientation among various homogenized matrices, we refer to their scheme in Fig. 2. The matrices \mathbf{A}_{eff} , $\underline{\mathbf{A}}_{\text{eff},\mathbf{N}}$, $\overline{\mathbf{A}}_{\text{eff},\mathbf{N}}$, $\mathbf{A}_{\text{eff},\mathbf{N}}$, and \mathbf{D}_N already introduced in Sec. 4.3 are in no relation to matrices $\mathbf{A}_{\text{eff},\mathbf{N}}^{\text{FFTH}}$, $\mathbf{A}_{\text{eff},\mathbf{D},\mathbf{N}}^{\text{FFTH}}$ from fully discrete formulation (4.14) because of the variational crime caused by numerical integration in Def. 4.17.

$$\begin{array}{ccccccc}
& & & \tilde{\mathbf{A}}_{\text{eff},D,N}^{\text{lin},M} & & \tilde{\mathbf{A}}_{\text{eff},N}^{\text{lin},M} & \\
& & & \Downarrow & & \Downarrow & \\
0 & \preceq & \underline{\underline{\mathbf{A}}}_{\text{eff},N}^{\text{con},M} & \preceq & \underline{\mathbf{A}}_{\text{eff},N} & \preceq & \overline{\mathbf{A}}_{\text{eff},N} & \preceq & \overline{\overline{\mathbf{A}}}_{\text{eff},N}^{\text{con},M} \\
& & & \parallel & & \parallel & \\
& & & \mathbf{A}_{\text{eff},N} - \mathbf{D}_N & \preceq & \mathbf{A}_{\text{eff},N} & \preceq & \mathbf{A}_{\text{eff},N} + \mathbf{D}_N \\
& & & \not\preceq & & \not\preceq & \\
& & & \mathbf{A}_{\text{eff},D,N}^{\text{FFTH}} & \underset{\preceq}{\text{equality if } N \text{ is odd}} & \mathbf{A}_{\text{eff},N}^{\text{FFTH}} &
\end{array}$$

Figure 2: The overview of homogenized material bounds

In Paper 6 in Sec. 3.5, we introduce the approximation of upper-lower bounds $\tilde{\mathbf{A}}_{\text{eff},D,N}^{\text{lin},M}$, $\tilde{\mathbf{A}}_{\text{eff},N}^{\text{lin},M}$ based on piecewise bilinear material coefficients. Next, piecewise constant material coefficients defined in a way to guaranty bounds produce upper-lower bounds of homogenized matrix, i.e. $\underline{\underline{\mathbf{A}}}_{\text{eff},N}^{\text{con},M}$, $\overline{\overline{\mathbf{A}}}_{\text{eff},N}^{\text{con},M}$. The next lemma shows sufficient condition to guaranty bounds with adjusted material coefficients.

Lemma 4.23 (Sufficient condition for the upper-lower bound structure, Lem. 3.31 in Paper 6). *Let $\mathbf{A} \in L_{\text{per}}^{\infty}(\mathcal{Y}; \mathbb{R}^{d \times d})$ be material coefficients and $\overline{\mathbf{A}}, \underline{\mathbf{A}} \in L_{\text{per}}^{\infty}(\mathcal{Y}; \mathbb{R}^{d \times d})$ its upper and lower approximations satisfying*

$$\underline{\mathbf{A}}(\mathbf{x}) \preceq \mathbf{A}(\mathbf{x}) \preceq \overline{\mathbf{A}}(\mathbf{x}), \quad \text{for almost all } \mathbf{x} \in \mathcal{Y}. \quad (4.16)$$

Let $\tilde{\mathbf{e}}_N^{(\alpha)} \in \mathcal{E}_N$ and $\tilde{\mathbf{j}}_N^{(\alpha)} \in \mathcal{J}_N$ be unitary minimizers for material coefficients \mathbf{A} , cf. Def. 2.6. Then matrices $\overline{\overline{\mathbf{A}}}_{\text{eff}}, \underline{\underline{\mathbf{A}}}_{\text{eff}} \in \mathbb{R}^{d \times d}$, defined as

$$(\overline{\overline{\mathbf{A}}}_{\text{eff},N})_{\alpha\beta} = (\overline{\mathbf{A}}(\boldsymbol{\epsilon}_{\beta} + \tilde{\mathbf{e}}_N^{(\beta)}), \boldsymbol{\epsilon}_{\alpha} + \tilde{\mathbf{e}}_N^{(\alpha)})_{L_{\text{per}}^2}, \quad (4.17a)$$

$$(\underline{\underline{\mathbf{A}}}_{\text{eff},N}^{-1})_{\alpha\beta} = (\underline{\mathbf{A}}^{-1}(\boldsymbol{\epsilon}_{\beta} + \tilde{\mathbf{j}}_N^{(\beta)}), \boldsymbol{\epsilon}_{\alpha} + \tilde{\mathbf{j}}_N^{(\alpha)})_{L_{\text{per}}^2}, \quad (4.17b)$$

comply with the upper-lower bound structure, i.e.

$$\underline{\underline{\mathbf{A}}}_{\text{eff},N} \preceq \underline{\mathbf{A}}_{\text{eff},N} \preceq \mathbf{A}_{\text{eff}} \preceq \overline{\mathbf{A}}_{\text{eff},N} \preceq \overline{\overline{\mathbf{A}}}_{\text{eff},N}.$$

A special situation occurs when the material coefficients are expressed as the linear combination of some functions placed at nodal points of regular grid. Then the integrals required for evaluating upper-lower bounds can be calculated by FFT algorithm, see Eq. (4.19) in next lemma and Lem. 3.33 in Paper 6.

Lemma 4.24 (Calculation of homogenized matrices, Lem. 3.32 in Paper 6). *Let $\mathbf{u}_N, \mathbf{v}_N \in \mathcal{T}_N^d$ be trigonometric polynomials and $\mathbf{A}_M \in L_{\text{per}}^{\infty}(\mathcal{Y}; \mathbb{R}_{\text{spd}}^{d \times d})$ for $M \in \mathbb{N}^d$ be function explicitly expressed as*

$$\mathbf{A}_M(\mathbf{x}) = \sum_{\mathbf{n} \in \mathbb{Z}_M^d} \psi(\mathbf{x} + \mathbf{x}_M^{\mathbf{n}}) \mathbf{A}^{\mathbf{n}}, \quad \mathbf{x} \in \mathcal{Y}$$

where $\psi \in L_{\text{per}}^{\infty}(\mathcal{Y}; \mathbb{R})$ is some basis function and $\mathbf{A} \in \mathbb{R}^{d \times d \times M}$. Then the integrals of the type occurring in Eq. (4.17) can be calculated as

$$(\mathbf{A}_M \mathbf{u}_N, \mathbf{v}_N)_{L_{\text{per}}^2} = \frac{1}{|\mathcal{Y}|_d} \sum_{\alpha, \beta} \sum_{\mathbf{m} \in \mathbb{Z}_{2N}^d} w(\mathbf{m}) \hat{u}_{N, \alpha, \beta}^{\mathbf{m}} \hat{v}_{\alpha, \beta}^{\mathbf{m}} \quad (4.18)$$

where the integration weight $w(\mathbf{m})$ is defined as $w(\mathbf{m}) := \int_{\mathcal{Y}} \psi(\mathbf{x}) \varphi_{\mathbf{m}}(\mathbf{x})$ and factors $\widehat{\mathbf{u}}_{N,\beta,\alpha}^{\mathbf{m}}$, $\widehat{\mathbf{A}}_{\alpha\beta}^{\mathbf{m}}$ are defined as

$$\widehat{\mathbf{u}}_{N,\beta,\alpha}^{\mathbf{m}} = \frac{1}{2|\mathbf{N}|_{\Pi}} \sum_{\mathbf{k} \in \mathbb{Z}_{2N}^d} \mathbf{u}_{N,\beta}(\mathbf{x}_{2N}^{\mathbf{k}}) \mathbf{v}_{N,\alpha}(\mathbf{x}_{2N}^{\mathbf{k}}) \omega_{2N}^{-\mathbf{m}\mathbf{k}} \quad (4.19a)$$

$$\widehat{\mathbf{A}}_{\alpha\beta}^{\mathbf{m}} = \sum_{\mathbf{n} \in \mathbb{Z}_M^d} \mathbf{A}_{\alpha\beta}^{\mathbf{n}} \omega_M^{-\mathbf{m}\mathbf{n}} \quad (4.19b)$$

Investigation of previous lemma reveals that functions ψ are chosen to have analytical expression of integral weights $w(\mathbf{m}) := \int_{\mathcal{Y}} \psi(\mathbf{x}) \varphi_{\mathbf{m}}(\mathbf{x})$. A reasonable choice is a constant and a bilinear function

$$\chi_N(\mathbf{x}) = \begin{cases} 1, & |x_{\alpha}| < \frac{Y_{\alpha}}{M_{\alpha}} \text{ for all } \alpha \\ 0, & \text{otherwise} \end{cases}, \quad \text{tri}_N(\mathbf{x}) = \prod_{\alpha} \max\{1 - |\frac{x_{\alpha} M_{\alpha}}{2Y_{\alpha}}|, 0\},$$

leading to the weights

$$w_N^0(\mathbf{m}) := \prod_{\alpha} \frac{2Y_{\alpha}}{M_{\alpha}} \text{sinc}\left(\frac{m_{\alpha}}{M_{\alpha}}\right) \quad w_N^1(\mathbf{m}) := \prod_{\alpha} \frac{2Y_{\alpha}}{M_{\alpha}} \text{sinc}^2\left(\frac{m_{\alpha}}{M_{\alpha}}\right)$$

with function $\text{sinc}(x) := \begin{cases} 1, & x = 0 \\ \frac{\sin(\pi x)}{\pi x}, & x \neq 0 \end{cases}$. Another example of the function with the analytical expression of the weight is a circle function with the weight based on the Bessel function.

More details about calculation of upper-lower bounds can be found in Sec. 3.5 in Paper 6.

4.4 Numerical experiments

In this section, numerical experiments presented in Papers 1,2, and 6 are summarized. Sec. 4.4.1 is dedicated to the acceleration of the FFT-based homogenization by Conjugate gradients, while Sec. 4.4.2 is dedicated to the guaranteed bounds of the homogenized matrix.

4.4.1 Acceleration by Conjugate gradients

In this section, the theoretical result about the acceleration of the original FFT-based homogenization by Conjugate gradients is validated. We show that Conjugate gradients are independent on reference conductivity \mathbf{A}^0 , the parameter of Lippmann-Schwinger equation. Moreover, both methods, the original and accelerated one, are compared in terms of computational times.

From Paper 2, a three-dimensional electric conduction in a cubic periodic unit cell $\mathcal{Y} = \prod_{\alpha=1}^3 (-\frac{1}{2}, \frac{1}{2})$ is chosen as a model problem. The conductivity parameters are defined as

$$\mathbf{A}(\mathbf{x}) = \begin{cases} \rho \mathbf{I}, & \|\mathbf{x}\|_2 < (\frac{3}{16\pi})^{\frac{1}{3}} \\ \begin{pmatrix} 1 & 0.2 & 0.2 \\ 0.2 & 1 & 0.2 \\ 0.2 & 0.2 & 1 \end{pmatrix}, & \text{otherwise} \end{cases}$$

where $\rho > 0$ denotes the contrast of phase conductivities; it represents a two-phase medium with spherical inclusions of 25% volume fraction. We consider the macroscopic field $\mathbf{E} := \boldsymbol{\epsilon}_1$ and discretize the unit cell with $\mathbf{N} = [n, n, n]$ nodes¹⁰. The conductivity of the homogeneous reference medium $\mathbf{A}^0 \in \mathbb{R}^{d \times d}$ is parametrized as

$$\mathbf{A}^0 = \lambda \mathbf{I}, \quad \lambda = 1 - \omega + \rho\omega, \quad (4.20)$$

where $\omega \approx 0.5$ delivers the optimal convergence of the original Moulinec-Suquet Fast-Fourier Transform-based Homogenization (FFTH) algorithm [32].

We first investigate the sensitivity of Conjugate Gradients (CG) to the choice of the reference medium. The results appear in Fig. 3(a) plotting the relative number of iterations for CG against the conductivity of the reference medium parametrized by ω , recall Eq. (4.20). As expected, CG solver achieve a significant improvement over FFTH method as it requires about 40% iterations of FFTH for a mildly-contrasted composite down to 4% for $\rho = 10^3$. The minor differences visible especially for $\rho = 10^3$ can be therefore attributed to accumulation of round-off errors.

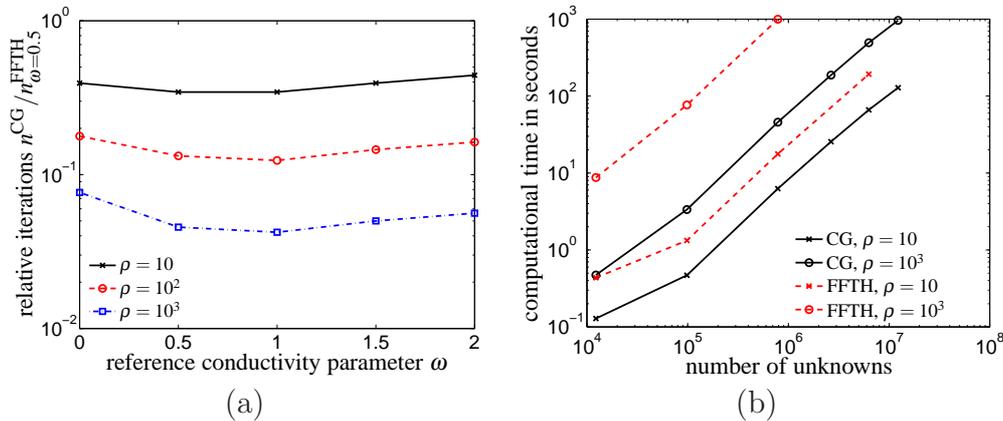


Figure 3: (a) Relative number of iterations as a function of the reference medium parameter ω and (b) computational time as a function of the number of unknowns.

In Fig. 3(b), we present the total computational time¹¹ as a function of the number of degrees of freedom and the phase ratio ρ . The results confirm that the computational times scales linearly with the increasing number of degrees of freedom for both schemes for fixed phase ratio ρ [57]. The ratio of the computational time for CG and FFTH algorithms remains almost constant, which indicates that the cost of a single iteration of CG and FFTH method is comparable.

In addition, the memory requirements of both schemes are also comparable. This aspect emphasized the major advantage of the short-recurrence CG-based scheme over alternative schemes for non-symmetric systems, such as GMRES. Finally, we note that finer discretizations can be treated by a straightforward parallel implementation.

¹⁰In particular, n was taken consequently as 16, 32, 64, 128 and 160 leading up to $3 \cdot 160^3 \doteq 12.2 \times 10^6$ unknowns

¹¹The problem was solved with a MATLAB[®] in-house code on a machine Intel[®] Core[™]2 Duo 3 GHz CPU, 3.28 GB computing memory with Debian linux 5.0 operating system.

4.4.2 Guaranteed bounds

In this section, the properties of the upper-lower bounds of the homogenized matrices are validated. We consider 2-dimensional model problem from Paper 6 with material coefficients defined on a periodic unit cell $\mathcal{Y} = \Pi_{\alpha=1}^2(-1, 1) \subset \mathbb{R}^2$ as

$$\mathbf{A}(\mathbf{x}) = \mathbf{I}[1 + 10f(\mathbf{x})], \quad \mathbf{x} \in \mathcal{Y},$$

where $\mathbf{I} \in \mathbb{R}^{d \times d}$ is identity matrix and $f : \mathcal{Y} \rightarrow \mathbb{R}$ is a scalar nonnegative function defined explicitly as

$$f(\mathbf{x}) = \begin{cases} 1, & \|\mathbf{x}\|_{\infty} < \frac{3}{4} \\ 0, & \text{otherwise} \end{cases}.$$

The problem is discretized with odd number of discretization points $\mathbf{N} = (n, n)$ where $n \in \{5, 15, 45, 135, 405, 1215\}$.

Fig. 4(a) shows the periodic unit cell with the interface between phases and the nodal points sets, $\{\mathbf{x}_{\mathbf{N}}^n \in \mathcal{Y} : \mathbf{n} \in \mathbb{Z}_{\mathbf{N}}^d\}$, for particular \mathbf{N} . Next in Fig. 4(b), we demonstrate the properties of the homogenized matrices for their particular diagonal component. The inequality $\underline{\mathbf{A}}_{\text{eff}, \mathbf{N}} \preceq \overline{\mathbf{A}}_{\text{eff}, \mathbf{N}}$ stated in (2.11) is satisfied and the error, difference between them, is approaching zero supporting Lem. 2.14.

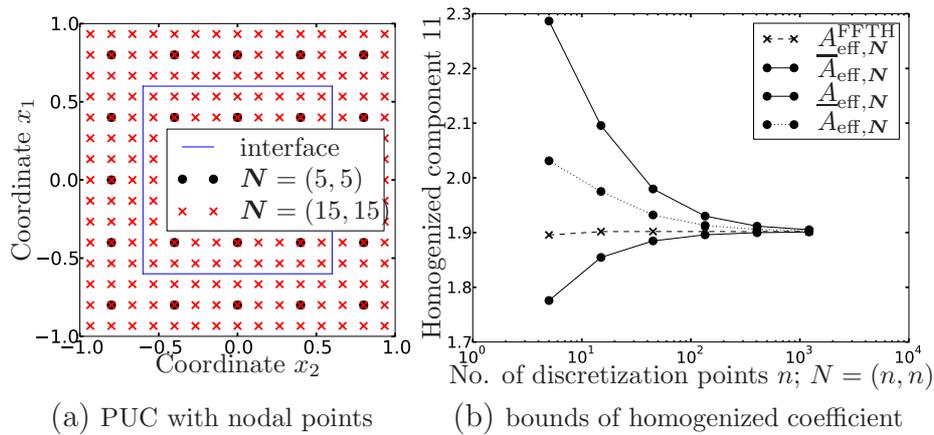


Figure 4: Periodic unit cell with nodal points and guaranteed bounds of homogenized material coefficient

In this case, approximate homogenized matrix $\underline{\mathbf{A}}_{\text{eff}, \mathbf{N}}^{\text{FFTH}}$ resembles the real homogenized coefficients $\underline{\mathbf{A}}_{\text{eff}}$ properly even for small \mathbf{N} compared to the mean value $\overline{\mathbf{A}}_{\text{eff}, \mathbf{N}} = \frac{1}{2}(\underline{\mathbf{A}}_{\text{eff}, \mathbf{N}} + \overline{\mathbf{A}}_{\text{eff}, \mathbf{N}})$ that overestimates. It is the consequence of good approximation of inclusion topology by nodal points, see Fig. 4(a); the interface lies exactly between the nodes. Generally, the approximate homogenized matrix $\underline{\mathbf{A}}_{\text{eff}, \mathbf{N}}^{\text{FFTH}}$ can be located either over or under the bounds — for a more detail, see Sec. 4.1 in Paper 6.

4.5 Applications

In this section, applications of FFT-based homogenization to linear elasticity are discussed according to Papers 3 and 4.

Linear elasticity is certainly more complicated than the scalar problem presented so far, however, the similar structure is observed. The complexity consists of more complicated material law between second order tensors, strain and stress, through fourth order stiffness tensor. Moreover, Green function and derived projection operator to admissible fields — analogue to (2.9) and (4.1) — are more complicated as it naturally splits into two parts, one corresponding to volumetric fields and second to deviatoric fields.

Paper 3 deals with FFT-based homogenization in the framework of representative volume element to produce homogenized stiffness matrix of cement paste, gypsum, and aluminum alloy. A grid nanoindentation is used for the determination of phase properties in grid points at microscale; the direct utilization of grid data simplifies the numerical evaluation and contributes to its efficiency. The method is compared to some simple analytical homogenization procedures with material phases obtained by a statistical deconvolution.

Paper 4 describes a technique to homogenize highly porous aluminium foam — porosity causes the violation of uniform ellipticity. Although, some modifications of FFT-based homogenization were proposed [27, 29] to overcome this difficulty, no rigorous convergence theory has been provided yet. Thus, the FFT-based homogenization is used as a part of two step homogenization, particularly, it is incorporated at a lower scale composed of aluminum melt with admixtures. A higher scale containing significant volume fraction of air voids exceeding 90% is homogenized using two-dimensional Finite element method.

5 Conclusions

The most important results addressing the questions in Sec. 3 are summarized:

1. It has been shown that Lippmann-Schwinger equation is equivalent, in the sense the unique solution coincide, to the corresponding weak formulation for an arbitrary parameter $\mathbf{A}^0 \in \mathbb{R}_{\text{spd}}^{d \times d}$.
2. The discretization of weak formulation (2.5) has been proposed to reproduce the original Moulinec-Suquet numerical algorithm, Rem 2.11; it can be newly considered as the Finite element method with trigonometric polynomials basis functions. Convergence of approximate solutions to continuous one has been proven. Moreover, the successful application of Conjugate gradients to non-symmetric linear system has been explained.
3. The method for guaranteed bounds of the homogenized matrix has been used for FFT-based homogenization; FFT algorithm can be utilized for the evaluation. The solution of the dual formulation, that is required for lower bound, can be avoided for the odd number of discretization points.

The presented results clarify FFT-based homogenization in a way to be further analyzed by standard mathematical instruments. Main areas of investigation and employment of results are:

- preconditioning of linear system coming from discretization,

- providing theory for linear elasticity,
- providing theory for porous materials,
- generalization to another physical problems (Stokes problem),
- validation of other homogenization methods with the help of reliable guaranteed bounds.

6 References

- [1] N. S. Bakhvalov and A. V. Knyazev, *Efficient computation of averaged characteristics of composites of a periodic structure of essentially different materials*, Soviet mathematics - Doklady **42** (1991), no. 1, 57–62.
- [2] T. Belytschko, T.J.R. Hughes, N. Patankar, C.T. Herakovich, and E.C. Bakis, *Research directions in computational and composite mechanics*, Tech. report, United States National Committee on Theoretical and Applied Mechanics, 2007.
- [3] A. Bensoussan, G. Papanicolau, and J.L. Lions, *Asymptotic analysis for periodic structures*, vol. 5, North Holland, 1978.
- [4] D. Braess, *Finite elements: Theory, fast solvers, and applications in solid mechanics*, Cambridge University Press, 2001.
- [5] H. Brézis, *Functional analysis, Sobolev spaces and partial differential equations*, Universitext Series, Springer, 2010.
- [6] S. Brisard and L. Dormieux, *FFT-based methods for the mechanics of composites: A general variational framework*, Computational Materials Science **49** (2010), no. 3, 663–671.
- [7] ———, *Combining Galerkin approximation techniques with the principle of Hashin and Shtrikman to derive a new FFT-based numerical method for the homogenization of composites*, Computer Methods in Applied Mechanics and Engineering (2012).
- [8] B. Budiansky, *On the elastic moduli of some heterogeneous materials*, Journal of the Mechanics and Physics of Solids **13** (1965), no. 4, 223–227.
- [9] C. S. Burrus and T. W. Parks, *DFT/FFT and convolution algorithms*, Wiley, New York, 1985.
- [10] P.G. Ciarlet, *Handbook of numerical analysis: Finite element methods*, vol. II, Elsevier Science Publishers B.V. (North-Holland), 1991.
- [11] D. Cioranescu and P. Donato, *An introduction to homogenization*, Oxford Lecture Series in Mathematics and Its Applications, Oxford University Press, 1999.
- [12] J. Dvořák, *Optimization of composite materials*, Master’s thesis, The Charles University in Prague, June 1993.

- [13] ———, *A reliable numerical method for computing homogenized coefficients*, Tech. report, available at CiteSeerX <http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.45.1190>, 1995.
- [14] I. Ekeland and R. Temam, *Convex analysis and variational problems*, SIAM, 1976.
- [15] 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** (1999), no. 1, 41–47.
- [16] R. Hill, *A self-consistent mechanics of composite materials*, Journal of the Mechanics and Physics of Solids **13** (1965), no. 4, 213–222.
- [17] H. Hoang-Duc and G. Bonnet, *Effective properties of viscoelastic heterogeneous periodic media: an approximate solution accounting for the distribution of heterogeneities.*, Mechanics of Materials (2012).
- [18] V.V. Jikov, S.M. Kozlov, and O.A. Oleinik, *Homogenization of differential operators and integral functionals*, Springer-Verlag, 1994.
- [19] L.C.Evans, *Partial differential equations*, vol. 19, American Mathematical Society, 2000.
- [20] H. Le Quang, T.L. Phan, and G. Bonnet, *Effective thermal conductivity of periodic composites with highly conducting imperfect interfaces*, International Journal of Thermal Sciences **50** (2011), no. 8, 1428–1444.
- [21] R.A. Lebensohn, A.K. Kanjarla, and P. Eisenlohr, *An elasto-viscoplastic formulation based on fast fourier transforms for the prediction of micromechanical fields in polycrystalline materials*, International Journal of Plasticity (2011).
- [22] R.A. Lebensohn, A.D. Rollett, and P. Suquet, *Fast fourier transform-based modeling for the determination of micromechanical fields in polycrystals*, JOM Journal of the Minerals, Metals and Materials Society **63** (2011), no. 3, 13–18.
- [23] S.B. Lee, RA Lebensohn, and A.D. Rollett, *Modeling the viscoplastic micromechanical response of two-phase materials using Fast Fourier Transforms*, International Journal of Plasticity **27** (2011), no. 5, 707–727.
- [24] J. Li, S. Meng, X. Tian, F. Song, and C. Jiang, *A non-local fracture model for composite laminates and numerical simulations by using the FFT method*, Composites Part B: Engineering (2011).
- [25] J. Li, X.X. Tian, and R. Abdelmoula, *A damage model for crack prediction in brittle and quasi-brittle materials solved by the FFT method*, International journal of fracture (2012), 1–12.
- [26] J. Lukeš and J. Malý, *Measure and integral*, Matfyzpress Prague, 1995.

- [27] J. C. Michel, H. Moulinec, and P. Suquet, *A computational method based on augmented Lagrangians and fast Fourier transforms for composites with high contrast*, CMES-Computer Modeling in Engineering & Sciences **1** (2000), no. 2, 79–88.
- [28] ———, *A computational scheme for linear and non-linear composites with arbitrary phase contrast*, International Journal for Numerical Methods in Engineering **52** (2001), no. 1–2, 139–160.
- [29] V. Monchiet and G. Bonnet, *A polarization-based FFT iterative scheme for computing the effective properties of elastic composites with arbitrary contrast*, International Journal for Numerical Methods in Engineering **89** (2012), no. 11, 1419–1436.
- [30] V. Monchiet, G. Bonnet, and G. Lauriat, *A FFT-based method to compute the permeability induced by a Stokes slip flow through a porous medium*, Comptes Rendus Mécanique **337** (2009), no. 4, 192–197.
- [31] T. Mori and K. Tanaka, *Average stress in matrix and average elastic energy of materials with misfitting inclusions*, Acta metallurgica **21** (1973), no. 5, 571–574.
- [32] H. Moulinec and P. Suquet, *A fast numerical method for computing the linear and nonlinear mechanical properties of composites*, Comptes rendus de l’Académie des sciences. Série II, Mécanique, physique, chimie, astronomie **318** (1994), no. 11, 1417–1423.
- [33] ———, *A numerical method for computing the overall response of nonlinear composites with complex microstructure*, Computer Methods in Applied Mechanics and Engineering **157** (1997), no. 1–2, 69–94.
- [34] ———, *Comparison of FFT-based methods for computing the response of composites with highly contrasted mechanical properties*, Physica B: Condensed Matter **338** (2003), no. 1–4, 58–60.
- [35] J. Nečas, *Direct methods in the theory of elliptic equations*, Springer, 2012.
- [36] J. Novák, A. Kučerová, and J. Zeman, *Microstructural enrichment functions based on stochastic Wang tilings*, arXiv preprint arXiv:1110.4183 (2011).
- [37] J. Němeček, V. Králík, and J. Vondřejc, *Micromechanical analysis of heterogeneous structural materials*, Cement and Concrete Composites (2012).
- [38] J. Němeček, V. Králík, and J. Vondřejc, *A two-scale micromechanical model for aluminium foam based on results from nanoindentation*, (2012), Submitted.
- [39] J.T. Oden, T. Belytschko, I. Babuška, and T.J.R. Hughes, *Research directions in computational mechanics*, Computer Methods in Applied Mechanics and Engineering **192** (2003), no. 7, 913–922.

-
- [40] O.A. Oleinik, A.S. Shamaev, and G.A. Yosifian, *Mathematical problems in elasticity and homogenization*, Studies in Mathematics and Its Applications, North-Holland, 1992.
- [41] W. Rudin, *Functional analysis*, McGraw-Hill, Inc., New York, 1991.
- [42] ———, *Real and complex analysis*, Tata McGraw-Hill Education, 2006.
- [43] J. Saranen and G. Vainikko, *Periodic integral and pseudodifferential equations with numerical approximation*, Springer Monographs Mathematics, 2000.
- [44] P. Suquet, H. Moulinec, O. Castelnau, M. Montagnat, N. Lahellec, F. Grennerat, P. Duval, and R. Brenner, *Multi-scale modeling of the mechanical behavior of polycrystalline ice under transient creep*, *Procedia IUTAM* **3** (2012), 64–78.
- [45] J. Sýkora, J. Zeman, and M. Šejnoha, *Selected topics in homogenization of transport processes in historical masonry structures*, arXiv preprint arXiv:1204.6199 (2012).
- [46] L. Tartar, *The general theory of homogenization: a personalized introduction*, vol. 7, Springer, 2009.
- [47] G. Vainikko, *Fast solvers of the Lippmann-Schwinger equation*, Direct and Inverse Problems of Mathematical Physics (R. P. Gilbert, J. Kajiwara, and Y. S. Xu, eds.), International Society for Analysis, Applications and Computation, vol. 5, Kluwer Academic Publishers, Dordrecht, The Netherlands, 2000, pp. 423–440.
- [48] R. Valenta and M. Šejnoha, *Two-step homogenization of asphalt mixtures*, *Bulletin of Applied Mechanics* **4** (2008), no. 14, 61–66.
- [49] V. Vinogradov and G. W. Milton, *An accelerated FFT algorithm for thermoelastic and non-linear composites*, *International Journal for Numerical Methods in Engineering* **76** (2008), no. 11, 1678–1695.
- [50] J. Vondřejc, J. Zeman, and I. Marek, *Analysis of a Fast Fourier Transform based method for modeling of heterogeneous materials*, *Lecture Notes in Computer Science* **7116** (2012), 512–522.
- [51] ———, *FFT-based finite element method for homogenization*, (2013), In preparation.
- [52] ———, *Guaranteed bounds of effective material properties using FFT-based FEM*, (2013), In preparation.
- [53] V. Šmilauer and Z.P. Bažant, *Identification of viscoelastic C-S-H behavior in mature cement paste by FFT-based homogenization method*, *Cement and Concrete Research* **40** (2010), no. 2, 197–207.
- [54] V. Šmilauer and Z. Bittnar, *Microstructure-based micromechanical prediction of elastic properties in hydrating cement paste*, *Cement and Concrete Research* **36** (2006), no. 9, 1708–1718.

-
- [55] Z. Wiećkowski, *Dual finite element methods in mechanics of composite materials*, Journal of Theoretical and Applied Mechanics **2** (1995), no. 33, 233–252.
- [56] J. Yvonnet, *A fast method for solving microstructural problems defined by digital images: a space lippmann–schwinger scheme*, International Journal for Numerical Methods in Engineering (2012).
- [57] 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** (2010), no. 21, 8065–8071.
- [58] J. Zeman and M. Šejnoha, *Numerical evaluation of effective elastic properties of graphite fiber tow impregnated by polymer matrix*, Journal of the Mechanics and Physics of Solids **49** (2001), no. 1, 69–90.

7 List of thesis papers

Paper 1 : ISI journal paper

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** (2010), no. 21, 8065–8071.

Paper 2 : Conference paper in SCOPUS

J. Vondřejc, J. Zeman, and I. Marek, *Analysis of a Fast Fourier Transform based method for modeling of heterogeneous materials*, Lecture Notes in Computer Science **7116** (2012), 512–522.

Paper 3 : ISI journal paper (In Press):

J. Němeček, V. Králík, and J. Vondřejc, *Micromechanical analysis of heterogeneous structural materials*, Cement and Concrete Composites (2012).

Paper 4 : Peer review in ISI journal

J. Němeček, V. Králík, and J. Vondřejc, *A two-scale micromechanical model for aluminium foam based on results from nanoindentation*, (2012).

Paper 5 : In preparation for ISI journal

J. Vondřejc, J. Zeman, and I. Marek, *FFT-based finite element method homogenization*, (2013), In preparation.

Paper 6 : In preparation for ISI journal

J. Vondřejc, J. Zeman, and I. Marek, *Guaranteed bounds of effective material properties using FFT-based FEM*, (2013), In preparation.

Part II

Paper 1

Authors:

Jan Zeman, Jaroslav Vondřejc, Jan Novák, and Ivo Marek

Title:

Accelerating a FFT-based solver for numerical homogenization of periodic media by conjugate gradients

Source:

JOURNAL OF COMPUTATIONAL PHYSICS

Volume:

229

Issue:

21

Pages:

8065–8071

DOI:

10.1016/j.jcp.2010.07.010

Accession number:

WOS:000282118500001

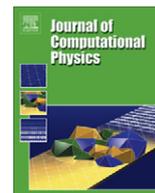
Published:

20th October 2010



Contents lists available at ScienceDirect

Journal of Computational Physics

journal homepage: www.elsevier.com/locate/jcp

Short note

Accelerating a FFT-based solver for numerical homogenization of periodic media by conjugate gradients

Jan Zeman^{a,*}, Jaroslav Vondřejc^a, Jan Novák^b, Ivo Marek^c^a Department of Mechanics, Faculty of Civil Engineering, Czech Technical University in Prague, Thákurova 7, 166 29 Prague 6, Czech Republic^b Centre for Integrated Design of Advanced Structures, Faculty of Civil Engineering, Czech Technical University in Prague, Thákurova 7, 166 29 Prague 6, Czech Republic^c Department of Mathematics, Faculty of Civil Engineering, Czech Technical University in Prague, Thákurova 7, 166 29 Prague 6, Czech Republic

ARTICLE INFO

Article history:

Received 7 April 2010

Received in revised form 30 June 2010

Accepted 9 July 2010

Available online 17 July 2010

Keywords:

Numerical homogenization

FFT-based solvers

Trigonometric collocation method

Conjugate gradient solvers

ABSTRACT

In this short note, we present a new technique to accelerate the convergence of a FFT-based solver for numerical homogenization of complex periodic media proposed by Moulinec and Suquet [1]. The approach proceeds from discretization of the governing integral equation by the trigonometric collocation method due to Vainikko [2], to give a linear system which can be efficiently solved by conjugate gradient methods. Computational experiments confirm robustness of the algorithm with respect to its internal parameters and demonstrate significant increase of the convergence rate for problems with high-contrast coefficients at a low overhead per iteration.

© 2010 Elsevier Inc. All rights reserved.

1. Introduction

A majority of computational homogenization algorithms rely on the solution of the unit cell problem, which concerns the determination of local fields in a representative sample of a heterogeneous material under periodic boundary conditions. Currently, the most efficient numerical solvers of this problem are based on discretization of integral equations. In the case of particulate composites with smooth bounded inclusions embedded in a matrix phase, the problem can be reduced to internal interfaces and solved with remarkable accuracy and efficiency by the fast multipole method, see [3, and references therein]. An alternative method has been proposed by Moulinec and Suquet [1] to treat problems with general microstructures supplied in the form of digital images. The algorithm is based on the Neumann series expansion of the inverse to an operator arising in the associated Lippmann–Schwinger equation and exploits the Fast Fourier Transform (FFT) to evaluate the action of the operator efficiently.

The major disadvantage of the FFT-based method consists in its poor convergence for composites exhibiting large jumps in material coefficients. To overcome this difficulty, Eyre and Milton proposed in [4] an accelerated scheme derived from a modified integral equation treated by means of the series expansion approach. In addition, Michel et al. [5] introduced an equivalent saddle-point formulation solved by the Augmented Lagrangian method. As clearly demonstrated in a numerical study by Moulinec and Suquet [6], both methods converge considerably faster than the original variant; the number of iterations is proportional to the square root of the phase contrast instead of the linear increase for the basic scheme. However,

* Corresponding author. Tel.: +420 2 2435 4482; fax: +420 2 2431 0775.

E-mail addresses: zemanj@cml.fsv.cvut.cz (J. Zeman), vondrej@gmail.com (J. Vondřejc), novakj@cml.fsv.cvut.cz (J. Novák), marek@mbox.ms.cuni.cz (I. Marek).

URL: <http://mech.fsv.cvut.cz/~zemanj> (J. Zeman).

this comes at the expense of increased computational cost per iteration and the sensitivity of the Augmented Lagrangian algorithm to the setting of its internal parameters.

In this short note, we introduce yet another approach to improve the convergence of the original FFT-based scheme [1] based on the trigonometric collocation method [7] and its application to the Helmholtz equation as introduced by Vainikko [2]. We observe that the discretization results in a system of linear equations with a structured dense matrix, for which a matrix–vector product can be computed efficiently using FFT (cf. Section 2). It is then natural to treat the resulting system by standard iterative solvers, such as the Krylov subspace methods, instead of the series expansion technique. In Section 3, the potential of such approach is demonstrated by means of a numerical study comparing the performance of the original scheme and the conjugate- and biconjugate-gradient methods for two-dimensional scalar electrostatics.

2. Methodology

In this section, we briefly summarize the essential steps of the trigonometric collocation-based solution to the unit cell problem by adapting the original exposition of Vainikko [2] to the setting of electrical conduction in periodic composites. In what follows, α , \mathbf{a} and \mathbf{A} denote scalar, vector and second-order tensor quantities with Greek subscripts used when referring to the corresponding components, e.g. $A_{\alpha\beta}$. Matrices are denoted by a serif font (e.g. \mathbf{A}) and a multi-index notation is employed, in which $\mathbb{R}^{\mathbf{N}}$ with $\mathbf{N} = (N_1, \dots, N_d)$ represents $\mathbb{R}^{N_1 \times \dots \times N_d}$ and $\mathbf{A}^{\mathbf{k}}$ stands for the (k_1, \dots, k_d) -th element of the matrix $\mathbf{A} \in \mathbb{R}^{\mathbf{N}}$.

2.1. Problem setting

We consider a composite material represented by a periodic unit cell $\mathcal{Y} = \prod_{\alpha=1}^d (-Y_\alpha, Y_\alpha) \subset \mathbb{R}^d$. In the context of linear electrostatics, the associated unit cell problem reads as

$$\nabla \times \mathbf{e}(\mathbf{x}) = \mathbf{0}, \quad \nabla \cdot \mathbf{j}(\mathbf{x}) = \mathbf{0}, \quad \mathbf{j}(\mathbf{x}) = \mathbf{L}(\mathbf{x}) \cdot \mathbf{e}(\mathbf{x}), \quad \mathbf{x} \in \mathcal{Y} \tag{1}$$

where \mathbf{e} is a \mathcal{Y} -periodic vectorial electric field, \mathbf{j} denotes the corresponding vector of electric current and \mathbf{L} is a second-order positive-definite tensor of electric conductivity. In addition, the field \mathbf{e} is subject to a constraint

$$\mathbf{e}^0 = \frac{1}{|\mathcal{Y}|} \int_{\mathcal{Y}} \mathbf{e}(\mathbf{x}) \, d\mathbf{x}, \tag{2}$$

where \mathbf{e}^0 denotes a prescribed macroscopic electric field and $|\mathcal{Y}|$ represents the d -dimensional measure of \mathcal{Y} .

Next, we introduce a homogeneous reference medium with constant conductivity \mathbf{L}^0 , leading to a decomposition of the electric current field in the form

$$\mathbf{j}(\mathbf{x}) = \mathbf{L}^0 \cdot \mathbf{e}(\mathbf{x}) + \delta \mathbf{L}(\mathbf{x}) \cdot \mathbf{e}(\mathbf{x}), \quad \delta \mathbf{L}(\mathbf{x}) = \mathbf{L}(\mathbf{x}) - \mathbf{L}^0. \tag{3}$$

The original problem (1)–(2) is then equivalent to the periodic Lippmann–Schwinger integral equation, formally written as

$$\mathbf{e}(\mathbf{x}) + \int_{\mathcal{Y}} \mathbf{r}^0(\mathbf{x} - \mathbf{y}) \cdot (\delta \mathbf{L}(\mathbf{y}) \cdot \mathbf{e}(\mathbf{y})) \, d\mathbf{y} = \mathbf{e}^0, \quad \mathbf{x} \in \mathcal{Y}, \tag{4}$$

where the \mathbf{r}^0 operator is derived from the Green's function of the problem (1)–(2) with $\mathbf{L}(\mathbf{x}) = \mathbf{L}^0$ and $\mathbf{e}^0 = \mathbf{0}$. Making use of the convolution theorem, Eq. (4) attains a local form in the Fourier space:

$$\hat{\mathbf{e}}(\mathbf{k}) = \begin{cases} |\mathcal{Y}|^{\frac{1}{2}} \mathbf{e}^0, & \mathbf{k} = \mathbf{0}, \\ -\widehat{\mathbf{r}}^0(\mathbf{k}) \cdot (\delta \widehat{\mathbf{L}} \cdot \hat{\mathbf{e}})(\mathbf{k}), & \mathbf{k} \in \mathbb{Z}^d \setminus \{\mathbf{0}\}, \end{cases} \tag{5}$$

where $\hat{f}(\mathbf{k})$ denotes the Fourier coefficient of $f(\mathbf{x})$ for the \mathbf{k} th frequency given by

$$\hat{f}(\mathbf{k}) = \int_{\mathcal{Y}} f(\mathbf{x}) \varphi_{-\mathbf{k}}(\mathbf{x}) \, d\mathbf{x}, \quad \varphi_{\mathbf{k}}(\mathbf{x}) = |\mathcal{Y}|^{-\frac{1}{2}} \exp\left(i\pi \sum_{\alpha=1}^d \frac{x_\alpha k_\alpha}{Y_\alpha}\right), \tag{6}$$

“ i ” is the imaginary unit and

$$\widehat{\mathbf{r}}^0(\mathbf{k}) = \begin{cases} \mathbf{0}, & \mathbf{k} = \mathbf{0}, \\ \frac{\mathbf{k} \otimes \mathbf{k}}{\mathbf{k} \cdot \mathbf{L}^0 \cdot \mathbf{k}}, & \mathbf{k} \in \mathbb{Z}^d \setminus \{\mathbf{0}\}, \end{cases} \tag{7}$$

Here, we refer to [4,8] for additional details.

2.2. Discretization via trigonometric collocation

Numerical solution of the Lippmann–Schwinger equation is based on a discretization of a unit cell \mathcal{Y} into a regular periodic grid with $N_1 \times \dots \times N_d$ nodal points and grid spacings $\mathbf{h} = (2Y_1/N_1, \dots, 2Y_d/N_d)$. The searched field \mathbf{e} in (4) is approximated by a trigonometric polynomial \mathbf{e}^N in the form (cf. [2])

$$\mathbf{e}(\mathbf{x}) \approx \mathbf{e}^N(\mathbf{x}) = \sum_{\mathbf{k} \in \mathbb{Z}^N} \hat{\mathbf{e}}(\mathbf{k}) \varphi_{\mathbf{k}}(\mathbf{x}), \quad \mathbf{x} \in \mathcal{Y}, \quad (8)$$

where $\mathbf{N} = (N_1, \dots, N_d)$, $\hat{\mathbf{e}}$ designates the Fourier coefficients defined in (6) and

$$\mathbb{Z}^N = \left\{ \mathbf{k} \in \mathbb{Z}^d : -\frac{N_\alpha}{2} < k_\alpha \leq \frac{N_\alpha}{2}, \alpha = 1, \dots, d \right\}. \quad (9)$$

We recall, e.g. from [2], that the α th component of the trigonometric polynomial expansion e_α^N admits two equivalent finite-dimensional representations. The first one is based on a matrix $\hat{e}_\alpha \in \mathbb{C}^N$ of the Fourier coefficients of the α th component and equation (8) with $\hat{e}_\alpha(\mathbf{k}) = \hat{e}_\alpha^{\mathbf{k}}$. Second, the data can be entirely determined by interpolation of nodal values

$$e_\alpha^N(\mathbf{x}) = \sum_{\mathbf{k} \in \mathbb{Z}^N} e_\alpha^{\mathbf{k}} \varphi_{\mathbf{k}}^N(\mathbf{x}), \quad \alpha = 1, \dots, d \quad (10)$$

where $e_\alpha \in \mathbb{R}^N$ is a matrix storing electric field values at grid points, $e_\alpha^{\mathbf{k}} = e_\alpha^N(\mathbf{x}^{\mathbf{k}})$ is the corresponding value at the \mathbf{k} th node with coordinates $\mathbf{x}^{\mathbf{k}} = (k_1 h_1, \dots, k_d h_d)$ and basis functions

$$\varphi_{\mathbf{k}}^N(\mathbf{x}) = |\mathbf{N}|^{-1} \sum_{\mathbf{m} \in \mathbb{Z}^N} \exp \left\{ i\pi \sum_{\alpha=1}^d m_\alpha \left(\frac{x_\alpha}{Y_\alpha} - \frac{2k_\alpha}{N_\alpha} \right) \right\} \quad (11)$$

satisfy the Dirac delta property $\varphi_{\mathbf{k}}^N(\mathbf{x}^{\mathbf{m}}) = \delta_{\mathbf{m}\mathbf{k}}$ with $|\mathbf{N}| = \prod_{\alpha=1}^d N_\alpha$. Both representations can be directly related to each other by

$$\hat{e}_\alpha = \mathbf{F} e_\alpha, \quad e_\alpha = \mathbf{F}^{-1} \hat{e}_\alpha, \quad (12)$$

where the Vandermonde matrices $\mathbf{F} \in \mathbb{C}^{N \times N}$ and $\mathbf{F}^{-1} \in \mathbb{C}^{N \times N}$

$$\mathbf{F}^{\mathbf{k}\mathbf{m}} = |\mathcal{Y}|^{-\frac{1}{2}} \exp \left(- \sum_{\alpha=1}^d 2\pi i \frac{k_\alpha m_\alpha}{N_\alpha} \right), \quad (13)$$

$$(\mathbf{F}^{-1})^{\mathbf{k}\mathbf{m}} = |\mathcal{Y}|^{\frac{1}{2}} |\mathbf{N}|^{-1} \exp \left(\sum_{\alpha=1}^d 2\pi i \frac{k_\alpha m_\alpha}{N_\alpha} \right), \quad (14)$$

implement the forward and inverse Fourier transform, respectively, e.g. [9, Section 4.6].

The trigonometric collocation method is based on the projection of the Lippmann–Schwinger equation (4) to the space of the trigonometric polynomials of the form $\{\sum_{\mathbf{k} \in \mathbb{Z}^N} c_{\mathbf{k}} \varphi_{\mathbf{k}}, c_{\mathbf{k}} \in \mathbb{C}\}$ (cf. [7,2]). In view of Eq. (10), this is equivalent to the collocation at grid points, with the action of \mathbf{L}^0 operator evaluated from the Fourier space expression (5) converted to the nodal representation by (12)₂. The resulting system of collocation equations reads

$$(\mathbf{I} + \mathbf{B})\mathbf{e} = \mathbf{e}^0, \quad (15)$$

where $\mathbf{e} \in \mathbb{R}^{d \times N}$ and $\mathbf{e}^0 \in \mathbb{R}^{d \times N}$ store the corresponding components of the solution and of the macroscopic field, respectively. Furthermore, \mathbf{I} is the $d \times d \times \mathbf{N} \times \mathbf{N}$ unit matrix and the non-symmetric matrix \mathbf{B} can be expressed, for the two-dimensional setting, in the partitioned format as

$$\mathbf{B} = \begin{bmatrix} \mathbf{F}^{-1} & 0 \\ 0 & \mathbf{F}^{-1} \end{bmatrix} \begin{bmatrix} \hat{\Gamma}_{11}^0 & \hat{\Gamma}_{12}^0 \\ \hat{\Gamma}_{21}^0 & \hat{\Gamma}_{22}^0 \end{bmatrix} \begin{bmatrix} \mathbf{F} & 0 \\ 0 & \mathbf{F} \end{bmatrix} \begin{bmatrix} \delta L_{11} & \delta L_{12} \\ \delta L_{21} & \delta L_{22} \end{bmatrix}, \quad (16)$$

with an obvious generalization to an arbitrary dimension. Here, $\hat{\Gamma}_{\alpha\beta}^0 \in \mathbb{R}^{N \times N}$ and $\delta L_{\alpha\beta} \in \mathbb{R}^{N \times N}$ are diagonal matrices storing the corresponding grid values, for which it holds

$$\left(\hat{\Gamma}_{\alpha\beta}^0 \right)^{\mathbf{k}\mathbf{k}} = \hat{\Gamma}_{\alpha\beta}^0(\mathbf{k}), \quad \delta L_{\alpha\beta}^{\mathbf{k}\mathbf{k}} = \delta L_{\alpha\beta}(\mathbf{x}^{\mathbf{k}}), \quad \alpha, \beta = 1, \dots, d \quad \text{and} \quad \mathbf{k} \in \mathbb{Z}^N. \quad (17)$$

2.3. Iterative solution of collocation equations

It follows from Eq. (16) that the cost of the multiplication by \mathbf{B} or by \mathbf{B}^T is driven by the forward and inverse Fourier transforms, which can be performed in $O(|\mathbf{N}| \log |\mathbf{N}|)$ operations by FFT techniques. This makes the resulting system (15) ideally suited for iterative solvers.

In particular, the original Fast Fourier Transform-based Homogenization (FFTH) scheme formulated by Moulinec and Suquet in [1] is based on the Neumann expansion of the matrix inverse $(\mathbf{I} + \mathbf{B})^{-1}$, so as to yield the m th iterate in the form

$$\mathbf{e}^{(m)} = \sum_{j=0}^m (-B)^j \mathbf{e}^0. \quad (18)$$

Convergence of the series (18) was comprehensively studied in [4,8], where it was shown that the optimal rate of convergence is achieved for

$$\mathbf{L}^0 = \frac{\lambda_{\min} + \lambda_{\max}}{2} \mathbf{I}, \quad (19)$$

with λ_{\min} and λ_{\max} denoting the minimum and maximum eigenvalues of $\mathbf{L}(\mathbf{x})$ on \mathcal{Y} and \mathbf{I} being the identity tensor.

Here, we propose to solve the non-symmetric system (15) by well-established Krylov subspace methods, in particular, exploiting the classical Conjugate Gradient (CG) method [10] and the biconjugate gradient (BiCG) algorithm [11]. Even though that CG algorithm is generally applicable to symmetric and positive-definite systems only, its convergence in the one-dimensional setting has been proven by Vondřejc [12, Section 6.2]. A successful application of CG method to a generalized Eshelby inhomogeneity problem has also been recently reported by Novák [13] and Kanaun [14].

3. Results

To assess the performance of the conjugate gradient algorithms, we consider a model problem of the transverse electric conduction in a square array of identical circular particles with 50% volume fraction. A uniform macroscopic field $\mathbf{e}^0 = (1, 0)$ is imposed on the corresponding single-particle unit cell, discretized by $\mathbf{N} = (255, 255)$ nodes¹ and the phases are considered to be isotropic with the conductivities set to $\mathbf{L} = \mathbf{I}$ for the matrix phase and to $\mathbf{L} = \varrho \mathbf{I}$ for the particle.

The conductivity of the homogeneous reference medium is parameterized as

$$\mathbf{L}^0(\omega) = (1 - \omega + \varrho\omega) \mathbf{I}, \quad (20)$$

where $\omega = 0.5$ corresponds to the optimal convergence of FFTH algorithm (19). All conjugate gradient-related results have been obtained using the implementations according to [16] and referred to as Algorithm 6.18 (CG method) and Algorithm 7.3 (BiCG scheme). Two termination criteria are considered. The first one is defined for the m th iteration as [15]

$$(\eta_e^{(m)})^2 = \frac{\sum_{\mathbf{k} \in \mathbb{Z}^N} (\mathbf{k} \cdot \hat{\mathbf{j}}^{\mathbf{k}(m)})^2}{\|\hat{\mathbf{j}}^{\mathbf{0}(m)}\|_2^2} \leq \varepsilon^2, \quad (21)$$

and provides the test of the equilibrium condition (1)₂ in the Fourier space. An alternative expression, related to the standard residual norm for iterative solvers, has been proposed by Vinogradov and Milton in [8] and admits the form

$$\eta_r^{(m)} = \frac{\|\mathbf{L}^0(\mathbf{e}^{(m+1)} - \mathbf{e}^{(m)})\|_2}{\|\mathbf{e}^0\|_2} \leq \varepsilon, \quad (22)$$

with the additional \mathbf{L}^0 term ensuring the proportionality to (21) at convergence. From the numerical point of view, the latter criterion is more efficient than the equilibrium variant, which requires additional operations per iteration. From the theoretical point of view, its usage is justified only when supported by a convergence result for the iterative algorithm. In the opposite case, the equilibrium norm appears to be more appropriate, in order to avoid spurious non-physical solutions.

3.1. Choice of reference medium and norm

Since no results for the optimal choice of the reference medium are known for (Bi)CG-based solvers, we first estimate their sensitivity to this aspect numerically. The results appear in Fig. 1(a), plotting the relative number of iterations for CG and BiCG solvers against the conductivity of the reference medium parameterized by ω , recall Eq. (20).

As expected, both CG and BiCG solvers achieve a significant improvement over FFTH method in terms of the number of iterations, ranging from 50% for a mildly-contrasted composite down to 2% for $\varrho = 10^4$. Moreover, contrary to all other available methods, the number of iterations is almost independent of the choice of the reference medium. We also observe, in agreement with results in [12, Section 6.2] for the one-dimensional setting, that CG and BiCG algorithms generate identical sequences of iterates; the minor differences visible for $\omega > 1$ or $\varrho = 10^4$ can be therefore attributed to accumulation of round-off errors. These conclusions hold for both equilibrium- and residual-based norms, which appear to be roughly proportional for the considered range of the phase contrasts (cf. Fig. 1(b)). Therefore, the residual criterion (22) will mostly be used in what follows.

In Fig. 2, we supplement the comparison by considering the total CPU time required to achieve a convergence. The data indicate that the cost of one iteration is governed by the matrix-vector multiplication, recall Eq. (16): the overhead of CG scheme is about 10% with respect to FFTH method, while the application of BiCG algorithm, which involves \mathbf{B} and \mathbf{B}^T products per iteration [11], is about twice as demanding. As a result, CG algorithm significantly reduces the overall computational

¹ Note that the odd number of discretization points is used to eliminate artificial high-frequency oscillations of the solution in the Fourier space, as reported in [15, Section 2.4].

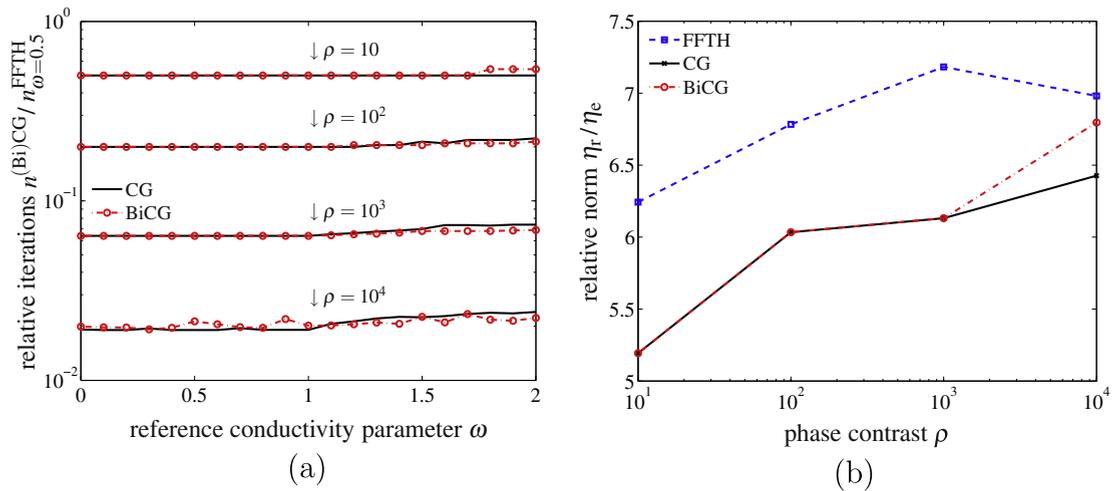


Fig. 1. (a) Relative number of iterations as a function of the reference medium parameter ω and (b) ratio between residual- and equilibrium-based norms at convergence for η_r termination condition with tolerance $\varepsilon = 10^{-4}$.

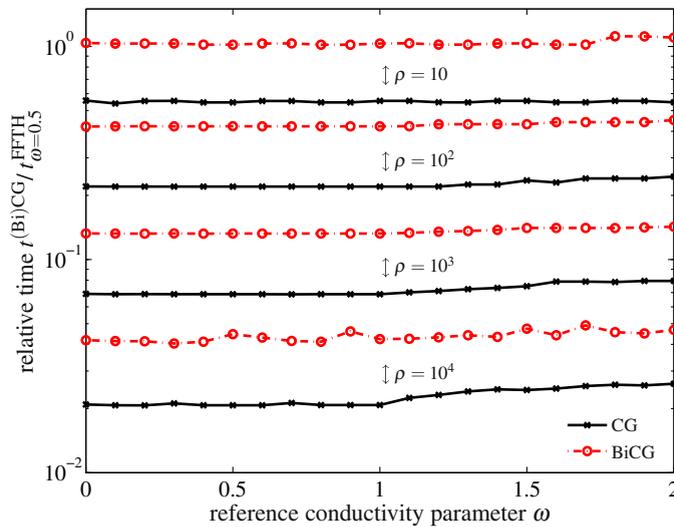


Fig. 2. Relative CPU time t of CG and BiCG solvers plotted against the conductivity parameter ω for η_r -based termination condition with tolerance $\varepsilon = 10^{-4}$.

time in the whole range of contrasts, whereas a similar effect has been reported for the candidate schemes only for $q \geq 10^3$ (cf. [6]).

3.2. Influence of phase contrast

As confirmed by all previous works, the phase contrast q is the critical parameter influencing the convergence of FFT-based iterative solvers. In Fig. 3, we compare the scaling of the total number of iterations with respect to phase contrast for CG and FFTH methods, respectively. The results clearly show that the number of iterations grows as \sqrt{q} instead of the linear increase for FFTH method. This follows from error bounds

$$\eta_r^{(m)} \leq \gamma^m \eta_r^{(0)}, \quad \gamma^{\text{FFTH}} = \frac{q-1}{q+1}, \quad \gamma^{\text{CG}} = \frac{\sqrt{q}-1}{\sqrt{q}+1}. \tag{23}$$

The first estimate was proven in [4], whereas the second expression is a direct consequence of the condition number of matrix B being proportional to q and a well-known result for the conjugate gradient method, e.g. [16, Section 6.11.3]. The CG-based method, however, failed to converge for the infinite contrast limit. Such behavior is equivalent to the Eyre-Milton

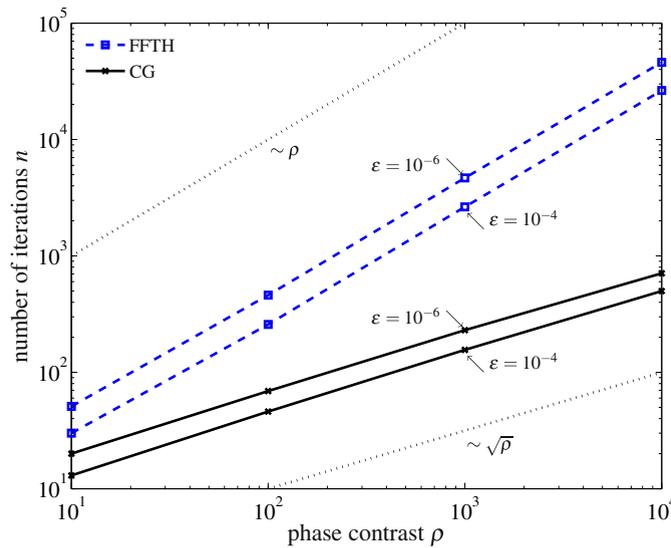


Fig. 3. Total number of iterations n plotted against phase contrast ρ ; the reference medium corresponds to for $\omega = 0.5$ and tolerance ε is related to η_r norm.

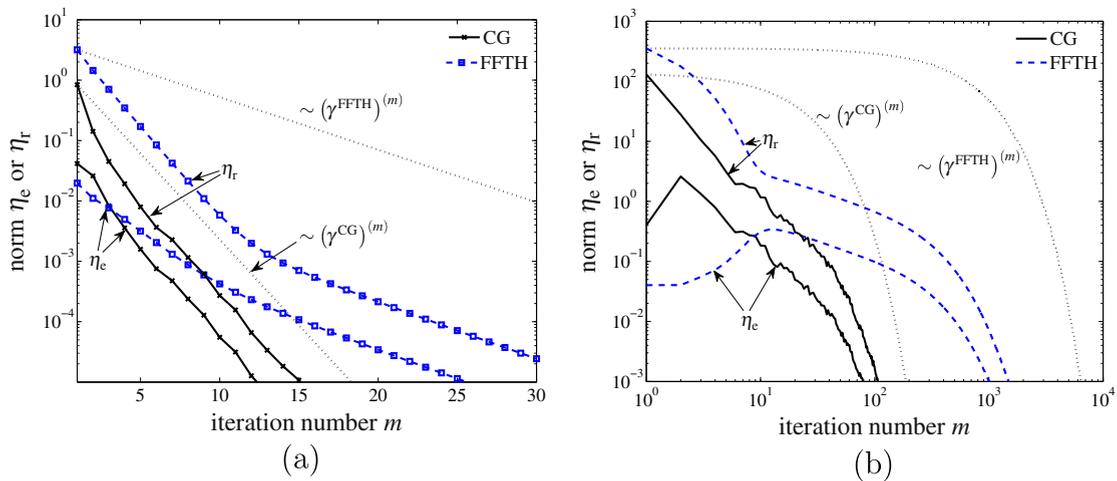


Fig. 4. Convergence progress of CG and FFTH methods for (a) $\rho = 10^1$ and (b) $\rho = 10^3$ as quantified by η_e and η_r norms; reference medium corresponds to $\omega = 0.5$ and the dot-and-dashed curves indicate the convergence rates (23).

scheme [4]. It is, however, inferior to the Augmented Lagrangian algorithm, for which the convergence rate improves with increasing ρ and the method converges even as $\rho \rightarrow \infty$. Nonetheless, such results are obtained for optimal, but not always straightforward, choice of the parameters [5].

3.3. Convergence progress

The final illustration of the CG-based algorithm is provided by Fig. 4, displaying a detailed convergence behavior for both low- and high-contrast cases. The results in Fig. 4(a) correspond well with estimates (23) for both residual and equilibrium-based norms. Influence of a higher phase contrast is visible from Fig. 4(b), plotted in the full logarithmic scale. For FFTH algorithm, two regimes can be clearly distinguished. In the first few iterations, the residual error rapidly decreases, but the iterates tend to deviate from equilibrium. Then, both residuals are simultaneously reduced. For CG scheme, the increase of the equilibrium residual appears only in the first iteration and then the method rapidly converges to the correct solution. However, its convergence curve is irregular and the algorithm repeatedly stagnates in two consecutive iterations. Further analysis of this phenomenon remains a subject of future work.

4. Conclusions

In this short note, we have presented a conjugate gradient-based acceleration of the FFT-based homogenization solver originally proposed by Moulinec and Suquet [1] and illustrated its performance on a problem of electric conduction in a periodic two-phase composite with isotropic phases. On the basis of obtained results, we conjecture that:

- (i) the non-symmetric system of linear equations (15), arising from discretization by the trigonometric collocation method [2], can be solved using the standard conjugate gradient algorithm,
- (ii) the convergence rate of the method is proportional to the square root of the phase contrast,
- (iii) the methods fails to converge in the infinite contrast limit,
- (iv) contrary to available improvements of the original FFT-solver [4,5], the cost of one iteration remains comparable to the basic scheme and the method is insensitive to the choice of auxiliary reference medium.

The presented computational experiments provide the first step towards further improvements of the method, including a rigorous analysis of its convergence properties, acceleration by multi-grid solvers and preconditioning and the extension to non-linear problems.

Acknowledgements

The authors thank Milan Jirásek (Czech Technical University in Prague) and Christopher Quince (University of Glasgow) for helpful comments on the manuscript. This research was supported by the Czech Science Foundation, through projects No. GAČR 103/09/1748, No. GAČR 103/09/P490 and No. GAČR 201/09/1544, and by the Grant Agency of the Czech Technical University in Prague through project No. SGS10/124/OHK1/2T/11.

References

- [1] H. Moulinec, P. Suquet, A fast numerical method for computing the linear and nonlinear mechanical properties of composites, *Comptes rendus de l'Académie des sciences. Série II, Mécanique, physique, chimie, astronomie* 318 (11) (1994) 1417–1423.
- [2] G. Vainikko, Fast solvers of the Lippmann–Schwinger equation, in: R.P. Gilbert, J. Kajiwara, Y.S. Xu (Eds.), *Direct and Inverse Problems of Mathematical Physics*, International Society for Analysis Applications and Computation, vol. 5, Kluwer Academic Publishers., Dordrecht, The Netherlands, 2000, pp. 423–440.
- [3] L. Greengard, J. Lee, Electrostatics and heat conduction in high contrast composite materials, *Journal of Computational Physics* 211 (1) (2006) 64–76.
- [4] D.J. Eyre, G.W. Milton, A fast numerical scheme for computing the response of composites using grid refinement, *The European Physical Journal Applied Physics* 6 (1) (1999) 41–47.
- [5] J.C. Michel, H. Moulinec, P. Suquet, A computational method based on augmented Lagrangians and fast Fourier transforms for composites with high contrast, *CMES-Computer Modeling in Engineering and Sciences* 1 (2) (2000) 79–88.
- [6] H. Moulinec, P. Suquet, Comparison of FFT-based methods for computing the response of composites with highly contrasted mechanical properties, *Physica B: Condensed Matter* 338 (1–4) (2003) 58–60.
- [7] J. Saranen, G. Vainikko, Trigonometric collocation methods with product integration for boundary integral equations on closed curves, *SIAM Journal on Numerical Analysis* 33 (4) (1996) 1577–1596.
- [8] V. Vinogradov, G.W. Milton, An accelerated FFT algorithm for thermoelastic and non-linear composites, *International Journal for Numerical Methods in Engineering* 76 (11) (2008) 1678–1695.
- [9] G. Golub, C.F. Van Loan, *Matrix Computations*, The Johns Hopkins University Press, Baltimore and London, 1996.
- [10] M.R. Hestenes, E. Stiefel, Methods of conjugate gradients for solving linear systems, *Journal of Research of the National Bureau of Standards* 49 (6) (1952) 409–463.
- [11] R. Fletcher, Conjugate gradient methods for indefinite systems, in: G. Watson (Ed.), *Numerical Analysis, Proceedings of the Dundee Conference on Numerical Analysis*, 1975, Lecture Notes in Mathematics, vol. 506, Springer-Verlag, New York, 1976, pp. 73–89.
- [12] J. Vondřejc, Analysis of heterogeneous materials using efficient meshless algorithms: one-dimensional study, Master's thesis, Czech Technical University in Prague (2009). <<http://mech.fsv.cvut.cz/~vondrej/download/ING.pdf>>.
- [13] J. Novák, Calculation of elastic stresses and strains inside a medium with multiple isolated inclusions, in: M. Papadarakakis, B. Topping (Eds.), *Proceedings of the Sixth International Conference on Engineering Computational Technology*, Stirlingshire, UK, 2008, pp. 16, paper 127. doi:10.4203/ccp.89.127.
- [14] S. Kanaun, Fast calculation of elastic fields in a homogeneous medium with isolated heterogeneous inclusions, *International Journal of Multiscale Computational Engineering* 7 (4) (2009) 263–276.
- [15] H. Moulinec, P. Suquet, A numerical method for computing the overall response of nonlinear composites with complex microstructure, *Computer Methods in Applied Mechanics and Engineering* 157 (1–2) (1998) 69–94.
- [16] Y. Saad, *Iterative Methods for Sparse Linear Systems*, second ed., Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 2003.

Part III

Paper 2

Authors:

Jaroslav Vondřejc, Jan Zeman, and Ivo Marek

Title:

Analysis of a Fast Fourier Transform Based Method for Modeling of Heterogeneous Materials

Conference:

8th International Conference on Large-Scale Scientific Computations, LSSC 2011; Sozopol; 6th June 2011–10th June 2011

Source:

Lecture Notes in Computer Science

Volume:

7116

Year:

2012

Pages:

515–522

DOI:

10.1007/978-3-642-29843-1_58

Analysis of a Fast Fourier Transform Based Method for Modeling of Heterogeneous Materials

Jaroslav Vondřejc¹, Jan Zeman¹, and Ivo Marek²

¹ Czech Technical University in Prague,
Faculty of Civil Engineering, Department of Mechanics
vondrej@gmail.com,

² Czech Technical University in Prague,
Faculty of Civil Engineering, Department of Mathematics

Abstract. The focus of this paper is on the analysis of the Conjugate Gradient method applied to a non-symmetric system of linear equations, arising from a Fast Fourier Transform-based homogenization method due to Moulinec and Suquet [1]. Convergence of the method is proven by exploiting a certain projection operator reflecting physics of the underlying problem. These results are supported by a numerical example, demonstrating significant improvement of the Conjugate Gradient-based scheme over the original Moulinec-Suquet algorithm.

Keywords: Homogenization, Fast Fourier Transform, Conjugate Gradients

1 Introduction

The last decade has witnessed a rapid development in advanced experimental techniques and modeling tools for microstructural characterization, typically provided in the form of pixel- or voxel-based geometry. Such data now allow for the design of bottom-up predictive models of the overall behavior for a wide range of engineering materials. Of course, such step necessitates the development of specialized algorithms, capable of handling large-scale voxel-based data in an efficient manner. In the engineering community, perhaps the most successful solver meeting these criteria was proposed by Moulinec and Suquet in [1]. The algorithm is based on the Neumann series expansion of the inverse of an operator arising in the associated Lippmann-Schwinger equation and exploits the Fast Fourier Transform to evaluate the action of the operator efficiently for voxel-based data. In our recent work [2], we have offered a new approach to the Moulinec-Suquet scheme, by exploiting the trigonometric collocation method due to Saranen and Vainikko [3]. Here, the Lippman-Schwinger equation is projected to a space of trigonometric polynomials to yield a non-symmetric system of linear equations, see Section 2 below. Quite surprisingly, numerical experiments revealed that the system can be efficiently solved using the standard Conjugate

Gradient algorithm. The analysis of this phenomenon, as presented in Section 3, is at the heart of this contribution. The obtained results are further supported by a numerical example in Section 4 and summarized in Section 5.

The following notation is used throughout the paper. Symbols a , \mathbf{a} and \mathbf{A} denote scalar, vector and second-order tensor quantities, respectively, with Greek subscripts used when referring to the corresponding components, e.g. $A_{\alpha\beta}$. The outer product of two vectors is denoted as $\mathbf{a} \otimes \mathbf{a}$, whereas $\mathbf{a} \cdot \mathbf{b}$ or $\mathbf{A} \cdot \mathbf{b}$ represents the single contraction between vectors (or tensors). A multi-index notation is employed, in which $\mathbb{R}^{\mathbf{N}}$ with $\mathbf{N} = (N_1, \dots, N_d)$ represents $\mathbb{R}^{N_1 \times \dots \times N_d}$ and $|\mathbf{N}|$ abbreviates $\prod_{\alpha=1}^d N_\alpha$. Block matrices are denoted by capital letters typeset in a bold serif font, e.g. $\mathbf{A} \in \mathbb{R}^{d \times d \times \mathbf{N} \times \mathbf{N}}$, and the superscript and subscript indexes are used to refer to the components, such that $\mathbf{A} = [A_{\alpha\beta}^{\mathbf{k}\mathbf{m}}]_{\alpha,\beta=1,\dots,d}^{\mathbf{k},\mathbf{m} \in \bar{\mathbb{Z}}^{\mathbf{N}}}$ with

$$\bar{\mathbb{Z}}^{\mathbf{N}} = \left\{ \mathbf{k} \in \mathbb{Z}^d : -\frac{N_\alpha}{2} < k_\alpha \leq \frac{N_\alpha}{2}, \alpha = 1, \dots, d \right\}.$$

Sub-matrices of \mathbf{A} are denoted as

$$\mathbf{A}_{\alpha\beta} = [A_{\alpha\beta}^{\mathbf{k}\mathbf{m}}]_{\mathbf{k},\mathbf{m} \in \bar{\mathbb{Z}}^{\mathbf{N}}} \in \mathbb{R}^{\mathbf{N} \times \mathbf{N}}, \quad \mathbf{A}^{\mathbf{k}\mathbf{m}} = [A_{\alpha\beta}^{\mathbf{k}\mathbf{m}}]_{\alpha,\beta=1,\dots,d} \in \mathbb{R}^{d \times d}$$

for $\alpha, \beta = 1, \dots, d$ and $\mathbf{k}, \mathbf{m} \in \bar{\mathbb{Z}}^{\mathbf{N}}$. Analogously, the block vectors are denoted by lower case letters, e.g. $\mathbf{e} \in \mathbb{R}^{d \times \mathbf{N}}$ and the matrix-by-vector multiplication is defined as

$$[\mathbf{A}\mathbf{e}]_\alpha^{\mathbf{k}} = \sum_{\beta=1}^d \sum_{\mathbf{m} \in \bar{\mathbb{Z}}^{\mathbf{N}}} A_{\alpha\beta}^{\mathbf{k}\mathbf{m}} \mathbf{e}_\beta^{\mathbf{m}} \in \mathbb{R}^{d \times \mathbf{N}}, \quad (1)$$

with $\alpha = 1, \dots, d$ and $\mathbf{k} \in \bar{\mathbb{Z}}^{\mathbf{N}}$.

2 Problem setting

Consider a composite material represented by a periodic unit cell

$$\mathcal{Y} = \prod_{\alpha=1}^d (-Y_\alpha, Y_\alpha) \subset \mathbb{R}^d.$$

In the context of linear electrostatics, the associated unit cell problem reads as

$$\nabla \times \mathbf{e}(\mathbf{x}) = \mathbf{0}, \quad \nabla \cdot \mathbf{j}(\mathbf{x}) = \mathbf{0}, \quad \mathbf{j}(\mathbf{x}) = \mathbf{L}(\mathbf{x}) \cdot \mathbf{e}(\mathbf{x}), \quad \mathbf{x} \in \mathcal{Y} \quad (2)$$

where \mathbf{e} is a \mathcal{Y} -periodic vectorial electric field, \mathbf{j} denotes the corresponding vector of electric current and \mathbf{L} is a second-order positive-definite tensor of electric conductivity. In addition, the field \mathbf{e} is subject to a constraint

$$\langle \mathbf{e}(\mathbf{x}) \rangle := \frac{1}{|\mathcal{Y}|} \int_{\mathcal{Y}} \mathbf{e}(\mathbf{x}) \, d\mathbf{x} = \mathbf{e}^0, \quad (3)$$

where $|\mathcal{Y}|$ denotes the d -dimensional measure of \mathcal{Y} and $\mathbf{e}^0 \neq \mathbf{0}$ a prescribed macroscopic electric field.

The original problem (2)–(3) is then equivalent to the periodic Lippmann-Schwinger integral equation, formally written as

$$\mathbf{e}(\mathbf{x}) + \int_{\mathcal{Y}} \Gamma(\mathbf{x} - \mathbf{y}; \mathbf{L}^0) \cdot (\mathbf{L}(\mathbf{y}) - \mathbf{L}^0) \cdot \mathbf{e}(\mathbf{y}) d\mathbf{y} = \mathbf{e}^0, \quad \mathbf{x} \in \mathcal{Y}, \quad (4)$$

where $\mathbf{L}^0 \in \mathbb{R}^{d \times d}$ denotes a homogeneous reference medium. The operator $\Gamma(\mathbf{x}, \mathbf{L}^0)$ is derived from the Green's function of the problem (2)–(3) with $\mathbf{L}(\mathbf{x}) = \mathbf{L}^0$ and can be simply expressed in the Fourier space

$$\hat{\Gamma}(\mathbf{k}; \mathbf{L}^0) = \begin{cases} \mathbf{0} & \mathbf{k} = \mathbf{0} \\ \frac{\boldsymbol{\xi} \otimes \boldsymbol{\xi}}{\boldsymbol{\xi} \cdot \mathbf{L}^0 \cdot \boldsymbol{\xi}} & \boldsymbol{\xi}(\mathbf{k}) = \left(\frac{k_\alpha}{Y_\alpha} \right)_{\alpha=1}^d; \mathbf{k} \in \mathbb{Z}^d \setminus \mathbf{0}. \end{cases} \quad (5)$$

Operator $\hat{f} = \hat{f}(\mathbf{k})$ stands for the Fourier coefficient of $f(\mathbf{x})$ for the \mathbf{k} -th frequency given by

$$\hat{f}(\mathbf{k}) = \int_{\mathcal{Y}} f(\mathbf{x}) \varphi_{-\mathbf{k}}(\mathbf{x}) d\mathbf{x}, \quad \varphi_{\mathbf{k}}(\mathbf{x}) = |\mathcal{Y}|^{-\frac{1}{2}} \exp\left(i\pi \sum_{\alpha=1}^d \frac{x_\alpha k_\alpha}{Y_\alpha}\right), \quad (6)$$

”i” is the imaginary unit ($i^2 = -1$). We refer to [2,4] for additional details. Note that the linear electrostatics serves here as a model problem; the framework can be directly extended to e.g. elasticity [5], (visco-)plasticity [6] or to multiferroics [7].

2.1 Discretization via trigonometric collocation

The numerical solution of the Lippmann-Schwinger equation is based on a discretization of a unit cell \mathcal{Y} into a regular periodic grid with $N_1 \times \dots \times N_d$ nodal points and grid spacings $\mathbf{h} = (2Y_1/N_1, \dots, 2Y_d/N_d)$. The searched field \mathbf{e} in (4) is approximated by a trigonometric polynomial \mathbf{e}^N in the form (cf. [3, Chapter 10])

$$\mathbf{e}(\mathbf{x}) \approx \mathbf{e}^N(\mathbf{x}) = \sum_{\mathbf{k} \in \bar{\mathbb{Z}}^N} \hat{\mathbf{e}}^{\mathbf{k}} \varphi_{\mathbf{k}}(\mathbf{x}), \quad \mathbf{x} \in \mathcal{Y}, \quad (7)$$

where $\hat{\mathbf{e}}^{\mathbf{k}} = (\hat{e}_\alpha^{\mathbf{k}})_{\alpha=1, \dots, d}$ designates the Fourier coefficients defined in (6). Notice that the trigonometrical polynomials are uniquely determined by a regular grid data, which makes them well-suited to problems with pixel- or voxel-based computations.

The trigonometric collocation method is based on the projection of the Lippmann-Schwinger equation (4) onto the space of the trigonometric polynomials

$$\mathcal{T}^N = \left\{ \sum_{\mathbf{k} \in \bar{\mathbb{Z}}^N} c_{\mathbf{k}} \varphi_{\mathbf{k}}, c_{\mathbf{k}} \in \mathbb{C} \right\}, \quad (8)$$

leading to a linear system in the form, cf. [2]

$$(\mathbf{I} + \mathbf{B})\mathbf{e} = \mathbf{e}^0, \quad \mathbf{B} = \mathbf{F}^{-1}\hat{\mathbf{F}}\mathbf{F}(\mathbf{L} - \mathbf{L}^0), \quad (9)$$

where $\mathbf{e} = (\mathbf{e}_\alpha^{\mathbf{k}})_{\alpha=1,\dots,d}^{\mathbf{k} \in \bar{\mathbb{Z}}^N} \in \mathbb{R}^{d \times N}$ is the unknown vector, $\mathbf{I} = [\delta_{\alpha\beta} \delta_{\mathbf{k}\mathbf{m}}]_{\alpha,\beta=1,\dots,d}^{\mathbf{k},\mathbf{m} \in \bar{\mathbb{Z}}^N} \in \mathbb{R}^{d \times d \times N \times N}$ is the identity matrix, expressed as the product of the Kronecker delta functions $\delta_{\alpha\beta}$ and $\delta_{\mathbf{k}\mathbf{m}}$, and $\mathbf{e}^0 = (e_\alpha^0)_{\alpha=1,\dots,d}^{\mathbf{k} \in \bar{\mathbb{Z}}^N} \in \mathbb{R}^{d \times N}$.

All the matrices in (9) exhibit a block-diagonal structure. In particular,

$$\hat{\mathbf{F}} = [\delta_{\mathbf{k}\mathbf{m}} \hat{\Gamma}_{\alpha\beta}^{\mathbf{k}\mathbf{m}}]_{\alpha,\beta=1,\dots,d}^{\mathbf{k},\mathbf{m} \in \bar{\mathbb{Z}}^N}, \quad \mathbf{L} = [\delta_{\mathbf{k}\mathbf{m}} \mathbf{L}_{\alpha\beta}^{\mathbf{k}\mathbf{m}}]_{\alpha,\beta=1,\dots,d}^{\mathbf{k},\mathbf{m} \in \bar{\mathbb{Z}}^N}, \quad \mathbf{L}^0 = [\delta_{\mathbf{k}\mathbf{m}} \mathbf{L}_{\alpha\beta}^0]_{\alpha,\beta=1,\dots,d}^{\mathbf{k},\mathbf{m} \in \bar{\mathbb{Z}}^N},$$

with $\hat{\Gamma}_{\alpha\beta}^{\mathbf{k}\mathbf{k}} = \hat{\Gamma}_{\alpha\beta}(\mathbf{k}; \mathbf{L}^0)$, $\mathbf{L}_{\alpha\beta}^{\mathbf{k}\mathbf{k}} = L_{\alpha\beta}(\mathbf{k})$ and $(\mathbf{L}^0)_{\alpha\beta} = L_{\alpha\beta}^0$. The matrix \mathbf{F} implements the Discrete Fourier Transform and is defined as

$$\mathbf{F} = [\delta_{\alpha\beta} \mathbf{F}^{\mathbf{k}\mathbf{m}}]_{\alpha,\beta=1,\dots,d}^{\mathbf{k},\mathbf{m} \in \bar{\mathbb{Z}}^N}, \quad \mathbf{F}^{\mathbf{k}\mathbf{m}} = \frac{|\mathcal{Y}|^{\frac{1}{2}}}{\prod_{\alpha=1}^d N_\alpha} \exp\left(-\sum_{\alpha=1}^d 2\pi i \frac{k_\alpha m_\alpha}{N_\alpha}\right), \quad (10)$$

with \mathbf{F}^{-1} representing the inverse transform.

It follows from Eq. (1) that the cost of multiplication by \mathbf{B} is dominated by the action of \mathbf{F} and \mathbf{F}^{-1} , which can be performed in $O(|N| \log |N|)$ operations by the Fast Fourier Transform techniques. This makes the system (9) well-suited for applying some iterative solution technique. In particular, the original Fast Fourier Transform-based Homogenization scheme formulated by Moulinec and Suquet in [1] is based on the Neumann expansion of the matrix inverse $(\mathbf{I} + \mathbf{B})^{-1}$, so as to yield the m -th iterate in the form

$$\mathbf{e}^{(m)} = \sum_{j=0}^m (-\mathbf{B})^j \mathbf{e}^0. \quad (11)$$

As indicated earlier, our numerical experiments [2] suggest that the system can be efficiently solved using the Conjugate Gradient method, despite the non-symmetry of \mathbf{B} evident from (9). This observation is studied in more detail in the next Section.

3 Solution by the Conjugate Gradient method

We start our analysis with recasting the system (9) into a more convenient form, by employing a certain operator and the associated sub-space introduced later. Note that for simplicity, the reference conductivity is taken as $\mathbf{L}^0 = \lambda \mathbf{I}$.

Definition 1. Given $\lambda > 0$, we define operator $\mathbf{P}_\mathcal{E} = \lambda \mathbf{F}^{-1} \hat{\mathbf{F}} \mathbf{F}$ and associated sub-space as

$$\mathcal{E} = \{\mathbf{P}_\mathcal{E} \mathbf{x} \text{ for } \mathbf{x} \in \mathbb{R}^{d \times N}\} \subset \mathbb{R}^{d \times N}.$$

Lemma 1. The operator $\mathbf{P}_\mathcal{E}$ is an orthogonal projection.

Proof. First, we will prove that $\mathbf{P}_{\mathcal{E}}$ is projection, i.e. $\mathbf{P}_{\mathcal{E}}^2 = \mathbf{P}_{\mathcal{E}}$. Since \mathbf{F} is a unitary matrix, it is easy to see that

$$\mathbf{P}_{\mathcal{E}}^2 = (\lambda\mathbf{F}^{-1}\hat{\mathbf{F}}\mathbf{F})(\lambda\mathbf{F}^{-1}\hat{\mathbf{F}}\mathbf{F}) = \mathbf{F}^{-1}(\lambda\hat{\mathbf{F}})^2\mathbf{F}. \quad (12)$$

Hence, in view of the block-diagonal character of $\hat{\mathbf{F}}$, it is sufficient to prove the projection property of sub-matrices $(\lambda\hat{\mathbf{F}})^{kk}$ only. This follows using a simple algebra, recall Eq. (5):

$$(\lambda\hat{\mathbf{F}})^{kk}(\lambda\hat{\mathbf{F}})^{kk} = \frac{\boldsymbol{\xi}(\mathbf{k}) \otimes \boldsymbol{\xi}(\mathbf{k})}{\boldsymbol{\xi}(\mathbf{k}) \cdot \boldsymbol{\xi}(\mathbf{k})} \cdot \frac{\boldsymbol{\xi}(\mathbf{k}) \otimes \boldsymbol{\xi}(\mathbf{k})}{\boldsymbol{\xi}(\mathbf{k}) \cdot \boldsymbol{\xi}(\mathbf{k})} = \frac{\boldsymbol{\xi}(\mathbf{k}) \otimes \boldsymbol{\xi}(\mathbf{k})}{\boldsymbol{\xi}(\mathbf{k}) \cdot \boldsymbol{\xi}(\mathbf{k})} = (\lambda\hat{\mathbf{F}})^{kk}.$$

The orthogonality of $\mathbf{P}_{\mathcal{E}}$ now follows from

$$\mathbf{P}_{\mathcal{E}}^* = \left(\lambda\mathbf{F}^{-1}\hat{\mathbf{F}}\mathbf{F}\right)^* = \lambda\mathbf{F}^*\hat{\mathbf{F}}^*(\mathbf{F}^{-1})^* = \lambda\mathbf{F}^{-1}\hat{\mathbf{F}}\mathbf{F} = \mathbf{P}_{\mathcal{E}},$$

according to a well-known result of linear algebra, e.g. Proposition 1.8 in [8]. \square

Remark 1. It follows from the previous results that the subspace \mathcal{E} collects the non-zero coefficients of trigonometric polynomials \mathcal{T}^N with zero rotation, which represent admissible solutions to the unit cell problem defined by (2). Note that the orthogonal space \mathcal{E}^\perp contains the trigonometric representation of constant fields, cf. [4, Section 12.7].

Lemma 2. *The solution \mathbf{e} to the linear system (9) admits the decomposition $\mathbf{e} = \mathbf{e}^0 + \mathbf{e}_{\mathcal{E}}$, with $\mathbf{e}_{\mathcal{E}} \in \mathcal{E}$ satisfying*

$$\mathbf{P}_{\mathcal{E}}\mathbf{L}\mathbf{e}_{\mathcal{E}} + \mathbf{P}_{\mathcal{E}}\mathbf{L}\mathbf{e}^0 = \mathbf{0}. \quad (13)$$

Proof. As $\mathbf{e} \in \mathbb{R}^{d \times N}$, Lemma 1 ensures that it can be decomposed into two orthogonal parts $\mathbf{e}_{\mathcal{E}} = \mathbf{P}_{\mathcal{E}}\mathbf{e}$ and $\mathbf{e}_{\mathcal{E}^\perp} = (\mathbf{I} - \mathbf{P}_{\mathcal{E}})\mathbf{e}$. Substituting this expression into (9), and using the identity $\mathbf{B} = \lambda\mathbf{F}^{-1}\hat{\mathbf{F}}\mathbf{F}\left(\frac{1}{\lambda}\mathbf{I} - \mathbf{I}\right)$, we arrive at

$$\frac{1}{\lambda}\mathbf{P}_{\mathcal{E}}\mathbf{L}\mathbf{e}_{\mathcal{E}} + \mathbf{e}_{\mathcal{E}^\perp} + \frac{1}{\lambda}\mathbf{P}_{\mathcal{E}}\mathbf{L}\mathbf{e}_{\mathcal{E}^\perp} = \mathbf{e}^0. \quad (14)$$

Since $\mathbf{e}^0 \in \mathcal{E}^\perp$, we have $\mathbf{e}_{\mathcal{E}^\perp} = \mathbf{e}^0$ and the proof is complete. \square

With these auxiliary results in hand, we are in the position to present our main result.

Proposition 1. *The non-symmetric system of linear equations (9) is solvable by the Conjugate Gradient method for an initial vector $\mathbf{e}_{(0)} = \mathbf{e}^0 + \tilde{\mathbf{e}}$ with $\tilde{\mathbf{e}} \in \mathcal{E}$. Moreover, the sequence of iterates is independent of the parameter λ .*

Proof (outline). It follows from Lemma 2 that the solution to (9) admits yet another, optimization-based, characterization in the form

$$\mathbf{e} = \mathbf{e}^0 + \arg \min_{\bar{\mathbf{e}} \in \mathcal{E}} \left[\frac{1}{2} (\mathbf{L}\bar{\mathbf{e}}, \bar{\mathbf{e}})_{\mathbb{R}^{d \times N}} + (\mathbf{L}\mathbf{e}^0, \bar{\mathbf{e}})_{\mathbb{R}^{d \times N}} \right]. \quad (15)$$

The residual corresponding to the initial vector $\mathbf{e}_{(0)}$ equals to

$$\mathbf{r}_{(0)} = \mathbf{e}^0 - (\mathbf{I} + \mathbf{B})(\mathbf{e}^0 + \tilde{\mathbf{e}}) = -\frac{1}{\lambda} \mathbf{P}_{\mathcal{E}} \mathbf{L} \mathbf{e}^0 - \frac{1}{\lambda} \mathbf{P}_{\mathcal{E}} \mathbf{L} \tilde{\mathbf{e}} \in \mathcal{E}.$$

It can be verified that the subspace \mathcal{E} is \mathbf{B} -invariant, thus $(\mathbf{I} + \mathbf{B})\mathcal{E} \subset \mathcal{E}$. Therefore, the Krylov subspace

$$\mathcal{K}_m(\mathbf{I} + \mathbf{B}, \mathbf{r}_{(0)}) = \text{span} \{ \mathbf{r}_{(0)}, (\mathbf{I} + \mathbf{B})\mathbf{r}_{(0)}, \dots, (\mathbf{I} + \mathbf{B})^m \mathbf{r}_{(0)} \} \subset \mathcal{E}$$

for arbitrary $m \in \mathbb{N}$. This implies that the residual $\mathbf{r}_{(m)}$ and the Conjugate Gradient search direction $\mathbf{p}_{(m)}$ at the m -th iteration satisfy $\mathbf{r}_{(m)} \in \mathcal{E}$ and $\mathbf{p}_{(m)} \in \mathcal{E}$. Since \mathbf{B} is symmetric and positive-definite on \mathcal{E} , the convergence of CG algorithm now follows from standard arguments, e.g. Theorem 6.6 in [8]. Observe that different choices of λ generate identical Krylov subspaces, thus the sequence of iterates is independent of λ . \square

Remark 2. Note that it is possible to show, using direct calculations based on the projection properties of $\mathbf{P}_{\mathcal{E}}$, that the Biconjugate Gradient algorithm produces exactly the same sequence of vectors as the Conjugate Gradient method, see [9].

4 Numerical example

To support our theoretical results, we consider a three-dimensional model problem of electric conduction in a cubic periodic unit cell $\mathcal{Y} = \prod_{\alpha=1}^3 (-\frac{1}{2}, \frac{1}{2})$, representing a two-phase medium with spherical inclusions of 25% volume fraction. The conductivity parameters are defined as

$$\mathbf{L}(\mathbf{x}) = \begin{cases} \rho \mathbf{I}, & \|\mathbf{x}\|_2 < (\frac{3}{16\pi})^{\frac{1}{3}} \\ \begin{pmatrix} 1 & 0.2 & 0.2 \\ 0.2 & 1 & 0.2 \\ 0.2 & 0.2 & 1 \end{pmatrix}, & \text{otherwise} \end{cases}$$

where $\rho > 0$ denotes the contrast of phase conductivities. We consider the macroscopic field $\mathbf{e}^0 = [1, 0, 0]$ and discretize the unit cell with $\mathbf{N} = [n, n, n]$ nodes³. The conductivity of the homogeneous reference medium $\mathbf{L}^0 \in \mathbb{R}^{d \times d}$ is parametrized as

$$\mathbf{L}^0 = \lambda \mathbf{I}, \quad \lambda = 1 - \omega + \rho\omega, \quad (16)$$

where $\omega \approx 0.5$ delivers the optimal convergence of the original Moulinec-Suquet Fast-Fourier Transform-based Homogenization (FFTH) algorithm [1].

We first investigate the sensitivity of Conjugate Gradient (CG) algorithm to the choice of reference medium. The results appear in Fig. 1(a), plotting the

³ In particular, n was taken consequently as 16, 32, 64, 128 and 160 leading up to $3 \cdot 160^3 \doteq 12.2 \times 10^6$ unknowns

relative number of iterations for CG against the conductivity of the reference medium parametrized by ω , recall Eq. (16). As expected, CG solver achieve a significant improvement over FFTH method as it requires about 40% iterations of FFTH for a mildly-contrasted composite down to 4% for $\rho = 10^3$. The minor differences visible especially for $\rho = 10^3$ can be therefore attributed to accumulation of round-off errors. These observations fully confirm our theoretical results presented earlier in Section 3.

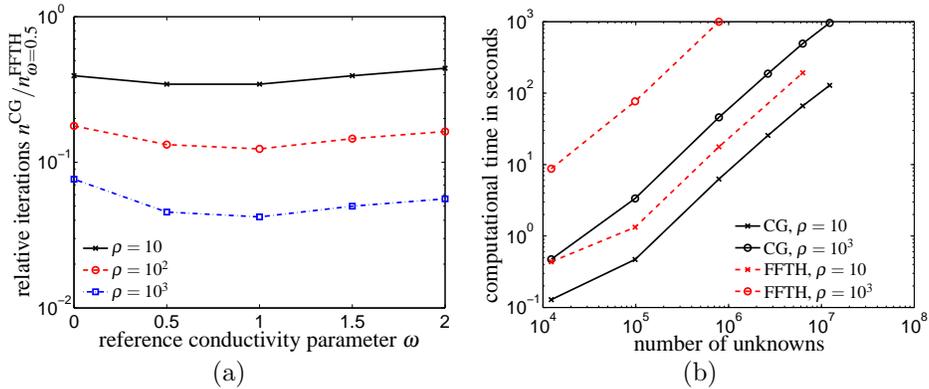


Fig. 1. (a) Relative number of iterations as a function of the reference medium parameter ω and (b) computational time as a function of the number of unknowns.

In Fig. 1(b), we present the total computational time⁴ as a function of the number of degrees of freedom and the phase ratio ρ . The results confirm that the computational times scales linearly with the increasing number of degrees of freedom for both schemes for a fixed ρ [2]. The ratio of the computational time for CG and FFTH algorithms remains almost constant, which indicates that the cost of a single iteration of CG and FFTH method is comparable.

In addition, the memory requirements of both schemes are also comparable. This aspect represents the major advantage of the short-recurrence CG-based scheme over alternative schemes for non-symmetric systems, such as GMRES. Finally note that finer discretizations can be treated by a straightforward parallel implementation.

5 Conclusions

In this work, we have proven the convergence of Conjugate Gradient method for a non-symmetric system of linear equations arising from periodic unit cell ho-

⁴ The problem was solved with a MATLAB[®] in-house code on a machine Intel[®] Core[™]2 Duo 3 GHz CPU, 3.28 GB computing memory with Debian linux 5.0 operating system.

mogenization problem and confirmed it by numerical experiment. The important conclusions to be pointed out are as follows:

1. The success of the Conjugate Gradient method follows from the projection properties of operator $\mathbf{P}_\mathcal{E}$ introduced in Definition 1, which reflect the structure of the underlying physical problem.
2. Contrary to all available extensions of the FFTH scheme, the performance of the Conjugate Gradient-based method is independent of the choice of reference medium. This offers an important starting point for further improvements of the method.

Apart from the already mentioned parallelization, performance of the scheme can further be improved by a suitable preconditioning procedure. This topic is currently under investigation.

Acknowledgments This work was supported by the Czech Science Foundation, through projects No. GAČR 103/09/1748, No. GAČR 103/09/P490, No. GAČR 201/09/1544, and by the Grant Agency of the Czech Technical University in Prague through project No. SGS10/124/OHK1/2T/11.

References

1. H. Moulinec, P. Suquet, A fast numerical method for computing the linear and nonlinear mechanical properties of composites, *Comptes rendus de l'Académie des sciences. Série II, Mécanique, physique, chimie, astronomie* 318 (11) (1994) 1417–1423.
2. J. Zeman, J. Vondřejc, J. Novák, I. Marek, Accelerating a FFT-based solver for numerical homogenization of periodic media by conjugate gradients, *Journal of Computational Physics* 229 (21) (2010) 8065–8071. [arXiv:1004.1122](https://arxiv.org/abs/1004.1122).
3. J. Saranen, G. Vainikko, *Periodic Integral and Pseudodifferential Equations with Numerical Approximation*, Springer Monographs in Mathematics, Springer-Verlag, Berlin, Heidelberg, 2002.
4. G. W. Milton, *The Theory of Composites*, Vol. 6 of Cambridge Monographs on Applied and Computational Mathematics, Cambridge University Press, Cambridge, UK, 2002.
5. V. Šmilauer, Z. Bittnar, Microstructure-based micromechanical prediction of elastic properties in hydrating cement paste, *Cement and Concrete Research* 36 (9) (2006) 1708–1718.
6. A. Prakash, R. A. Lebensohn, Simulation of micromechanical behavior of polycrystals: finite elements versus fast Fourier transforms, *Modelling and Simulation in Materials Science and Engineering* 17 (6) (2009) 064010+.
7. R. Brenner, J. Bravo-Castillero, Response of multiferroic composites inferred from a fast-Fourier-transform-based numerical scheme, *Smart Materials and Structures* 19 (11) (2010) 115004+.
8. Y. Saad, *Iterative Methods for Sparse Linear Systems*, second edition with corrections Edition, Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 2003.
9. J. Vondřejc, Analysis of heterogeneous materials using efficient meshless algorithms: One-dimensional study, Diploma thesis, Czech Technical University in Prague, URL: <http://mech.fsv.cvut.cz/~vondrejck/download/ING.pdf> (2009).

Part IV

Paper 3

Authors:

Jiří Němeček, Vlastimil Králík, and Jaroslav Vondřejc

Title:

Micromechanical analysis of heterogeneous structural materials

Source:

Cement and Concrete Composites

DOI:

10.1016/j.cemconcomp.2012.06.015

Published:

In Press, Available online 4th July 2012



Contents lists available at SciVerse ScienceDirect

Cement & Concrete Composites

journal homepage: www.elsevier.com/locate/cemconcomp

Micromechanical analysis of heterogeneous structural materials

Jiří Němeček*, Vlastimil Králík, Jaroslav Vondřejc

Czech Technical University in Prague, Department of Mechanics, Faculty of Civil Engineering, Thákurova 7, 166 29 Prague 6, Czech Republic

ARTICLE INFO

Article history:

Received 31 January 2012

Received in revised form 25 June 2012

Accepted 26 June 2012

Available online xxx

Keywords:

Micromechanics

Nanoindentation

Heterogeneous materials

Grid indentation

Deconvolution

Homogenization

FFT

ABSTRACT

This paper shows an efficient methodology based on micromechanical framework and grid nanoindentation for the assessment of effective elastic properties on several types of microscopically heterogeneous structural materials. Such task is a prerequisite for successful nano- and micro-structural material characterization, development and optimization. The grid nanoindentation and statistical deconvolution methods previously described in the literature e.g. for cementitious materials [1,2], alkali activated materials [3] or high-performance concretes [4] have been employed. In this paper we demonstrate their utilization also for other types of structural composites with crystalline nature and we validate the results by using enhanced numerical method based on fast Fourier transform (FFT). The direct procedure of using grid nanoindentation data in the FFT method simplifies the evaluation of effective composite properties and leads to the assemblage of the full stiffness matrix compared to simple analytical approaches.

The paper deals namely with cement paste, gypsum and aluminum alloy. Nanoindentation is used for the determination of phase properties in grid points at the scale below one micrometer. Statistical approach and deconvolution methods are applied to assess intrinsic phase properties. Elastic properties obtained by nanoindentation are homogenized in the frame of the representative volume element (RVE) by means of analytical and numerical FFT-based schemes. Good correlation of the results from all methods was found for the tested materials due to the close-to-isotropic nature of the composites in the RVE having dimensions $\sim 100\text{--}200\ \mu\text{m}$. Results were also verified against macroscopic experimental results. The proposed and validated numerical approach can be successively used for the material modeling in finite element software or for optimization of materials with inhomogeneous microstructures.

© 2012 Elsevier Ltd. All rights reserved.

1. Introduction

Structural composites such as concrete, gypsum, metals and others are often characterized by a heterogeneous microstructure at different length scales (nm to m). Traditionally, their mechanical properties are assessed from macroscopic tests on samples with cm to m dimensions that can only describe overall (averaged) properties like overall Young's modulus or strength. Nowadays, nanoindentation [5] can be successfully applied to access the nanometer scale and to assess individual phase properties like C–S–H gels, Portlandite or clinker. However, the properties extracted from nanoindentation are measured for small material volumes (nm to μm). The large gap between the scales can be crossed by using multiscale models and micromechanical framework which uses the concept of the representative volume element (RVE) [6] defined for each material level. Homogenization of individual contributions of the RVE microstructural components is provided by multiple micromechanical approaches that search for effective properties by solving

matrix-inclusion problems. There is a variety of analytical methods and estimates (Voigt, Reuss or Hashin–Strikmann bounds, Mori–Tanaka method, self-consistent scheme and others [6]) that usually need to assess phase properties and their volume fractions prior to the analysis. Such assessment is not straightforward in the case of structural composites whose microstructure develops in space and time during their lifetime. Therefore, statistical estimates obtained from grid nanoindentation need to be employed. The grid nanoindentation and statistical deconvolution methods have been described and used e.g. by Ulm and coworkers [1,2] for cement based materials, Němeček et al. [3] for alkali activated materials or Sorelli et al. [4] for high performance concrete.

In the case of numerical methods (e.g. finite elements or FFT based methods), homogenization can be much easier due to the direct use of grid point mechanical data as will be demonstrated later in the paper.

2. Methods

In this paper, we first deal with the evaluation of nanoindentation data received from large statistical sets (hundreds of indents) on the scale of several hundreds of micrometers which is a scale

* Corresponding author. Tel.: +420 224 354 309.

E-mail addresses: jiri.nemecek@fsv.cvut.cz (J. Němeček), vlastimil.kralik@fsv.cvut.cz (V. Králík), jaroslav.vondrej@fsv.cvut.cz (J. Vondřejc).

that includes all material phases within RVE in a sufficient content. Since the microstructure of the composites is very complex in this scale and the determination of pure individual micromechanically distinct phases is not straightforward, we assess the individual properties by using grid indentation technique [2] with subsequent statistical deconvolution method [2–4]. Mathematically, the deconvolution is an ill-posed problem that can be regularized by a prior setting of the number of mechanically different phases that are determined. Therefore, we link this number with the number of chemically different phases or groups of mechanically similar constituents as described in Section 6. We also adapt the originally proposed deconvolution method [2] by using different minimizing criteria and modified Monte Carlo simulations as described in Němeček et al. [3]. Such methodology gives us mean phase properties together with the estimation of their volume fractions based on the experimental dataset from the whole grid.

After setting the RVE size and receiving phase properties, effective elastic properties are determined by both analytical and numerical homogenization schemes. The comparison of the methods is provided by comparing the differences between the output stiffness matrices. As mentioned earlier, the application of the numerical scheme does not require the knowledge of intrinsic phase properties and the direct use of grid data is utilized.

3. Tested materials and test setup

3.1. Cement paste

Selected heterogeneous structural materials were chosen for this study. At first, cement paste samples were prepared from Portland cement CEM-I 42,5 R (locality Mokrý, CZ) with water to cement weight ratio equal to 0.5 [7]. Samples were stored in water for two years. Therefore, high degree of hydration (over 90%) can be anticipated in the samples. The microstructure of cement paste in the tested volume includes several chemical phases known from cement chemistry, namely low- and high- density calcium–silica hydrates (LD and HD C–S–H), calcium hydroxide $\text{Ca}(\text{OH})_2$, residual clinker, porosity and some other minor phases. The cement paste microstructure is shown in Fig. 1a. Very light areas in Fig. 1a can be attributed to the residual clinker, light grey areas are rich of $\text{Ca}(\text{OH})_2$, dark grey zone belongs to C–S–H gels and black color represents very low density regions or capillary porosity. Note, that C–S–H gel and $\text{Ca}(\text{OH})_2$ zones are spatially intermixed in small volumes ($\ll 10 \mu\text{m}$) and the resolution of SEM–BSE images does not allow for a direct separation of these phases from the image. The majority of the material volume mostly consists of poorly crystalline or amorphous phases (C–S–H) and partly of crystalline phases ($\text{Ca}(\text{OH})_2$). Portlandite crystals are known for their anisotropy. Since their size and volume is not large in the sample and they can be mixed with C–S–H, all phases will be supposed to be mechanically isotropic for simplification in the analysis.

Cement paste includes also wide distribution of pores. Majority of pores lies in the nanometer range ($< 100 \text{ nm}$, as checked with He/Hg-porosimetry) and, on the other hand, large capillary pores are present in the scale above the indentation level (i.e. $\gg 1 \mu\text{m}$). Therefore, the indentation depth was chosen so that the nanoporosity was included in the tested volume but the large capillary porosity was not. The depth range $\sim 100\text{--}300 \text{ nm}$ was suitable for the analysis.

Cement paste was indented by a grid consisting of $20 \times 20 = 400$ indents with $10 \mu\text{m}$ spacing which yields the RVE size $\sim 200 \mu\text{m}$. The indents were prescribed as load controlled (maximum force 2 mN, loading/unloading rate 12 mN/min, holding for 30 s). Examples of load-penetration diagrams for different

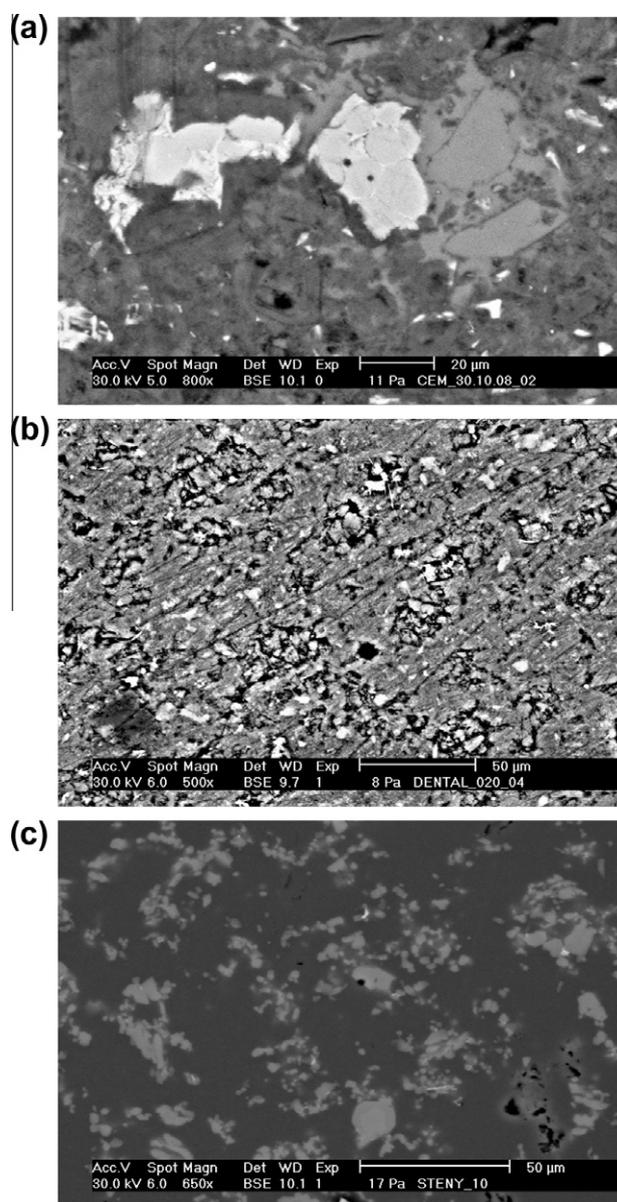


Fig. 1. Microstructures of (a) cement paste, (b) gypsum and (c) Al-alloy.

constituents are shown in Fig. 2a. The final penetration depths vary for the phases depending on their stiffness.

3.2. Gypsum

Secondly, dental gypsum (Interdent[®]) was chosen as a model representative for gypsum based materials. Samples were prepared with water to gypsum ratio 0.2 and matured in ambient conditions for two months. From the chemistry point of view, every gypsum binder is composed of three main components – calcium sulfate anhydrite (CaSO_4), calcium sulfate hemihydrate ($\text{CaSO}_4 \cdot \frac{1}{2}\text{H}_2\text{O}$) in two modifications: α - or β -hemihydrate, and calcium sulfate dihydrate ($\text{CaSO}_4 \cdot 2\text{H}_2\text{O}$). The gypsum binder includes also some impurities and additives in the case of natural sources. The Interdent gypsum is a low-porosity purified α -hemihydrate used for dental purposes.

The hardened gypsum mass is a porous material with a relatively large internal surface consisting of interlocking crystals in the form of plates and needles (Singh and Middendorf [8]). Note

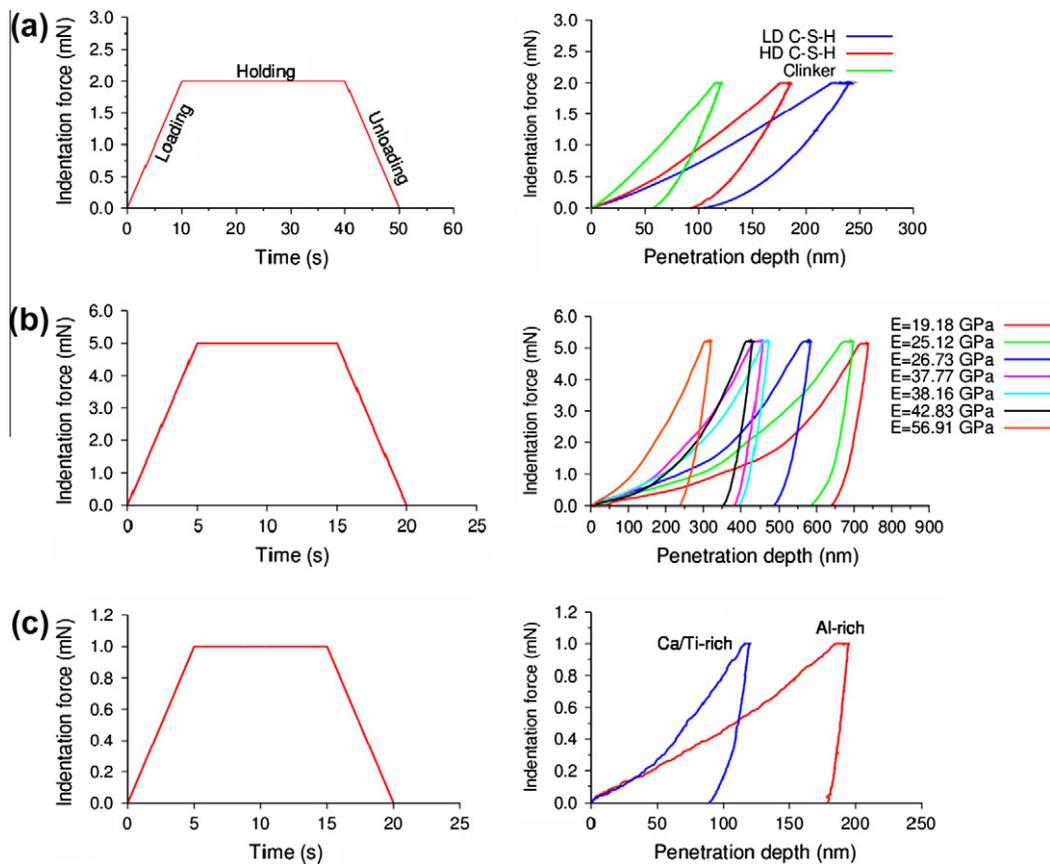


Fig. 2. Nanoindentation load–time and load–depth diagrams for (a) cement paste, (b) gypsum and (c) Al-alloy.

that in the case of β -hemihydrate hydration, the resulting sample porosity is typically very large (more than 50% for higher water to binder ratios) and crystals are interlocked very weakly. Therefore, ordinary gypsum systems used for building purposes which are based on β -hemihydrate are characterized with relatively low strengths (<10 MPa in compression). In contrast, hydration of our samples based on α -hemihydrate produced a dense matrix with total sample porosity just 19%. The majority of pores lay in the nano-range 0–300 nm (0–100 nm 7%, 100–300 nm 4%, 300–1000 nm 1%) and virtually no pores appeared between 1 and 100 μm (<0.5%). Due to the very low porosity, the strength of this material is much higher (>50 MPa in compression). The gypsum microstructure is depicted in Fig. 1b in which dark areas can be attributed to the porosity, very light areas belong to low hydrated CaSO_4 grains or carbonates. The majority of the sample volume in Fig. 1b composes of hydrated crystalline mass.

Two locations were tested on gypsum samples. Each place was covered by $15 \times 12 = 180$ indents with 15 μm spacing. Similar loading as in the case of cement was used (load controlled test to maximum force 5 mN). Typical loading diagrams are depicted in Fig. 2b. A bit wider range of final depths on indented phases (200–800 nm) was obtained due to larger differences in the phase stiffness. However, the majority of indents were performed to the mean final depths around 400–500 nm. The RVE size defined by the tested area is again $\sim 200 \mu\text{m}$ in this case.

3.3. Aluminum alloy

For the sake of comparison with different kind of material, an aluminum alloy used for the production of lightweight aluminum foams Alporas[®] was also studied [9,10]. The material is produced from an aluminum intermixed with 1.5 wt.% of Ca and 1.6 wt.%

TiH_2 . Ca/Ti-rich discrete precipitates and diffuse Al_4Ca areas develop in the metal solid [11] that can be seen as lighter areas in Fig. 1c. Therefore, two distinct phases denoted as Al-rich and Ca/Ti-rich were separated in this study.

Nanoindentation was applied to the cell walls of the foam. Loading to maximum force 1 mN was used. Final depths arrived at ~ 100 –200 nm. Typical differences between the loading diagrams of different phases obtained from nanoindentation are shown in Fig. 2c. Results from 200 indents (two locations 10×10 indents) with 10 μm spacing were evaluated. The RVE size related to the tested region is $\sim 100 \mu\text{m}$ in this case.

4. Nanoindentation, sample preparation and evaluation of phase properties

As mentioned above, nanoindentation has been applied to receive elastic constants of individual material phases. Nanoindenter (Nanohardness tester, CSM Instruments) located in Prague's laboratory at the Czech Technical University was employed in our measurements. The apparatus was equipped with a diamond pyramidal Berkovich tip with the apex radius ~ 100 nm.

The already well-known principle of nanoindentation lies in bringing a very small tip (Berkovich in our case) to the surface of the material to make an imprint. Material constants are deduced from the measured load–displacement curves performed on flat surfaces.

For our measurements, the depth of penetration was kept around ~ 300 nm for cement paste, ~ 500 nm for gypsum and ~ 200 nm for aluminum in order to capture each material phase on one hand and to minimize phase interactions on the other hand. The depth of the affected volume under the indenter tip can be estimated as $3 \times$ the penetration depth [2], i.e. around 0.6^3 –

$1.5^3 \mu\text{m}^3$ for the studied cases. Such size roughly corresponds to 1/10 of most of the grains or single phase areas which justifies the use of phase devolution [1,2]. The indentation volume contains also a part of nanoporosity that is naturally included in phase results.

All samples were mechanically polished prior to the testing in order to achieve smooth and flat surface with substantially smaller roughness compared to indentation depths. The surface roughness (evaluated as root-mean-square on the scanned area of $10 \times 10 \mu\text{m}$) was checked with AFM. It was found to be $\sim 25 \text{ nm}$ on cement paste, $\sim 38 \text{ nm}$ on gypsum and $\sim 12 \text{ nm}$ on Al-alloy. Therefore, the sample roughness was acceptable in relation to the awaited indentation depths.

The indentation loading history contained three segments: loading, holding and unloading periods. The holding period was included in order to minimize creep effects on the elastic unloading [7]. Elastic properties were evaluated for individual indents using analytical formulae derived by Oliver and Pharr [12], which account for an elasto-plastic contact of a conical indenter with an isotropic half-space. The reduced (combined) elastic modulus is then defined as:

$$E_r = \frac{1}{2\beta} \frac{\sqrt{\pi} dP}{\sqrt{A} dh} \quad (1)$$

in which A is the projected contact area of the indenter at the peak load, β is geometrical constant ($\beta = 1.034$ for the used Berkovich tip) and dP/dh is a slope of the unloading branch evaluated at the peak. Elastic modulus E of the measured sample can be found using contact mechanics which accounts for the effect of non-rigid indenter as:

$$\frac{1}{E_r} = \frac{(1 - \nu^2)}{E} + \frac{(1 - \nu_i^2)}{E_i} \quad (2)$$

in which ν is the Poisson's ratio of the tested material, E_i a ν_i are known elastic modulus and Poisson's ratio of the indenter.

The solution of the contact problem for anisotropic materials can be found in [13,14]. In this work, all material phases were treated as elastically isotropic. Such simplification was adopted due to the following reasons. In cement paste, the degree of crystallinity is poor in the majority of specimen volume (e.g. in C-S-H gel). The content of crystalline $\text{Ca}(\text{OH})_2$ phases is low and due to the limited space for the crystal growth the degree of crystallinity decreases.

On the other hand, gypsum is composed of a polycrystalline matter with locally anisotropic character. However, the response in grid nanoindentation is measured on differently oriented crystals and also on a combination of differently oriented crystals located under the indenter in the affected volume $\sim 1.5^3 \mu\text{m}^3$. The tested location can be viewed as a set of mechanically different phases that are physically averaged by an indenter. Apparent isotropic elasticity constants associated with the tested indentation volume can be derived in this case. Similarly, isotropic estimates were derived for the measured volume in case of Al-alloy disregarding the local anisotropy on a crystalline level.

The distinction of the chemically and/or mechanically different material phases is often not possible on the microlevel ($< 1 \mu\text{m}$) even with the use of SEM-EDX images. In order to receive statistically relevant data from all material phases, we applied grid indentation over the tested RVE (Fig. 1). Large matrices containing hundreds of indents have been performed on tested samples (see Section 3). To assess individual phase properties, statistical deconvolution was employed [2,3]. In this method, experimental data are analyzed from the frequency plots. Mean elastic properties as well as phase volume fraction are estimated based on the best fit of the experimental data with a limited number of Gauss distributions (Fig. 3).

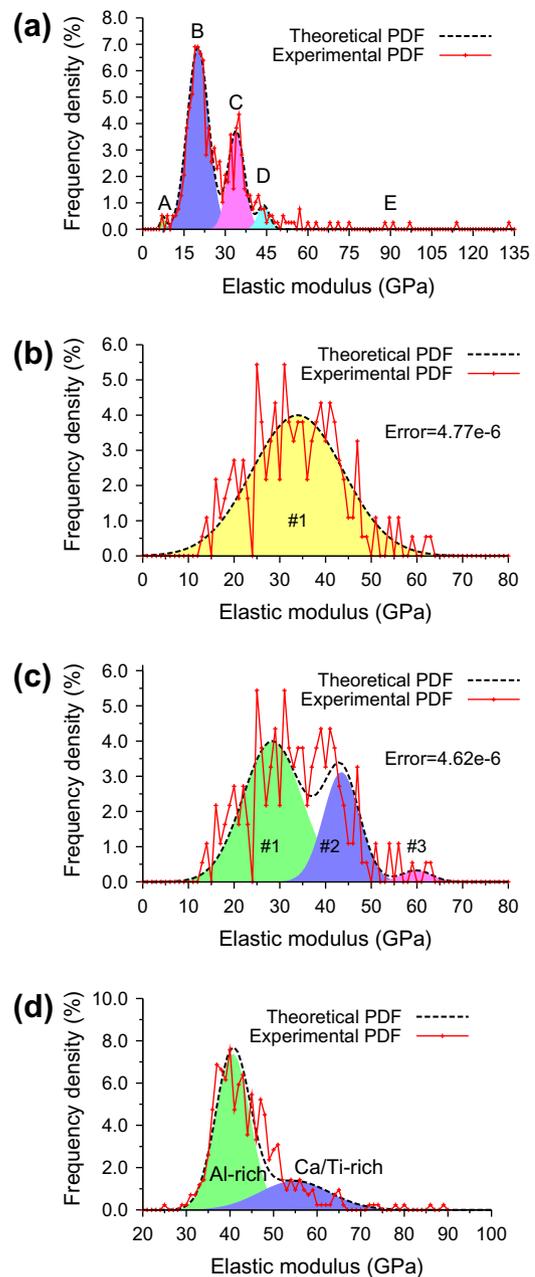


Fig. 3. Deconvolution of modulus of elasticity frequency plots into mechanical phases on (a) cement paste, (b) gypsum three phases fit (c) gypsum single phase fit and (d) Al-alloy.

5. Micromechanical homogenization

5.1. Analytical and numerical schemes

In general, homogenization methods search for effective material properties. The previously mentioned concept of RVE which includes all microstructural inhomogeneities that should be substantially smaller than the RVE size is utilized. The homogenization problem can be solved either by analytical methods or by numerical approximations. Analytical schemes often rely on simplified assumptions concerning inclusion geometry, boundary conditions or isotropy. More complex results can be obtained from numerical methods that are based on finite element solution or fast Fourier transformation (Moulinec and Suquet [17]), for instance.

The classical analytical solution based on constant stress/strain fields in individual microscale components for an ellipsoidal inclusion embedded in an infinite body was derived by Eshelby [15]. Effective elastic properties are then obtained through averaging over the local contributions. Various estimates considering different geometrical constraints or special choices of the reference medium such as the Mori–Tanaka method or the self-consistent scheme [6,16] can be used. For the case of a composite material with prevailing matrix and spherical inclusions the Mori–Tanaka method [16] was previously found to be a simple but powerful tool to estimate effective composite properties also for structural materials [2] and, therefore, it was used in this work. In the Mori–Tanaka method, the homogenized isotropic bulk and shear moduli of an r -phase composite are assessed as follows:

$$k_{\text{hom}} = \frac{\sum_r f_r k_r \left(1 + \alpha_0 \left(\frac{k_r}{k_0} - 1\right)\right)^{-1}}{\sum_r f_r \left(1 + \alpha_0 \left(\frac{k_r}{k_0} - 1\right)\right)^{-1}} \quad (3)$$

$$\mu_{\text{hom}} = \frac{\sum_r f_r \mu_r \left(1 + \beta_0 \left(\frac{\mu_r}{\mu_0} - 1\right)\right)^{-1}}{\sum_r f_r \left(1 + \beta_0 \left(\frac{\mu_r}{\mu_0} - 1\right)\right)^{-1}} \quad (4)$$

$$\alpha_0 = \frac{3k_0}{3k_0 + 4\mu_0}, \quad \beta_0 = \frac{6k_0 + 12\mu_0}{15k_0 + 20\mu_0} \quad (5)$$

where the subscript 0 corresponds to the reference medium and r corresponds to a particulate inclusion. Thus, k_0 and μ_0 are the bulk and shear moduli of the reference medium, while k_r and μ_r refer to the inclusion phases. Further, bulk and shear moduli can be recomputed to engineering values of elastic modulus and Poisson's ratio as:

$$E = \frac{9k\mu}{3k + \mu}, \quad \nu = \frac{3k - 2\mu}{6k + 2\mu} \quad (6)$$

Materials with no preference of matrix phase (e.g. polycrystalline metals) are usually modeled with the self-consistent scheme [6]. It is an implicit scheme, similar to the Mori–Tanaka method, in which the reference medium points back to the homogenized medium itself.

Local strain and stress fields in a RVE can also be found by numerical methods like finite element method or a method based on fast Fourier transformation (FFT). The later one was proven to be a reliable and computationally inexpensive method which only utilizes mechanical data in the regular grid (i.e. equidistant discretization points). Such a concept perfectly matches with the concept of the grid nanoindentation. Therefore, the FFT method was chosen for our purposes. The numerical scheme used here solves the problem of finding the effective elasticity tensor with a periodically repeating RVE by using discretization of an integral Lippmann–Schwinger equation:

$$\varepsilon(\mathbf{x}) = \varepsilon^0 - \int_{\Omega} \Gamma^0(\mathbf{x} - \mathbf{y}) : (L(\mathbf{y}) - L^0) : \varepsilon(\mathbf{y}) d\mathbf{y} \quad (7)$$

in which ε and L stand for the local strain and stiffness tensor, respectively, and ε^0 is the homogenized strain defined as a spatial average over RVE domain Ω as

$$\varepsilon^0 = \langle \varepsilon \rangle = \frac{1}{\Omega} \int_{\Omega} \varepsilon(\mathbf{x}) d\mathbf{x} \quad (8)$$

Γ^0 is the periodic Green operator associated with the reference elasticity tensor L^0 which is a parameter of the method [17,18]. The problem is further discretized using trigonometric collocation method [19,20] which leads to the assemblage of a nonsymmetrical linear system of equations:

$$[I + F^{-1} \hat{\Gamma} F (L - L^0)] e = e^0 \quad (9)$$

where the vector e stores a strain field at discretization points and e^0 the macroscopic strain, L and L^0 stores the material coefficients at discretization points and reference elasticity tensor respectively, I denotes the identity matrix, $\hat{\Gamma}$ stores the values corresponding to the integral kernel in the Fourier space, and F (F^{-1}) stores the (inverse) discrete Fourier transform matrices that can be provided by fast Fourier transform algorithm. The possibility to solve the nonsymmetric linear system by the conjugate gradient method (CG) is proposed by Zeman et al. in [21], where also the particular expression of individual matrices can be found for the problem of electric conductivity or heat transfer. The linear system (Eq. (9)) depends only on stiffness coefficients at grid points that can be obtained using nanoindentation and thus the homogenized (effective) tensor (further denoted as $L_{\text{eff}}^{\text{FFT}}$) can be calculated.

In practice, the homogenization procedure includes several steps:

- (1) Definition of a periodic unit cell (PUC) with discretization points corresponding to indents' locations (a regular grid).
- (2) Assessment of Young's moduli E and Poisson's ratio ν with the help of nanoindentation in grid points (Oliver and Pharr method [12] was used for the extraction of Young's moduli from load–displacement indentation curves).
- (3) Assemblage of local elastic stiffness tensors in grid points (plane strain assumption used) which in Mandel's notation reads:

$$L = \frac{E}{(1 + \nu)(1 - 2\nu)} \begin{bmatrix} 1 - \nu & \nu & 0 \\ \nu & 1 - \nu & 0 \\ 0 & 0 & 1 - 2\nu \end{bmatrix} \quad (10)$$

- (4) Calculation of local strain (from a linear system, Eq. (9), using CG algorithm [21]) and stress fields ($\sigma = L:e$) in grid points when applying homogeneous macroscopic strain (unit loads e^0) to the PUC domain.
- (5) Calculation of an average stress in the PUC by integration over its volume

$$\langle \sigma \rangle = \frac{1}{\Omega} \int_{\Omega} \sigma d\mathbf{x} \quad (11)$$

- (6) Calculation of the homogenized elasticity tensor for PUC from average stress and prescribed macroscopic strain

$$L_{\text{eff}}^{\text{FFT}} : e^0 = \langle \sigma \rangle \quad (12)$$

The resulting homogenized stiffness matrix for PUC must be symmetric, positive definite, but generally anisotropic. The resulting anisotropy of the matrix depends on the topology of inclusions in PUC regardless of the fact that the individual points are treated as locally isotropic. Note also, that the FFT homogenization takes no assumptions on the morphology of the phases as in the case of analytical schemes. It works only with the stiffness coefficients distributed within the PUC and its accuracy depends only on the density of the grid points.

5.2. Comparison of analytical and numerical schemes

The simple analytical methods used in this work (Mori–Tanaka, self-consistent) operate with the assumption of isotropic effective properties. Such assumption is usually acceptable for disordered structural materials. In this case, the isotropic stiffness matrix and plane strain conditions takes the form (in Mandel's notation):

$$L_{eff}^A = \frac{E_{eff}}{(1 + \nu_{eff})(1 - 2\nu_{eff})} \begin{bmatrix} 1 - \nu_{eff} & \nu & 0 \\ \nu & 1 - \nu_{eff} & 0 \\ 0 & 0 & 1 - 2\nu_{eff} \end{bmatrix}$$

$$= \begin{bmatrix} k + \frac{4}{3}\mu & k - \frac{2}{3}\mu & 0 \\ k - \frac{2}{3}\mu & k + \frac{4}{3}\mu & 0 \\ 0 & 0 & 2\mu \end{bmatrix} \quad (13)$$

in which E_{eff} and ν_{eff} are analytically computed effective Young's modulus and Poisson's ratio, respectively. Alternatively, effective bulk and shear moduli k and μ can be used for the calculation. The difference between this stiffness matrix and that received from FFT homogenization can be expressed using a stiffness error norm:

$$\delta = \sqrt{\frac{\left(L_{eff}^{FFT} - L_{eff}^A \right) :: \left(L_{eff}^{FFT} - L_{eff}^A \right)}{\left(L_{eff}^{FFT} :: L_{eff}^{FFT} \right)}} \quad (14)$$

in which L_{eff}^{FFT} is the (anisotropic) effective stiffness matrix computed by the FFT method.

To assess the degree of anisotropy of the L_{eff}^{FFT} matrix, one can use different measures. Here, we define the degree of anisotropy as:

$$\delta_{ISO} = \inf_{L_{ISO}} \sqrt{\frac{\left(L_{eff}^{FFT} - L_{ISO} \right) :: \left(L_{eff}^{FFT} - L_{ISO} \right)}{\left(L_{eff}^{FFT} :: L_{ISO} \right)}} \quad (15)$$

where the infimum is taken over all isotropic positive definite matrices. We simply calculate the upper estimate $\delta_{ISO}^{FFT} \geq \delta_{ISO}$:

$$\delta_{ISO}^{FFT} = \sqrt{\frac{\left(L_{eff}^{FFT} - L_{ISO}^{FFT} \right) :: \left(L_{eff}^{FFT} - L_{ISO}^{FFT} \right)}{\left(L_{eff}^{FFT} :: L_{ISO}^{FFT} \right)}} \quad (16)$$

by a particular choice of an isotropic matrix:

$$L_{ISO}^{FFT} = \begin{bmatrix} k_{ISO} + \frac{4}{3}\mu_{ISO} & k_{ISO} - \frac{2}{3}\mu_{ISO} & 0 \\ k_{ISO} - \frac{2}{3}\mu_{ISO} & k_{ISO} + \frac{4}{3}\mu_{ISO} & 0 \\ 0 & 0 & 2\mu_{ISO} \end{bmatrix} \quad (17)$$

with

$$\mu_{ISO} = \frac{L_{eff,33}^{FFT}}{2}, \quad k_{ISO} = \frac{L_{eff,11}^{FFT} + L_{eff,22}^{FFT}}{2} - \frac{4}{3}\mu_{ISO}$$

6. Results and discussion

The resulting frequency plot of elastic moduli measured on cement paste merged from all positions (400 indents) was deconvoluted into five mechanical phases (that correspond to chemical ones) as specified in Table 1. Note, that the values in Table 1 (and similarly in Tables 2 and 3) were found as the best fit in the minimization problem solved by the deconvolution algorithm. The bin size was set to 1 GPa in the construction of probability density functions (PDFs). Ulm et al. [1] suggested the use of cumulative density function (CDF) in the deconvolution rather than PDF. Using CDF does not require the choice of a bin size. On the other hand, using PDF is more physically intuitive and in the case of large dataset leads to similar results.

The deconvoluted phases on cement paste correspond to the peaks shown in Fig. 3a. They are denoted as A = low stiffness phase, B = low density C–S–H, C = high density C–S–H, D = Ca(OH)₂, E = clinker. In this case, the notation of mechanically distinct phases matches well with the cement chemistry. Note, that the stiffest microstructural component, the clinker, is not captured well by nanoindentation since the stiffness contrast with respect to other components is too high [2,7]. However, the content of residual clinker is very low in the case of matured paste and it does

Table 1

Data received from statistical deconvolution and homogenized values on cement paste.

Deconvoluted phase	E (GPa)	Poisson's ratio	Volume fraction
Low stiffness phase (A)	7.45	0.2	0.011
Low density C–S–H (B)	20.09	0.2	0.632
High density C–S–H (C)	33.93	0.2	0.263
Portlandite (D)	43.88	0.3	0.046
Clinker (E)	121.0 ^a	0.3	0.048
<i>Homogenization</i>			
C–S–H level (B + C) by M–T	23.36	0.2	
C–S–H level (B + C) by SCS	23.41	0.2	
Cement paste level (B + C) + A + D + E by M–T	25.39	0.207	1.0
Cement paste level (B + C) + A + D + E by SCS	25.44	0.208	1.0

M–T stands for the Mori–Tanaka scheme; SCS stands for the self-consistent scheme.

^a Note: Clinker value was adjusted to 121 GPa according to [7].

Table 2

Data received from statistical deconvolution to the three phases and homogenized values on gypsum.

Deconvoluted phase	E (GPa)	Poisson's ratio	Volume fraction
#1	28.36	0.32	0.663
#2	43.46	0.32	0.310
#3	59.89	0.32	0.027
<i>Homogenization method</i>			
M–T	32.96	0.32	1.0
SCS	33.02	0.32	1.0

Note: M–T stands for the Mori–Tanaka scheme; SCS stands for the self-consistent scheme.

Table 3

Data received from statistical deconvolution and homogenized values on Al-alloy.

Deconvoluted phase	E (GPa)	Poisson's ratio	Volume fraction
Al-rich zone	61.88	0.35	0.64
Ca/Ti-rich zone	87.40	0.35	0.36
<i>Homogenization method</i>			
M–T	70.09	0.35	1.0
SCS	70.15	0.35	1.0

Note: M–T stands for the Mori–Tanaka scheme, SCS stands for the self-consistent scheme.

not significantly influence the rest of the results. Nevertheless, the proper value of elastic modulus for homogenization was taken from ex situ measurements of clinker [7,22].

Two-step homogenization was used in the case of cement paste. Firstly, homogenized properties for the C–S–H level were obtained from low- and high-density C–S–H phases (RVE ~1 μm). Upper level homogenization for RVE (~200 μm) was performed in the second step in which homogenized C–S–H properties were considered together with the rest of the phases (i.e. low stiffness phase, Portlandite and clinker). Results for cement paste are summarized in Table 1. Very similar estimates have been obtained with the Mori–Tanaka and the self-consistent schemes.

Nanoindentation data received on gypsum samples (two locations with 180 indents each) revealed the polycrystalline nature of the composite with an anisotropic character. Since the gypsum crystals are dispersed in the sample volume in a random manner, surface measurements by nanoindentation show high scatter. As mentioned earlier, apparent isotropic moduli associated with the indentation volume ~1.5³ μm³ were assessed. The scatter in received results (Fig. 3b) can be treated as a set of mechanically different responses from different crystal orientations. As such, we

can either use deconvolution to separate mechanically significant groups of these orientations (further denoted as phases) or compute apparent elastic moduli of isotropic solid from all responses in an ensemble (i.e. compute average value from all results). Both approaches have been tested.

The physical motivation for identifying the mechanical phases lies in the fact that gypsum crystallizes in the monoclinic system which is characterized with three significant crystallographic orientations. Therefore, derivation of the three significant peaks by deconvolution of frequency plot was tested (Fig. 3b). Numerical results from this deconvolution are summarized in Table 2.

On the other hand, further simplification based on the assessment of only a single apparent isotropic phase is possible. Then, only one Gaussian distribution is assumed in the calculation. Such fit is depicted in Fig. 3c. The mean value derived from the histogram ($E = 33.90$ GPa) can be interpreted as an effective gypsum Young's modulus valid for the RVE ($\sim 100 \mu\text{m}$) which includes also intrinsic nanoporosity.

The difference between the two solutions in terms of an error computed as a sum of squared differences between the experimental and theoretical curves in the deconvolution analysis [3] is very small ($\sim 3\%$). Thus, both fits are almost equally good as indicated in Fig. 3b and c. Also, the comparison of the resulting effective Young's moduli computed by the self-consistent scheme or the Mori–Tanaka method in the case of the three phase medium (Table 2) with an apparent Young's modulus in the case of a single phase ($E = 33.90$ GPa) shows small differences ($\sim 2.7\%$).

Two mechanically distinct phases were found by the statistical deconvolution (from 200 indents) on Al-alloy sample (Fig. 3c and Table 3). According to the SEM–EDX studies, the dominant phase was denoted as Al-rich zone, whereas the lower stiffness phase was Ca/Ti-rich area. The bin size in the frequency plots was set again to 1 GPa in both cases of gypsum and Al-alloy.

Based on the nanoindentation data analytical homogenization were employed for the assessment of effective RVE elastic properties at first (Tables 1–3). Very similar results have been produced by the Mori–Tanaka method or the self-consistent scheme.

At second, the comparison of stiffness matrices (Eq. (13)) derived from analytical results (Mori–Tanaka scheme was considered for cement pastes and Al-alloy; self-consistent scheme for gypsum), and those from FFT homogenization was performed. Results are specified in the following equations 18–20. The stiffness values are given in GPa. Respective error norms are computed in Eq. (21).

$$\text{cement : } L_{eff}^A = \begin{bmatrix} 28.44 & 7.43 & 0 \\ 7.43 & 28.44 & 0 \\ 0 & 0 & 21.02 \end{bmatrix}$$

$$L_{eff}^{FFT} = \begin{bmatrix} 26.177 & 6.778 & 0.068 \\ 6.778 & 26.224 & 0.014 \\ 0.068 & 0.014 & 19.818 \end{bmatrix} \quad (18)$$

$$\text{Gypsum : 3 phase fit : } L_{eff}^A = \begin{bmatrix} 47.25 & 22.24 & 0 \\ 22.24 & 47.25 & 0 \\ 0 & 0 & 25.02 \end{bmatrix}$$

$$\text{1phase fit : } L_{eff}^A = \begin{bmatrix} 48.51 & 22.84 & 0 \\ 22.84 & 48.51 & 0 \\ 0 & 0 & 25.69 \end{bmatrix}$$

$$L_{eff}^{FFT} = \begin{bmatrix} 45.302 & 21.185 & 0.101 \\ 21.185 & 45.497 & -0.008 \\ 0.101 & -0.008 & 24.396 \end{bmatrix} \quad (19)$$

$$\text{Al-alloy : } L_{eff}^A = \begin{bmatrix} 112.479 & 60.566 & 0 \\ 60.566 & 112.479 & 0 \\ 0 & 0 & 51.913 \end{bmatrix}$$

$$L_{eff}^{FFT} = \begin{bmatrix} 117.130 & 62.741 & -0.163 \\ 62.741 & 117.106 & -0.143 \\ -0.163 & -0.143 & 54.313 \end{bmatrix} \quad (20)$$

$$\text{Errors : } \delta_{\text{cement}} = 0.08, \delta_{\text{gypsum}} = 0.07, \delta_{\text{Al-alloy}} = 0.04 \quad (21)$$

It is clear from the above equations that both simple analytical and advanced FFT-based method give comparable results in our case. The differences given by error norms for cement and gypsum (7–8%) are acceptable and show good agreement of the results received from different methods. The best agreement of the methods was reached on Al-alloy (error 4%) which can be attributed to the fact that both material phases (Al-rich, and Ca/Ti-rich zones) are even more homogeneously dispersed at microscale RVE compared to the phases that appear in cement paste or gypsum.

The upper bound of the degree of anisotropy for the FFT-based stiffness matrices was assessed by an index defined in Eq. (16) with the following results:

$$\delta_{ISO}^{\text{cement}} = 0.0132, \delta_{ISO}^{\text{gypsum}} = 0.0043, \delta_{ISO}^{\text{Al-alloy}} = 0.0016 \quad (22)$$

Low values in Eq. (22) (0.1–1.3%) show the close-to-isotropic nature of the tested materials within the specified RVE. In other words, microstructural inhomogeneities are uniformly dispersed in the RVE and consequently it also justifies the usage of analytical methods producing isotropic effective (homogenized) properties.

It must be emphasized again that although both analytical and numerical methods give similar results, there is a clear advantage of the FFT method which works directly with the grid indentation data compared to analytical Mori–Tanaka method which needs the assessment of phase properties and volume fractions. Moreover, the full stiffness matrix including possible anisotropy is captured by using the FFT method.

Comparison with macroscopic experimental values of elastic moduli for the given materials also shows good agreement with model predictions. Hydrated compound of cement paste was studied e.g. by Němeček [7] ($E = 26.4 \pm 1.8$ GPa), Constantinides and Ulm [23,24] ($E = 22.8 \pm 0.5$ GPa) or Hughes and Trtik [25] ($E = 26.5$ GPa). The values correspond well with our results ($E = 25.4$ GPa).

Gypsum elastic properties were studied e.g. by Meille and Garboczi [26,27] who estimated the plane strain values of the Young's modulus (computed as an angular average from anisotropic crystal elastic moduli tensor) as ~ 45.7 GPa. Such value was also reported for zero crystal porosity by Sanahuja et al. [28]. If one takes into account an intercrystalline porosity 12% (i.e. the gypsum nanoporosity measured for our specific case; see Section 3.2) the Young's modulus drops down to ~ 34 GPa [28] which is in excellent agreement with our homogenized value ($E = 32\text{--}33.90$ GPa).

Homogenized Al-alloy properties ($E = 70.1$ GPa) agree very well with experimental values reported e.g. by Jeon et al. [29] or Ashby et al. [30] ($E = 70$ GPa).

7. Conclusions

Nanoindentation was successfully used for the assessment of elastic parameters of intrinsic material constituents at the scale below one micrometer and effective composite properties were evaluated with analytical Mori–Tanaka, self-consistent and FFT numerical schemes for three typical structural composites with heterogeneous microstructure. Based on the micromechanical approaches and proposed methodologies we can draw the following conclusions.

- (1) It has been shown that the use of grid indentation gives access to both phase properties as well as volume fractions in the case of testing highly heterogeneous microstructures of cement paste, gypsum and Al-alloy.
- (2) Effective elastic properties of their microstructural RVEs (100–200 μm) were successfully determined with analytical Mori–Tanaka or self-consistent schemes. However, such approach assumes isotropic nature of the composite with spherical inclusions and several assumptions concerning mainly the number of mechanically different phases and bin size need to be made in the deconvolution algorithm. Therefore, an additional knowledge about the composite microstructure and its microstructural composition is necessary in this case.
- (3) Further, numerical FFT-based method was used for the assessment of effective elastic composite properties. The direct use of grid indentation data is employed in this method. The method provides effective stiffness matrix and can capture also possible anisotropy.
- (4) The performance of both analytical and numerical approaches was in good agreement for the tested materials mainly due to the close-to-isotropic nature in their RVEs.
- (5) Comparison with macroscopic experimental data also shows good correlation of measured effective values and the predicted ones.
- (6) The proposed numerical procedure for the estimation of effective elastic properties can be further applied also to other nano- or micro-heterogeneous structural composites in order to assess their anisotropic stiffness matrices or to optimize their composition.

Acknowledgment

Support of the Czech Science Foundation (P105/12/0824 and P105/12/0331) and the Grant Agency of the Czech Technical University in Prague (SGS12/116/OHK1/2T/11) is gratefully acknowledged.

References

- [1] Ulm F-J, Vandamme M, Bobko C, Ortega JA. Statistical indentation techniques for hydrated nanocomposites: concrete, bone, and shale. *J Am Ceram Soc* 2007;90(9):2677–92.
- [2] Constantinides G, Chandran KR, Ulm F-J, Vliet KV. Grid indentation analysis of composite microstructure and mechanics: principles and validation. *Mater Sci Eng: A* 2006;430(1–2):189–202.
- [3] Němeček J, Šmilauer V, Kopecký L. Nanoindentation characteristics of alkali-activated aluminosilicate materials. *Cem Concr Compos* 2011;33(2):163–70.
- [4] Sorelli L, Constantinides G, Ulm F-J, Toutlemonde F. The nano-mechanical signature of ultra high performance concrete by statistical nanoindentation techniques. *Cem Concr Res* 2008;38(12):1447–56.
- [5] Fischer-Cripps AC. *Nanoindentation*. New York: Springer Verlag; 2002.
- [6] Zaoui A. Continuum micromechanics: survey. *J Eng Mech* 2002;128:808–16.
- [7] Němeček J. Creep effects in nanoindentation of hydrated phases of cement pastes. *Mater Characterization* 2009;60(9):1028–34.
- [8] Singh NB, Middendorf B. Calcium sulfate hemihydrate hydration leading to gypsum crystallization. *Prog Crystal Growth Characterization Mater* 2007;53(1):57–77.
- [9] Miyoshi T, Itoh M, Akiyama S, Kitahara A. Aluminium foam “ALPORAS”: the production process, properties and application. *Mater Res Soc Sympos Proc* 1998;521.
- [10] Němeček J, Králík V, Vondřejc J, Němečková J. Identification of micromechanical properties on metal foams using nanoindentation. In: *Proceedings of the thirteenth international conference on civil, structural and environmental engineering computing [CD-ROM]*. Edinburgh: Civil-Comp Press; 2011. p. 1–12 [ISBN 978-1-905088-46-1].
- [11] Simone AE, Gibson LJ. Aluminum foams produced by liquid-state processes. *Acta Mater* 1998;46(9):3109–23.
- [12] Oliver WC, Pharr GM. An improved technique for determining hardness and elastic modulus using load and displacement sensing indentation experiments. *J Mater Res* 1992;7(6):1564–83.
- [13] Swadener JG, Pharr GM. Indentation of elastically anisotropic half-spaces by cones and parabolae of revolution. *Phil Mag A* 2001;81(2):447–66.
- [14] Vlassak JJ, Ciavarella M, Barber JR, Wang X. The indentation modulus of elastically anisotropic materials for indenters of arbitrary shape. *J Mech Phys Solids* 2003;51:1701–21.
- [15] Eshelby JD. The determination of the elastic field of an ellipsoidal inclusion and related problems. *Proc R Soc London A* 1957;241:376–96.
- [16] Mori T, Tanaka K. Average stress in the matrix and average elastic energy of materials with misfitting inclusions. *Acta Metall* 1973;21(5):571–4.
- [17] Moulinec H, Suquet P. A fast numerical method for computing the linear and nonlinear mechanical properties of composites. *Comptes rendus de l'Académie des sciences. Série II, Mécanique, physique, chimie, astronomie* 1994; 318(11): 1417–1423.
- [18] Moulinec H, Suquet P. A numerical method for computing the overall response of nonlinear composites with complex microstructure. *Comput Methods Appl Mech Eng* 1998;157(1–2):69–94.
- [19] Saranen J, Vainikko G. *Periodic integral and pseudodifferential equations with numerical approximation*. Berlin: Springer; 2002.
- [20] Vainikko G. Fast solvers of the Lippmann–Schwinger equation. In: Gilbert RP, Kajiwara J, Xu YS, editors. *Direct and inverse problems of mathematical physics, international society for analysis applications and computation, vol. 5*. Dordrecht (The Netherlands): Kluwer Academic Publishers; 2000. p. 423–40.
- [21] Zeman J, Vondřejc J, Novák J, Marek I. Accelerating a FFT-based solver for numerical homogenization of periodic media by conjugate gradients. *J Comput Phys* 2010;229(21):8065–71.
- [22] Velez K et al. Determination of nanoindentation of elastic modulus and hardness of pure constituents of Portland cement clinker. *Cem Concr Res* 2001;31:555–61.
- [23] Constantinides G, Ulm F-J. The nanogranular nature of C–S–H. *J Mech Phys Solids* 2007;55:64–90.
- [24] Constantinides G, Ulm F-J. The effect of two types of C–S–H on the elasticity of cement-based materials: results from nanoindentation and micromechanical modeling. *Cem Concr Res* 2004;34(1):67–80.
- [25] Hughes JJ, Trtik P. Micro-mechanical properties of cement paste measured by depth-sensing nanoindentation: a preliminary correlation of physical properties with phase type. *Mater Characterization* 2004;53:223–31.
- [26] Meille S, Garboczi EJ. Linear elastic properties of 2D and 3D models of porous materials made from elongated objects. *Modell Simul Mater Sci Eng* 2001;9(5):371–90.
- [27] Garboczi et al. *Modeling and measuring the structure and properties of cement-based materials*, National Institute of Standards and Technology; 1989–2011. <<http://ciks.cbt.nist.gov/garbocz/monograph>>.
- [28] Sanahuja J, Dormieux L, Meille S, Hellmich C, Fritsch A. Micromechanical explanation of elasticity and strength of gypsum: from elongated anisotropic crystals to isotropic porous polycrystals. *J Eng Mech* 2010;136(2):239–53.
- [29] Jeon I et al. Cell wall mechanical properties of closed-cell Al foam. *Mech Mater* 2009;41:60–73.
- [30] Ashby MF, Evans AG, Fleck NA, Gibson LJ, Hutchinson JW, Wadley HNG. *Metal Foams: A Design Guide*. Butterworth-Heinemann; 2000.

Part V

Paper 4

Authors:

Jiří Němeček, Vlastimil Králík, and Jaroslav Vondřejc

Title:

A Two-scale micromechanical model for aluminium foam based on results from nanoindentation

A Two-scale micromechanical model for aluminium foam based on results from nanoindentation

Jiří Němeček¹, Vlastimil Králík¹, Jaroslav Vondřejc¹

¹Department of Mechanics, Faculty of Civil Engineering

Czech Technical University in Prague, Czech Republic

Abstract

The main aim of this paper is to develop and verify simple but effective model for assessing elastic properties of a porous aluminium foam system and to compare results received from experimental micromechanics with solutions given by simple analytical or more advanced numerical methods. The material is characterized by a closed pore system with very thin but microscopically inhomogeneous pore walls (~0.1 mm) and large air pores (~2.9 mm). Therefore, two material levels can be distinguished.

The lower level of the proposed model contains inhomogeneous solid matter of the foam cell walls produced from aluminium melt with admixtures. Elastic parameters as well as volume fractions of microstructural material phases at this level are assessed with nanoindentation and effective properties computed via analytical and numerical homogenization schemes. The effective Young's modulus of the cell walls was found close to 70 GPa irrespective to the used homogenization procedure.

The higher foam scale contains homogenized cell wall properties and a significant volume fraction of air voids (91.4%). Since analytical schemes fail to predict effective properties of this highly porous structure, numerical homogenization based on simple 2-D finite element model is utilized. The model geometry is based on foam optical images and the beam structure is produced using Voronoi tessellation. Effective foam Young's modulus was found to be 1.36-1.38 GPa which is in relation with ~1.45 GPa obtained from uniaxial compression experiments. The stiffness underestimation in the 2-D model is caused likely by the lack of the real 3-D confinement that can not be fully captured in the simplified model.

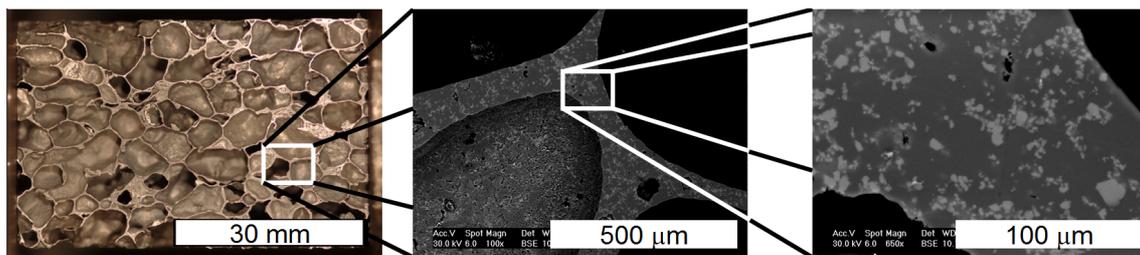
Keywords: aluminium foam, multi-scale model, nanoindentation, statistical deconvolution, elastic properties, image analysis, homogenization.

1 Introduction

Metal foams and especially lightweight aluminium foams belong to the group of up-to-date engineering materials with high potential to many applications. Metal foam is a highly porous hierarchical material with a cellular microstructure. Macroscopically, it can be characterized by attractive mechanical and physical properties such as high stiffness and strength in conjunction with very low weight, excellent impact energy absorption, high damping capacity and good sound absorption capability. The usual source material for the production of metal foams are aluminium and aluminium alloys because of low specific density ($\sim 2700 \text{ kg/m}^3$), low melting point ($\sim 660 \text{ }^\circ\text{C}$), non-flammability, possibility of recycling and excellent corrosion resistance. Metal foams

are used in applications ranging from automotive and aerospace industries (e.g. bumpers, car body sills, motorcycle helmets) to building industry (e.g. sound proofing panels). Our aim was to characterize and to model a commercially available foam Alporas[®] produced by Shinko Wire Company, Ltd. .

Alporas is characterized with a hierarchical system of pores containing different cell morphologies (in shape and size) in dependence on the foam density and inhomogeneous material properties of the cell walls . A typical cross section of the foam can be seen in Figure 1 in which large pores (having typically 1-13 mm in diameter) with detailed view on thin walls ($\sim 100 \mu\text{m}$ thick) is shown.



(a)

(b)

(c)

Figure 1: (a) Overall view on a foam structure (further denoted as Level II); (b) ESEM image of a cell wall; (c) detailed ESEM image of a cell wall showing Al-rich (dark grey) and Ca/Ti-rich areas (light zones; denoted as Level I).

It follows from its hierarchical microstructure that the mechanical properties of metal foams are governed by two major factors:

- (i) cell morphology (shape, size and distribution of cells) and
- (ii) material properties of the cell walls .

Traditionally, mechanical properties of metal foams are obtained using conventional macroscopic testing techniques on large samples that can give overall (effective) properties, e.g. -. However, conventional measurements face significant obstacles in the form of very small dimensions of cell walls, low local bearing capacity, local yielding and bending of the cell walls. These problems can be overcome using micromechanical experimental methods in which the load–displacement curve is obtained in the sub-micrometer range. A few attempts have been carried out in the past, e.g. , .

In this paper, we focused on the prediction of overall foam elastic properties from microscopic measurements and on the model validation against experimental results. For accessing the cell wall properties we employed statistical nanoindentation and deconvolution technique for the phase separation ,. Compared to traditional macroscopic techniques nanoindentation can distinguish between individual inhomogeneous microstructural entities. The effective cell wall properties have been obtained through analytical and numerical up-scaling techniques . Finally, simple 2-D finite element model for the upper composite scale has been proposed and results validated by full-scale experiments.

2 Experimental part

2.1 Materials and sample preparation

Commercial aluminium foam Alporas[®] (Shinko Wire Company, Ltd) was used in this study. The manufacturing process of the Alporas is a batch casting process in which 1.5 wt.% of calcium is added to the aluminium molten at 680 °C. Calcium serves as a

thickening agent which increases viscosity and stabilizes the air bubbles. The alloy is poured into a casting mold and stirred with an admixture of 1.6 wt.% TiH₂ that is used as a blowing agent. Then, the foamed molten material is cooled down. A typical resulting internal structure of the aluminium foam is shown in Figure 1a.

Firstly, a large panel of Alporas (160 × 100 × 60 mm) was polished and scanned with a high resolution scanner. Acquired images were segmented to binary ones and further used in image analyses. Then, a smaller Alporas block was cut into thin slices (~5 mm) and embedded into epoxy resin to fill the pores. The surface was mechanically grinded and polished to reach minimum surface roughness suitable for nanoindentation. Very low roughness $R_q \approx 10$ nm was achieved on cell walls. The sample was investigated with electron microscopy (ESEM) and nanoindentation.

2.2 ESEM and microstructural analysis

The microstructure of cell walls was firstly studied in electron microscope (ESEM). It was found that a significant inhomogeneity of the microstructural material phases exists on the level of tens of micrometers (Figure 1b,c). Two distinct phases, that exhibit different color in back-scattered electron (BSE) images, can be distinguished. The chemical composition of the two phases was checked with EDX element analysis in ESEM. It was found that the majority of the volume (dark zone in Figure 1c, 2a) consists of aluminium (~67 wt.%), oxygen (~32 wt.%) and further trace elements (Mg, Ti, Fe, Co, Ni, Cu, Si <2 wt.%). Lighter zones in Figure 2 consist of Al (~60 wt.%), O (~30 wt.%), Ca (~5 wt.%), Ti (~5 wt.%) and other elements (<1 wt.%). As expected, the majority of the volume (dark zone) is composed of aluminum and aluminium oxide

Al_2O_3 (further denoted as Al-rich area). Lighter zones contain significant amount of calcium and titanium (further denoted as Ca/Ti-rich area). The non-uniform distribution of these zones shows inhomogeneous mixing of the admixtures that are added during the production process.

2.3 Image analysis and porosity

In order to estimate the volume fractions of Al-rich and Ca/Ti-rich areas image analysis based on previously taken ESEM images was employed. Ten arbitrarily chosen areas on wall cross sections were explored. Images were segmented to two phases using a common threshold value of a grey level for all images (Figure 2). The Ca/Ti-rich area was estimated to cover $22\pm 4\%$ of the whole area.

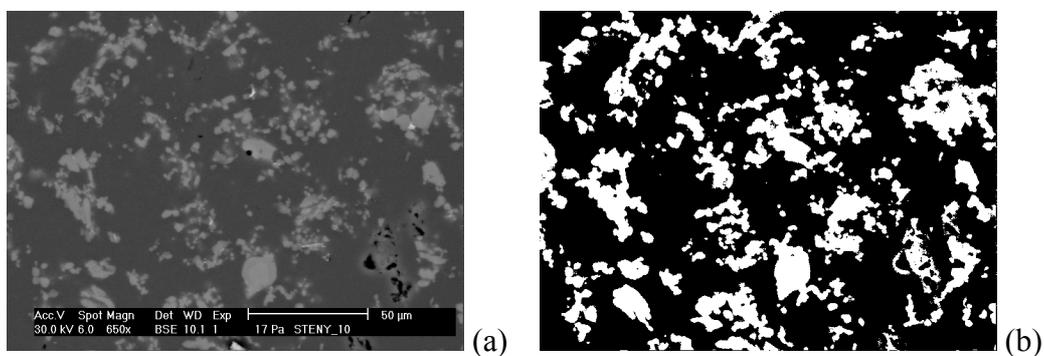


Figure 2: An example of (a) ESEM image of the cell wall and (b) processed image segmented to two phases (white=Ca/Ti-rich, black=Al-rich area).

The overall porosity of the sample was assessed by weighing of a large Alporas panel (knowing the sample dimensions and solid mass density $2700 \text{ kg}\cdot\text{m}^{-3}$). The

porosity reached 91.4% which corresponds to e.g. . In other words, solid mass (i.e. cell walls) occupied only 8.6% of the total volume in the specimen.

Further, distribution of cell wall thicknesses and the distribution of the pore sizes were studied by means of pore contour detection in the Matlab environment. At first, the contours were generated for every pore in the image and areal characteristics (centroid, area, second moment of inertia) were computed (Figure 3). The wall thicknesses were calculated as the minimum distance between the neighboring contours. The distribution of the thicknesses is shown in Figure 4. It can be seen in Figure 4 that a significant peak occurs around $\sim 60 \mu\text{m}$ which can be understood as a characteristic cell wall thickness.

Then, equivalent ellipses were constructed from contour areal characteristics under the condition that they have the same area and the same principal second moment of inertia. Such assumption led to the evaluation of two main half axes (a_i and b_i) for each equivalent ellipse. In order to characterize the shape of pores, an equivalent ellipse

shape factor was defined as the ratio $e_i = \frac{a_i}{b_i}$. The distribution of the shape factor is depicted in Figure 5. It can be concluded that pores have typically a round shape with the shape factor lying mostly between 1 and 2. The peak with the highest occurrence in Figure 5 appears around $e_i=1.15$.

Due to the round shape of pores it makes sense to compute also an equivalent pore diameter using circular pore replacement. The distribution of equivalent circular pores is depicted in Figure 6. Wide distribution of pores with diameters 0-6 mm was found. The mean equivalent diameter was found to be 2.9 ± 1.5 mm for the specific specimen.

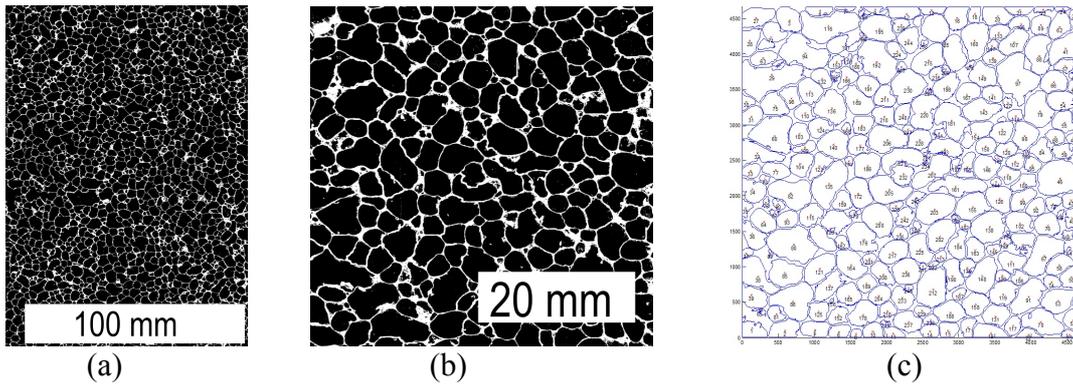


Figure 3: (a) Binary image of the polished foam panel. (b) Binary image of $\sim 50 \times 50$ mm foam cut. (c) Cell contours in the cut (prepared in Matlab).

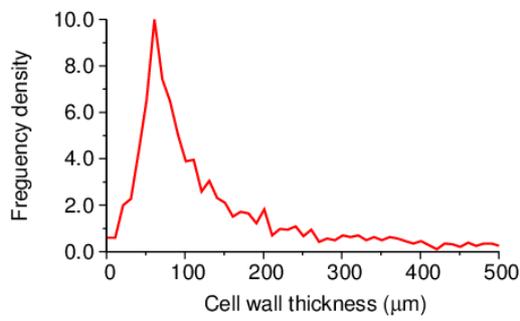


Figure 4: Distribution of cell wall thicknesses.

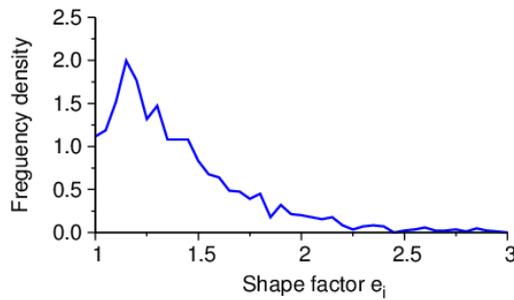


Figure 5: Distribution of equivalent ellipse shape factor.

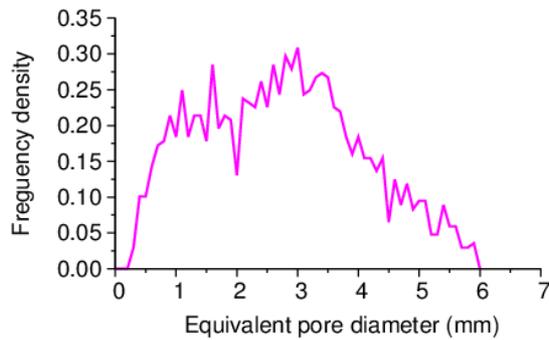


Figure 6: Distribution of equivalent circular pores.

2.4 Nanoindentation

Micromechanical properties of the cell walls were measured by means of nanoindentation. The tests were performed using a Hysitron Tribolab system[®] at the Czech Technical University in Prague. This system consists of in-situ SPM imaging which was used for scanning the sample surface. Three-sided pyramidal diamond tip (Berkovich type) was used for all measurements. Two distant locations were chosen on the sample to capture its heterogeneity. Each location was covered by a series of 10×10 indents with 10 μm spacing (Figure 7). It yields 200 indents in total which was considered to give sufficiently large statistical set of data. Standard load controlled test of an individual indent consisted of three segments: loading, holding at the peak and unloading. Loading and unloading of this trapezoidal loading function lasted for 5 seconds, the holding part lasted for 10 seconds. Maximum applied load was 1 mN. Maximum indentation depths were ranging between 100 and 300 nm depending on the stiffness of the indented phase. Elastic modulus was evaluated for individual indents using standard Oliver and Pharr methodology which accounts for elasto-plastic contact of a conical indenter with an isotropic half-space as

$$E_r = \frac{1}{2\beta} \frac{\sqrt{\pi}}{\sqrt{A}} \frac{dP}{dh} \quad (1)$$

in which E_r is the reduced modulus measured in an experiment, A is the projected contact area of the indenter at the peak load, β is geometrical constant ($\beta=1.034$ for the

used Berkovich tip) and $\frac{dP}{dh}$ is a slope of the unloading branch evaluated at the peak.

Elastic modulus E of the measured media can be found using contact mechanics which accounts for the effect of non-rigid indenter as

$$\frac{1}{E_r} = \frac{(1-\nu^2)}{E} + \frac{(1-\nu_i^2)}{E_i} \quad (2)$$

in which ν is the Poisson's ratio of the tested material, E_i a ν_i are known elastic modulus and Poisson's ratio of the indenter. In our case, $\nu=0.35$ was taken as an estimate for all indents.

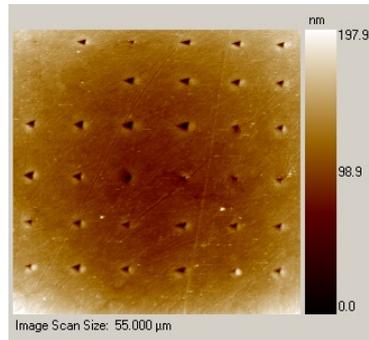


Figure 7: Part of the indentation matrix showing 6×6 indents with 10 μm spacing as scanned with Hysitron Tribolab.

Results of elastic moduli in the form of histograms have been further analyzed with the statistical deconvolution technique , . The technique seeks for parameters of individual phases included in overall results. It searches for n-Gauss distributions in an experimental Probability Density Function - PDF (Figure 9). Random seed and

minimizing criteria of the differences between the experimental and theoretical overall PDFs (particularly quadratic norm of the differences) are computed in the algorithm to find the best fit. Details on the deconvolution technique can be found in . Two-phase system (one dominant Al-rich phase and one minor Ca/Ti-rich phase) was assumed in our deconvolution.

3 Numerical part

3.1 Scale separation

In order to describe heterogeneous systems and their effective properties in a statistical sense, representative volume element (RVE) have been previously introduced . RVE statistically represents a higher structural level of the material and serves for evaluation of the effective (homogenized) properties within the defined volume. It includes all microstructural inhomogeneities that should be substantially smaller than the RVE size. The definition of the material scales can be defined through the scale separation inequality:

$$d \ll L \ll D \quad (3)$$

in which d is the characteristic size of the largest microstructural inhomogeneity, L is the RVE size and D is a characteristic structural length scale. Knowing the material and geometrical properties of the microstructural material phases a homogenization can be performed.

Nanoindentation is able to access intrinsic material properties of individual micro-scale phases provided the dimension of an indent (h) is small enough, i.e. $h \ll d$. As a rule of thumb $h < d/10$ is usually used to access material properties of individual constituents without any dependence on the length scale.

As mentioned above, the metal foam material has a hierarchical microstructure. At least two levels need to be considered:

- **Level I** (the cell wall level) has a characteristic dimension defined by the mean midspan wall thickness $L_I \sim 60 \mu\text{m}$. This level consists of prevailing aluminium matrix (Al-rich area) with embedded heterogeneities in the form of Ti/Ca-rich areas. Intrinsic elastic properties of the constituents were assessed by nanoindentation at this level. Individual indent size was prescribed to be considerably smaller ($h \sim 100\text{-}300 \text{ nm}$) than a characteristic size of Ca/Ti inhomogeneities ($\sim 4 \mu\text{m}$).
- **Level II** (the foam level) has a characteristic dimension of $L_{II} \sim 50 \text{ mm}$. At this level, large pores with an average equivalent diameter $\sim 2.9 \text{ mm}$ (assuming circular pores) occur in the total volume of 91.4%. At level II, cell walls are considered as homogeneous having the properties that come from the Level I homogenization.

3.2 Analytical homogenizations of Level I

The RVE with substantially smaller dimensions than the macroscale body allows imposing homogeneous boundary conditions over the RVE. Continuum micromechanics provides a framework, in which elastic properties of heterogeneous

microscale phases are homogenized to give overall effective properties on the upper scale . A significant group of analytical homogenization methods relies on the Eshelby's solution that is derived for ellipsoidal inclusions embedded in an infinite body. Then, uniform stress field appears in inclusions when macroscopic load is applied in infinity. Effective elastic properties are obtained through averaging over the local contributions.

From the material point of view, composite materials are usually characterized by a prevailing matrix phase, which serves as a reference medium in homogenization methods, reinforced with geometrically distinguishable inclusions. For example, the Mori-Tanaka method can be appropriate for these cases. In this method, the effective bulk k_{eff} and shear μ_{eff} moduli of the composite are computed as follows

$$k_{eff} = \frac{\sum_r f_r k_r (1 + \alpha_0 (\frac{k_r}{k_0} - 1))^{-1}}{\sum_r f_r (1 + \alpha_0 (\frac{k_r}{k_0} - 1))^{-1}}, \quad (4)$$

$$\mu_{eff} = \frac{\sum_r f_r \mu_r (1 + \beta_0 (\frac{\mu_r}{\mu_0} - 1))^{-1}}{\sum_r f_r (1 + \beta_0 (\frac{\mu_r}{\mu_0} - 1))^{-1}}, \quad (5)$$

$$\alpha_0 = \frac{3k_0}{3k_0 + 4\mu_0}, \beta_0 = \frac{6k_0 + 12\mu_0}{15k_0 + 20\mu_0} \quad (6)$$

where f_r is the volume fraction of the r^{th} phase, k_r its bulk modulus and the coefficients α_0 and β_0 describe bulk and shear properties of the 0th phase, i.e. the reference medium , . The bulk and shear moduli can be directly linked with Young's modulus E and Poisson's ratio ν used in engineering computations as

$$E = \frac{9k\mu}{3k + \mu}, \quad (7)$$

$$\nu = \frac{3k - 2\mu}{6k + 2\mu}. \quad (8)$$

Polycrystalline metals, in which no preference of matrix phase exists, are usually modeled with the self-consistent scheme in which the reference medium refers back to the homogenized medium itself. Regardless the most suitable homogenization technique, which would be probably the Mori-Tanaka method in our case, we used multiple estimates assuming spherical inclusions. Namely, the Mori-Tanaka method, self-consistent scheme, Voigt and Reuss bounds (parallel or serial configuration of phases with perfect bonding). Results from nanoindentation have been used as input parameters to the methods.

3.3 Numerical homogenization of Level I based on FFT

In order to verify results from simple analytical schemes advanced homogenization method based on fast Fourier transformation (FFT) was used. The behavior of any heterogeneous materials consisting of periodically repeating RVE (occupying domain $\Omega = \prod_{i=1}^d (-Y_i, Y_i)$, where Y_i is the axial size and d denotes the space dimension) can be described with differential equations with periodic boundary conditions and prescribed macroscopic load as

$$\langle \varepsilon \rangle := \frac{1}{|\Omega^*|} \int_{\Omega^*} \varepsilon(\mathbf{x}) d\mathbf{x} = \varepsilon^0 \quad (9)$$

$$\boldsymbol{\sigma}(\mathbf{x}) = \mathbf{L}(\mathbf{x}) : \boldsymbol{\varepsilon}(\mathbf{x}) \quad \text{div}\boldsymbol{\sigma}(\mathbf{x}) = \mathbf{0} \quad \mathbf{x} \in \Omega \quad (10)$$

where $\boldsymbol{\sigma}$ denotes symmetric second order stress tensor, $\boldsymbol{\varepsilon}$ symmetric second order strain tensor and $\mathbf{L}(\mathbf{x})$ the fourth order tensor of elastic stiffness at individual locations \mathbf{x} . The effective (homogenized) material tensor \mathbf{L}_{eff} is such a tensor satisfying

$$\langle \boldsymbol{\sigma} \rangle = \mathbf{L}_{\text{eff}} \langle \boldsymbol{\varepsilon} \rangle \quad (11)$$

for an arbitrary macroscopic strain $\boldsymbol{\varepsilon}^0$. Thus the problem of finding effective material tensor is composed of finding corresponding strain field $\boldsymbol{\varepsilon}$ and associated stress field $\boldsymbol{\sigma}$ for known elastic properties \mathbf{L} and prescribed strain $\boldsymbol{\varepsilon}^0$ using differential Equation 10.

In addition to discretization of the weak formulation leading to classical finite element method, the problem can be solved by method based on fast Fourier transform, proposed by Moulinec and Suquet in , based on an integral (Lippmann–Schwinger) equation

$$\boldsymbol{\varepsilon}(\mathbf{x}) + \int_{\Omega} \Gamma^0(\mathbf{x} - \mathbf{y}) : (\mathbf{L}(\mathbf{y}) - \mathbf{L}^0) : \boldsymbol{\varepsilon}(\mathbf{y}) d\mathbf{y} = \boldsymbol{\varepsilon}^0 \quad (12)$$

where Γ^0 is the periodic Green's operator associated with the reference elasticity tensor \mathbf{L}^0 which is a parameter of the method. The operator is expressed in the Fourier space as

$$\hat{\Gamma}_{ijkl}^0(\boldsymbol{\xi}) = \frac{1}{4\mu|\boldsymbol{\xi}|^2} (\delta_{ki}\xi_l\xi_j + \delta_{li}\xi_k\xi_j + \delta_{kj}\xi_l\xi_i + \delta_{lj}\xi_k\xi_i) - \frac{\lambda + \mu}{\mu(\lambda + 2\mu)} \frac{\xi_i\xi_j\xi_k\xi_l}{|\boldsymbol{\xi}|^4} \quad (13)$$

Numerical solution of Equation 13 is based on the discretization of a unit cell Ω into a regular periodic grid with $N_1 \times \dots \times N_d$ nodal points and grid spacings

$\mathbf{h} = (2 \frac{Y_1}{N_1}, \dots, 2 \frac{Y_d}{N_d})$. The searched field $\boldsymbol{\varepsilon}$ is approximated by a trigonometric

polynomial \mathbf{e}_N in the form

$$\mathbf{e}(\mathbf{x}) \approx \mathbf{e}_N(\mathbf{x}) = \sum_{\mathbf{k} \in \mathbf{Z}_N^d} \hat{\mathbf{e}}(\mathbf{k}) \varphi_{\mathbf{k}}(\mathbf{x}), \mathbf{x} \in \Omega, \quad (14)$$

where $\mathbf{N} = (N_1, \dots, N_d)$, $\hat{\mathbf{e}}$ designates the Fourier coefficients and

$$\mathbf{Z}_N^d = \left\{ \mathbf{k} \in \mathbf{Z}^d : -\frac{N_\alpha}{2} < k_\alpha \leq \frac{N_\alpha}{2}, \alpha = 1, \dots, d \right\}. \quad (15)$$

The discretization leads to a nonsymmetrical linear system of equations

$$[\mathbf{I} + \mathbf{F}^{-1} \hat{\Gamma} \mathbf{F} (\mathbf{L} - \mathbf{L}^0)] \mathbf{e} = \mathbf{e}^0 \quad (16)$$

where the vector \mathbf{e} stores a strain field at discretization points and \mathbf{e}^0 the macroscopic strain, \mathbf{L} and \mathbf{L}^0 stores the material coefficients at discretization points and reference elasticity tensor respectively, \mathbf{I} denotes the identity matrix, $\hat{\Gamma}$ stores the values corresponding to the integral kernel in the Fourier space, and \mathbf{F} (\mathbf{F}^{-1}) stores the

(inverse) discrete Fourier transform matrices that can be provided by fast Fourier transform algorithm. The possibility to solve the nonsymmetric linear system by the conjugate gradient method is proposed by Zeman et.al. in and justified in Vondřejc et al. , where also the particular expression of individual matrices can be found for the problem of electric conductivity or heat transfer. The linear system (Equation 16) depends only on stiffness coefficients at grid points that can be obtained using nanoindentation and thus the homogenized (effective) tensor $\mathbf{L}_{\text{eff}}^{\text{FFT}}$ can be calculated from Equation 11. The particular case of homogenization of elastic properties received from nanoindentation on a sample surface (half-space) also requires an assumption of plane strain conditions.

3.4 Comparison of analytical and numerical schemes

The comparison of analytical and FFT schemes includes an assessment of the stiffness matrix (here in Mandel's notation) for isotropic material assuming plane strain conditions (equally with the FFT scheme) as

$$\mathbf{L}_{\text{eff}}^{\text{A}} = \frac{E_{\text{eff}}}{(1 + \nu_{\text{eff}})(1 - 2\nu_{\text{eff}})} \begin{bmatrix} 1 - \nu_{\text{eff}} & \nu & 0 \\ \nu & 1 - \nu_{\text{eff}} & 0 \\ 0 & 0 & 1 - 2\nu_{\text{eff}} \end{bmatrix}. \quad (17)$$

The difference between the analytical results ($\mathbf{L}_{\text{eff}}^{\text{A}}$) and numerically computed stiffness matrix ($\mathbf{L}_{\text{eff}}^{\text{FFT}}$) can be expressed using a stiffness error norm as

$$\delta = \sqrt{\frac{\left(\mathbf{L}_{\text{eff}}^{\text{FFT}} - \mathbf{L}_{\text{eff}}^{\text{A}} \right) \therefore \left(\mathbf{L}_{\text{eff}}^{\text{FFT}} - \mathbf{L}_{\text{eff}}^{\text{A}} \right)}{\left(\mathbf{L}_{\text{eff}}^{\text{FFT}} \therefore \mathbf{L}_{\text{eff}}^{\text{FFT}} \right)}}. \quad (18)$$

4 Results and discussion

4.1 Nanoindentation

Results from nanoindentation clearly showed heterogeneity of the cell walls, i.e. the presence of mechanically different inclusions. An example of typical loading diagrams gained from nanoindentation at Al-rich area (dark zone in Figure 1b,c) and Ca/Ti-rich area (light zone in Figure 1b,c) are shown in Figure 8. Due to the load controlled nanoindentation test, the final penetration depth varied for differently stiff phases. An average maximum depth of penetration reached by the indenter was around ~180 nm. Higher values for more compliant Al-rich zone were reached (~190 nm) whereas the indentation depths to harder but less frequent Ca/Ti-rich areas were around 100 nm.

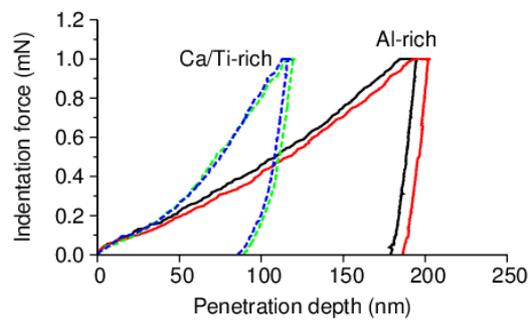


Figure 8: Typical loading diagrams for Al-rich and Ca/Ti-rich zones.

Elastic moduli were evaluated for each individual indent. Overall results are depicted in Figure 9a in which histogram of all elastic moduli from two different positions and results merged from both positions are shown. No significant differences between the positions were found. Therefore, merged results were further used in the deconvolution of phase elastic properties.

Two-phase system (one dominant Al-rich phase and one minor Ca/Ti-rich phase) was assumed in the deconvolution algorithm (Figure 9b). It can be seen in Figure 9b that a significant peak appears around 62 GPa. This value can be considered as a dominant characteristic of the prevailing phase (Al-rich). The rest of results can be attributed to the minor Ca/Ti-rich phase. Table 1 contains numerical results from the deconvolution with the estimated volume fractions of the phases.

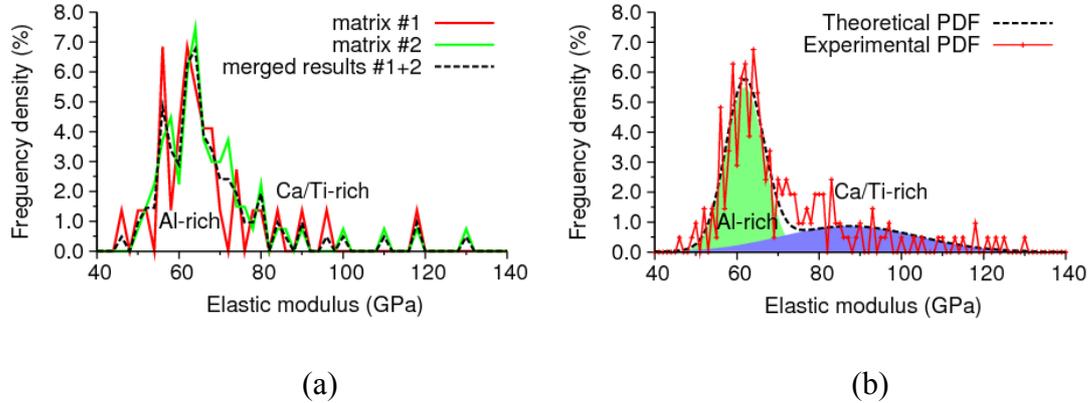


Figure 9: (a) Probability density functions of elastic moduli from two measured positions and (b) merged results with deconvoluted phases.

Table 1: Elastic moduli and volume fractions from deconvolution.

Phase	Mean (GPa)	St. dev. (GPa)	Volume fraction (-)
1 (Al-rich zone)	61.88	4.6	0.638
2 (Ca/Ti-rich zone)	87.40	16.7	0.362

The characteristic value for the first phase roughly corresponds to the elastic modulus of pure aluminium (70 GPa, ref.). The lower value obtained from nanoindentation suggests that probably some small-scale porosity or impurities (Ca) added to the molten are intrinsically included in the results of this mechanically dominant phase. The determined elastic modulus value of Al-rich zone is also in excellent agreement with the value 61.7 GPa measured by Jeon et al. on melted Al-1.5 wt.%Ca alloy.

4.2 Level I homogenization

It is clear from ESEM images (Figure 1 and 2) that the Ca/Ti-rich areas occupy much larger space of the solid compared to the initial batch volume fractions (Ca and TiH₂ content is less than 1 vol.%). Chemical reactions and precipitation during hardening form new compounds in the Al matrix. It follows from other studies that Ca/Ti-rich discrete precipitates and diffuse Al₄Ca areas develop in the metal solid. These areas are denoted as Ca/Ti-rich areas in this study. Based on the color in ESEM images the area is estimated as 22±4% by image analysis. Results from statistical nanoindentation (36.2%) suggest that a substantially larger part of the matrix is mechanically influenced by the Ca/Ti addition and a higher fraction of the volume belongs to this mechanically distinct phase.

The homogenized elastic modulus for the two considered microscale phases in the cell wall (i.e. Level I) is summarized in Table 2 for individual homogenization techniques. Very close bounds and insignificant differences in the elastic moduli estimated by the schemes were found.

Table 2: Values of the Level I effective Young's modulus computed by different homogenization schemes.

Scheme	Mori-Tanaka	Self-consist. scheme	Voigt bound	Reuss bound
E (GPa)	70.083	70.135	71.118	69.195

In the following considerations, we use the result received from the Mori-Tanaka scheme, i.e. we take the homogenized isotropic elastic constants (Young's modulus and Poisson's ratio) of the Level I as $E_{eff,I}=70.083$ GPa, $\nu_{eff,I}=0.35$. The stiffness matrices computed in Mandel's notation from analytical Mori-Tanaka results (using Equations 4-8) and from FFT homogenization are:

$$\mathbf{L}_{eff}^A = \begin{bmatrix} 112.479 & 60.566 & 0 \\ 60.566 & 112.479 & 0 \\ 0 & 0 & 51.913 \end{bmatrix} \text{ (GPa)} \quad (19)$$

$$\mathbf{L}_{eff}^{FFT} = \begin{bmatrix} 117.1300 & 62.7413 & -0.1625 \\ 62.7413 & 117.1060 & -0.1430 \\ -0.1625 & -0.1430 & 54.3132 \end{bmatrix} \text{ (GPa)}. \quad (20)$$

It is worth noting that the analytical form of the stiffness matrix \mathbf{L}_{eff}^A as well as \mathbf{L}_{eff}^{FFT} contains perfect symmetry by definition. Further, low values of off-axis components associated with shear strains in \mathbf{L}_{eff}^{FFT} (that are zero in case of \mathbf{L}_{eff}^A) show close to isotropic nature of the material. In other words, microstructural inhomogeneities are uniformly dispersed in the RVE. Consequently, this finding also

justifies the usage of analytical methods producing isotropic effective (homogenized) properties. Also, the stiffness error evaluated by Equation 18 is $\delta=0.0393$. The difference in schemes less than 4% shows a very good agreement of the methods for the studied case.

4.3 Level II homogenization

At this level, cell walls are considered as a homogeneous phase having the properties that come from the Level I homogenization. The solid phase is very sparse in the sample volume due to its porosity (91.4% of air). The walls create a matrix phase and the large air pores can be considered as inclusions in this homogenization.

Since analytical approaches are often used also for extreme cases of large stiffness contrast of phases or for large sample porosities (e.g. Šejnoha et al.), we firstly tried to estimate effective elastic properties with the same analytical schemes used at Level I. The result is summarized in Table 3. Voigt and Reuss bounds are quite distant in this case. Unfortunately, simple analytical schemes also fail to predict correctly the composite stiffness due to the extreme sample porosity. The Mori-Tanaka method approaches the arithmetic mean between the bounds, whereas the self-consistent scheme tends to reach the stiffness of the phase with higher occurrence (i.e. the air).

Table 3: Values of the Level II effective Young's modulus computed by different analytical homogenization schemes.

Scheme	Mori-Tanaka	Self-consist. scheme	Voigt bound	Reuss bound
E (GPa)	3.1510	0.0012	6.0200	0.0011

Most of analytical studies on the homogenization of foams are based on models with a regular periodic microstructure. Nevertheless, real foam microstructures are characterized with different sizes and shapes and sizes of pores rather than with periodic structures as shown in Section 2 of this paper. The solution can be to solve the problem of irregular microstructures by an analysis of a large representative volume element containing large enough number of pores. Such model can be solved in two or three dimensions.

Therefore, more appropriate (but still simple) two dimensional microstructure based FEM model was proposed. The model geometry was generated from high resolution optical image of polished foam cross-section (Figure 10a) converted to binary one. Square domain with 106×106 mm size (i.e. being much larger than average pore size ~ 2.9 mm) was extracted from the image. At this domain, pore centroids were detected, Delaunay triangulation applied and Voronoi cells created. Then, an equivalent 2D-beam structure was generated from Voronoi cell boundaries (Figure 10b). Based on several numerical studies performed for this purpose (but not shown here in details), it was found that the distribution of cross sectional areas and bending stiffness of individual beams do not play a significant role in the evaluation of the homogenized properties.

The overall stiffness is influenced mainly by the sum of the beam cross sectional areas and by the beam inclination to the load direction. The contribution of the beam bending

stiffness is diminished due to the very large beam length compared to its small cross sectional dimensions. Therefore, as an approximate but sufficient estimate, uniform cross-sectional area and uniform second moment of inertia were prescribed to all beams. The beam cross sectional area (A_{beam}) was computed from the total sample porosity ($\varphi=0.086$) and the total length of all beams (l_{total}) in the RVE with rectangular dimensions $a \times b$ as

$$A_{beam} = \frac{ab\varphi}{l_{total}}. \quad (21)$$

Taking into account 2-D case (i.e. unit thickness of the plane) beam height can be set as

$$h_{z,beam} = \frac{A_{beam}}{1} = A_{beam}. \quad (22)$$

Assuming rectangular shape of a cross section one can readily obtain the second moment of inertia as

$$I_{y,beam} = \frac{1}{12} 1 h_{z,beam}^3 = \frac{A_{beam}^3}{12}. \quad (23)$$

In the analysis, macroscopic strain $\boldsymbol{\varepsilon}^0$ is prescribed to the RVE and microscopic strains and stresses are solved. Volumetric averaging of microscopic stresses leads to the assessment of an average macroscopic stress and finally estimation of effective stiffness parameters. The key issue of the computation is the size of RVE and application of boundary conditions around the domain. Since the domain size is always smaller than an infinite body, any constraints can strongly influence the results. Application of the kinematic boundary conditions leads to the overestimation of effective stiffness and it can give an upper bound, whereas the static boundary conditions give a lower bound .

The best solution is usually provided by applying periodic boundary conditions to RVE which are, however, difficult to implement into commercial codes.

Nevertheless, the influence of the boundary conditions on microscopic strains and stresses in the domain decrease in distant points from the boundary. The size of our domain (106×106 mm) allowed us to solve the problem with kinematic boundary conditions. For homogenization, considerably smaller region (later found optimum 35-50 mm) in the central part was used. Microscopic strains and stresses were computed inside this smaller area which was assumed to be still sufficiently large to describe the material inhomogeneities and to serve as material RVE.

Kinematic constrains were applied on all domain sides. Free beams located around the boundary and not connected to any cell were deleted and supports put on the nodes located on the closest cell. Such arrangement of beams and supports prevented the structure from unreasonably large deformations of these free boundary beams. The whole domain (106×106 mm) was subjected to homogeneous macroscopic strain in one axial direction ($\boldsymbol{\varepsilon}^0 = \{1,0,0\}^T$) by imposing prescribed displacement to one domain side (Figure 10c). The test was performed using commercial Ansys FEM software and microscopic strains and stresses solved in the domain. Strains and stresses (structural forces for the case of beams, respectively) inside the smaller area (35-50 mm) were averaged and used for computation of the homogenized stiffness matrix (one column in the matrix, respectively). Assuming material isotropy, the first component (1,1) at the material stiffness matrix is given by:

$$L_{11} = E \frac{(1-\nu)}{(1+\nu)(1-2\nu)} \quad (24)$$

in which E is the Young's modulus and ν Poisson's ratio, respectively. Since the Poisson's ratio of the whole foam is close to zero (as confirmed by experimental measurements) the L_{11} member coincides with the Young's modulus E .

For the tension test in x-direction (Figure 10c), the homogenized Young's modulus was found to be RVE size dependent. Experimental investigations of the dependence of sample size on apparent elastic modulus and strength were conducted e.g. by Ashby et al. . They found, the modulus and strength become independent of size when the sample dimensions exceeded about seven cell diameters. This would imply minimum RVE size 20.3 mm for our typical cell size (2.9 mm). On the other hand, the RVE size should not exceed roughly 1/3 to 1/2 of the whole domain size not to be influenced by boundary conditions which implies maximum RVE width about 35 to 53 mm for our 106 mm wide domain. To find an optimum RVE size a numerical study was conducted for different RVE sizes in the range 20 to 90 mm (see Figure 11). An optimum RVE was confirmed to be between 35 and 50 mm for our specific domain. Results for smaller RVEs (<35 mm) are influenced by the beam inhomogeneity inside the RVE (in other words, such small RVE is not representative enough) whereas larger RVEs (>50 mm) are already influenced by the vicinity of boundary conditions. For optimum RVE sizes, the effective Young's modulus varied in the range $E_{eff,II}=1.36-1.38$ GPa.

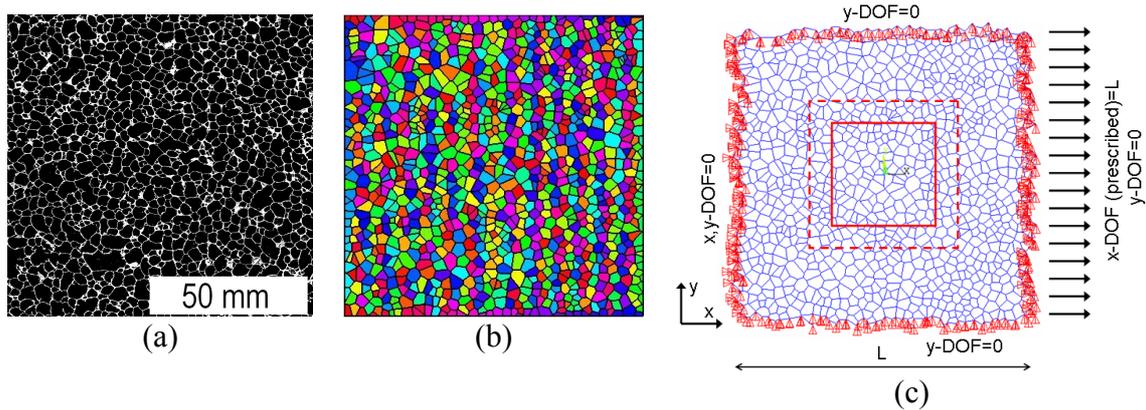


Figure 10: (a) Binary image of the foam (106×106 mm). (b) Voronoi tessellation. (c) 2-D beam model with boundary constraints (red squares indicate optimum RVE sizes from which homogenized properties have been obtained; solid line=35×35 mm, dashed line=50×50 mm).

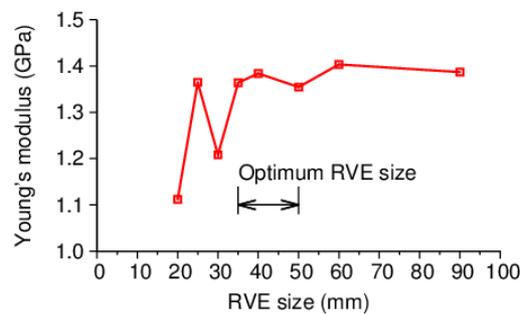


Figure 11: Dependence of effective Level II Young's modulus on RVE size.

The resulting homogenized Young's modulus is comparable with the range of experimental values (0.4–1 GPa) reported for Alporas[®] e.g. by Ashby et. al. . Experimental measurements in uniaxial compression performed on our samples (30×30×60/90 mm Alporas blocks) indicate $E=1.45\pm 0.15$ GPa (see Section 4.4). The slightly lower stiffness obtained from two-dimensional model can be explained by the lack of additional confinement appearing in the three-dimensional case. The influence

of the RVE size can also play a role as described above. However, the obtained difference is small ($\sim 5\%$), probably also due to the almost zero foam Poisson's ratio. Anyway, results of the simplified 2-D model have to be treated as a relatively close but only the first estimate of the Level II material properties which should be refined e.g. by using more precise three-dimensional model.

4.4 Results from macroscopic measurements

Uniaxial compression tests on $30\times 30\times 60$ and $30\times 30\times 90$ mm Alporas blocks (Figure 12a) were performed in an electromechanical press to verify numerical results on the Level II. Specimens were loaded-unloaded by five to ten cycles at very low strains and then fully compressed up to $\sim 5\%$ longitudinal strain (Figure 12b,c). Longitudinal and transversal (engineering) strains were evaluated by means of digital image correlation from CCD camera images taken during the test. Negligible differences have been found between the slopes of loading/unloading cycles (Figure 12c) which justifies evaluation of elastic properties from this part of the loading diagram. Young's modulus was finally computed as the average slope from all relevant cycles (i.e. all cycles except the first and the last one that both can be influenced by non-linear effects). Young's modulus was determined as $E=1.45\pm 0.15$ GPa on six foam samples. Poisson's ratio was found to be $\nu\approx 0$ in the elastic regime.

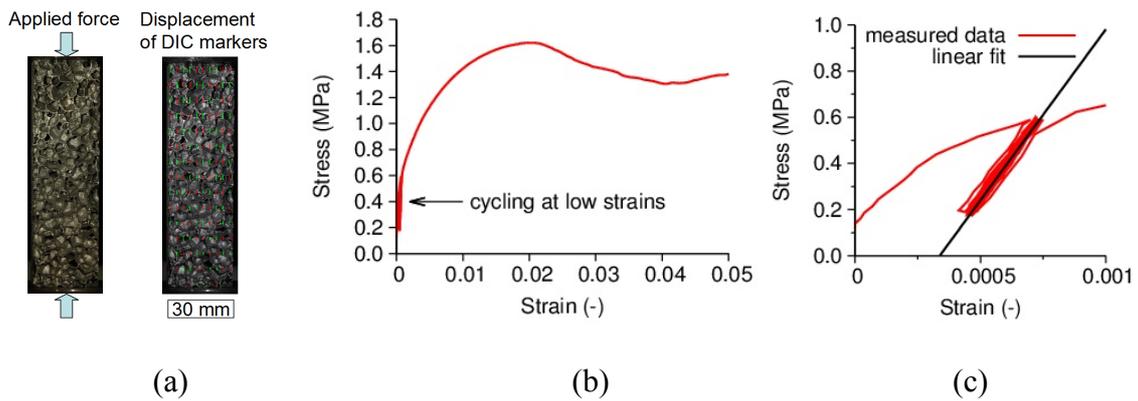


Figure 12: Compression test. (a) Foam sample and an image with digital image correlation markers, (b) stress strain diagram, (c) detail of loading-unloading cycles at low strains.

5 Conclusions

A simple but effective two-scale microstructure based model of closed-cell aluminium foam for the assessment of homogenized elastic properties was proposed in this paper. The first level characterized by thin cell walls ($\sim 60 \mu\text{m}$) was successfully homogenized with several analytical continuum mechanics schemes. Two different material phases (Al-rich and Ca/Ti-rich) were detected at this lower scale by ESEM and statistical grid nanoindentation. Effective Young's modulus $E_{eff,i} \approx 70 \text{ GPa}$ was received regardless the used scheme. The value was also justified by numerical FFT-based homogenization with a very good agreement (error less than 4%).

The upper foam level (Level II) contained homogenized walls and large air pores. Here, analytical tools were applied without success. Very poor estimates were given by the Mori-Tanaka or self-consistent due to extremely high air content in the foam and large stiffness contrast. To better describe the real foam microstructure, a FEM model

was proposed for the numerical homogenization at the second level. The model geometry was generated from large optical scan of polished foam cross section converted to binary image. Delaunay triangulation and Voronoi tessellation have been applied and equivalent beam structure generated. The dependence of RVE size was solved in a large domain (106×106 mm) supported by kinematic boundary conditions. An optimum RVE size was found to be in the range 35-50 mm (i.e. 33-47% of the domain size) for which effective elastic properties were assessed ($E_{eff,II}=1.36-1.38$ GPa).

The model has proven to realistically describe macroscopic elastic properties of the foam. Nevertheless, two-dimensional approximation slightly underestimated the experimentally obtained stiffness ($E\approx 1.45$ GPa). It is likely due to the inability to capture additional confinement coming from three-dimensional material microstructure or due to the RVE size inaccuracy. Further development of the numerical model and generation of the model geometry from micro-CT data (i.e. extension to 3-D) are planned as future developments.

Acknowledgements

Financial support of the Czech Science Foundation is gratefully acknowledged (projects P105/12/0824 and P105/12/0331).

References

- [1] Banhart J, Manufacture, characterisation and application of cellular metals and metal foams, Progress in Materials Science 46: 559-632, 2001.

- [2] Miyoshi T, Itoh M, Akiyama S, Kitahara A, Aluminium foam "ALPORAS": The production process, properties and application” Materials Research Society Symp. Proc. 521, 1998.
- [3] Hasan MA, Kim A, Lee H-J, Measuring the cell wall mechanical properties of Al-alloy foams using the nanoindentation method, Composite Structures 83: 180-188, 2008.
- [4] [Papadopoulos DP](#), [Konstantinidis IC](#), [Papanastasiou N](#), [Skolianos S](#), [Lefakis H](#) and [Tsipas DN](#), Mechanical properties of Al metal foams, [Materials Letters](#) 58 (21): 2574-2578, 2004.
- [5] Jeon I et al., Cell wall mechanical properties of closed-cell Al foam, Mechanics of Materials, 41 (1): 60-73, 2009.
- [6] Sugimura Y, Meyer J, He MY, Bart-Smith H, Grenstedt J, Evans AG, On the mechanical performance of closed cell Al alloy foams, Acta Mater. 45: 5245–5259, 1997.
- [7] Idris MI, Vodenitcharova, Hoffman M, Mechanical behaviour and energy absorption of closed-cell aluminium foam panels in uniaxial compression, Materials Science and Engineering A 517, 37–45, 2009.
- [8] Yongliang M, Guangchun Y, Hongjie L, Effect of cell shape anisotropy on the compressive behavior of closed-cellaluminum foams, Materials and Design 31, 1567–1569, 2010.
- [9] De Giorgi M, Carofalo A, Dattoma V, Nobile R, Palano F, Aluminium foams structural modelling, Computers and Structures 88, 25–35, 2010.

- [10] Constantinides G, Chandran FR, Ulm F-J, Vliet KV, Grid indentation analysis of composite microstructure and mechanics: Principles and validation, *Materials Science and Engineering: A*, [430 \(1-2\)](#):189-202, 2006.
- [11] Němeček J, Šmilauer V, Kopecký L, Nanoindentation characteristics of alkali-activated aluminosilicate materials, *Cement and Concrete Composites* [33 \(2\)](#): 163-170, 2011.
- [12] Zaoui A, Continuum Micromechanics: Survey, *Journal of Engineering Mechanics* 128, 2002.
- [13] ISO 4287-1997, “Geometrical Product Specifications (GPS) - Surface texture: Profile method - Terms, definitions and surface texture parameters”.
- [14] Ashby MF et al., *Metal Foams: A Design Guide*, Elsevier, 2000.
- [15] Oliver W, Pharr GM, An improved technique for determining hardness and elastic modulus using load and displacement sensing indentation experiments, *J. Mater. Res.* 7 (6): 1564–1583, 1992.
- [16] Eshelby JD, The determination of the elastic field of an ellipsoidal inclusion and related problem, *Proc. Roy. Soc. London A* 241: 376–396 1957.
- [17] Mori T, Tanaka K, Average stress in matrix and average elastic energy of materials with misfitting inclusions, *Acta Metallurgica* 21 (5): 571-574, 1973.
- [18] Webelements on-line library,
<http://www.webelements.com/aluminium/physics.html>
- [19] Moulinec H, Suquet P, A fast numerical method for computing the linear and nonlinear mechanical properties of composites, *Comptes rendus de l'Académie des sciences. Série II, Mécanique, physique, chimie, astronomie*, 318(11): 1417-1423, 1994.

- [20] Vainikko G, Fast solvers of the Lippmann-Schwinger equation, in R. P. Gilbert and J. Kajiwara and Y. S. Xu, editors, *Direct and Inverse Problems of Mathematical Physics in International Society for Analysis, Applications and Computation*, pages 423-440. Kluwer Academic Publishers, Dordrecht, The Netherlands, 2000.
- [21] Zeman J, Vondřejc J, Novák J, Marek I, Accelerating a FFT-based solver for numerical homogenization of periodic media by conjugate gradients, *Journal of Computational Physics*, 229(21): 8065-8071, 2010.
- [22] Vondřejc J, Zeman J, Marek I, Analysis of a Fast Fourier Transform Based Method for Modeling of Heterogeneous Materials, *Lecture Notes in Computer Science* 7116: 512-522, 2012.
- [23] Simone AE, Gibson LJ, Aluminum Foams Produced by Liquid-state Processes, *Acta Mater* 46 (9): 3109-3123, 1998.
- [24] Šejnoha M, Šmilauer V, Němeček J, Kopecký L, Application of Micromechanics in Engineering Practice, *Proceedings of the Eighth International Conference on Engineering Computational Technology*, B.H.V. Topping, (Editor), Civil-Comp Press, Stirlingshire, Scotland, 2012.
- [25] Gibson LJ, Ashby MF, *Cellular Solids: Structure and Properties*, Cambridge University Press, 1999.
- [26] Hohe J, Becker W, A probabilistic approach to the numerical homogenization of irregular solid foams in the finite strain regime, *International Journal of Solids and Structures* 42, 3549–3569, 2005.

-
- [27] Šmilauer V, Bittnar Z, Microstructure-based micromechanical prediction of elastic properties in hydrating cement paste, *Cement and Concrete Research* 36, 1708–1718, 2006.
- [28] Ashby MF, Evans A, Fleck NA, Gibson LJ, Hutchinson JW, Wadley HN, *Metal foams: a design guide*, Butterworth-Heinemann, Oxford, UK, 2000.
- [29] Jandajsek I, Jiroušek O, Vavřík D, Precise strain measurement in complex materials using Digital Volumetric Correlation and time lapse micro-CT data, *Procedia Engineering* 10, 1730–1735, 2011.

Figure captions

Figure 1: (a) Overall view on a foam structure (further denoted as Level II); (b) ESEM image of a cell wall; (c) detailed ESEM image of a cell wall showing Al-rich (dark grey) and Ca/Ti-rich areas (light zones; denoted as Level I).

Figure 2: An example of (a) ESEM image of the cell wall and (b) processed image segmented to two phases (white=Ca/Ti-rich, black=Al-rich area).

Figure 3: (a) Binary image of the polished foam panel. (b) Binary image of $\sim 50 \times 50$ mm foam cut. (c) Cell contours in the cut (prepared in Matlab).

Figure 4: Distribution of cell wall thicknesses.

Figure 5: Distribution of equivalent ellipse shape factor.

Figure 6: Distribution of equivalent circular pores.

Figure 7: Part of the indentation matrix showing 6×6 indents with $10 \mu\text{m}$ spacing as scanned with Hysitron Tribolab.

Figure 8: Typical loading diagrams for Al-rich and Ca/Ti-rich zones.

Figure 9: (a) Probability density functions of elastic moduli from two measured positions and (b) merged results with deconvoluted phases.

Figure 10: (a) Binary image of the foam (106×106 mm). (b) Voronoi tessellation. (c) 2-D beam model with boundary constraints (red squares indicate optimum RVE sizes from which homogenized properties have been obtained; solid line=35×35 mm, dashed line=50×50 mm).

Figure 11: Dependence of effective Level II Young's modulus on RVE size.

Figure 12: Compression test. (a) Foam sample and an image with digital image correlation markers, (b) stress strain diagram, (c) detail of loading-unloading cycles at low strains.

Tables

Table 1: Elastic moduli and volume fractions from deconvolution.

Phase	Mean (GPa)	St. dev. (GPa)	Volume fraction (-)
1 (Al-rich zone)	61.88	4.6	0.638
2 (Ca/Ti-rich zone)	87.40	16.7	0.362

Table 2: Values of the Level I effective Young's modulus computed by different homogenization schemes.

Scheme	Mori-Tanaka	Self-consist. scheme	Voigt bound	Reuss bound
E (GPa)	70.083	70.135	71.118	69.195

Table 3: Values of the Level II effective Young's modulus computed by different analytical homogenization schemes.

Scheme	Mori-Tanaka	Self-consist. scheme	Voigt bound	Reuss bound
E (GPa)	3.1510	0.0012	6.0200	0.0011

Part VI

Paper 5

Authors:

Jaroslav Vondřejc, Jan Zeman, and Ivo Marek

Title:

FFT-based Finite element method for homogenization

FFT-based Finite element method for homogenization

J.Vondřejc, J.Zeman and I.Marek

January, 2013

Abstract

We present a mathematical theory to a FFT-based homogenization, the numerical method introduced by Moulinec and Suquet [14], The method is based on the Lippmann-Schwinger integral equation including the Green function. We show its equivalence to weak formulation in the sense the unique solution coincide. Then we provide an unifying concept of the discretization with trigonometric polynomials [19] that is applicable for both formulations. Moreover, we explain the solution of the resulting non-symmetric linear system by Conjugate gradients as proposed in [24]. Finally, the convergence of discrete solutions to continuous one is provided.

1 Introduction

A majority of computational homogenization techniques rely on the solution to the unite cell problem, which concerns the determination of local fields in a representative sample of a heterogeneous material under periodic boundary conditions. In order to show the approach, a scalar problem modeling heat transfer, diffusion, or electric conductivity (our choice), is considered and formulated as

$$\operatorname{curl} \mathbf{e}(\mathbf{x}) = \mathbf{0}, \quad \operatorname{div} \mathbf{j}(\mathbf{x}) = \mathbf{0}, \quad \mathbf{j}(\mathbf{x}) = \mathbf{A}(\mathbf{x})\mathbf{e}(\mathbf{x}), \quad \mathbf{x} \in \mathcal{Y} \quad (1.1)$$

where $\mathcal{Y} = \prod_{\alpha=1}^d (-Y_\alpha, Y_\alpha) \subset \mathbb{R}^d$ denotes a periodic unit cell in d -dimensional space, $\mathbf{e} = (e_\alpha)_{\alpha=1}^d : \mathcal{Y} \rightarrow \mathbb{R}^d$ a vector valued electric field, $\mathbf{j} = (j_\alpha)_{\alpha=1}^d : \mathcal{Y} \rightarrow \mathbb{R}^d$ the corresponding vector of electric current, $\mathbf{A} = (A_{\alpha\beta})_{\alpha,\beta=1}^d : \mathcal{Y} \rightarrow \mathbb{R}^{d \times d}$ a second-order uniformly elliptic bounded tensor of electric conductivity, and $\mathbf{A}\mathbf{e}$ express the tensor by vector field multiplication, i.e. $\mathbf{A}\mathbf{e} = \left(\sum_{\beta=1}^d A_{\alpha\beta} e_\beta \right)_{\alpha=1}^d$. The differential equation (1.1) is constrained to a periodic boundary condition with period $\mathbf{Y} \in \mathbb{R}^d$ and the prescribed average load coming from the macroscopic level average electric field \mathbf{E} , i.e.

$$\langle \mathbf{e} \rangle := \frac{1}{|\mathcal{Y}|_d} \int_{\mathcal{Y}} \mathbf{e}(\mathbf{x}) \, d\mathbf{x} = \mathbf{E}, \quad (1.2)$$

where \mathbf{E} denotes a prescribed macroscopic electric field and $|\mathcal{Y}|_d$ represents the d -dimensional measure of \mathcal{Y} .

The solution of initial problem, (1.1) and (1.2), can be very demanding especially for complicated microstructures and high contrasts in conductivity coefficients in combination with higher-dimensional problems. There are various works dealing with discretization of weak formulation, Def. 2.20, especially by Finite Element Method (FEM) with polynomials as the basis functions [1, 6, 7, 22, 25].

Another method based on the Lippmann-Schwinger equation uses discretization with trigonometric collocation method proposed by Vainikko in [19] and leading to the non-symmetric non-sparse linear system. The numerical solution based on Neumann series expansion was introduced by Moulinec and Suquet in [14]. Nowadays, the proposed algorithm is in extensive investigation and some improvements and modifications were proposed [15, 8, 20, 8, 20, 13, 2, 3, 23], however no rigorous theory has been published yet.

Under the assumption of positive definiteness of material coefficients \mathbf{A} at periodic unit cell \mathcal{Y} , we show in Lemma 2.29 the connection between Lippmann-Schwinger integral equation, Def. (2.23), and weak formulation, Def. (2.20). The possibility to obtain the equivalent linear systems from both formulations has been already published in [21]. The discretization via trigonometric polynomials requires the values of material coefficients \mathbf{A} at regular grid, thus it is fitted for data obtained as digital images from e.g. X-ray tomography [4] or nanoindentation [16, 17].

Since the multiplication with the matrix from the linear system can be efficiently provided using Fast Fourier Transform (FFT), the solution is appropriate for Krylov subspace methods. Zeman et.al. in [24] proposed the solution using Conjugate gradients in spite of the absence of positive definiteness and symmetry of linear system and showed the behavior of the algorithm, its equivalence with Biconjugate gradient method and especially its independence on reference conductivity \mathbf{A}^0 , the parameter occurring in the Lippmann-Schwinger type integral formulation.

In Section 3.4, we show the convergence of discrete solution expressed as linear combination of trigonometric polynomials to the continuous solution and Section 3.6 validates the usage of Conjugate gradients for suitable initial approximation.

For the later use in this work, the following notation is introduced. The letter d denotes the dimension of the problem, assuming $d = 2, 3$; the Greek letters $\alpha, \beta, \gamma, \zeta, \theta$ are reserved to indices relating dimension, thus ranging $1, \dots, d$ (the range is for simplicity often omitted).

The set \mathbb{N}_0 represent natural numbers including zero. The sets \mathbb{C}^d and \mathbb{R}^d are spaces of complex and real vectors with canonical basis $\{\epsilon_\alpha\}$ and are equipped with the Lebesgue measure $d\mathbf{x}$. We denote by $|\Omega|_d$ the d -dimensional Lebesgue measure of a measurable set $\Omega \subset \mathbb{R}^d$. The norm $\|\cdot\|_2$ on \mathbb{C}^d is induced by scalar product $(\mathbf{u}, \mathbf{v})_{\mathbb{C}^d} = \sum_\alpha u_\alpha \bar{v}_\alpha$ for $\mathbf{u}, \mathbf{v} \in \mathbb{C}^d$.

The set $\mathbb{R}_{\text{spd}}^{d \times d}$ denotes the space of symmetric positive definite matrices of size $d \times d$ with norm $\|\mathbf{C}\|_2 = \max_{\mathbf{x} \in \mathbb{R}^d, \|\mathbf{x}\|=1} \|\mathbf{C}\mathbf{x}\|_2$ that equals to a largest eigenvalue.

A function $f : \mathbb{R}^d \mapsto \mathbb{R}$ is \mathbf{Y} -periodic (with period $\mathbf{Y} \in \mathbb{R}^d$) if $f(\mathbf{x} + \mathbf{Y} \odot \mathbf{k}) = f(\mathbf{x})$ for arbitrary $\mathbf{x} \in \mathbb{R}^d, \mathbf{k} \in \mathbb{Z}^d$, where operator \odot denotes element-wise multiplication. The \mathbf{Y} -periodic functions are sufficient to define only on a periodic unit cell (PUC), set to $\mathcal{Y} := (-Y_\alpha, Y_\alpha)_{\alpha=1}^d \subset \mathbb{R}^d$. Two integrable functions which are almost everywhere equal are identified.

Finally, operator \oplus^\perp denotes the direct sum of mutually orthogonal subspaces, e.g. $\mathbb{R}^d = \epsilon_1 \oplus^\perp \epsilon_2 \oplus^\perp \dots \oplus^\perp \epsilon_d$.

2 Continuous formulation

This section describing the continuous solution to initial problem (1.1) and (1.2) is split into two parts. Sec. 2.1 summarize the well known facts about function spaces and their properties that are used in Sec. 2.2 defining a weak formulation and describing the Lippmann-Schwinger integral equation and showing their equivalence.

2.1 Definitions

This section provides some facts about periodic distribution, the periodic L_{per}^2 space and its splitting into subspaces used thorough the work. The facts are well known and can be found in e.g. books of Saranen and Vainikko [18], Jikov, Kozlov, and Oleinik [9], or in other books dealing with function spaces and the equations of mathematical physics.

We begin with definitions of fundamental spaces and their properties.

Definition 2.1 (\mathbf{Y} -periodic L_{per}^p space). *Let $p \in \mathbb{R}$ with $1 \leq p \leq \infty$. Define a \mathbf{Y} -periodic space of square integrable functions as*

$$L_{\text{per}}^p(\mathcal{Y}) = \{f \in L_{\text{loc}}^p(\mathbb{R}^d) : f \text{ is } \mathbf{Y}\text{-periodic a.e.}\}$$

2 CONTINUOUS FORMULATION

3

where $L_{loc}^p(\mathcal{Y})$ denotes

$$\begin{cases} \{f : \mathbb{R}^d \rightarrow \mathbb{C}; f \text{ measurable} : \int_{\Omega} f^p(\mathbf{x}) \, d\mathbf{x} < \infty \text{ with } \Omega \subset \mathbb{R}^d \text{ bounded}\}, & 1 \leq p < \infty \\ \{f : \mathbb{R}^d \rightarrow \mathbb{C}; f \text{ measurable such that there exists } C \in \mathbb{R} \text{ with } |f| < C \text{ a.e.}\}, & p = \infty. \end{cases}$$

unifying the functions equaling one another almost everywhere.

The space $L_{\text{per}}^2(\mathcal{Y})$ is a Hilbert space with a scalar product of functions f and g defined as

$$(f, g)_{L_{\text{per}}^2(\mathcal{Y})} := |\mathcal{Y}|_d^{-1} \int_{\mathcal{Y}} f(\mathbf{x}) \overline{g(\mathbf{x})} \, d\mathbf{x} \quad (2.1)$$

and induced norm

$$\|f\|_{L_{\text{per}}^2(\mathcal{Y})} := \sqrt{(f, f)_{L_{\text{per}}^2(\mathcal{Y})}} \quad (2.2)$$

where $|\mathcal{Y}|_d$ denotes the d -dimensional measure of the PUC $\mathcal{Y} \subset \mathbb{R}^d$. The space $L_{\text{per}}^p(\mathcal{Y})$ for $1 \leq p \leq \infty$ is a Banach space with norm

$$\|f\|_{L_{\text{per}}^p(\mathcal{Y})} = \begin{cases} \left(\int_{\mathcal{Y}} |f(\mathbf{x})|^p \, d\mathbf{x} \right)^{\frac{1}{p}}, & 1 \leq p < \infty \\ \inf\{C \in \mathbb{R} : |f(\mathbf{x})| < C \text{ for almost all } \mathbf{x}\}, & p = \infty. \end{cases} \quad (2.3)$$

The powerful method to analyze functions is the technique of Fourier Transform and/or, in the periodic case, the Fourier series working with trigonometric orthonormal basis of space $L_{\text{per}}^2(\mathcal{Y})$ that is given by the set $\{\varphi_{\mathbf{k}}\}_{\mathbf{k} \in \mathbb{Z}^d}$ with $\varphi_{\mathbf{k}}$ defined as

$$\varphi_{\mathbf{k}}(\mathbf{x}) = \exp\left(i\pi \sum_{\alpha} \frac{k_{\alpha} x_{\alpha}}{Y_{\alpha}}\right), \quad \mathbf{k} \in \mathbb{Z}^d. \quad (2.4)$$

Definition 2.2 (Fourier representation of functions). *We define the Fourier representation $f_{\mathcal{F}} \in L_{\text{per}}^2(\mathcal{Y})$ of function $f \in L_{\text{per}}^2(\mathcal{Y})$ as*

$$f_{\mathcal{F}}(\mathbf{x}) = \sum_{\mathbf{k} \in \mathbb{Z}^d} \hat{f}(\mathbf{k}) \varphi_{\mathbf{k}}(\mathbf{x})$$

where the numbers $\hat{f}(\mathbf{k}) \in \mathbb{C}$, $\mathbf{k} \in \mathbb{Z}^d$ denotes the Fourier coefficients of function f defined as

$$\hat{f}(\mathbf{k}) := (f, \varphi_{\mathbf{k}})_{L_{\text{per}}^2(\mathcal{Y})} = |\mathcal{Y}|_d^{-1} \int_{\mathcal{Y}} f(\mathbf{x}) \varphi_{-\mathbf{k}}(\mathbf{x}) \, d\mathbf{x}.$$

where $\overline{\varphi_{\mathbf{k}}} = \varphi_{-\mathbf{k}}$.

It can be shown that Fourier representation $f_{\mathcal{F}}$ equals to the initial function f almost everywhere thus they can be identified. Moreover function $f \in C^1(\mathcal{Y})$, function with continuous partial derivatives, are represented by the Fourier series everywhere. Next, we can alternatively, due to Parseval's identity, express the scalar product of functions $f, g \in L_{\text{per}}^2(\mathcal{Y})$ as

$$(f, g)_{L_{\text{per}}^2(\mathcal{Y})} := \sum_{\mathbf{k} \in \mathbb{Z}^d} \hat{f}(\mathbf{k}) \overline{\hat{g}(\mathbf{k})} \quad (2.5)$$

and thus the norm as $\|f\|_{L_{\text{per}}^2(\mathcal{Y})} = \left(\sum_{\mathbf{k} \in \mathbb{Z}^d} |\hat{f}(\mathbf{k})|^2 \right)^{\frac{1}{2}}$.

In order to work with derivatives of L_{per}^2 functions, we have to weaken the notion of derivative using distributions, the generalized functions. First, it is necessary to define the periodic test functions being the linear function space and a proper convergence on it. Next we follow with the Fourier representation of a derivative.

2 CONTINUOUS FORMULATION

4

Lemma 2.3 (Fourier coefficients for a derivative). *For a function $f \in C_{\text{per}}^1(\mathcal{Y})$ we can express the Fourier coefficient for a partial derivative as*

$$\widehat{\frac{\partial f}{\partial x_\alpha}}(\mathbf{k}) = i\pi\xi_\alpha(\mathbf{k})\hat{f}(\mathbf{k})$$

where $\mathbf{k} \in \mathbb{Z}^d$ and $\xi_\alpha(\mathbf{k}) = \frac{k_\alpha}{Y_\alpha}$.

Definition 2.4 (Test functions). *For $d \in \mathbb{N}$ and $\mathbf{Y} \in \mathbb{R}^d$ regarding a dimension and a period, we define spaces $\mathcal{D}(\mathbb{R}^d)$ and $\mathcal{D}_{\text{per}}(\mathbb{R}^d)$ of test functions and \mathbf{Y} -periodic test functions on \mathbb{R}^d as*

$$\begin{aligned} \mathcal{D}(\mathbb{R}^d) &= \{\psi \in C^\infty(\mathbb{R}^d) : \text{supp } \psi \text{ is compact}\} \\ \mathcal{D}_{\text{per}}(\mathbb{R}^d) &= \{\psi \in C^\infty(\mathbb{R}^d) : \mathcal{T}_{\mathbf{Y}}^{\mathbf{k}}\psi(\mathbf{x}) = \psi(\mathbf{x}), \forall \mathbf{x} \in \mathbb{R}^d, \forall \mathbf{k} \in \mathbb{Z}^d\} \end{aligned}$$

where $C^\infty(\mathbb{R}^d)$ denotes the set of infinitely differentiable scalar functions on \mathbb{R}^d , operator $\text{supp } \psi$ denotes support of a function defined as a set $\text{supp } \psi = \overline{\{\mathbf{x} \in \mathbb{R}^d : \psi(\mathbf{x}) \neq 0\}}$, and operator $\mathcal{T}_{\mathbf{Y}}^{\mathbf{k}}\psi$ denotes the translation of function ψ , i.e.

$$(\mathcal{T}_{\mathbf{Y}}^{\mathbf{k}}\psi)(\mathbf{x}) = \psi\left(\mathbf{x} + \sum_{\alpha=1}^d k_\alpha Y_\alpha \epsilon_\alpha\right)$$

with $\epsilon_\alpha = (\delta_{\alpha\beta})_{\beta=1}^d$ for $\alpha = 1, \dots, d$ being the vectors of the canonical basis of \mathbb{R}^d and $\delta_{\alpha\beta}$ denoting the Kronecker delta.

Definition 2.5 (Convergence in $\mathcal{D}(\mathbb{R}^d)$). *Having the sequence of test functions, i.e. $\{\psi_n\}_{n=1}^\infty, \psi_n \in \mathcal{D}(\mathbb{R}^d), n \in \mathbb{N}_0$, we say that the sequence $\{\psi_n\}_{n=1}^\infty$ converges to φ in $\mathcal{D}(\mathbb{R}^d)$, and write $\psi_n \rightarrow \psi$ in $\mathcal{D}(\mathbb{R})$, if the following properties are valid*

- there exists a compact subset Ω of \mathbb{R}^d such that $\text{supp } \psi_n \subset \Omega$ for all $n \in \mathbb{N}$
- for the arbitrary $k_1, \dots, k_d \in \mathbb{N}_0$, the uniform convergence holds

$$\frac{\partial^k \psi_n(\mathbf{x})}{\partial x_1^{k_1} \partial x_2^{k_2} \dots \partial x_d^{k_d}} \rightarrow \frac{\partial^k \psi(\mathbf{x})}{\partial x_1^{k_1} \partial x_2^{k_2} \dots \partial x_d^{k_d}}, \quad \text{for } n \rightarrow \infty$$

with $k = \sum_{\alpha=1}^d k_\alpha$.

Definition 2.6 (distribution). *A linear functional $T : \mathcal{D}(\mathbb{R}^d) \rightarrow \mathbb{R}$, with $T(\cdot)$ denoted as $\langle T, \cdot \rangle$, is a distribution on \mathbb{R}^d if it is continuous with respect to the convergence of $\mathcal{D}(\mathbb{R}^d)$, i.e.*

$$\psi_n \rightarrow \psi \text{ in } \mathcal{D}(\mathbb{R}) \quad \Rightarrow \quad \langle T, \psi_n \rangle \rightarrow \langle T, \psi \rangle \quad (2.6)$$

The space of distributions is denoted by $\mathcal{D}'(\mathbb{R}^n)$.

Definition 2.7 (\mathbf{Y} -periodic distribution). *Distribution $T \in \mathcal{D}'(\mathbb{R}^n)$ is \mathbf{Y} -periodic distribution, if*

$$\langle T, \mathcal{T}_{\mathbf{Y}}^{\mathbf{k}}\psi \rangle = \langle T, \psi \rangle, \quad \forall \psi \in \mathcal{D}(\mathbb{R}^d)$$

where $\mathcal{T}_{\mathbf{Y}}^{\mathbf{k}}\psi$ denotes the translation of function ψ described in Def. (2.4). The space of \mathbf{Y} -periodic distributions, i.e. $\{T \in \mathcal{D}'(\mathbb{R}) : T \text{ is } \mathbf{Y}\text{-periodic}\}$, is denoted as $\mathcal{D}'_{\text{per}}(\mathbb{R}^d)$.

In order to define the derivative of integrable function without any requirement for regularity, the notions of the derivative of distribution and convergence of distribution are introduced.

Definition 2.8 (derivative of periodic distribution). *For a periodic distribution $f \in \mathcal{D}'_{\text{per}}(\mathbb{R}^d)$ we define its partial derivative as distribution $\frac{\partial f}{\partial x_\alpha} \in \mathcal{D}'_{\text{per}}(\mathbb{R}^d)$ satisfying for $\alpha = 1, \dots, d$*

$$\left\langle \frac{\partial f}{\partial x_\alpha}, \psi \right\rangle = -\left\langle f, \frac{\partial \psi}{\partial x_\alpha} \right\rangle, \quad \forall \psi \in \mathcal{D}_{\text{per}}(\mathbb{R}^d). \quad (2.7)$$

Definition 2.9 (convergence of periodic distribution). *We say that a periodic distributions $f_n \in \mathcal{D}'_{\text{per}}(\mathbb{R}^d)$ converge to $f \in \mathcal{D}'_{\text{per}}(\mathbb{R}^d)$, and write*

$$\lim_{n \rightarrow \infty} f_n = f,$$

if

$$\lim_{n \rightarrow \infty} \langle f_n, \psi \rangle = \langle f, \psi \rangle, \quad \forall \psi \in \mathcal{D}_{\text{per}}(\mathbb{R}^d). \quad (2.8)$$

After a definition of Fourier coefficients of periodic distribution we follows with a Lemma about the representation of periodic distribution by Fourier series.

Definition 2.10 (Fourier coefficients of periodic distribution). *We define the Fourier coefficients $\hat{f}(\mathbf{k}) \in \mathbb{C}^d$ for $\mathbf{k} \in \mathbb{Z}^d$ of a periodic distribution $f \in \mathcal{D}'_{\text{per}}(\mathbb{R}^d)$ as*

$$\hat{f}(\mathbf{k}) = \langle f, \varphi_{\mathbf{k}} \rangle \quad (2.9)$$

as $\varphi_{\mathbf{k}} \in \mathcal{D}_{\text{per}}(\mathbb{R}^d)$.

Lemma 2.11 (Fourier series of distributions). *Assume $f \in \mathcal{D}'_{\text{per}}(\mathbb{R}^d)$ and $\psi \in \mathcal{D}_{\text{per}}(\mathbb{R}^d)$. Then there holds*

(a) *there exists $p \in \mathbb{N}$ and $c_r \in \mathbb{R}$ such that $|\hat{f}(\mathbf{k})| \leq c_p \|\mathbf{k}\|_2^p, \forall \mathbf{k} \in \mathbb{Z}^d \setminus \mathbf{0}$*

(b) *there exists $c_r \in \mathbb{R}$ such that for arbitrary $r > 0$ we have $|\hat{\psi}(\mathbf{k})| \neq c_r \|\mathbf{k}\|_2^r, \forall \mathbf{k} \in \mathbb{Z}^d \setminus \mathbf{0}$*

(c) $\langle f, \psi \rangle = \sum_{\mathbf{k} \in \mathbb{Z}^d} \hat{f}(\mathbf{k}) \hat{\psi}(-\mathbf{k})$

(d) $\lim_{\min N_1, \dots, N_d \rightarrow \infty} \sum_{\mathbf{k} \in \underline{\mathbb{Z}}_N^d} \hat{f}(\mathbf{k}) \varphi_{\mathbf{k}}(\mathbf{x}) \rightarrow f(\mathbf{x})$ in $\mathcal{D}'_{\text{per}}(\mathbb{R}^d)$

where $\underline{\mathbb{Z}}_N^d = \left\{ \mathbf{k} \in \mathbb{Z}^d : -\frac{N_\alpha}{2} \leq k_\alpha < \frac{N_\alpha}{2}, \alpha = 1, \dots, d \right\}$.

The Lemma and the proof that can be founded in e.g. [18] says in (a) and (b) that the higher the regularity of a function/distribution the higher the decay of Fourier coefficients in infinity giving the sense of expression (c). The property (c) is important especially with the definition of derivative of distribution, cf. Def. (2.8). The property (d) legitimize the expression of the distribution $f \in \mathcal{D}'_{\text{per}}(\mathbb{R}^d)$ as a Fourier series.

Next notes are dedicated to Sobolev spaces.

Definition 2.12 (Derivative with multiindex notation). *Let $m \in \mathbb{N}_0$, $\mathbf{l} \in \mathbb{N}_0^d$, and $u(\mathbf{x}) \in L^2_{\text{per}}(\mathcal{Y})$, we define*

$$D^{\mathbf{l}}u(\mathbf{x}) = \frac{\partial^{|\mathbf{l}|} u(\mathbf{x})}{\prod_{\alpha=1}^d \partial^{l_\alpha} x_\alpha}$$

$$D^m u(\mathbf{x}) = \{D^{\mathbf{l}}u : \|\mathbf{l}\|_1 = m\}$$

with a partial derivative in the sense of distributions.

Definition 2.13 (Sobolev space). *Let $1 \leq p \leq \infty$ and $\mu \in \mathbb{N}$. Sobolev space $W^{\mu,p}_{\text{per}}(\mathcal{Y})$ of \mathbf{Y} -periodic functions (distributions) is defined as*

$$W^{\mu,p}_{\text{per}}(\mathcal{Y}) = \left\{ f \in L^p_{\text{per}}(\mathcal{Y}); D^{\mathbf{l}}u \in L^p_{\text{per}} \text{ for } \mathbf{l} \in \mathbb{N}_0^d : \|\mathbf{l}\|_1 \leq \mu \right\}.$$

For a special case $p = 2$, we define

$$H^{\mu}_{\text{per}}(\mathcal{Y}) := W^{\mu,2}_{\text{per}}(\mathcal{Y}).$$

We define a norm of Sobolev space $\|\cdot\|_{W^{\mu,p}_{\text{per}}(\mathcal{Y})}$ as

$$\|u\|_{W^{\mu,p}_{\text{per}}(\mathcal{Y})} = \begin{cases} \left(\sum_{\mathbf{l} \in \mathbb{N}_0^d : \|\mathbf{l}\|_1 \leq \mu} \|D^{\mathbf{l}}u\|_{L^p_{\text{per}}(\mathcal{Y})}^p \right)^{\frac{1}{p}}, & \text{for } 1 \leq p < \infty \\ \max_{\mathbf{l} \in \mathbb{N}_0^d : \|\mathbf{l}\|_1 \leq \mu} \|D^{\mathbf{l}}u\|_{L^\infty_{\text{per}}(\mathcal{Y})}, & \text{for } p = \infty. \end{cases}$$

where norm $\|\cdot\|_{L^p_{\text{per}}(\mathcal{Y})}$ is defined in Eq. (2.3).

Remark 2.14. For a Sobolev space $H_{\text{per}}^\mu(\mathcal{Y})$ we also use equivalent norm

$$\|u\|_{H_{\text{per}}^\mu(\mathcal{Y})} = \left(\sum_{\mathbf{k} \in \mathbb{Z}^d} \|\underline{\xi}(\mathbf{k})\|_2^{2\mu} |\hat{u}(\mathbf{k})|^2 \right)^{\frac{1}{2}}.$$

where $\hat{u}(\mathbf{k})$ denotes Fourier coefficients of function u and

$$\underline{\xi}(\mathbf{k}) = \begin{cases} 1, & \text{for } \mathbf{k} = \mathbf{0}, \\ \left(\frac{k_\alpha}{2Y_\alpha} \right)_{\alpha=1}^d, & \text{for } \mathbf{0} \neq \mathbf{k} \in \mathbb{Z}^d. \end{cases}$$

Remark 2.15 (vector valued functions). Let $\mu \in \mathbb{N}$ and $1 \leq p \leq \infty$. We define a space of vector valued functions as

$$\begin{aligned} L_{\text{per}}^p(\mathcal{Y}; \mathbb{C}^d) &= \{ \mathbf{f} : \mathcal{Y} \rightarrow \mathbb{C}^d \mid \mathbf{f} = (f_\alpha)_{\alpha=1}^d, f_\alpha \in L_{\text{per}}^p(\mathcal{Y}) \text{ for } \alpha = 1, \dots, d \} \\ W_{\text{per}}^{\mu,p}(\mathcal{Y}; \mathbb{C}^d) &= \{ \mathbf{f} : \mathcal{Y} \rightarrow \mathbb{C}^d \mid \mathbf{f} = (f_\alpha)_{\alpha=1}^d, f_\alpha \in W_{\text{per}}^{\mu,p}(\mathcal{Y}) \text{ for } \alpha = 1, \dots, d \} \end{aligned}$$

with norms

$$\|\mathbf{f}\|_{L_{\text{per}}^p(\mathcal{Y}; \mathbb{C}^d)} = \|(\|\mathbf{f}(\mathbf{x})\|_p)\|_{L_{\text{per}}^p(\mathcal{Y})} \quad \|\mathbf{f}\|_{W_{\text{per}}^{\mu,p}(\mathcal{Y}; \mathbb{C}^d)} = \|(\|\mathbf{f}(\mathbf{x})\|_p)\|_{W_{\text{per}}^{\mu,p}(\mathcal{Y})}$$

and scalar product

$$(\mathbf{f}, \mathbf{g})_{L_{\text{per}}^2(\mathcal{Y}; \mathbb{C}^d)} = \sum_{\alpha} (f_\alpha, g_\alpha)_{L_{\text{per}}^2(\mathcal{Y})}$$

where $\|\mathbf{v}\|_p := (\sum_{\alpha} v_\alpha^p)^{\frac{1}{p}}$ denotes the norm on space \mathbb{R}^d . The above space naturally inherits the properties stated in previous parts.

The Fourier representation $\mathbf{f}_{\mathcal{F}} \in L_{\text{per}}^2(\mathcal{Y}; \mathbb{C}^d)$ of the vector valued function $\mathbf{f} \in L_{\text{per}}^2(\mathcal{Y}; \mathbb{C}^d)$ can be defined as a Fourier coefficients of its components, i.e.

$$\mathbf{f}_{\mathcal{F}}(\mathbf{x}) = \sum_{\mathbf{k} \in \mathbb{Z}^d} \hat{\mathbf{f}}(\mathbf{k}) \varphi_{\mathbf{k}}(\mathbf{x}) \quad \hat{\mathbf{f}}(\mathbf{k}) = \left((f_\alpha, \varphi_{\mathbf{k}})_{L_{\text{per}}^2(\mathcal{Y})} \right)_{\alpha}^d.$$

The matrix valued functions or tensor, i.e. $L_{\text{per}}^2(\mathcal{Y}; \mathbb{C}^{d \times d})$ and $W_{\text{per}}^{\mu,p}(\mathcal{Y}; \mathbb{C}^{d \times d})$, are defined analogically with norms

$$\|\mathbf{f}\|_{L_{\text{per}}^p(\mathcal{Y}; \mathbb{C}^d)} = \left\| \left(\|f_{\alpha\beta}\|_{L_{\text{per}}^2(\mathcal{Y})} \right)_{\alpha,\beta=1}^d \right\|_p \quad \|\mathbf{f}\|_{W_{\text{per}}^{\mu,p}(\mathcal{Y}; \mathbb{C}^d)} = \left\| \left(\|f_{\alpha\beta}\|_{W_{\text{per}}^{\mu,p}(\mathcal{Y})} \right)_{\alpha,\beta=1}^d \right\|_p$$

where $\|\cdot\|_p$ denotes a matrix norm on $\mathbb{R}^{d \times d}$ defined as

$$\|\mathbf{A}\|_p = \sup_{\mathbf{x} \in \mathbb{R}^d, \mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{A}\mathbf{x}\|_p}{\|\mathbf{x}\|_p}.$$

Now, we will define the differential operators divergence and curl on the space L_{per}^2 , i.e. in the sense of distributions.

Definition 2.16 (divergence and curl). Let $\mathbf{f} \in L_{\text{per}}^2(\mathcal{Y}; \mathbb{R}^d)$ be vector valued function, then we define divergence and curl in distributional sense as

$$\begin{aligned} \text{div } \mathbf{f} &:= (\text{div } \mathbf{f}, \psi) = - \sum_{\alpha=1}^d \left\langle f_\alpha, \frac{\partial \psi}{\partial x_\alpha} \right\rangle, \quad \forall \psi \in \mathcal{D}_{\text{per}}(\mathbb{R}^d), \\ \text{curl } \mathbf{f} &:= (\text{curl } \mathbf{f}, \psi) = - \left(\left\langle f_\alpha, \frac{\partial \psi}{\partial x_\beta} \right\rangle - \left\langle f_\beta, \frac{\partial \psi}{\partial x_\alpha} \right\rangle \right)_{\alpha,\beta=1,\dots,d}, \quad \forall \psi \in \mathcal{D}_{\text{per}}(\mathbb{R}^d). \end{aligned}$$

2 CONTINUOUS FORMULATION

7

Remark 2.17 (Helmholtz decomposition). *Next, we introduce Helmholtz decomposition $L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d) = \mathcal{U} \oplus^\perp \mathcal{E} \oplus^\perp \mathcal{J}$ to the spaces of constant, curl-free with zero mean, and divergence free with zero mean fields*

$$\mathcal{U} = \{\mathbf{v} \in L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d) : \mathbf{v}(\mathbf{x}) = \text{const.}\}, \quad (2.10a)$$

$$\mathcal{E} = \{\mathbf{v} \in L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d) : \text{curl } \mathbf{v} = \mathbf{0}, \langle \mathbf{v} \rangle = \mathbf{0}\}, \quad (2.10b)$$

$$\mathcal{J} = \{\mathbf{v} \in L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d) : \text{div } \mathbf{v} = 0, \langle \mathbf{v} \rangle = \mathbf{0}\}, \quad (2.10c)$$

where $\langle \mathbf{v} \rangle := \frac{1}{|\mathcal{Y}|^d} \int_{\mathcal{Y}} \mathbf{v}(\mathbf{x}) \, d\mathbf{x} \in \mathbb{R}^d$ denotes the mean value of function \mathbf{v} over periodic unit cell \mathcal{Y} . Since the space \mathcal{U} consists of constant functions, we identify the space \mathcal{U} with \mathbb{R}^d ; this validates the operations such as $\mathbf{E} + \mathbf{v} \in L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)$ for $\mathbf{E} \in \mathbb{R}^d$, $\mathbf{v} \in L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)$ and $\mathbf{C}\mathbf{J} \in \mathbb{R}^d$ for $\mathbf{C} \in \mathbb{R}^{d \times d}$ and $\mathbf{J} \in \mathcal{U}$.

Lemma 2.18 (Existence of potential). *Vector valued function $\mathbf{e} \in \mathcal{E}$ has potential, i.e. there exist scalar function $u \in H^1_{\text{per}}(\mathcal{Y}) := \{u \in L^2_{\text{per}}(\mathcal{Y}); |\hat{u}(0)|^2 + \sum_{\mathbf{k} \in \mathbb{Z}^d} \|\boldsymbol{\xi}(\mathbf{k})\|_2^2 |\hat{u}(\mathbf{k})|^2 < \infty\}$. Moreover, space \mathcal{E} can be alternatively expressed as $\mathcal{E} = \{\nabla u; u \in H^1_{\text{per}}(\mathcal{Y})\}$.*

Proof. Vector valued function $\mathbf{e} \in \mathcal{E} \subset L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)$ can be expressed using Fourier series

$$\mathbf{e} = \sum_{\mathbf{k} \in \mathbb{Z}^d \setminus \{0\}} \hat{f}_\alpha(\mathbf{k}) \varphi_{\mathbf{k}}$$

noting $\langle \mathbf{e} \rangle = \mathbf{0}$. Then we define a scalar function $u \in L^2_{\text{per}}(\mathcal{Y})$ as

$$u(\mathbf{x}) := \sum_{\mathbf{k} \in \mathbb{Z}^d \setminus \{0\}} \frac{\hat{e}_1(\mathbf{k})}{i\pi\xi_1(\mathbf{k})} \varphi_{\mathbf{k}}(\mathbf{x})$$

and we show that it is a potential for \mathbf{e} .

Next, we show that $\frac{\partial u}{\partial x_\alpha} = e_\alpha$ in $L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)$ by testing with function $\frac{\partial \psi}{\partial x_1}$, i.e.

$$\begin{aligned} \left\langle \frac{\partial u}{\partial x_\alpha} - e_\alpha, \frac{\partial \psi}{\partial x_1} \right\rangle &= \left\langle \frac{\partial u}{\partial x_\alpha}, \frac{\partial \psi}{\partial x_1} \right\rangle - \left\langle e_1, \frac{\partial \psi}{\partial x_\alpha} \right\rangle \\ &= \sum_{\mathbf{k} \in \mathbb{Z}^d} \left(i\pi\xi_\alpha(\mathbf{k}) \frac{\hat{e}_1(\mathbf{k})}{i\pi\xi_1(\mathbf{k})} i\pi\xi_1(\mathbf{k}) \hat{\psi}(\mathbf{k}) - \hat{e}_1(\mathbf{k}) i\pi\xi_\alpha(\mathbf{k}) \hat{\psi}(\mathbf{k}) \right) = 0 \end{aligned}$$

where the identity $\langle e_\alpha, \frac{\partial \psi}{\partial x_1} \rangle = \langle e_1, \frac{\partial \psi}{\partial x_\alpha} \rangle$ was used as the function \mathbf{e} is curl free from definition.

The inclusion $\{\nabla u; u \in H^1_{\text{per}}(\mathcal{Y})\} \subseteq \mathcal{E}$ follows from

$$\begin{aligned} \text{curl } \nabla u &:= (\text{curl } \nabla u, \psi) = - \left(\left\langle \frac{\partial u}{\partial x_\alpha}, \frac{\partial \psi}{\partial x_\beta} \right\rangle - \left\langle \frac{\partial u}{\partial x_\beta}, \frac{\partial \psi}{\partial x_\alpha} \right\rangle \right)_{\alpha, \beta=1, \dots, d} \\ &= \left(\left\langle u, \frac{\partial^2 \psi}{\partial x_\alpha \partial x_\beta} \right\rangle - \left\langle u, \frac{\partial^2 \psi}{\partial x_\alpha \partial x_\beta} \right\rangle \right)_{\alpha, \beta=1, \dots, d} = (0)_{\alpha, \beta=1, \dots, d} \end{aligned}$$

for arbitrary $\psi \in \mathcal{D}_{\text{per}}(\mathbb{R})$.

Finally, we show that u is from Sobolev space, i.e. $u \in H^1_{\text{per}}(\mathcal{Y})$. It could be shown that functions $u_\alpha \in L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)$ for $\alpha = 2, \dots, d$ defined as

$$u_\alpha(\mathbf{x}) := \sum_{\mathbf{k} \in \mathbb{Z}^d} \frac{\hat{e}_\alpha(\mathbf{k})}{i\pi\xi_\alpha(\mathbf{k})} \varphi_{\mathbf{k}}(\mathbf{x})$$

2 CONTINUOUS FORMULATION

8

are also potentials of function \mathbf{e} . Clearly we have $\nabla u_\alpha = \mathbf{e}$ for $\alpha = 2, \dots, d$ and thus $\nabla u_\alpha = \nabla u$ and consequently even $u_\alpha = u$. Then $u \in H_{\text{per}}^1(\mathcal{Y})$ as

$$\begin{aligned} \|u\|_{H_{\text{per}}^1}^2 &= \sum_{\mathbf{k} \in \mathbb{Z}^d} \|\boldsymbol{\xi}(\mathbf{k})\|_2^2 |\hat{u}(\mathbf{k})|^2 = \sum_{\mathbf{k} \in \mathbb{Z}^d \setminus \{0\}} \|\boldsymbol{\xi}(\mathbf{k})\|_2^2 |\hat{u}(\mathbf{k})|^2 \\ &\leq \sum_{\mathbf{k} \in \mathbb{Z}^d \setminus \{0\}} |\xi_1(\mathbf{k})|^2 |\hat{u}(\mathbf{k})|^2 + |\xi_2(\mathbf{k})|^2 |\hat{u}(\mathbf{k})|^2 + \dots + |\xi_d(\mathbf{k})|^2 |\hat{u}(\mathbf{k})|^2 \\ &= \sum_{\mathbf{k} \in \mathbb{Z}^d \setminus \{0\}} |\xi_1(\mathbf{k})|^2 \left| \frac{\hat{e}_1(\mathbf{k})}{i\pi\xi_1(\mathbf{k})} \right|^2 + |\xi_2(\mathbf{k})|^2 \left| \frac{\hat{e}_2(\mathbf{k})}{i\pi\xi_2(\mathbf{k})} \right|^2 + \dots + |\xi_d(\mathbf{k})|^2 \left| \frac{\hat{e}_d(\mathbf{k})}{i\pi\xi_d(\mathbf{k})} \right|^2 \\ &\leq \|\mathbf{e}\|_{L_{\text{per}}^2(\mathcal{Y}; \mathbb{R}^d)}^2. \end{aligned}$$

□

2.2 Integral and weak formulation

In this section, we define weak formulation Def. 2.20 and the formulation based on Lippmann-Schwinger equation Def. 2.23. In Lemma 2.29 we show the equivalence of both formulations. It is based on projection operator that is derived from Green function occurring in the Lippmann-Schwinger equation, see Lem. 2.28.

Notation 2.19. Here and in the sequel, $\mathbf{A} \in L_{\text{per}}^\infty(\mathcal{Y}, \mathbb{R}_{\text{spd}}^{d \times d})$ denotes symmetric¹ and uniformly elliptic material coefficients of electric conductivity. It means that there exists positive constant $c_A > 0$ such that for almost all $\mathbf{x} \in \mathcal{Y}$ and all nonzero $\mathbf{u} \in \mathbb{R}^d$, inequality

$$c_A \|u\|_2^2 \leq (\mathbf{A}(\mathbf{x})\mathbf{u}, \mathbf{u})_{\mathbb{R}^d}$$

holds. For the next use we define constant of upper bound $C_A := \|\mathbf{A}\|_{L_{\text{per}}^\infty}$.

Moreover, function $\mathbf{e} \in \mathcal{E}$ denotes perturbation of electric field, and $\mathbf{E} \in \mathcal{U}$ its macroscopic counterparts. Then their summation $(\mathbf{E} + \mathbf{e}) \in \mathcal{U} \oplus^\perp \mathcal{E}$ represents microscopic field.

Weak formulation in the next definition is derived with the multiplication of differential equation (1.1) by a test function $\mathbf{v} \in \mathcal{D}_{\text{per}}(\mathcal{Y})$, the application of the Green theorem and some algebraic emendation.

Definition 2.20 (Weak formulation). Let $\mathbf{A} \in L_{\text{per}}^\infty(\mathcal{Y}, \mathbb{R}^{d \times d})$ be symmetric positive definite material coefficients. Then we define bilinear form $B[\cdot, \cdot] : L_{\text{per}}^2(\mathcal{Y}, \mathbb{R}^d) \times L_{\text{per}}^2(\mathcal{Y}, \mathbb{R}^d) \rightarrow \mathbb{R}$ as

$$B[\mathbf{e}, \mathbf{v}] := (\mathbf{A}\mathbf{e}, \mathbf{v})_{L_{\text{per}}^2(\mathcal{Y}, \mathbb{R}^d)} \quad (2.11)$$

and the linear functional $f[\cdot] : L_{\text{per}}^2(\mathcal{Y}, \mathbb{R}^d) \rightarrow \mathbb{R}$ as

$$f[\mathbf{v}] := -(\mathbf{A}\mathbf{E}, \mathbf{v})_{L_{\text{per}}^2(\mathcal{Y}, \mathbb{R}^d)}. \quad (2.12)$$

where $\mathbf{E} \in \mathcal{U}$ is a macroscopic load. Function $\tilde{\mathbf{e}} \in \mathcal{E}$ is the weak solution of the differential equation (1.1) with periodic boundary condition and constraint $\langle \tilde{\mathbf{e}} \rangle = \mathbf{E}$ if

$$B[\tilde{\mathbf{e}}, \mathbf{v}] = f[\mathbf{v}], \quad \forall \mathbf{v} \in \mathcal{E}. \quad (2.13)$$

We say that $\mathbf{e} = \mathbf{E} + \tilde{\mathbf{e}}$ is microscopic field for macroscopic load $\mathbf{E} \in \mathcal{U}$.

Remark 2.21 (Existence of the unique solution). Since the bilinear form from Def. (2.20) is positive definite $c_A \|\mathbf{v}\|_{L_{\text{per}}^2} \leq B[\mathbf{v}, \mathbf{v}]$ and bounded $B[\mathbf{u}, \mathbf{v}] \leq C_A \|\mathbf{u}\|_{L_{\text{per}}^2} \|\mathbf{v}\|_{L_{\text{per}}^2}$ for all $\mathbf{u}, \mathbf{v} \in \mathcal{E}$ and the linear functional f is bounded $\|f\| \leq C_A \|\mathbf{e}_0\|_{L_{\text{per}}^2}$, the well known Lax-Milgram theorem provide the existence and uniqueness of the solution.

¹For almost all $\mathbf{x} \in \mathcal{Y}$, equality $\mathbf{A}(\mathbf{x}) = \mathbf{A}(\mathbf{x})^T$ holds.

Remark 2.22 (The solution of weak formulation as a minimizer). *The bilinear form is symmetric, i.e. $B[\mathbf{u}, \mathbf{v}] = B[\mathbf{v}, \mathbf{u}]$ for all $\mathbf{u}, \mathbf{v} \in L_{\text{per}}^2$, as \mathbf{A} is symmetric, thus we can express the solution as a minimizer of quadratic functional*

$$\mathbf{e} = \mathbf{E} + \tilde{\mathbf{e}} = \mathbf{E} + \operatorname{argmin}_{\mathbf{v} \in \mathcal{E}} \frac{1}{2} B[\mathbf{v}, \mathbf{v}] - f[\mathbf{v}]$$

In order to formulate Lippmann-Schwinger equation, we introduce a homogeneous reference medium with constant conductivity $\mathbf{A}^0 \in \mathbb{R}^{d \times d}$ that is symmetric and positive definite. Then we can decompose the electric current from Eq. (1.1) into the form

$$\operatorname{div}[\mathbf{A}^0 \mathbf{e}(\mathbf{x})] = -\operatorname{div}[(\mathbf{A}(\mathbf{x}) - \mathbf{A}^0) \mathbf{e}(\mathbf{x})] \quad (2.14)$$

that is appropriate for the use of Fourier Transform technique for a homogeneous problem with conductivity \mathbf{A}^0 ; the right-hand side of Eq. (2.14) then represents the sources, the divergence of so called polarization field. The original problem (1.1)–(1.2) is then equivalent to the periodic Lippmann-Schwinger equation with a derivative of Green function as an integral kernel.

Definition 2.23 (Integral formulation). *Having the same assumption as in Def. (2.20) and in addition let $\mathbf{A}^0 \in \mathbb{R}_{\text{spd}}^{d \times d}$ be symmetric positive definite matrix. We say that $\mathbf{e} \in L_{\text{per}}^2(\mathcal{Y}; \mathbb{R}^d)$ is the solution of Lippmann-Schwinger equation if it satisfies*

$$\mathbf{e}(\mathbf{x}) + \int_{\mathcal{Y}} \Gamma(\mathbf{x} - \mathbf{y}; \mathbf{A}^0) (\mathbf{A}(\mathbf{y}) - \mathbf{A}^0) \mathbf{e}(\mathbf{y}) \, d\mathbf{y} = \mathbf{E}, \quad \mathbf{x} \in \mathcal{Y} \quad (2.15)$$

where the convolution integral $\mathbf{f} \in L_{\text{per}}^2(\mathcal{Y}; \mathbb{R}^d) \rightarrow \int_{\mathcal{Y}} \Gamma(\mathbf{x} - \mathbf{y}; \mathbf{A}^0) \mathbf{f}(\mathbf{y}) \, d\mathbf{y} \in L_{\text{per}}^2(\mathcal{Y}; \mathbb{R}^d)$ is defined with the help of Fourier series as

$$\int_{\mathcal{Y}} \Gamma(\mathbf{x} - \mathbf{y}; \mathbf{A}^0) \mathbf{f}(\mathbf{y}) \, d\mathbf{y} := \sum_{\mathbf{k} \in \mathbb{Z}^d \setminus \{\mathbf{0}\}} \frac{\boldsymbol{\xi}(\mathbf{k}) \otimes \boldsymbol{\xi}(\mathbf{k})}{(\mathbf{A}^0 \boldsymbol{\xi}(\mathbf{k}), \boldsymbol{\xi}(\mathbf{k}))_{\mathbb{R}^d}} \hat{\mathbf{f}}(\mathbf{k}) \varphi_{\mathbf{k}}(\mathbf{x}) \quad (2.16a)$$

$$= \sum_{\mathbf{k} \in \mathbb{Z}^d \setminus \{\mathbf{0}\}} \left(\sum_{\beta=1}^d \frac{k_{\alpha} k_{\beta}}{Y_{\alpha} Y_{\beta}} \hat{\mathbf{f}}_{\beta}(\mathbf{k}) \varphi_{\mathbf{k}}(\mathbf{x}) \right)_{\alpha=1, \dots, d} \quad (2.16b)$$

where $\boldsymbol{\xi}(\mathbf{k}) \in \mathbb{R}^d$ is a vector with components $\xi_{\alpha}(\mathbf{k}) = \frac{k_{\alpha}}{Y_{\alpha}}$, binary operator \otimes denotes tensor product and $(\mathbf{A}^0 \boldsymbol{\xi}(\mathbf{k}), \boldsymbol{\xi}(\mathbf{k}))_{\mathbb{R}^d}$ is a quadratic form in \mathbb{R}^d and $\hat{\mathbf{f}}(\mathbf{k})$ for $\mathbf{k} \in \mathbb{Z}^d$ are the Fourier coefficients.

Remark 2.24. *The definition of convolution integral in Eq. (2.16) is correct, it really maps the function to the space $L_{\text{per}}^2(\mathcal{Y}; \mathbb{R}^d)$ as*

$$\left\| \int_{\mathcal{Y}} \Gamma(\mathbf{x} - \mathbf{y}; \mathbf{A}^0) \mathbf{f}(\mathbf{y}) \, d\mathbf{y} \right\|_{L_{\text{per}}^2(\mathcal{Y}; \mathbb{R}^d)} \leq \frac{1}{c_A} \|\mathbf{f}\|_{L_{\text{per}}^2(\mathcal{Y}; \mathbb{R}^d)}$$

In order to show the equivalence of both formulations in Def. (2.20) and (2.23), we introduce an operator $\mathcal{G}[\cdot]$ based on convolution integral of Lippmann-Schwinger equation. In the following lemma, we show that it is a projection on space \mathcal{E} , the space of curl-free fields with zero mean.

Definition 2.25. *We define an operator $\mathcal{G}[\cdot] : L_{\text{per}}^2(\mathcal{Y}; \mathbb{R}^d) \rightarrow L_{\text{per}}^2(\mathcal{Y}; \mathbb{R}^d)$ as*

$$\mathcal{G}[\mathbf{f}](\mathbf{x}) := \int_{\mathcal{Y}} \Gamma(\mathbf{x} - \mathbf{y}; \mathbf{A}^0) \mathbf{A}^0 \mathbf{f}(\mathbf{y}) \, d\mathbf{y}$$

where the convolution integral is defined in Def. (2.23).

Remark 2.26. The operator \mathcal{G} can be explicitly expressed as

$$\mathcal{G}[\mathbf{f}](\mathbf{x}) = \sum_{\mathbf{k} \in \mathbb{Z}^d \setminus \{\mathbf{0}\}} \sum_{\beta, \gamma} \frac{\xi_\alpha(\mathbf{k}) \xi_\beta(\mathbf{k})}{(\mathbf{A}^0 \boldsymbol{\xi}(\mathbf{k}), \boldsymbol{\xi}(\mathbf{k}))_{\mathbb{R}^d}} A_{\beta\gamma}^0 \widehat{f}_\gamma(\mathbf{k}) \varphi_{\mathbf{k}}(\mathbf{x}) \quad (2.17)$$

where $\xi_\alpha(\mathbf{k}) = \frac{k_\alpha}{Y_\alpha} \in \mathbb{R}$ for $\alpha = 1, \dots, d$ as in Def. (2.23).

Remark 2.27 (The adjoint of \mathcal{G}). For the next use we also discuss the adjoint operator \mathcal{G}^* to \mathcal{G} , i.e. the operator satisfying $(\mathcal{G}\mathbf{v}, \mathbf{w}) = (\mathbf{v}, \mathcal{G}^*\mathbf{w})$ for all $\mathbf{v}, \mathbf{w} \in L^2_{\text{per}}$. With the help of previous Remark 2.26 we can easily express the adjoint operator \mathcal{G}^* as

$$\mathcal{G}^*[\mathbf{f}](\mathbf{x}) = \sum_{\mathbf{k} \in \mathbb{Z}^d \setminus \{\mathbf{0}\}} \sum_{\beta, \gamma} A_{\alpha\beta}^0 \frac{\xi_\beta(\mathbf{k}) \xi_\gamma(\mathbf{k})}{(\mathbf{A}^0 \boldsymbol{\xi}(\mathbf{k}), \boldsymbol{\xi}(\mathbf{k}))_{\mathbb{R}^d}} \widehat{f}_\gamma(\mathbf{k}) \varphi_{\mathbf{k}}(\mathbf{x}) \quad (2.18)$$

leading to a property $\mathbf{A}^0 \mathcal{G} = \mathcal{G}^* \mathbf{A}^0$.

The following lemma provides a projection on space \mathcal{E} and it is based on the well known results about projections of matrices in Fourier space $\frac{\boldsymbol{\xi}(\mathbf{k}) \otimes \boldsymbol{\xi}(\mathbf{k})}{(\boldsymbol{\xi}(\mathbf{k}), \boldsymbol{\xi}(\mathbf{k}))_{\mathbb{R}^d}}$ that are used in homogenization theory, see e.g. [12]. However, the operator \mathcal{G} , including parameter \mathbf{A}^0 and acting on L^2_{per} , has not been presented yet, to the best of our knowledge.

Lemma 2.28 (Projection on curl-free fields). Operator $\mathcal{G}[\cdot] : L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d) \rightarrow L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)$ defined in Def. (2.25) is a projection on the space \mathcal{E} .

Proof. First, we note that Def. (2.25) is correct, i.e. it maps the functions into the space $L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)$ as $\|\mathcal{G}[\mathbf{f}]\|_{L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)} \leq \frac{C_A}{c_A} \|\mathbf{f}\|_{L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)}$. The Fourier coefficients of the real functions satisfies the following symmetry $\widehat{\mathbf{f}}(\mathbf{k}) = \text{conj} \widehat{\mathbf{f}}(-\mathbf{k})$ for all $\mathbf{k} \in \mathbb{Z}^d$. The images of operator \mathcal{G} provides such symmetry

$$\begin{aligned} \text{conj}(\widehat{\mathcal{G}\mathbf{f}}(-\mathbf{k})) &= \text{conj} \left(\sum_{\beta, \gamma=1}^d \frac{\xi_\alpha \xi_\beta}{(\mathbf{A}^0 \boldsymbol{\xi}, \boldsymbol{\xi})_{\mathbb{R}^d}} A_{\beta\gamma}^0 f_\gamma(-\mathbf{k}) \right)_{\alpha=1}^d \\ &= \left(\sum_{\beta, \gamma=1}^d \frac{\xi_\alpha \xi_\beta}{(\mathbf{A}^0 \boldsymbol{\xi}, \boldsymbol{\xi})_{\mathbb{R}^d}} A_{\beta\gamma}^0 \text{conj}(f_\gamma(-\mathbf{k})) \right)_{\alpha=1}^d \\ &= \left(\sum_{\beta, \gamma=1}^d \frac{\xi_\alpha \xi_\beta}{(\mathbf{A}^0 \boldsymbol{\xi}, \boldsymbol{\xi})_{\mathbb{R}^d}} A_{\beta\gamma}^0 f_\gamma(\mathbf{k}) \right)_{\alpha=1}^d = \widehat{\mathcal{G}\mathbf{f}}(\mathbf{k}) \end{aligned}$$

We show that operator \mathcal{G} is the projection by showing $\mathcal{G}^2 = \mathcal{G}$, hence

$$\begin{aligned} \mathcal{G}^2[\mathbf{f}](\mathbf{x}) &= \mathcal{G}[\mathcal{G}[\mathbf{f}]](\mathbf{x}) = \mathcal{G} \left[\sum_{\mathbf{k} \in \mathbb{Z}^d \setminus \{\mathbf{0}\}} \frac{\boldsymbol{\xi}(\mathbf{k}) \otimes \boldsymbol{\xi}(\mathbf{k})}{(\mathbf{A}^0 \boldsymbol{\xi}(\mathbf{k}), \boldsymbol{\xi}(\mathbf{k}))_{\mathbb{R}^d}} \mathbf{A}^0 \widehat{\mathbf{f}}(\mathbf{k}) \varphi_{\mathbf{k}}(\mathbf{x}) \right] \\ &= \sum_{\mathbf{k} \in \mathbb{Z}^d \setminus \{\mathbf{0}\}} \frac{\boldsymbol{\xi}(\mathbf{k}) \otimes \boldsymbol{\xi}(\mathbf{k})}{(\mathbf{A}^0 \boldsymbol{\xi}(\mathbf{k}), \boldsymbol{\xi}(\mathbf{k}))_{\mathbb{R}^d}} \mathbf{A}^0 \frac{\boldsymbol{\xi}(\mathbf{k}) \otimes \boldsymbol{\xi}(\mathbf{k})}{(\mathbf{A}^0 \boldsymbol{\xi}(\mathbf{k}), \boldsymbol{\xi}(\mathbf{k}))_{\mathbb{R}^d}} \mathbf{A}^0 \widehat{\mathbf{f}}(\mathbf{k}) \varphi_{\mathbf{k}}(\mathbf{x}) \\ &= \sum_{\mathbf{k} \in \mathbb{Z}^d \setminus \{\mathbf{0}\}} \frac{\boldsymbol{\xi}(\mathbf{k}) \otimes \boldsymbol{\xi}(\mathbf{k})}{(\mathbf{A}^0 \boldsymbol{\xi}(\mathbf{k}), \boldsymbol{\xi}(\mathbf{k}))_{\mathbb{R}^d}} \mathbf{A}^0 \widehat{\mathbf{f}}(\mathbf{k}) \varphi_{\mathbf{k}}(\mathbf{x}) = \mathcal{G}[\mathbf{f}](\mathbf{x}) \end{aligned}$$

Next, we prove that it maps into the space \mathcal{E} by showing $\text{curl} \mathcal{G}[\mathbf{f}] = 0$ and $\langle \mathcal{G}[\mathbf{f}] \rangle = 0$ for all $\mathbf{f} \in L^2_{\text{per}}(\mathcal{Y}, \mathbb{R}^d)$. The later case is trivial as we sum in Eq. (2.17) over the set $\mathbb{Z}^d \setminus \{\mathbf{0}\}$, i.e. the constant

2 CONTINUOUS FORMULATION

11

term is omitted. Hence, the first case follows for arbitrary test function $\psi \in \mathcal{D}_{\text{per}}(\mathbb{R}^d)$ as

$$\text{curl } \mathcal{G}[\mathbf{f}] := (\text{curl } \mathcal{G}[\mathbf{f}], \psi) = - \left(\langle (\mathcal{G}[\mathbf{f}])_\alpha, \frac{\partial \psi}{\partial x_\beta} \rangle - \langle (\mathcal{G}[\mathbf{f}])_\beta, \frac{\partial \psi}{\partial x_\alpha} \rangle \right)_{\alpha, \beta=1}^d \quad (2.19a)$$

$$= - \left(\sum_{\mathbf{k} \in \mathbb{Z}} \widehat{(\mathcal{G}[\mathbf{f}])}_\alpha(\mathbf{k}) \frac{\widehat{\partial \psi}}{\partial x_\beta}(-\mathbf{k}) - \sum_{\mathbf{k} \in \mathbb{Z}} \widehat{(\mathcal{G}[\mathbf{f}])}_\beta(\mathbf{k}) \frac{\widehat{\partial \psi}}{\partial x_\alpha}(-\mathbf{k}) \right)_{\alpha, \beta=1}^d \quad (2.19b)$$

$$= \left(\sum_{\mathbf{k} \in \mathbb{Z} \setminus \{\mathbf{0}\}} \sum_{\zeta, \theta} \frac{\xi_\alpha \xi_\zeta A_{\zeta\theta}^0 \widehat{f}_\theta}{(\mathbf{A}^0 \boldsymbol{\xi}(\mathbf{k}), \boldsymbol{\xi}(\mathbf{k}))_{\mathbb{R}^d}} i\pi \xi_\beta \widehat{\psi} - \sum_{\mathbf{k} \in \mathbb{Z} \setminus \{\mathbf{0}\}} \sum_{\zeta, \theta} \frac{\xi_\beta \xi_\zeta A_{\zeta\theta}^0 \widehat{f}_\theta}{(\mathbf{A}^0 \boldsymbol{\xi}(\mathbf{k}), \boldsymbol{\xi}(\mathbf{k}))_{\mathbb{R}^d}} i\pi \xi_\alpha \widehat{\psi} \right)_{\alpha, \beta=1}^d \quad (2.19c)$$

$$= (0)_{\alpha, \beta=1}^d \quad (2.19d)$$

where Eq. (2.19a) is just from Def. (2.16), Eq. (2.19b) comes from a property (c) in Lemma 2.11, and Eq. (2.19c) follows from the definition of operator \mathcal{G} and Lemma 2.3.

We show that the projection maps on space \mathcal{E} by showing $\mathcal{G}[\mathbf{f}] = \mathbf{f}$ for all $\mathbf{f} \in \mathcal{E}$. Using again property (c) in Lemma 2.11, we can write for an arbitrary test function $\psi \in \mathcal{D}_{\text{per}}(\mathbb{R}^d)$ and arbitrary α

$$(f_\alpha - \mathcal{G}[\mathbf{f}]_\alpha, \psi)_{L^2_{\text{per}}(\mathcal{Y})} = \sum_{\mathbf{k} \in \mathbb{Z} \setminus \{\mathbf{0}\}} \left(\widehat{f}_\alpha(\mathbf{k}) - \sum_{\beta, \gamma} \frac{\xi_\alpha(\mathbf{k}) \xi_\beta(\mathbf{k})}{(\mathbf{A}^0 \boldsymbol{\xi}(\mathbf{k}), \boldsymbol{\xi}(\mathbf{k}))_{\mathbb{R}^d}} A_{\beta\gamma}^0 \widehat{f}_\gamma(\mathbf{k}) \right) \widehat{\psi}(\mathbf{k})$$

Using Lemma 2.18 giving a potential $u \in H^1_{\text{per}}(\mathcal{Y})$ such that $\nabla u = \mathbf{f}$, we can continue as

$$\begin{aligned} &= \sum_{\mathbf{k} \in \mathbb{Z} \setminus \{\mathbf{0}\}} \left(\frac{\widehat{\partial u}}{\partial x_\alpha}(\mathbf{k}) - \sum_{\beta, \gamma} \frac{\xi_\alpha(\mathbf{k}) \xi_\beta(\mathbf{k})}{(\mathbf{A}^0 \boldsymbol{\xi}(\mathbf{k}), \boldsymbol{\xi}(\mathbf{k}))_{\mathbb{R}^d}} A_{\beta\gamma}^0 \frac{\widehat{\partial u}}{\partial x_\gamma}(\mathbf{k}) \right) \widehat{\psi}(\mathbf{k}) \\ &= \sum_{\mathbf{k} \in \mathbb{Z} \setminus \{\mathbf{0}\}} \left(i\pi \xi_\alpha(\mathbf{k}) \widehat{u}(\mathbf{k}) - \sum_{\beta, \gamma} \frac{\xi_\alpha(\mathbf{k}) \xi_\beta(\mathbf{k})}{(\mathbf{A}^0 \boldsymbol{\xi}(\mathbf{k}), \boldsymbol{\xi}(\mathbf{k}))_{\mathbb{R}^d}} A_{\beta\gamma}^0 i\pi \xi_\gamma(\mathbf{k}) \widehat{u}(\mathbf{k}) \right) \widehat{\psi}(\mathbf{k}) = 0 \end{aligned}$$

where we have used Lemma 2.3 about the Fourier coefficients of a derivative. \square

Using previously proven Lemma, we show the connection between both formulations. We notice that the solutions of both formulations differ by a constant; the unique solution of weak formulation $\tilde{\mathbf{e}} \in \mathcal{E}$ has zero mean while the solution of Lippmann-Schwinger equation \mathbf{e}_{LS} , if exists, has the mean value equal to macroscopic load $\langle \mathbf{e}_{\text{LS}} \rangle = \mathbf{E}$.

Theorem 2.29. *Weak formulation, Def. (2.20), and Lippmann-Schwinger equation, Def. (2.23), are equivalent in the sense the unique solution coincide.*

Proof. The unique microscopic field $\mathbf{e} = \mathbf{E} + \tilde{\mathbf{e}}$ with $\tilde{\mathbf{e}} \in \mathcal{E}$ and $\mathbf{E} \in \mathcal{U}$ coming from a weak formulation satisfies

$$(\mathbf{A}\tilde{\mathbf{e}}, \mathbf{w})_{L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)} = -(\mathbf{A}\mathbf{E}, \mathbf{w})_{L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)}, \quad \forall \mathbf{w} \in \mathcal{E}.$$

Using the property of operator \mathcal{G} , i.e. $\mathcal{G}[\mathbf{w}] = \mathbf{w}$ for all $\mathbf{w} \in \mathcal{E}$, and the identity $\mathbf{A}^0(\mathbf{A}^0)^{-1}$ noting that \mathbf{A}^0 is positive definite and symmetric, we can rewrite the weak formulation into

$$\begin{aligned} (\mathbf{A}^0(\mathbf{A}^0)^{-1}\mathbf{A}(\mathbf{E} + \tilde{\mathbf{e}}), \mathcal{G}\mathbf{w})_{L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)} &= 0, \quad \forall \mathbf{w} \in L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d), \\ ((\mathbf{A}^0)^{-1}\mathbf{A}(\mathbf{E} + \tilde{\mathbf{e}}), \mathbf{A}^0\mathcal{G}\mathbf{w})_{L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)} &= 0, \quad \forall \mathbf{w} \in L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d), \end{aligned}$$

3 DISCRETE SOLUTIONS

12

where we have enlarged the space of test functions. Next, using the properties of adjoint operator \mathcal{G}^* stated in Remark 2.27 and the fact that $\mathbf{A}^0(\cdot) : \mathbf{w} \in L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d) \rightarrow \mathbf{A}^0 \mathbf{w} \in L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)$ is isomorphism, we write

$$\begin{aligned} ((\mathbf{A}^0)^{-1} \mathbf{A}(\mathbf{E} + \tilde{\mathbf{e}}), \mathcal{G}^* \mathbf{A}^0 \mathbf{w})_{L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)} &= 0, \quad \forall \mathbf{w} \in L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d) \\ (\mathcal{G}(\mathbf{A}^0)^{-1} \mathbf{A} \mathbf{e}, \mathbf{v})_{L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)} &= 0, \quad \forall \mathbf{v} \in L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d) \end{aligned}$$

where $\mathbf{v} = \mathbf{A}^0 \mathbf{w}$.

Adding and subtracting the terms $(\tilde{\mathbf{e}}, \mathbf{v})_{L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)}$ and $(\mathbf{E}, \mathbf{v})_{L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)}$ and using a property $\mathcal{G} \mathbf{E} = 0$ heads to

$$\begin{aligned} (\tilde{\mathbf{e}}, \mathbf{v})_{L^2_{\text{per}}} + (\mathbf{E}, \mathbf{v})_{L^2_{\text{per}}} + (\mathcal{G}(\mathbf{A}^0)^{-1} \mathbf{A} \mathbf{e}, \mathbf{v})_{L^2_{\text{per}}} - (\tilde{\mathbf{e}}, \mathbf{v})_{L^2_{\text{per}}} &= (\mathbf{E}, \mathbf{v})_{L^2_{\text{per}}}, \quad \forall \mathbf{v} \in L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d) \\ (\mathbf{e} + \mathcal{G}(\mathbf{A}^0)^{-1}(\mathbf{A} - \mathbf{A}^0)\mathbf{e}, \mathbf{v})_{L^2(\mathcal{Y}; \mathbb{R}^d)} &= (\mathbf{E}, \mathbf{v})_{L^2(\mathcal{Y}; \mathbb{R}^d)}, \quad \forall \mathbf{v} \in L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d) \end{aligned}$$

Since it is tested for all test function $\mathbf{v} \in L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)$, we abandon the scalar product to equation holding almost everywhere

$$\begin{aligned} \mathbf{e}(\mathbf{x}) + \mathcal{G}[(\mathbf{A}^0)^{-1}(\mathbf{A} - \mathbf{A}^0)\mathbf{e}](\mathbf{x}) &= \mathbf{E}, \quad \text{for a.a. } \mathbf{x} \in \mathcal{Y} \\ \mathbf{e}(\mathbf{x}) + \int_{\mathcal{Y}} \mathbf{\Gamma}(\mathbf{x} - \mathbf{y})(\mathbf{A}(\mathbf{y}) - \mathbf{A}^0)\mathbf{e}(\mathbf{y}) \, d\mathbf{y} &= \mathbf{E}, \quad \text{for a.a. } \mathbf{x} \in \mathcal{Y} \end{aligned}$$

The functions from $L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)$ equal to each other if they are the same almost everywhere, thus we can choose an appropriate representative in order to satisfy the equation everywhere and the first implication is done.

The contrary case can be done by multiplication with a test function $\mathbf{v} \in L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)$, by integration over \mathcal{Y} , by split of the solution with projection $\mathbf{e} = \mathcal{G} \mathbf{e} + (\mathbf{I} - \mathcal{G})\mathbf{e}$, and by the same arguments about the projections \mathcal{G} and its adjoint \mathcal{G}^* as in the first part of the proof. \square

3 Discrete solutions

This section is dedicated to the discretization of weak formulation in Def. 2.20. The trigonometric polynomials are used as the basis functions, see Sec. 3.1, and their approximation properties are stated in Sec. 3.2.

Next, in Sec. 3.3, we define Galerkin approximation (GA) and Galerkin approximation with numerical integration (GAwNI) and we show the fully discrete formulation of GAwNI that is appropriate for the solution by Conjugate gradients. In Sec. 3.4, we provide convergence of discrete solutions to the solution of weak formulation, then in Sec. 3.5 the convergence is shown for rough material coefficients by their regularization. Finally, Sec. 3.6 provides information about the numerical solution of a corresponding linear system by Conjugate gradients.

In the sequel, vector $\mathbf{N} \in \mathbb{N}^d$ is reserved for a number of discretization points, then scalar $|\mathbf{N}|_{\Pi} := \prod_{\alpha} N_{\alpha}$ denotes the number of degrees of freedom. If N_{α} is odd for all α we talk about the odd number of discretization points. We use a simplification in order to make the theory of real valued trigonometric polynomials uncomplicated — in the sequel, only the odd number of discretization points is considered, Sec. 3.1.

A multi-index notation is employed, in which $\mathbb{R}^{\mathbf{N}}$ represents $\mathbb{R}^{N_1 \times \dots \times N_d}$. Set $\mathbb{R}^{d \times \mathbf{N}}$ denotes the space of vectors \mathbf{v} with components $v_{\alpha}^{\mathbf{n}}$ and $\mathbb{R}^{d \times d \times \mathbf{N} \times \mathbf{N}}$ the space of matrices \mathbf{A} with components $A_{\alpha\beta}^{\mathbf{nm}}$ for α, β and $\mathbf{n}, \mathbf{m} \in \mathbb{Z}_{\mathbf{N}}^d$. Next, vectors $\mathbf{v}^{\mathbf{n}} \in \mathbb{R}^d$ for $\mathbf{n} \in \mathbb{Z}_{\mathbf{N}}^d$ and $\mathbf{v}_{\alpha} \in \mathbb{R}^{\mathbf{N}}$ for α represent subvectors of \mathbf{v} with components $v_{\alpha}^{\mathbf{n}}$. Analogically, submatrices $\mathbf{A}^{\mathbf{nm}} \in \mathbb{R}^{d \times d}$ and $\mathbf{A}_{\alpha\beta} \in \mathbb{R}^{\mathbf{N} \times \mathbf{N}}$ can be defined. A scalar product on set $\mathbb{R}^{d \times \mathbf{N}}$ is defined as $(\mathbf{u}, \mathbf{v})_{\mathbb{R}^{d \times \mathbf{N}}} := \sum_{\alpha} \sum_{\mathbf{n} \in \mathbb{Z}_{\mathbf{N}}^d} u_{\alpha}^{\mathbf{n}} v_{\alpha}^{\mathbf{n}}$ and matrix \mathbf{A} by vector \mathbf{v} multiplication as $(\mathbf{A}\mathbf{v})_{\alpha}^{\mathbf{n}} := \sum_{\beta} \sum_{\mathbf{m} \in \mathbb{Z}_{\mathbf{N}}^d} A_{\alpha\beta}^{\mathbf{nm}} v_{\beta}^{\mathbf{m}}$. Matrix \mathbf{A} is symmetric positive definite if relation

$A_{\alpha\beta}^{mn} = A_{\beta\alpha}^{nm}$ holds for all components and inequality $(\mathbf{A}\mathbf{v}, \mathbf{v})_{\mathbb{R}^{d \times N}} > 0$ applies for arbitrary $\mathbf{v} \in \mathbb{R}^{d \times N}$. We use the serif font for vectors $\mathbf{v} \in \mathbb{R}^{d \times N}$ and matrices $\mathbf{A} \in \mathbb{R}^{d \times d \times N \times N}$ to distinguish from vectors $\mathbf{E} \in \mathbb{R}^d$ and matrices $\mathbf{A}(\mathbf{x}) \in \mathbb{R}^{d \times d}$ and from vector valued functions $\mathbf{v} \in L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)$. In order to differentiate vectors and matrices for different number of discretization points N , we often write them with subscript N , i.e. \mathbf{v}_N and \mathbf{A}_N .

3.1 Finite dimensional space of trigonometric polynomials

In this section we define a finite dimensional space of trigonometric polynomials and provide its properties.

Definition 3.1. For odd N , we define the space of trigonometric polynomials as

$$\mathcal{T}_N = \left\{ \sum_{\mathbf{k} \in \underline{\mathbb{Z}}_N^d} \hat{\mathbf{c}}^{\mathbf{k}} \varphi_{\mathbf{k}}(\mathbf{x}); \hat{\mathbf{c}}^{\mathbf{k}} \in \mathbb{C}, \hat{\mathbf{c}}^{\mathbf{k}} = \overline{(\hat{\mathbf{c}}^{-\mathbf{k}})} \text{ for } \mathbf{k} \in \underline{\mathbb{Z}}_N^d \right\} \quad (3.1)$$

where the index set $\underline{\mathbb{Z}}_N^d$, already defined in Lemma 2.11, is expressed as

$$\underline{\mathbb{Z}}_N^d = \left\{ \mathbf{k} \in \mathbb{Z}^d : -\frac{N_\alpha}{2} \leq k_\alpha < \frac{N_\alpha}{2}, \right\} \quad (3.2)$$

and $\varphi_{\mathbf{k}}$ are basis functions defined already in Eq. (2.4), i.e. $\varphi_{\mathbf{k}}(\mathbf{x}) = \exp\left(i\pi \sum_{\alpha=1}^d \frac{k_\alpha x_\alpha}{Y_\alpha}\right)$. The scalar valued functions from the space \mathcal{T}_N are denoted with a subscript N , e.g. $f_N \in \mathcal{T}_N$. Analogically, we define the vector valued trigonometric polynomials as

$$\mathcal{T}_N^d = \{ \mathbf{f} \in L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d) : f_\alpha \in \mathcal{T}_N \text{ for } \alpha = 1, \dots, d \},$$

Remark 3.2. The condition $\hat{\mathbf{c}}^{\mathbf{k}} = \overline{(\hat{\mathbf{c}}^{-\mathbf{k}})}$ for $\mathbf{k} \in \underline{\mathbb{Z}}_N^d$ in previous definition satisfies that trigonometric polynomials are real valued. The fact that N is odd satisfies that the boundary frequencies $\frac{N_\alpha}{2}$ do not occur at all. This assures that the index set $\underline{\mathbb{Z}}_N^d$ is symmetric around zero, thus all frequencies $\mathbf{k} \in \underline{\mathbb{Z}}_N^d$ has the opposite counterpart $-\mathbf{k}$ as well as the Fourier coefficients.

Now, we show two possible representations of trigonometric polynomials and the connection between them. They can be expressed through their Fourier coefficients and through their function values at nodal points.

Remark 3.3 (Representation of $\mathbf{v}_N \in \mathcal{T}_N^d$ through its Fourier coefficients). The function of trigonometric polynomials, $v_N \in \mathcal{T}_N$, can be expressed as

$$v_N(\mathbf{x}) = \sum_{\mathbf{k} \in \underline{\mathbb{Z}}_N^d} \hat{v}_N(\mathbf{k}) \varphi_{\mathbf{k}}(\mathbf{x}),$$

where $\hat{v}_N(\mathbf{k}) \in \mathbb{C}^d$ are Fourier coefficients defined in Def. (2.2).

Definition 3.4 (nodal points). We define the nodal points of a periodic unit cell \mathcal{Y} as points $\mathbf{x}_N^{\mathbf{k}} \in \mathcal{Y}$ for $\mathbf{k} \in \underline{\mathbb{Z}}_N^d$ explicitly expressed as

$$\mathbf{x}_N^{\mathbf{k}} := \sum_{\alpha} \frac{2Y_\alpha k_\alpha}{N_\alpha} \boldsymbol{\epsilon}_\alpha = \left(\frac{2Y_1 k_1}{N_1}, \frac{2Y_2 k_2}{N_2}, \dots, \frac{2Y_d k_d}{N_d} \right).$$

Thus each vector $\mathbf{k} \in \underline{\mathbb{Z}}_N^d$ points out one point $\mathbf{x}_N^{\mathbf{k}} \in \mathcal{Y}$ from the regular grid $\{\mathbf{x}_N^{\mathbf{k}} \in \mathcal{Y}; \mathbf{k} \in \underline{\mathbb{Z}}_N^d\}$.

Remark 3.5 (Representation of $\mathbf{v}_N \in \mathcal{T}_N^d$ through its nodal values). The trigonometric polynomials can be represented through its nodal values as

$$\mathbf{v}_N(\mathbf{x}) = \sum_{\mathbf{m} \in \underline{\mathbb{Z}}_N^d} \mathbf{v}_N(\mathbf{x}_N^{\mathbf{m}}) \varphi_{N,\mathbf{m}}(\mathbf{x})$$

where the functions $\varphi_{N,\mathbf{k}}(\mathbf{x}) \in \mathcal{T}_N^d$ state for

$$\varphi_{N,\mathbf{m}}(\mathbf{x}) = \frac{1}{|N|_{\Pi}} \sum_{\mathbf{k} \in \mathbb{Z}_N^d} \varphi_{\mathbf{k}}(\mathbf{x}) \omega_N^{-\mathbf{m}\mathbf{k}} = \frac{1}{|N|_{\Pi}} \sum_{\mathbf{k} \in \mathbb{Z}_N^d} \exp \left\{ i\pi \sum_{\alpha=1}^d k_{\alpha} \left(\frac{x_{\alpha}}{Y_{\alpha}} - \frac{2m_{\alpha}}{N_{\alpha}} \right) \right\} \quad (3.3)$$

The are called shape basis functions and satisfy the Dirac delta property at nodal points

$$\varphi_{N,\mathbf{m}}(\mathbf{x}_N^{\mathbf{k}}) = \delta_{\mathbf{k}\mathbf{m}}, \quad (3.4)$$

i.e. it is such trigonometric polynomial having the value one in a particular nodal point of regular grid and value zero in the rest of the nodal points.

Remark 3.6 (Connection of representations). Let $v_N(\mathbf{x}) \in \mathcal{T}_N$. Then Fourier coefficients $\hat{v}_N(\mathbf{k})$, $\mathbf{k} \in \mathbb{Z}_N^d$ can be obtained with the Discrete Fourier Transform from the values at nodal points $(v_N(\mathbf{x}_N^{\mathbf{m}}))_{\mathbf{m} \in \mathbb{Z}_N^d}$. Hence

$$\begin{aligned} \hat{v}(\mathbf{k}) &= \frac{1}{|N|_{\Pi}} \sum_{\mathbf{m} \in \mathbb{Z}_N^d} v(\mathbf{x}_N^{\mathbf{m}}) \omega_N^{-\mathbf{k}\mathbf{m}}, \quad \mathbf{k} \in \mathbb{Z}_N^d \\ v(\mathbf{x}_N^{\mathbf{m}}) &= \sum_{\mathbf{k} \in \mathbb{Z}_N^d} \hat{v}(\mathbf{k}) \omega_N^{\mathbf{k}\mathbf{m}}, \quad \mathbf{m} \in \mathbb{Z}_N^d \end{aligned}$$

where $\omega_N^{\mathbf{k}\mathbf{m}} = \exp \left(2\pi i \sum_{\alpha} \frac{k_{\alpha} m_{\alpha}}{N_{\alpha}} \right)$ are the coefficients of DFT with the property

$$\sum_{\mathbf{m} \in \mathbb{Z}_N^d} \omega_N^{\mathbf{k}\mathbf{m}} \overline{\omega_N^{\mathbf{l}\mathbf{m}}} = \delta_{\mathbf{k}\mathbf{l}} |N|_{\Pi} \quad (3.5)$$

Remark 3.7 (Orthogonality of $\varphi_{N,\mathbf{k}}$). For shape basis functions, the following holds

$$(\varphi_{N,\mathbf{k}}, \varphi_{N,\mathbf{l}}) = \frac{\delta_{\mathbf{k}\mathbf{l}}}{|N|_{\Pi}}$$

Proof.

$$\begin{aligned} (\varphi_{N,\mathbf{k}}, \varphi_{N,\mathbf{l}})_{L_{\text{per}}^2} &= \frac{1}{|N|_{\Pi}} \sum_{\mathbf{m} \in \mathbb{Z}_N^d} \omega_N^{\mathbf{m}\mathbf{k}} (\varphi_{\mathbf{m}}, \varphi_{N,\mathbf{l}})_{L_{\text{per}}^2} \\ &= \left(\frac{1}{|N|_{\Pi}} \right)^2 \sum_{\mathbf{m} \in \mathbb{Z}_N^d} \sum_{\mathbf{n} \in \mathbb{Z}_N^d} \omega_N^{\mathbf{m}\mathbf{k}} \overline{\omega_N^{\mathbf{n}\mathbf{l}}} (\varphi_{\mathbf{m}}, \varphi_{\mathbf{n}})_{L_{\text{per}}^2} \\ &= \left(\frac{1}{|N|_{\Pi}} \right)^2 \sum_{\mathbf{m} \in \mathbb{Z}_N^d} \sum_{\mathbf{n} \in \mathbb{Z}_N^d} \omega_N^{\mathbf{m}\mathbf{k}} \overline{\omega_N^{\mathbf{n}\mathbf{l}}} \delta_{\mathbf{m}\mathbf{n}} \\ &= \left(\frac{1}{|N|_{\Pi}} \right)^2 \sum_{\mathbf{m} \in \mathbb{Z}_N^d} \omega_N^{\mathbf{m}\mathbf{k}} \omega_N^{-\mathbf{m}\mathbf{l}} = \frac{\delta_{\mathbf{k}\mathbf{l}}}{|N|_{\Pi}} \end{aligned}$$

□

Definition 3.8. Operator $\mathcal{I}_N : \mathcal{T}_N^d \rightarrow \mathbb{R}^{d \times N}$ stocks the values of the trigonometric polynomials at the nodal points to a vector $\mathcal{I}_N[\mathbf{v}_N] = (\mathbf{v}_{N,\alpha}(\mathbf{x}_N^{\mathbf{n}}))_{\alpha=1, \dots, d}^{\mathbf{n} \in \mathbb{Z}_N^d}$.

Lemma 3.9. Operator \mathcal{I}_N from the previous definition is an isomorphism.

Proof. The proof is the consequence of its Def. 3.8 and the Dirac delta property (3.4). □

3.2 Orthogonal and Interpolation Projection

This section provides lemmas about trigonometric approximations, i.e. about the estimates of orthogonal projection and interpolation projection. In [18, Section 8], this topic is well described for a one-dimensional setting and partially for a two-dimensional setting with equal number of discretization points in each direction. Thus we present here the generalization of the estimates for arbitrary dimension and for arbitrary number of discretization points. However, the proofs follows the same ideas as in [18, Section 8].

Definition 3.10 (Orthogonal projection). *For odd $N \in \mathbb{N}$, we define an operator $P_N[\cdot] : L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d) \rightarrow \mathcal{T}_N^d$ as*

$$P_N[\mathbf{v}](\mathbf{x}) = \sum_{\mathbf{k} \in \mathbb{Z}_N^d} \hat{\mathbf{v}}(\mathbf{k}) \varphi_{\mathbf{k}}(\mathbf{x})$$

where $\hat{\mathbf{v}}(\mathbf{k})$ denotes Fourier coefficients of \mathbf{v} .

Remark 3.11. *It can be easily shown that operator $P_N[\cdot]$ is the orthogonal projection.*

The possibility of expression functions as a Fourier series is stated in the following lemma with a proof that is based on the density of the set of trigonometric polynomials $\{\varphi_{\mathbf{k}}\}_{\mathbf{k} \in \mathbb{Z}^d}$ in space L^2_{per} , e.g. [11].

Lemma 3.12 (Approximation by orthogonal projection). *For a function $\mathbf{v} \in L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)$, the following holds*

$$\lim_{N \rightarrow \infty} \|\mathbf{v} - P_N[\mathbf{v}]\|_{L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)} = 0$$

where $N \rightarrow \infty$ means $\min_{\alpha} N_{\alpha} \rightarrow \infty$.

Lemma 3.13 (Approximation by orthogonal projection). *For a function $\mathbf{v} \in H^{\mu}_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)$ with $\mu \in \mathbb{R}$ we have*

$$\|\mathbf{v} - P_N[\mathbf{v}]\|_{H^{\lambda}_{\text{per}}} \leq \left(\min_{\alpha=1, \dots, d} \frac{N_{\alpha}}{2Y_{\alpha}} \right)^{(\lambda-\mu)} \|\mathbf{v}\|_{H^{\mu}_{\text{per}}} \quad (3.6)$$

where $\lambda \leq \mu$ and $N \neq \mathbf{0}$.

Proof. The proof is a consequence of direct calculation.

$$\begin{aligned} \|\mathbf{v} - P_N[\mathbf{v}]\|_{H^{\lambda}_{\text{per}}}^2 &= \sum_{\mathbf{Z}^d \setminus \mathbb{Z}_N^d} \|\underline{\xi}(\mathbf{k})\|_2^{2\lambda} \|\hat{\mathbf{v}}(\mathbf{k})\|_2^2 \\ &= \sum_{\mathbf{Z}^d \setminus \mathbb{Z}_N^d} \|\underline{\xi}(\mathbf{k})\|_2^{2(\lambda-\mu)} \|\underline{\xi}(\mathbf{k})\|_2^{2\mu} \|\hat{\mathbf{v}}(\mathbf{k})\|_2^2 \\ &\leq \left(\min_{\alpha=1, \dots, d} \frac{N_{\alpha}}{2Y_{\alpha}} \right)^{2(\lambda-\mu)} \sum_{\mathbf{Z}^d \setminus \mathbb{Z}_N^d} \|\underline{\xi}(\mathbf{k})\|_2^{2\mu} \|\hat{\mathbf{v}}(\mathbf{k})\|_2^2 \\ &\leq \left(\min_{\alpha=1, \dots, d} \frac{N_{\alpha}}{2Y_{\alpha}} \right)^{2(\lambda-\mu)} \|\mathbf{v}\|_{H^{\mu}_{\text{per}}}^2 \end{aligned}$$

□

Definition 3.14 (Interpolation projection). *For \mathbf{Y} -periodic continuous functions $C_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)$ and odd N , we define an operator $Q_N[\cdot] : C_{\text{per}} \rightarrow \mathcal{T}_N^d$ with conditions*

$$Q_N[\mathbf{v}] \in \mathcal{T}_N^d \quad Q_N[\mathbf{v}](\mathbf{x}_N^{\mathbf{k}}) = \mathbf{v}(\mathbf{x}_N^{\mathbf{k}}) \quad \text{for } \mathbf{k} \in \mathbb{Z}_N^d.$$

Alternatively, the interpolation operator can be defined as

$$Q_N[v] = \sum_{\mathbf{k} \in \mathbb{Z}_N^d} v(\mathbf{x}_N^{\mathbf{k}}) \varphi_{N,\mathbf{k}}.$$

This is based on two possibilities of trigonometric polynomial representation and their connection Rem. 3.6.

Remark 3.15. *The interpolation projection is well defined for $H_{\text{per}}^\mu(\mathcal{Y})$ with $\mu > \frac{d}{2}$ as the space $H_{\text{per}}^\mu(\mathcal{Y})$ is embedded into the space of continuous functions C_{per} .*

Remark 3.16. *The mapping Q_N from previous definition is clearly a projection, nevertheless not orthogonal anymore.*

Next, we state the lemma about interpolation operator in Fourier space; the lemma is used for the facts about approximation by interpolation projection.

Lemma 3.17 (Interpolation operator Q_N in Fourier space). *For $v \in H_{\text{per}}^\mu(\mathcal{Y})$, $\mu > \frac{d}{2}$*

$$Q_N[v](\mathbf{x}) = \sum_{\mathbf{k} \in \mathbb{Z}_N^d} \left[\sum_{\mathbf{l} \in \mathbb{Z}^d} \hat{v}(\mathbf{k} + \mathbf{l} \odot \mathbf{N}) \right] \varphi_{\mathbf{k}}(\mathbf{x})$$

where \odot denotes element-wise multiplication.

Proof. Using the property of operator $Q_N[\cdot]$ to being a projection on \mathcal{T}_N , we have $Q_N[\varphi_{\mathbf{k}}] = \varphi_{\mathbf{k}}$ as $\varphi_{\mathbf{k}} \in \mathcal{T}_N$. Moreover for all $\mathbf{k}, \mathbf{m} \in \mathbb{Z}_N^d$ and $\mathbf{l} \in \mathbb{Z}^d$, the basis functions $\varphi_{\mathbf{k}}(\mathbf{x})$ and $\varphi_{\mathbf{k} + \mathbf{l} \odot \mathbf{N}}(\mathbf{x})$ equals one another² at nodal points $\mathbf{x}_N^{\mathbf{m}} = \sum_{\alpha=1}^d \frac{2Y_\alpha k_\alpha}{N_\alpha} \epsilon_\alpha$ as

$$\begin{aligned} \varphi_{\mathbf{k} + \mathbf{l} \odot \mathbf{N}}(\mathbf{x}_N^{\mathbf{m}}) &= \varphi_{\mathbf{k}}(\mathbf{x}_N^{\mathbf{m}}) \cdot \exp\left(i\pi \sum_{\alpha=1}^d \frac{l_\alpha N_\alpha 2Y_\alpha m_\alpha}{Y_\alpha N_\alpha}\right) \\ &= \varphi_{\mathbf{k}}(\mathbf{x}_N^{\mathbf{m}}) \cdot \exp\left(2i\pi \sum_{\alpha=1}^d l_\alpha m_\alpha\right) = \varphi_{\mathbf{k}}(\mathbf{x}_N^{\mathbf{m}}) \end{aligned}$$

and thus we even have

$$Q_N[\varphi_{\mathbf{k} + \mathbf{l} \odot \mathbf{N}}] = Q_N[\varphi_{\mathbf{k}}] = \varphi_{\mathbf{k}}, \quad (\mathbf{k} \in \mathbb{Z}_N^d, \mathbf{l} \in \mathbb{Z}^d).$$

As the operator Q_N is obviously linear, we can finish the proof

$$\begin{aligned} Q_N[v](\mathbf{x}) &= Q_N \left[\sum_{\mathbf{m} \in \mathbb{Z}^d} \hat{v}(\mathbf{m}) \varphi_{\mathbf{m}}(\mathbf{x}) \right] (\mathbf{x}) = Q_N \left[\sum_{\mathbf{k} \in \mathbb{Z}_N^d} \sum_{\mathbf{l} \in \mathbb{Z}^d} \hat{v}(\mathbf{k} + \mathbf{l} \odot \mathbf{N}) \varphi_{\mathbf{k} + \mathbf{l} \odot \mathbf{N}}(\mathbf{x}) \right] (\mathbf{x}) \\ &= \sum_{\mathbf{k} \in \mathbb{Z}_N^d} \left[\sum_{\mathbf{l} \in \mathbb{Z}^d} \hat{v}(\mathbf{k} + \mathbf{l} \odot \mathbf{N}) \right] \varphi_{\mathbf{k}}(\mathbf{x}) \end{aligned}$$

□

Lemma 3.18 (Approximation by interpolation projection). *Let $v \in H_{\text{per}}^\mu(\mathcal{Y})$ with $\mu > \frac{d}{2}$, then*

$$\|v - Q_N[v]\|_{H_{\text{per}}^\lambda(\mathcal{Y})} \leq c_{\lambda,\mu} \left(\min_{\alpha=1,\dots,d} \frac{N_\alpha}{2Y_\alpha} \right)^{\lambda-\mu} \|v\|_{H_{\text{per}}^\mu(\mathcal{Y})} \quad (3.7)$$

²operator \odot denotes element-wise multiplication

for $0 \leq \lambda \leq \mu$ and

$$c_{\lambda, \mu} = \left(1 + d^\lambda \rho^{2\lambda} \sum_{\mathbf{l} \in \mathbb{N}_0^d \setminus \{\mathbf{0}\}} \|\mathbf{l}\|_2^{-2\mu} \right)^{\frac{1}{2}}$$

$$\rho = \frac{\max_{\alpha=1, \dots, d} \frac{N_\alpha}{2Y_\alpha}}{\min_{\alpha=1, \dots, d} \frac{N_\alpha}{2Y_\alpha}}.$$

Proof. The proof generalizes the proof of Theorems 8.2.1 and 8.5.3 in [18]. Using the identity

$$\|v - Q_{\mathbf{N}}[v]\|_{H_{\text{per}}^\lambda(\mathcal{Y})}^2 = \|v - P_{\mathbf{N}}[v]\|_{H_{\text{per}}^\lambda(\mathcal{Y})}^2 + \|P_{\mathbf{N}}[v] - Q_{\mathbf{N}}[v]\|_{H_{\text{per}}^\lambda(\mathcal{Y})}^2$$

and Lem. 3.13 about approximation with orthogonal projection, we have to estimate the second term $\|P_{\mathbf{N}}[v] - Q_{\mathbf{N}}[v]\|_{H_{\text{per}}^\lambda(\mathcal{Y})}^2$.

Using Lem. 3.17, we can express Fourier coefficients of undergoing term

$$(P_{\mathbf{N}}[v] - Q_{\mathbf{N}}[v])(\mathbf{k}) = \begin{cases} \sum_{\mathbf{l} \in \mathbb{Z}^d \setminus \{\mathbf{0}\}} \hat{v}(\mathbf{k} + \mathbf{l} \odot \mathbf{N}), & \text{for } \mathbf{k} \in \mathbb{Z}_{\mathbf{N}}^d \\ 0, & \text{for } \mathbf{k} \in \mathbb{Z}^d \setminus \mathbb{Z}_{\mathbf{N}}^d. \end{cases}$$

Then some algebra with Cauchy inequality yield

$$\begin{aligned} \|P_{\mathbf{N}}[v] - Q_{\mathbf{N}}[v]\|_{H_{\text{per}}^\lambda(\mathcal{Y})}^2 &= \sum_{\mathbf{k} \in \mathbb{Z}_{\mathbf{N}}^d} \frac{\|\underline{\xi}(\mathbf{k})\|_2^{2\lambda}}{\|\underline{\xi}(\mathbf{k} + \mathbf{l} \odot \mathbf{N})\|_2^{2\mu}} \left| \sum_{\mathbf{l} \in \mathbb{Z}^d \setminus \{\mathbf{0}\}} \hat{v}(\mathbf{k} + \mathbf{l} \odot \mathbf{N}) \right|^2 \\ &\leq \sum_{\mathbf{k} \in \mathbb{Z}_{\mathbf{N}}^d} \left(\sum_{\mathbf{l} \in \mathbb{Z}^d \setminus \{\mathbf{0}\}} \frac{\|\underline{\xi}(\mathbf{k})\|_2^{2\lambda}}{\|\underline{\xi}(\mathbf{k} + \mathbf{l} \odot \mathbf{N})\|_2^{2\mu}} \cdot \|\underline{\xi}(\mathbf{k} + \mathbf{l} \odot \mathbf{N})\|_2^{2\mu} |\hat{v}(\mathbf{k} + \mathbf{l} \odot \mathbf{N})| \right)^2 \\ &\leq \sum_{\mathbf{k} \in \mathbb{Z}_{\mathbf{N}}^d} \left(\sum_{\mathbf{l} \in \mathbb{Z}^d \setminus \{\mathbf{0}\}} \frac{\|\underline{\xi}(\mathbf{k})\|_2^{2\lambda}}{\|\underline{\xi}(\mathbf{k} + \mathbf{l} \odot \mathbf{N})\|_2^{2\mu}} \right) \left(\sum_{\mathbf{l} \in \mathbb{Z}^d \setminus \{\mathbf{0}\}} \|\underline{\xi}(\mathbf{k} + \mathbf{l} \odot \mathbf{N})\|_2^{2\mu} |\hat{v}(\mathbf{k} + \mathbf{l} \odot \mathbf{N})|^2 \right) \\ &\leq \varepsilon_{\mathbf{N}}^2 \|v\|_{H_{\text{per}}^\mu(\mathcal{Y})}^2 \end{aligned}$$

where

$$\begin{aligned} \varepsilon_{\mathbf{N}}^2 &= \max_{\mathbf{k} \in \mathbb{Z}_{\mathbf{N}}^d} \left(\sum_{\mathbf{l} \in \mathbb{Z}^d \setminus \{\mathbf{0}\}} \frac{\|\underline{\xi}(\mathbf{k})\|_2^{2\lambda}}{\|\underline{\xi}(\mathbf{k} + \mathbf{l} \odot \mathbf{N})\|_2^{2\mu}} \right) \\ &\leq \|\underline{\xi}(\mathbf{N}/2)\|_2^{2\lambda} \max_{\mathbf{k} \in \mathbb{Z}_{\mathbf{N}}^d} \left(\sum_{\mathbf{l} \in \mathbb{Z}^d \setminus \{\mathbf{0}\}} \frac{1}{\left[\sum_{\alpha=1}^d \left(\frac{k_\alpha + l_\alpha N_\alpha}{Y_\alpha} \right)^2 \right]^\mu} \right) \\ &= \|\underline{\xi}(\mathbf{N}/2)\|_2^{2\lambda} \max_{\mathbf{k} \in \mathbb{Z}_{\mathbf{N}}^d} \left(\sum_{\mathbf{l} \in \mathbb{Z}^d \setminus \{\mathbf{0}\}} \frac{1}{\left[\sum_{\alpha=1}^d \left(\frac{N_\alpha}{2Y_\alpha} \right)^2 \left(\frac{k_\alpha}{2N_\alpha} + 2l_\alpha \right)^2 \right]^\mu} \right) \\ &= \|\underline{\xi}(\mathbf{N}/2)\|_2^{2\lambda} \left(\sum_{\mathbf{l} \in \mathbb{Z}^d \setminus \{\mathbf{0}\}} \frac{1}{\left[\sum_{\alpha=1}^d \left(\frac{N_\alpha}{2Y_\alpha} \right)^2 (1 + 2l_\alpha)^2 \right]^\mu} \right) \\ &\leq \|\underline{\xi}(\mathbf{N}/2)\|_2^{2\lambda} \left(\min_{\alpha=1, \dots, d} \frac{N_\alpha}{2Y_\alpha} \right)^{-2\mu} \left(\sum_{\mathbf{l} \in \mathbb{Z}^d \setminus \{\mathbf{0}\}} \frac{1}{\left[\sum_{\alpha=1}^d (1 + 2l_\alpha)^2 \right]^\mu} \right) \\ &\leq d^\lambda \left(\max_{\alpha=1, \dots, d} \frac{N_\alpha}{2Y_\alpha} \right)^{2\lambda} \left(\min_{\alpha=1, \dots, d} \frac{N_\alpha}{2Y_\alpha} \right)^{-2\mu} \left(\sum_{\mathbf{l} \in \mathbb{N}_0^d \setminus \{\mathbf{0}\}} \frac{1}{\left[\sum_{\alpha=1}^d l_\alpha^2 \right]^\mu} \right). \end{aligned}$$

Hence

$$\|P_{\mathbf{N}}[v] - Q_{\mathbf{N}}[v]\|_{H_{\text{per}}^\lambda(\mathcal{Y})}^2 \leq d^\lambda \rho^{2\lambda} \sum_{\mathbf{l} \in \mathbb{N}_0^d \setminus \{\mathbf{0}\}} \frac{1}{\|\mathbf{l}\|_2^{2\mu}} \left(\min_{\alpha=1, \dots, d} \frac{N_\alpha}{2Y_\alpha} \right)^{2(\lambda-\mu)} \|v\|_{H_{\text{per}}^\mu(\mathcal{Y})}^2$$

and together with Eq. (3.6), estimate in Eq. (3.7) follows. \square

3.3 Galerkin approximation with numerical integration

This section presents Galerkin approximation (GA) and Galerkin approximation with numerical integration (GAwNI) that leads to a linear system equivalent with Moulinec-Suquet algorithm, cf. [21].

Definition 3.19 (Finite dimensional space). *We define the finite dimensional space used for discretization $\mathcal{E}_{\mathbf{N}} = \mathcal{T}_{\mathbf{N}}^d \cap \mathcal{E}$, the space of curl-free trigonometric polynomials with zero mean.*

There is no doubt that the space $\mathcal{E}_{\mathbf{N}}$ is empty. Since we have projection \mathcal{G} from Def. 2.25, the Helmholtz decomposition for the space of trigonometric polynomials can be written down

$$\mathcal{T}_{\mathbf{N}}^d = \mathcal{U} \oplus^\perp \mathcal{E}_{\mathbf{N}} \oplus^\perp \mathcal{J}_{\mathbf{N}}. \quad (3.8)$$

Definition 3.20 (Galerkin approximation). *Let $\mathcal{E}_{\mathbf{N}}$ be a space from Def. 3.19. The function $\tilde{e}_{\mathbf{N}} \in \mathcal{E}_{\mathbf{N}}$ is called Galerkin approximation if it satisfies*

$$B[\tilde{e}_{\mathbf{N}}, \mathbf{v}_{\mathbf{N}}] = f[\mathbf{v}_{\mathbf{N}}], \quad \forall \mathbf{v}_{\mathbf{N}} \in \mathcal{E}_{\mathbf{N}}$$

Remark 3.21 (Existence of the unique solution). *The same arguments (positive definiteness and boundedness) with Lax-Milgram lemma as in Remark 2.21 provides the existence of the unique solution.*

Lemma 3.22 (Convergence of Galerkin approximation). *The solution of Galerkin approximation stated in Def. 3.20 converges to the solution of weak formulation Def. 2.20.*

Proof. Using Cea's lemma with constants c_A and C_A from 2.20, we can write

$$\begin{aligned} \|e - e_{\mathbf{N}}\|_{L_{\text{per}}^2} &\leq \frac{C_A}{c_A} \inf_{\mathbf{v}_{\mathbf{N}} \in \mathcal{E}_{\mathbf{N}}} \|e - \mathbf{v}_{\mathbf{N}}\|_{L_{\text{per}}^2} \\ &\leq \frac{C_A}{c_A} \|e - P_{\mathbf{N}}[e]\|_{L_{\text{per}}^2}. \end{aligned}$$

The limit passage with Lemma 3.12 provides the converge result. \square

Remark 3.23. *The better result than in previous lemma can be obtained with Lemma 3.13 for higher regularity of conductivity coefficients and consequently the higher regularity of the solution.*

Definition 3.24 (Galerkin approximation with numerical integration). *Let $\mathbf{A} \in L_{\text{per}}^2(\mathcal{Y}; \mathbb{R}_{\text{spd}}^{d \times d}) \cap C_{\text{per}}$ be material coefficients that are additionally continuous. Then the function $\tilde{e}_{\mathbf{N}} \in \mathcal{E}_{\mathbf{N}}$ is the solution of GAwNI if it satisfies*

$$B_{\mathbf{N}}[\tilde{e}_{\mathbf{N}}, \mathbf{v}_{\mathbf{N}}] = f_{\mathbf{N}}[\mathbf{v}_{\mathbf{N}}], \quad \forall \mathbf{v}_{\mathbf{N}} \in \mathcal{E}_{\mathbf{N}} \quad (3.9)$$

where

$$B_{\mathbf{N}}[e_{\mathbf{N}}, \mathbf{v}_{\mathbf{N}}] := (Q_{\mathbf{N}}[\mathbf{A}e_{\mathbf{N}}], \mathbf{v}_{\mathbf{N}})_{L_{\text{per}}^2(\mathcal{Y}, \mathbb{R}^d)} \quad (3.10a)$$

$$f_{\mathbf{N}}[\mathbf{v}_{\mathbf{N}}] := (Q_{\mathbf{N}}[\mathbf{A}\mathbf{E}], \mathbf{v}_{\mathbf{N}})_{L_{\text{per}}^2(\mathcal{Y}, \mathbb{R}^d)}. \quad (3.10b)$$

Remark 3.25. *The regularity requirement of conductivity coefficients \mathbf{A} in previous lemma satisfies that the function $\mathbf{A}e_{\mathbf{N}}$ is continuous. Hence, the interpolation operator $Q_{\mathbf{N}}[\cdot]$ is well defined.*

Next, we define the fully discrete space, the space of function values at regular grid.

Definition 3.26 (Curl-free vectors with zero mean). *We define the space of curl-free vectors with zero mean as $\mathbb{E}_N := \mathcal{I}_N[\mathcal{E}_N] \subset \mathbb{R}^{d \times N}$ where operator \mathcal{I}_N is defined in Def. 3.8.*

With the help of isomorphism \mathcal{I}_N , space $\mathbb{R}^{d \times N}$ can be split with Helmholtz decomposition into $\mathbb{R}^{d \times N} = \mathbb{U}_N \oplus^\perp \mathbb{E}_N \oplus^\perp \mathbb{J}_N$, cf. (3.8). The following definition provides the fully discrete formulation of GAwNI and the next lemma shows its correctness.

Definition 3.27 (Fully discrete representation of Galerkin approximation with numerical integration). *Let $\mathbf{A} \in C_{\text{per}}$ be the conductivity coefficients and*

$$\mathbf{A}_N = [\delta_{km} A_{\alpha\beta}(\mathbf{x}_N^k)]_{\alpha,\beta=1,\dots,d}^{km \in \underline{\mathbb{Z}}_N^d} \in \mathbb{R}^{d \times d \times N \times N} \quad (3.11)$$

be the matrix composed of its values at regular grid. Analogically let \mathbf{E} be the prescribed macroscopic electric field and

$$\mathbf{E}_N = \mathcal{I}_N[\mathbf{E}] \in \mathbb{U}_N \quad (3.12)$$

its fully discrete counterpart. Then we define the matrix representation of Galerkin approximation with numerical integration as the following problem:

find $\tilde{\mathbf{e}}_N \in \mathbb{E}_N$ such that

$$(\mathbf{A}_N \tilde{\mathbf{e}}_N, \mathbf{v}_N)_{\mathbb{R}^{d \times N}} = -(\mathbf{A}_N \mathbf{E}_N, \mathbf{v}_N)_{\mathbb{R}^{d \times N}}, \quad \forall \mathbf{v}_N \in \mathcal{E}_N.$$

Lemma 3.28 (Fully discrete representation of Galerkin approximation with numerical integration). *Let $\mathbf{A} \in C_{\text{per}}$ be the conductivity coefficients. Then the solution of Galerkin approximation with numerical integration $\tilde{\mathbf{e}}_N \in \mathcal{E}_N$ from Def. 3.24 is equivalent to the solution of fully discrete representation of Galerkin approximation with numerical integration $\tilde{\mathbf{e}}_N \in \mathbb{E}_N$ defined in Def. 3.27 as the solutions are related with isomorphism*

$$\tilde{\mathbf{e}}_N = \mathcal{I}_N[\tilde{e}_N].$$

Proof. We substitute a function \mathbf{v}_N expressed as the Fourier series

$$\mathbf{v}_N(\mathbf{x}) = \sum_{\mathbf{k} \in \underline{\mathbb{Z}}_N^d} \mathbf{v}_N(\mathbf{x}_N^{\mathbf{k}}) \varphi_{N,\mathbf{k}}(\mathbf{x})$$

into the bilinear form to obtain

$$B_N[e_N, \mathbf{v}_N] = \sum_{\mathbf{k} \in \underline{\mathbb{Z}}_N^d} \mathbf{v}_N(\mathbf{x}_N^{\mathbf{k}}) \cdot (Q_N[\mathbf{A}e_N], \varphi_{N,\mathbf{k}})_{L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)} \quad (3.13a)$$

$$= \sum_{\mathbf{k} \in \underline{\mathbb{Z}}_N^d} \mathbf{v}^{\mathbf{k}} \cdot (Q_N[\mathbf{A}e_N], \varphi_{N,\mathbf{k}})_{L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)} \quad (3.13b)$$

where with a term $(Q_N[\mathbf{A}e_N], \varphi_{N,\mathbf{k}})_{L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)}$ we understand $\left((Q_N[\mathbf{A}e_N]_{\alpha}, \varphi_{N,\mathbf{k}})_{L^2_{\text{per}}(\mathcal{Y})} \right)_{\alpha=1}^d$, operator ‘ \cdot ’ thus denotes the scalar product on \mathbb{R}^d , and vector $\mathbf{v}^{\mathbf{k}} \in \mathbb{R}^d$ for $\mathbf{k} \in \underline{\mathbb{Z}}_N^d$ is a subvector of vector $\mathbf{v}_N := \mathcal{I}_N[\mathbf{v}_N] \in \mathbb{R}^{d \times N}$.

Next, we analogically express the term with interpolation projection $Q_N[\mathbf{A}e_N]$ through its nodal points

$$Q_N[\mathbf{A}e_N](\mathbf{x}) = \sum_{\mathbf{m} \in \underline{\mathbb{Z}}_N^d} \mathbf{A}(\mathbf{x}_N^{\mathbf{m}}) e_N(\mathbf{x}_N^{\mathbf{m}}) \varphi_{N,\mathbf{m}}(\mathbf{x}). \quad (3.14)$$

It can be observed that nodal values $\mathbf{A}(\mathbf{x}_N^{\mathbf{m}}) e_N(\mathbf{x}_N^{\mathbf{m}})$ for $\mathbf{m} \in \underline{\mathbb{Z}}_N^d$ can be expressed with matrix

$$\mathbf{A}_N = [\delta_{km} A_{\alpha\beta}(\mathbf{x}_N^k)]_{\alpha,\beta=1,\dots,d}^{km \in \underline{\mathbb{Z}}_N^d} \in \mathbb{R}^{d \times d \times N \times N}$$

and vector $\mathbf{e}_N = \mathcal{I}_N[e_N]$ as

$$\mathbf{A}(\mathbf{x}_N^m) \mathbf{e}_N(\mathbf{x}_N^m) = (\mathbf{A}_N \mathbf{e}_N)^m \in \mathbb{R}^d. \quad (3.15)$$

Substitution of Eq. (3.14) with (3.15) into bilinear form (3.13), we deduce

$$B_N[\mathbf{e}_N, \mathbf{v}_N] = \sum_{\mathbf{k} \in \underline{\mathbb{Z}}_N^d} \sum_{\mathbf{m} \in \underline{\mathbb{Z}}_N^d} (\mathbf{A}_N \mathbf{e}_N)^m \cdot \mathbf{v}_N^{\mathbf{k}}(\varphi_{N,\mathbf{m}}, \varphi_{N,\mathbf{k}})_{L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)} \quad (3.16a)$$

$$= \sum_{\mathbf{k} \in \underline{\mathbb{Z}}_N^d} \sum_{\mathbf{m} \in \underline{\mathbb{Z}}_N^d} (\mathbf{A}_N \mathbf{e}_N)^m \cdot \mathbf{v}_N^{\mathbf{k}} \frac{\delta_{\mathbf{m}\mathbf{k}}}{|\mathbf{N}|_{\Pi}} \quad (3.16b)$$

$$= \sum_{\mathbf{k} \in \underline{\mathbb{Z}}_N^d} (\mathbf{A}_N \mathbf{e}_N)^{\mathbf{k}} \cdot \mathbf{v}_N^{\mathbf{k}} \frac{1}{|\mathbf{N}|_{\Pi}} = \frac{(\mathbf{A}_N \mathbf{e}_N, \mathbf{v}_N)_{\mathbb{R}^{d \times N}}}{|\mathbf{N}|_{\Pi}}. \quad (3.16c)$$

where we have used the orthogonality property of functions $\varphi_{N,\mathbf{k}}$ stated in Remark 3.7. The formula for a linear functional

$$f_N[\mathbf{v}_N] = \frac{(\mathbf{A}_N \mathbf{E}_N, \mathbf{v}_N)_{\mathbb{R}^{d \times N}}}{|\mathbf{N}|_{\Pi}}. \quad (3.17)$$

with $\mathbf{E}_N := \mathcal{I}_N[\mathbf{E}]$ can be deduced accordingly. Using Eq. (3.16) and (3.17), the bilinear form

$$B_N[\mathbf{e}_N, \mathbf{v}_N] = f_N[\mathbf{v}_N], \quad \forall \mathbf{v}_N \in \mathcal{E}_N$$

can be transformed to

$$\begin{aligned} \frac{(\mathbf{A}_N \mathbf{e}_N, \mathbf{v}_N)_{\mathbb{R}^{d \times N}}}{|\mathbf{N}|_{\Pi}} &= \frac{(\mathbf{A}_N \mathbf{E}_N, \mathbf{v}_N)_{\mathbb{R}^{d \times N}}}{|\mathbf{N}|_{\Pi}}, \quad \forall \mathbf{v}_N \in \mathbb{E}_N \\ (\mathbf{A}_N \mathbf{e}_N, \mathbf{v}_N)_{\mathbb{R}^{d \times N}} &= (\mathbf{A}_N \mathbf{E}_N, \mathbf{v}_N)_{\mathbb{R}^{d \times N}}, \quad \forall \mathbf{v}_N \in \mathbb{E}_N \end{aligned}$$

where \mathbf{e}_N and \mathbf{e}_N are related with isomorphism $\mathbf{e}_N = \mathcal{I}_N[e_N]$. \square

Lemma 3.29 (Existence of the unique solution of GAwNI). *The approximated bilinear form (3.10a) is uniformly elliptic, i.e.*

$$\frac{1}{c_A} \|\mathbf{v}_N\|_{L^2_{\text{per}}}^2 \leq B_N[\mathbf{v}_N, \mathbf{v}_N], \quad (3.18)$$

and bounded

$$B_N[\mathbf{e}_N, \mathbf{v}_N] \leq C_A \|\mathbf{e}_N\|_{L^2_{\text{per}}} \|\mathbf{v}_N\|_{L^2_{\text{per}}}. \quad (3.19)$$

The linear function is bounded with

$$f_N[\mathbf{v}_N] \leq C_A \|\mathbf{e}_0\|_{L^2_{\text{per}}} \|\mathbf{v}_N\|_{L^2_{\text{per}}} \quad (3.20)$$

Moreover, there exists the unique solution of Galerkin approximation with numerical integration.

Proof. The formulas (3.18), (3.19), and (3.20) follow easily from expressions in Remark 3.28.

$$\begin{aligned} \frac{1}{c_A} \|\mathbf{v}_N\|_{L^2_{\text{per}}}^2 &= \frac{1}{c_A} \|\mathbf{v}\|_2^2 \leq (\mathbf{A}\mathbf{v}, \mathbf{v})_{\mathbb{R}^{d \times N}} = B_N[\mathbf{v}_N, \mathbf{v}_N] \\ B_N[\mathbf{e}_N, \mathbf{v}_N] &= (\mathbf{A}\mathbf{e}, \mathbf{v})_{\mathbb{R}^{d \times N}} \leq \|\mathbf{A}\|_{\infty} \|\mathbf{e}\|_2 \|\mathbf{v}\|_2 = C_A \|\mathbf{e}_N\|_{L^2_{\text{per}}} \|\mathbf{v}_N\|_{L^2_{\text{per}}} \\ f_N[\mathbf{v}_N] &= (\mathbf{A}\mathbf{E}, \mathbf{v})_{\mathbb{R}^{d \times N}} \leq \|\mathbf{A}\|_{\infty} \|\mathbf{E}\|_2 \|\mathbf{v}\|_2 = C_A \|\mathbf{E}\|_{L^2_{\text{per}}} \|\mathbf{v}_N\|_{L^2_{\text{per}}} \end{aligned}$$

Now, the existence and uniqueness of the solution is provided using Lax-Milgram theorem. \square

3.4 Convergence to continuous solution

This section provides convergence of discrete solutions obtained with GAwNI to the solution of weak formulation. The rate of converge is primarily based on the regularity of the solution; for completeness, the basic regularity results are also provided.

We start with a definition of difference quotients and Lemma 3.31 that characterizes the Sobolev space H_{per}^1 .

Definition 3.30 (difference quotient). *For a function $\mathbf{u} \in L_{\text{per}}^2(\mathcal{Y}; \mathbb{R}^d)$, the i -th difference quotient of size h is $D_i^h \mathbf{u}(\mathbf{x}) \in \mathbb{R}^d$ defined as*

$$\bar{D}_i^h \mathbf{u}(\mathbf{x}) := \frac{\mathbf{u}(\mathbf{x} + h\boldsymbol{\epsilon}_i) - \mathbf{u}(\mathbf{x})}{h}, \quad i = 1, \dots, d$$

for $\mathbf{x} \in \mathbb{R}^d$, $h \in \mathbb{R}$ and $\boldsymbol{\epsilon}_\alpha = (\delta_{\alpha\beta})_{\beta=1}^d$ being a vector of canonical basis of space \mathbb{R}^d . Moreover we note

$$\bar{D}^h \mathbf{u}(\mathbf{x}) := (\bar{D}_\alpha^h \mathbf{u}(\mathbf{x}))_{\alpha=1}^d.$$

Lemma 3.31 (Difference quotient and Sobolev space). *The following two statements hold:*

- Assume $1 < p < \infty$, $\mathbf{u} \in L_{\text{per}}^p(\mathcal{Y}; \mathbb{R}^d)$, and there exists constant C independent of h such that

$$\|\bar{D}^h \mathbf{u}\|_{L_{\text{per}}^p} \leq C.$$

Then

$$\mathbf{u} \in W_{\text{per}}^{1,p}(\mathcal{Y}; \mathbb{R}^d)$$

with $\|D^1 \mathbf{u}\|_{L_{\text{per}}^p} \leq C$.

- Assume $1 < p < \infty$ and $\mathbf{u} \in W_{\text{per}}^{1,p}(\mathcal{Y}; \mathbb{R}^d)$. Then

$$\|\bar{D}^h \mathbf{u}\|_{L_{\text{per}}^p} \leq \|D^1 \mathbf{u}\|_{L_{\text{per}}^p}$$

for $h \in \mathbb{R}$.

Proof. For proof see Theorem 3 in Section 5.8.2 in [10]. □

Lemma 3.32 (Regularity result). *Let $\mathbf{A} \in W_{\text{per}}^{1,\infty}(\mathcal{Y}; \mathbb{R}^{d \times d})$ and $\mathbf{e} \in L_{\text{per}}^2(\mathcal{Y}; \mathbb{R}^d)$ be a solution of weak formulation, c.f. Def. 2.20. Then $\mathbf{e} \in H_{\text{per}}^1(\mathcal{Y}; \mathbb{R}^d)$ and it satisfies*

$$\begin{aligned} \|\mathbf{e}\|_{H_{\text{per}}^1} &\leq \frac{1}{c_A} \left(\|\mathbf{A}\mathbf{E}\|_{H_{\text{per}}^1} + \|\mathbf{A}\|_{W_{\text{per}}^{1,\infty}} \|\mathbf{e}\|_{L_{\text{per}}^2} \right) \\ &\leq \frac{1}{c_A} \|\mathbf{A}\|_{W_{\text{per}}^{1,\infty}} \left(\|\mathbf{E}\|_{L_{\text{per}}^2} + \|\mathbf{e}\|_{L_{\text{per}}^2} \right) \\ &\leq \frac{1}{c_A} \|\mathbf{A}\|_{W_{\text{per}}^{1,\infty}} \|\mathbf{E}\|_{L_{\text{per}}^2} \left(1 + \frac{c_A}{c_A} \right). \end{aligned}$$

Proof. The proof based on Lem. 3.31 can be found for a more general case in Section 6.3.1 in [10] as Theorem 1. In order to work with difference quotients, we use simple formulas

$$(\mathbf{v}, \bar{D}_k^{-h} \mathbf{w})_{L_{\text{per}}^2} = -(\bar{D}_k^h \mathbf{v}, \mathbf{w})_{L_{\text{per}}^2} \quad (3.21)$$

$$\bar{D}_k^h (\mathbf{v}\mathbf{w}) = \mathbf{v}^h \bar{D}_k^h \mathbf{w} + \mathbf{w} \bar{D}_k^h \mathbf{v} \quad (3.22)$$

where $\mathbf{v}^h(\mathbf{x}) = \mathbf{v}(\mathbf{x} + h\boldsymbol{\epsilon}_k)$. Then we start with weak formulation, c.f. Def. 2.20, i.e.

$$(\mathbf{A}\mathbf{e}, \mathbf{v})_{L_{\text{per}}^2(\mathcal{Y}; \mathbb{R}^d)} = -(\mathbf{f}, \mathbf{v})_{L_{\text{per}}^2(\mathcal{Y}; \mathbb{R}^d)}$$

where $\mathbf{f} := \mathbf{A}\mathbf{E} \in H_{\text{per}}^1(\mathcal{Y}; \mathbb{R}^d)$, with a special choice of test function $\mathbf{v} = -\bar{D}_k^{-h}(\bar{D}_k^h \mathbf{e})$ leading to

$$\begin{aligned} -(\mathbf{A}\mathbf{e}, -\bar{D}_k^{-h}(\bar{D}_k^h \mathbf{e}))_{L_{\text{per}}^2} &= (\mathbf{f}, \bar{D}_k^{-h}(\bar{D}_k^h \mathbf{e}))_{L_{\text{per}}^2} \\ (\bar{D}_k^h \mathbf{A}\mathbf{e}, \bar{D}_k^h \mathbf{e})_{L_{\text{per}}^2} &= -(\bar{D}_k^h \mathbf{f}, \bar{D}_k^h \mathbf{e})_{L_{\text{per}}^2} \\ (\mathbf{A}^h(\mathbf{x})\bar{D}_k^h \mathbf{e}, \bar{D}_k^h \mathbf{e})_{L_{\text{per}}^2} &= -(\bar{D}_k^h \mathbf{f}, \bar{D}_k^h \mathbf{e})_{L_{\text{per}}^2} - (\bar{D}_k^h \mathbf{A}\mathbf{e}, \bar{D}_k^h \mathbf{e})_{L_{\text{per}}^2}. \end{aligned}$$

The positive definiteness and boundedness of \mathbf{A} implies

$$\begin{aligned} c_A \|\bar{D}_k^h \mathbf{e}\|_{L_{\text{per}}^2}^2 &\leq \|\bar{D}_k^h \mathbf{f}\|_{L_{\text{per}}^2} \|\bar{D}_k^h \mathbf{e}\|_{L_{\text{per}}^2} + \|\bar{D}_k^h \mathbf{A}\|_{L_{\text{per}}^\infty} \|\mathbf{e}\|_{L_{\text{per}}^2} \|\bar{D}_k^h \mathbf{e}\|_{L_{\text{per}}^2} \\ \|\bar{D}_k^h \mathbf{e}\|_{L_{\text{per}}^2} &\leq \frac{1}{c_A} \|\bar{D}_k^h \mathbf{f}\|_{L_{\text{per}}^2} + \|\bar{D}_k^h \mathbf{A}\|_{L_{\text{per}}^\infty} \|\mathbf{e}\|_{L_{\text{per}}^2} \\ &\leq \frac{1}{c_A} \|\mathbf{f}\|_{H_{\text{per}}^1} + \|\mathbf{A}\|_{W_{\text{per}}^{1,\infty}} \|\mathbf{e}\|_{L_{\text{per}}^2} \\ &\leq \frac{1}{c_A} \|\mathbf{A}\|_{W_{\text{per}}^{1,\infty}} \left(\|\mathbf{E}\|_{L_{\text{per}}^2} + \|\mathbf{e}\|_{L_{\text{per}}^2} \right). \end{aligned}$$

The solution of weak formulation \mathbf{e} is bounded in L_{per}^2 , i.e.

$$c_A \|\mathbf{e}\|_{L_{\text{per}}^2}^2 < B[\mathbf{e}, \mathbf{e}] = f[\mathbf{e}] < C_A \|\mathbf{E}\|_{L_{\text{per}}^2} \|\mathbf{e}\|_{L_{\text{per}}^2},$$

the difference quotient $\|\bar{D}_k^h L\|_{L_{\text{per}}^\infty}$ is bounded due to Lem. 3.31 and thus the required Eq. (3.21) follows. \square

Lemma 3.33 (Higher regularity result). *Let $\mathbf{A} \in W_{\text{per}}^{m,\infty}(\mathcal{Y}; \mathbb{R}^{d \times d})$ for $m \in \mathbb{N}$ and $\mathbf{e} \in L_{\text{per}}^2(\mathcal{Y}; \mathbb{R}^d)$ be a solution of weak formulation, Def. 2.20. Then*

$$\mathbf{e} \in H_{\text{per}}^m(\mathcal{Y}; \mathbb{R}^d) \tag{3.23}$$

and it satisfies

$$\|\mathbf{e}\|_{H_{\text{per}}^m} \leq C(\mathbf{E}) \|\mathbf{A}\|_{W_{\text{per}}^{m,\infty}} \|\mathbf{E}\|_{L_{\text{per}}^2} \tag{3.24}$$

where $C(\mathbf{E})$ depends only on \mathbf{E} .

Proof. The proof based on Lem. 3.31 can be found for a more general case in Section 6.3.1 in [10] as Theorem 2.

We will establish (3.23) and estimate (3.24) by induction hypothesis.

Base case: The case $m = 1$ is proven in Lem. 3.32.

Induction hypothesis: Assume now (3.23) and estimate (3.24) be valid for $\mathbf{A} \in W_{\text{per}}^{m,\infty}$.

Induction step: Suppose additionally $\mathbf{A} \in W_{\text{per}}^{(m+1),\infty}$ to prove $\mathbf{e} \in H_{\text{per}}^{m+1}(\mathcal{Y}; \mathbb{R}^d)$ with estimate

$$\|\mathbf{e}\|_{H_{\text{per}}^{m+1}} < C \|\mathbf{A}\|_{W_{\text{per}}^{m+1,\infty}} \|\mathbf{E}\|_{L_{\text{per}}^2}.$$

Let \mathbf{l} be any multiindex with $\|\mathbf{l}\|_1 = m$ and take any test function $\tilde{\mathbf{v}} \in \mathcal{D}_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)$ and put $\mathbf{v} := (-1)^{\|\mathbf{l}\|_1} D^{\mathbf{l}} \tilde{\mathbf{v}}$ into the formula of weak formulation

$$\begin{aligned} B[\mathbf{e}, \mathbf{v}] &= f[\mathbf{v}] \\ (\mathbf{A}\mathbf{e}, \mathbf{v})_{L_{\text{per}}^2} &= (\mathbf{A}\mathbf{E}, \mathbf{v})_{L_{\text{per}}^2} \end{aligned}$$

Integration by parts leads to

$$B[\tilde{\mathbf{e}}, \tilde{\mathbf{v}}] = \tilde{f}[\tilde{\mathbf{v}}] \tag{3.25}$$

where

$$\begin{aligned}\tilde{\mathbf{e}}(\mathbf{x}) &= D^{\mathbf{l}} \mathbf{e}(\mathbf{x}) \in L^2_{\text{per}} \\ \tilde{f}[\tilde{\mathbf{v}}] &= (\tilde{f}, \tilde{\mathbf{v}})_{L^2_{\text{per}}} \\ \tilde{f}(\mathbf{x}) &= D^{\mathbf{l}} \mathbf{A}(\mathbf{x}) \mathbf{E} - \sum_{\mathbf{k}: \mathbf{k} \leq \mathbf{l}, \mathbf{l} \neq \mathbf{k}} \binom{\mathbf{l}}{\mathbf{k}} D^{\mathbf{l}-\mathbf{k}} \mathbf{A}(\mathbf{x}) D^{\mathbf{k}} \mathbf{e}(\mathbf{x}).\end{aligned}$$

and $\mathbf{x} \in \mathbb{R}^d$. Since the test function $\tilde{\mathbf{v}}$ is arbitrary, $\tilde{\mathbf{e}}(\mathbf{x})$ is a weak solution of

$$B[\tilde{\mathbf{e}}, \tilde{\mathbf{v}}] = \tilde{f}[\tilde{\mathbf{v}}], \quad \forall \tilde{\mathbf{v}} \in L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d),$$

Next, we can observe that $\tilde{\mathbf{f}} \in H^1_{\text{per}}$ as

$$\begin{aligned}\|\tilde{\mathbf{f}}\|_{H^1_{\text{per}}} &\leq \|D^{\mathbf{l}} \mathbf{A}\|_{W^{1,\infty}_{\text{per}}} \|\mathbf{E}\|_{L^2_{\text{per}}} + \sum_{\mathbf{k}: \mathbf{k} \leq \mathbf{l}, \mathbf{l} \neq \mathbf{k}} \binom{\mathbf{l}}{\mathbf{k}} \|D^{\mathbf{l}-\mathbf{k}} \mathbf{A}\|_{W^{1,\infty}_{\text{per}}} \|D^{\mathbf{k}} \mathbf{e}\|_{H^1_{\text{per}}} \\ &\leq \|\mathbf{A}\|_{W^{m+1,\infty}_{\text{per}}} \|\mathbf{E}\|_{L^2_{\text{per}}} + C \|\mathbf{A}\|_{W^{m+1,\infty}_{\text{per}}} \|\mathbf{e}\|_{H^m_{\text{per}}} \\ &\leq C \|\mathbf{A}\|_{W^{m+1,\infty}_{\text{per}}} \left(\|\mathbf{E}\|_{L^2_{\text{per}}} + \|\mathbf{e}\|_{H^m_{\text{per}}} \right) \\ &\leq C \|\mathbf{A}(\mathbf{x})\|_{W^{m+1,\infty}_{\text{per}}} \|\mathbf{E}\|_{L^2_{\text{per}}}\end{aligned}$$

Using Lem. 3.32 on problem stated in (3.25), we estimate

$$\begin{aligned}\|\tilde{\mathbf{e}}\|_{H^1_{\text{per}}} &\leq C \left(\|\tilde{\mathbf{f}}\|_{H^1_{\text{per}}} + \|\mathbf{A}\|_{W^{1,\infty}_{\text{per}}} \|\tilde{\mathbf{e}}\|_{L^2_{\text{per}}} \right) \\ &\leq C \left(\|\mathbf{A}\|_{W^{m+1,\infty}_{\text{per}}} \|\mathbf{E}\|_{L^2_{\text{per}}} + \|\mathbf{A}\|_{W^{1,\infty}_{\text{per}}} \|\mathbf{E}\|_{L^2_{\text{per}}} \right) \\ &\leq C \|\mathbf{A}\|_{W^{m+1,\infty}_{\text{per}}} \|\mathbf{E}\|_{L^2_{\text{per}}}.\end{aligned}$$

This inequality holds for each multiindex \mathbf{l} such that $\|\mathbf{l}\|_1 = m$ and $\tilde{\mathbf{e}}(\mathbf{x}) = D^{\mathbf{l}} \mathbf{e}(\mathbf{x})$, thus we conclude $\mathbf{e} \in H^{m+1}_{\text{per}}$ with

$$\|\mathbf{e}\|_{H^{m+1}_{\text{per}}} \leq C \|\mathbf{A}\|_{W^{m+1,\infty}_{\text{per}}} \|\mathbf{E}\|_{L^2_{\text{per}}}.$$

□

Lemma 3.34 (Strang). *Consider a family of discrete problems whose associated approximate bilinear forms are uniformly elliptic. Then there exists a constant C independent of the space $\mathcal{T}_{\mathbf{N}}$ such that*

$$\begin{aligned}\|\mathbf{e} - \mathbf{e}_{\mathbf{N}}\|_{L^2_{\text{per}}} &\leq \inf_{\mathbf{v}_{\mathbf{N}} \in \mathcal{E}_{\mathbf{N}}} \left\{ \left(\frac{C_A}{c_A} + 1 \right) \|\mathbf{e} - \mathbf{v}_{\mathbf{N}}\|_{L^2_{\text{per}}} + \frac{1}{c_A} \sup_{\mathbf{w}_{\mathbf{N}} \in \mathcal{E}_{\mathbf{N}}} \frac{|B[\mathbf{v}_{\mathbf{N}}, \mathbf{w}_{\mathbf{N}}] - B_{\mathbf{N}}[\mathbf{v}_{\mathbf{N}}, \mathbf{w}_{\mathbf{N}}]|}{\|\mathbf{w}_{\mathbf{N}}\|_{L^2_{\text{per}}}} \right\} \\ &\quad + \frac{1}{c_A} \sup_{\mathbf{w}_{\mathbf{N}} \in \mathcal{E}_{\mathbf{N}}} \frac{|f[\mathbf{w}_{\mathbf{N}}] - f_{\mathbf{N}}[\mathbf{w}_{\mathbf{N}}]|}{\|\mathbf{w}_{\mathbf{N}}\|_{L^2_{\text{per}}}}\end{aligned}$$

where the constants are from Def. 2.20.

Proof. The proof can be found as Theorem 26.1 in [5].

□

Lemma 3.35 (Convergence of discrete solutions to continuous one). *Let the conductivity coefficients be $\mathbf{A} \in W^{\mu,\infty}_{\text{per}}(\mathcal{Y}; \mathbb{R}^{d \times d})$ with $\mu \in \mathbb{N}$ such that $W^{\mu,\infty}_{\text{per}}$ is embedded into continuous functions. Then the sequence of solutions of Galerkin approximation with numerical integration, c.f. Def. 3.24, converge to the solution of weak formulation, Def. 2.20, in $\|\cdot\|_{L^2_{\text{per}}}$ norm, i.e.*

$$\|\mathbf{e} - \mathbf{e}_{\mathbf{N}}\|_{L^2_{\text{per}}} \leq C \left(\min_{\alpha} \frac{N_{\alpha}}{2Y_{\alpha}} \right)^{-\mu} \rightarrow 0 \quad \text{for } \min_{\alpha} \frac{N_{\alpha}}{2Y_{\alpha}} \rightarrow \infty$$

where

$$C = \left(\frac{C_A}{c_A} + 1 \right) \|e\|_{H_{\text{per}}^\mu} + \frac{c_{0,\mu}}{c_A} \|\mathbf{A}\|_{W_{\text{per}}^{\mu,\infty}} \|e\|_{H_{\text{per}}^\mu} + \frac{c_{0,\mu}}{c_A} \|\mathbf{A}\|_{W_{\text{per}}^{\mu,\infty}} \|\mathbf{E}\|_{L_{\text{per}}^2}.$$

Proof. We use Lem. 3.34 due to Strang, results about regularization Lem. 3.32 and 3.33, and approximation by orthogonal projection Lem. 3.13 and by interpolation projection Lem. 3.18.

$$\begin{aligned} \|e - e_N\|_{L_{\text{per}}^2} &\leq \inf_{\mathbf{v}_N \in \mathcal{E}_N} \left\{ \left(\frac{C_A}{c_A} + 1 \right) \|e - \mathbf{v}_N\|_{L_{\text{per}}^2} + \frac{1}{c_A} \sup_{\mathbf{w}_N \in \mathcal{E}_N} \frac{|(\mathbf{A}\mathbf{v}_N, \mathbf{w}_N)_{L_{\text{per}}^2} - (Q_N[\mathbf{A}\mathbf{v}_N], \mathbf{w}_N)_{L_{\text{per}}^2}|}{\|\mathbf{w}_N\|_{L_{\text{per}}^2}} \right\} \\ &\quad + \frac{1}{c_A} \sup_{\mathbf{w}_N \in \mathcal{E}_N} \frac{|(\mathbf{A}\mathbf{E}, \mathbf{w}_N)_{L_{\text{per}}^2} - (Q_N[\mathbf{A}\mathbf{E}], \mathbf{w}_N)_{L_{\text{per}}^2}|}{\|\mathbf{w}_N\|_{L_{\text{per}}^2}} \\ &\leq \inf_{\mathbf{v}_N \in \mathcal{E}_N} \left\{ \left(\frac{C_A}{c_A} + 1 \right) \|e - \mathbf{v}_N\|_{L_{\text{per}}^2} + \frac{1}{c_A} \|\mathbf{A}\mathbf{v}_N - Q_N[\mathbf{A}\mathbf{v}_N]\|_{L_{\text{per}}^2} \right\} + \frac{1}{c_A} \|\mathbf{A}\mathbf{E} - Q_N[\mathbf{A}\mathbf{E}]\|_{L_{\text{per}}^2} \\ &\leq \left(\frac{C_A}{c_A} + 1 \right) \|e - P_N e\|_{L_{\text{per}}^2} + \frac{1}{c_A} \|\mathbf{A}P_N e - Q_N[\mathbf{A}P_N e]\|_{L_{\text{per}}^2} + \frac{1}{c_A} \|\mathbf{A} - Q_N[\mathbf{A}]\|_{L_{\text{per}}^2} \|\mathbf{E}\|_{L^\infty} \\ &\leq \left(\frac{C_A}{c_A} + 1 \right) \left(\min_{\alpha} \frac{N_\alpha}{2Y_\alpha} \right)^{-\mu} \|e\|_{H_{\text{per}}^\mu} + \frac{c_{0,\mu}}{c_A} \left(\min_{\alpha} \frac{N_\alpha}{2Y_\alpha} \right)^{-\mu} \|\mathbf{A}\|_{W_{\text{per}}^{\mu,\infty}} \|e\|_{H_{\text{per}}^\mu} \\ &\quad + \frac{c_{0,\mu}}{c_A} \left(\min_{\alpha} \frac{N_\alpha}{2Y_\alpha} \right)^{-\mu} \|\mathbf{A}\|_{W_{\text{per}}^{\mu,\infty}} \|\mathbf{E}\|_{L_{\text{per}}^2} \\ &\leq \left(\min_{\alpha} \frac{N_\alpha}{2Y_\alpha} \right)^{-\mu} \left[\left(\frac{C_A}{c_A} + 1 \right) \|e\|_{H_{\text{per}}^\mu} + \frac{c_{0,\mu}}{c_A} \|\mathbf{A}\|_{W_{\text{per}}^{\mu,\infty}} \|e\|_{H_{\text{per}}^\mu} + \frac{c_{0,\mu}}{c_A} \|\mathbf{A}\|_{W_{\text{per}}^{\mu,\infty}} \|\mathbf{E}\|_{L_{\text{per}}^2} \right] \end{aligned}$$

Now, the limit passage reveals the proof. \square

3.5 Regularization of rough material coefficients

This section provide a method for calculation of discrete solutions if the material coefficients \mathbf{A} are rough, do not possess sufficient regularity. The method is based on regularization of material coefficient. Although the convergence is provided, the convergence can be arbitrarily slow — we cannot expect any order of convergence.

Lemma 3.36. *Let $\mathbf{A} \in L_{\text{per}}^\infty(\mathcal{Y}; \mathbb{R}^{d \times d})$ be a symmetric and uniformly elliptic material coefficients and \mathbf{A}_h be their regularization such that*

$$\begin{aligned} \mathbf{A}_h &\in L_{\text{per}}^\infty(\mathcal{Y}; \mathbb{R}^{d \times d}) \\ \|\mathbf{A} - \mathbf{A}_h\|_{L_{\text{per}}^2(\mathcal{Y}; \mathbb{R}^{d \times d})} &\rightarrow 0 \quad \text{for } h \rightarrow 0 \end{aligned}$$

and be positive definite and bounded with constants $c_{A,h}$ and $C_{A,h}$, i.e.

$$c_{A,h} \leq \frac{\sum_{\alpha,\beta=1}^d (A_h)_{\alpha\beta}(\mathbf{x}) \xi_\alpha \xi_\beta}{\|\xi\|^2} \leq C_{A,h}, \quad \forall \xi \in \mathbb{R}^d, \quad \text{and for a.a. } \mathbf{x} \in \mathcal{Y}.$$

Suppose $e, e_h \in L_{\text{per}}^2(\mathcal{Y}; \mathbb{R}^d)$ be solutions of weak formulation, c.f. Def. 2.20, for material coefficients \mathbf{A} and \mathbf{A}_h resp. Then

$$\lim_{h \rightarrow 0} \|e - e_h\|_{L_{\text{per}}^2(\mathcal{Y}; \mathbb{R}^d)} = 0 \quad (3.26)$$

Proof. First, define a bilinear form and linear functional for regularized coefficients \mathbf{A}_h as

$$B_h[e_h, \mathbf{v}] = (\mathbf{A}_h e_h, \mathbf{v})_{L_{\text{per}}^2(\mathcal{Y}; \mathbb{R}^d)} \quad f_h[\mathbf{v}] = (\mathbf{A}_h \mathbf{E}, \mathbf{v})_{L_{\text{per}}^2(\mathcal{Y}; \mathbb{R}^d)}$$

and a weak formulation

$$B_h[e_h, \mathbf{v}] = f_h[\mathbf{v}], \quad \forall \mathbf{v} \in \mathcal{E}$$

where \mathbf{e}_h denotes its solution.

We show from positive definiteness and boundedness of \mathbf{A} that solutions of regularized problems \mathbf{e}_h are uniformly bounded

$$\|\mathbf{e}_h\|_{L^2}^2 \leq \frac{1}{c_{A,h}} (\mathbf{A}_h \mathbf{e}_h, \mathbf{e}_h) = -\frac{1}{c_{A,h}} (\mathbf{A}_h \mathbf{E}, \mathbf{e}_h) \leq \frac{C_{A,h}}{c_{A,h}} \|\mathbf{E}\|_{L^2} \|\mathbf{e}_h\|_{L^2}.$$

Thus using Banach selection principle for sequence $\{\mathbf{e}_h\}$ under a reflexive and separable space such as L^2_{per} , we conclude there exists some subsequence converging weakly to some function $\bar{\mathbf{e}}$, i.e.

$$\mathbf{e}_h \rightharpoonup \bar{\mathbf{e}} \Leftrightarrow \lim_{h \rightarrow 0} (\mathbf{e}_h, \mathbf{v})_{L^2_{\text{per}}} = (\bar{\mathbf{e}}, \mathbf{v})_{L^2_{\text{per}}}, \quad \forall \mathbf{v} \in L^2_{\text{per}}.$$

Next, we show that the function $\bar{\mathbf{e}}$ satisfies the weak formulation with coefficients \mathbf{A} , hence for arbitrary $\mathbf{v} \in \mathcal{E}$

$$(\mathbf{A}_h \mathbf{e}_h, \mathbf{v})_{L^2_{\text{per}}} = (\mathbf{A}_h \mathbf{E}, \mathbf{v})_{L^2_{\text{per}}} \xrightarrow{h \rightarrow 0} (\mathbf{A} \mathbf{E}, \mathbf{v})_{L^2_{\text{per}}} = (\mathbf{A} \bar{\mathbf{e}}, \mathbf{v})_{L^2_{\text{per}}}$$

and since $(\mathbf{A}_h \mathbf{e}_h, \mathbf{v})_{L^2_{\text{per}}} \xrightarrow{h \rightarrow 0} (\mathbf{A} \bar{\mathbf{e}}, \mathbf{v})_{L^2_{\text{per}}}$, we can conclude that $\bar{\mathbf{e}} = \mathbf{e}$.

Next we show that the subsequence \mathbf{e}_h converge even strongly

$$\begin{aligned} c_{A,h} \|\mathbf{e} - \mathbf{e}_h\|_{L^2_{\text{per}}}^2 &\leq (\mathbf{A}_h (\mathbf{e} - \mathbf{e}_h), \mathbf{e} - \mathbf{e}_h)_{L^2_{\text{per}}} \\ &= (\mathbf{A}_h \mathbf{e}, \mathbf{e} - \mathbf{e}_h)_{L^2_{\text{per}}} - (\mathbf{A}_h \mathbf{e}_h, \mathbf{e} - \mathbf{e}_h)_{L^2_{\text{per}}} \pm (\mathbf{A} \mathbf{e}, \mathbf{e} - \mathbf{e}_h)_{L^2_{\text{per}}} \\ &= (\mathbf{A} \mathbf{e}, \mathbf{e} - \mathbf{e}_h)_{L^2_{\text{per}}} - (\mathbf{A}_h \mathbf{e}_h, \mathbf{e} - \mathbf{e}_h)_{L^2_{\text{per}}} + (\mathbf{A}_h \mathbf{e}, \mathbf{e} - \mathbf{e}_h)_{L^2_{\text{per}}} - (\mathbf{A} \mathbf{e}, \mathbf{e} - \mathbf{e}_h)_{L^2_{\text{per}}} \\ &= -((\mathbf{A} - \mathbf{A}_h) \mathbf{E}, \mathbf{e} - \mathbf{e}_h)_{L^2_{\text{per}}} + (\mathbf{A}_h \mathbf{e}, \mathbf{e} - \mathbf{e}_h)_{L^2_{\text{per}}} - (\mathbf{A} \mathbf{e}, \mathbf{e} - \mathbf{e}_h)_{L^2_{\text{per}}} \\ &\leq \|(\mathbf{A} - \mathbf{A}_h)\|_{L^\infty_{\text{per}}} (\mathbf{E}, \mathbf{e} - \mathbf{e}_h)_{L^2_{\text{per}}} + \|\mathbf{A}_h - \mathbf{A}\|_{L^\infty_{\text{per}}} (\mathbf{e}, \mathbf{e} - \mathbf{e}_h) \end{aligned}$$

concluding that there is a subsequence converging strongly $\mathbf{e}_h \rightarrow \mathbf{e}$.

Finally, since every convergent subsequence satisfies the same weak formulation with an unique solution, all limit points have to equal each other. Thus not only the subsequence but the whole sequence converge to the unique solution of weak formulation. \square

Lemma 3.37. *Let the assumptions from previous Lem. 3.36 be satisfied. Moreover let $\mathbf{e}_{h,\mathbf{N}}$ be solution of Galerkin approximation with numerical integration for regularized conductivity coefficients \mathbf{A}_h , c.f. Def. 3.24. Then for arbitrary $\varepsilon > 0$ there exists $h \in \mathbb{R}$ and $\mathbf{N} \in \mathbb{N}^d$ (not the other way around) such that*

$$\|\mathbf{e} - \mathbf{e}_{h,\mathbf{N}}\|_{L^2_{\text{per}}} \leq \varepsilon$$

Proof. From the previous Lem. 3.36, we can choose fixed $h > 0$ sufficiently small such that $\|\mathbf{e} - \mathbf{e}_h\| < \frac{\varepsilon}{2}$. Next, we can choose \mathbf{N} such that

$$\min_{\alpha} \frac{N_{\alpha}}{2Y_{\alpha}} \geq \left(\frac{2C}{\varepsilon} \right)^{\frac{1}{\mu}}$$

with constant C from Lem. 3.35. Then Lem. 3.35 implies

$$\|\mathbf{e}_h - \mathbf{e}_{h,\mathbf{N}}\|_{L^2_{\text{per}}} \leq \frac{\varepsilon}{2}$$

and the triangle inequality finishes the proof

$$\|\mathbf{e} - \mathbf{e}_{h,\mathbf{N}}\|_{L^2_{\text{per}}} \leq \|\mathbf{e} - \mathbf{e}_h\|_{L^2_{\text{per}}} + \|\mathbf{e}_h - \mathbf{e}_{h,\mathbf{N}}\|_{L^2_{\text{per}}} \leq \varepsilon.$$

\square

Remark 3.38 (Example of regularization of conductivity coefficients). *The example of regularized data \mathbf{A}_h satisfying conditions in Lem. 3.36 is a standard mollification, i.e.*

$$\mathbf{A}_h(\mathbf{x}) = \int_{\mathbb{R}^d} \eta_h(\mathbf{x} - \mathbf{y}) \mathbf{A}(\mathbf{y}) \, d\mathbf{y}$$

where $\eta_h(\mathbf{x}) = \frac{1}{h^d} \eta\left(\frac{\mathbf{x}}{h}\right)$,

$$\eta(\mathbf{x}) = \begin{cases} C \exp\left(\frac{1}{|\mathbf{x}|^2 - 1}\right), & \text{for } |\mathbf{x}| < 1, \\ 0, & \text{otherwise} \end{cases}$$

with constant C chosen so that $\int_{\mathbb{R}^d} \eta(\mathbf{x}) \, d\mathbf{x} = 1$.

3.6 Solution of Linear systems using Conjugate gradients

In this section, we describe the solution of Galerkin approximation with numerical integration by Conjugate gradients. The next lemma shows that the solution by Conjugate gradients corresponds to minimization of quadratic functional over subspace $\mathbb{E}_N \subset \mathbb{R}^{d \times N}$. The further definition and lemma provide an projection on \mathbb{E}_N . This satisfies Conjugate gradients are working in appropriate space. Moreover in [21], the relation of fully discrete formulation of GAwNI to Moulinec-Suquet algorithm [14] is provided.

Lemma 3.39. *The Galerkin approximation with numerical integration, defined in Def. 3.24, is solvable by Conjugate gradients.*

Proof. In Lemma 3.28, we have shown that Galerkin approximation with numerical integration can be characterized with the problem: find $\tilde{\mathbf{e}}_N \in \mathbb{E}_N$ such that

$$(\mathbf{A}_N \tilde{\mathbf{e}}_N, \mathbf{v}_N)_{\mathbb{R}^{d \times N}} = -(\mathbf{A}_N \mathbf{E}_N, \mathbf{v}_N)_{\mathbb{R}^{d \times N}}, \quad \forall \mathbf{v}_N \in \mathbb{E}_N. \quad (3.27)$$

where

$$\mathbf{A}_N = [\delta_{\mathbf{k}\mathbf{m}} A_{\alpha\beta}(\mathbf{x}_N^{\mathbf{k}})]_{\alpha,\beta=1,\dots,d}^{\mathbf{k},\mathbf{m} \in \mathbb{Z}_N^d} \in \mathbb{R}^{d \times d \times N \times N} \quad \mathbf{E}_N = \mathcal{I}_N[\mathbf{E}] \in \mathbb{R}^{d \times N}$$

being dependent on conductivity \mathbf{A} and macroscopic electric field \mathbf{E} . Since \mathbf{A} is symmetric and positive definite, the matrix \mathbf{A}_N is. Then the solution $\tilde{\mathbf{e}}_N$ of Eq. (3.27) can be expressed as

$$\tilde{\mathbf{e}}_N = \operatorname{argmin}_{\mathbf{v}_N \in \mathbb{E}_N} \frac{1}{2} (\mathbf{A}_N \mathbf{v}_N, \mathbf{v}_N)_{\mathbb{R}^{d \times N}} + (\mathbf{A}_N \mathbf{E}_N, \mathbf{v}_N)_{\mathbb{R}^{d \times N}} \quad (3.28)$$

and thus it is solvable by Conjugate gradients. \square

However, the previous lemma does not show how to effectively minimize the quadratic functional, Eq. (3.28), on a subspace $\mathbb{E}_N \subset \mathbb{R}^{d \times N}$.

Definition 3.40. *We define a matrix composed of integral kernel in Fourier space, c.f. Def. 2.23 and especially Eq. (2.16), as*

$$\hat{\mathbf{G}} = \left[\delta_{\mathbf{k}\mathbf{m}} \frac{\xi_\alpha(\mathbf{k}) \xi_\beta(\mathbf{k})}{(\boldsymbol{\xi}(\mathbf{k}), \boldsymbol{\xi}(\mathbf{k}))_{\mathbb{R}^d}} \right]_{\alpha,\beta=1,\dots,d}^{\mathbf{k},\mathbf{m} \in \mathbb{Z}_N^d} \in \mathbb{R}^{d \times d \times N \times N}$$

where $\boldsymbol{\xi}(\mathbf{k}) = \left(\frac{k_\alpha}{Y_\alpha}\right)_{\alpha=1,\dots,d} \in \mathbb{R}^d$. Next, we define matrices composed of (inverse) Discrete Fourier Transform

$$\mathbf{F} = \left[\frac{\delta_{\alpha\beta} \omega_N^{-\mathbf{k}\mathbf{m}}}{|\mathbf{N}|_\Pi} \right]_{\alpha,\beta=1,\dots,d}^{\mathbf{k},\mathbf{m} \in \mathbb{Z}_N^d} \in \mathbb{C}^{d \times d \times N \times N} \quad \mathbf{F}^{-1} = [\delta_{\alpha\beta} \omega_N^{\mathbf{k}\mathbf{m}}]_{\alpha,\beta=1,\dots,d}^{\mathbf{k},\mathbf{m} \in \mathbb{Z}_N^d} \in \mathbb{C}^{d \times d \times N \times N}.$$

where

$$\omega_N^{km} = \exp\left(2\pi i \sum_{\alpha=1}^d \frac{k_\alpha m_\alpha}{N_\alpha}\right), \quad \mathbf{m}, \mathbf{k} \in \mathbb{Z}_N^d.$$

Finally, we define matrix

$$\mathbf{G} = \mathbf{F}^{-1} \hat{\mathbf{G}} \mathbf{F} \in \mathbb{C}^{d \times d \times N \times N} \quad (3.29)$$

Remark 3.41. Submatrix $\mathbf{F}_{\alpha\alpha}$ is exactly the matrix of Discrete Fourier Transform. Thus the linear operator $\mathbf{F} : \mathbf{v} \in \mathbb{R}^{d \times N} \rightarrow \mathbf{Fv} \in \mathbb{C}^{d \times N}$ can be treated as the d -dimensional fast Fourier Transform routine applied on each submatrix \mathbf{v}_α for $\alpha = 1, \dots, d$.

Lemma 3.42. The matrix \mathbf{G} from Def. 3.40 defining a linear operator $\mathbf{G} : \mathbb{R}^{d \times N} \rightarrow \mathbb{C}^{d \times N}$ is a projection on \mathbb{E}_N .

Proof. The proof is a consequence of the fact that matrix \mathbf{G} is defined via continuous projection \mathcal{G} Def. 2.25 and that operator \mathcal{I}_N from Def. 3.8 is isomorphism. \square

Acknowledgments This work was supported by the Czech Science Foundation, through projects No. GAČR 103/09/1748, No. GAČR 103/09/P490, No. GAČR P105/12/0331 and by the Grant Agency of the Czech Technical University in Prague through project No. SGS10/124/OHK1/2T/11 and No. SGS 12/027/OHK1/1T/11.

References

- [1] N. S. Bakhvalov and A. V. Knyazev, *Efficient computation of averaged characteristics of composites of a periodic structure of essentially different materials*, Soviet mathematics - Doklady **42** (1991), no. 1, 57–62.
- [2] S. Brisard and L. Dormieux, *FFT-based methods for the mechanics of composites: A general variational framework*, Computational Materials Science **49** (2010), no. 3, 663–671.
- [3] ———, *Combining Galerkin approximation techniques with the principle of Hashin and Shtrikman to derive a new FFT-based numerical method for the homogenization of composites*, Computer Methods in Applied Mechanics and Engineering (2012).
- [4] J. Y. Buffière, P. Cloetens, W. Ludwig, E. Maire, and L. Salvo, *In situ X-ray tomography studies of microstructural evolution combined with 3D modeling*, MRS Bulletin **33** (2008), no. 6, 611–619.
- [5] P.G. Ciarlet, *Handbook of numerical analysis: Finite element methods*, vol. II, Elsevier Science Publishers B.V. (North-Holland), 1991.
- [6] J. Dvořák, *Optimization of composite materials*, Master's thesis, The Charles University in Prague, June 1993.
- [7] ———, *A reliable numerical method for computing homogenized coefficients*, Tech. report, available at CiteSeerX <http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.45.1190>, 1995.
- [8] 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** (1999), no. 1, 41–47.
- [9] V.V. Jikov, S.M. Kozlov, and O.A. Oleinik, *Homogenization of differential operators and integral functionals*, Springer-Verlag, 1994.
- [10] L.C.Evans, *Partial differential equations*, vol. 19, American Mathematical Society, 2000.
- [11] J. Lukeš and J. Malý, *Measure and integral*, Matfyzpress Prague, 1995.

REFERENCES

28

- [12] G. W. Milton and R. V. Kohn, *Variational bounds on the effective moduli of anisotropic composites*, Journal of Mechanics and Physics of Solids **36** (1988), no. 6, 597–629.
- [13] V. Monchiet and G. Bonnet, *A polarization-based FFT iterative scheme for computing the effective properties of elastic composites with arbitrary contrast*, International Journal for Numerical Methods in Engineering **89** (2012), no. 11, 1419–1436.
- [14] H. Moulinec and P. Suquet, *A fast numerical method for computing the linear and nonlinear mechanical properties of composites*, Comptes rendus de l'Académie des sciences. Série II, Mécanique, physique, chimie, astronomie **318** (1994), no. 11, 1417–1423.
- [15] ———, *A numerical method for computing the overall response of nonlinear composites with complex microstructure*, Computer Methods in Applied Mechanics and Engineering **157** (1997), no. 1–2, 69–94.
- [16] J. Němeček, V. Králík, and J. Vondřejc, *Micromechanical analysis of heterogeneous structural materials*, Cement and Concrete Composites (2012).
- [17] J. Němeček, V. Králík, and J. Vondřejc, *A two-scale micromechanical model for aluminium foam based on results from nanoindentation*, (2012), Submitted.
- [18] J. Saranen and G. Vainikko, *Periodic integral and pseudodifferential equations with numerical approximation*, Springer Monographs Mathematics, 2000.
- [19] G. Vainikko, *Fast solvers of the Lippmann-Schwinger equation*, Direct and Inverse Problems of Mathematical Physics (R. P. Gilbert, J. Kajiwara, and Y. S. Xu, eds.), International Society for Analysis, Applications and Computation, vol. 5, Kluwer Academic Publishers, Dordrecht, The Netherlands, 2000, pp. 423–440.
- [20] V. Vinogradov and G. W. Milton, *An accelerated FFT algorithm for thermoelastic and non-linear composites*, International Journal for Numerical Methods in Engineering **76** (2008), no. 11, 1678–1695.
- [21] J. Vondřejc, J. Zeman, and I. Marek, *Analysis of a Fast Fourier Transform based method for modeling of heterogeneous materials*, Lecture Notes in Computer Science **7116** (2012), 512–522.
- [22] Z. Wieçkowski, *Dual finite element methods in mechanics of composite materials*, Journal of Theoretical and Applied Mechanics **2** (1995), no. 33, 233–252.
- [23] J. Yvonnet, *A fast method for solving microstructural problems defined by digital images: a space lippmann-schwinger scheme*, International Journal for Numerical Methods in Engineering (2012).
- [24] 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** (2010), no. 21, 8065–8071.
- [25] J. Zeman and M. Šejnoha, *Numerical evaluation of effective elastic properties of graphite fiber tow impregnated by polymer matrix*, Journal of the Mechanics and Physics of Solids **49** (2001), no. 1, 69–90.

Part VII

Paper 6

Authors:

Jaroslav Vondřejc, Jan Zeman, and Ivo Marek

Title:

Arbitrary precise guaranteed bounds of homogenized material coefficients by FFT-based finite element method

**ARBITRARY PRECISE GUARANTEED BOUNDS OF
HOMOGENIZED MATERIAL COEFFICIENTS BY FFT-BASED
FINITE ELEMENT METHOD ***

JAROSLAV VONDŘEJC[†], JAN ZEMAN^{‡§}, AND IVO MAREK[¶]

Abstract. In this work, the guaranteed bounds of homogenized material coefficients are calculated with an arbitrary precision; primal and dual variational formulations are evaluated with approximate microscopic (local) solutions to produce the upper-lower bounds. Contrary to Dvořák [5, 6] and Więckowski [32] employing the approach with local solutions from h and p-version of Finite Element Method (FEM), we utilize the FFT-based FEM using trigonometric polynomials as the basis functions [31]. The connection between the primal and the dual formulations is investigated in discretized form; the solution of the dual formulation can be avoided, however, the theory substantially differs for even and odd number of discretization points. Numerical examples confirm the theoretical results about the rates of convergence and about the upper-lower bounds of the homogenized matrix.

Key words. FFT, homogenization, guaranteed bounds, Finite Element Method

AMS subject classifications. 65N15, 74Q20

1. Introduction. Guaranteed bounds of homogenized (effective) material properties, with some confidence in its accuracy, is essential for a reliable design [24]; a special attention is addressed to the upper-lower bounds for linear periodic elliptic problems, a scalar one or elasticity. The majority of bounds — e.g. Voigt and Reuss bounds [29, 25], Hashin and Shtrikman bounds [9, 10, 11] — are derived for arbitrary geometry of material phases, only the knowledge of volume fractions is assumed. However, these a priori estimates are rather wide especially for materials with highly oscillating coefficients.

Besides, assuming material properties are well known, there are a posteriori estimates producing reliable and guaranteed homogenized bounds with an arbitrary precision independently introduced in [6, 32] for a scalar problem and elasticity resp., see Sec. 2.2 for a summary.

This approach is based on the primal and the dual variational formulations. Approximated microscopic fields, satisfying linear second order elliptic partial differential equations with periodic boundary conditions and prescribed macroscopic (averaged) values, are calculated and their a posteriori evaluation in both formulations produces the guaranteed bounds of homogenized properties (matrix).

Contrary to already mentioned works [6, 32] incorporating classical p- and h-version resp. of Finite Element Method (FEM) for computation of local fields, we rather employ the FFT-based FEM originated from [19] and theoretically described in [31], see Sec. 3.3 for its overview. Noting, the method — a widely used in engineering problems, recently in [22, 23, 14, 28, 12, 15, 21] — has obtained miscellaneous variants

*This work was supported by the Czech Science Foundation through project No. GAČR P105/12/0331 and by the Grant Agency of the Czech Technical University in Prague through project No. SGS 12/027/OHK1/1T/11.

[†]Department of Mechanics, Faculty of Civil Engineering, Czech Technical University in Prague, Thákurova 7, 166 29 Prague 6, Czech Republic (vondrej@gmail.com).

[‡]Department of Mechanics, Faculty of Civil Engineering, Czech Technical University in Prague, Thákurova 7, 166 29 Prague 6, Czech Republic and Centre of Excellence IT4Innovations, VŠB-TU Ostrava, 17. listopadu 15/2172 708 33 Ostrava-Poruba, Czech Republic (zemanj@cml.fsv.cvut.cz).

[§]Department of Mathematics, Faculty of Civil Engineering, Czech Technical University in Prague, Thákurova 7, 166 29 Prague 6, Czech Republic (marekivo@mat.fsv.cvut.cz).

and modifications [8, 16, 2, 3, 18].

In this work, for simplicity, the scalar problem of electric conduction is considered, however, an extension to elasticity is feasible. Numerical homogenization of a periodic media with its upper and lower bounds, according to [5, 6], is summarized in Section 2.

Section 3.3 deals with the FFT-based FEM based on the discretization with trigonometric polynomials. In our previous works [30, 31], only the discretization with the odd number of points in each direction is considered; the problem with even number of discretization points was identified and partially solved in [20, section 2.4.2]. We describe the theory for the even number of discretization points in a way to satisfy the conformity of the method; finite dimensional spaces are subspaces of the trial space — it is the natural requirement for guaranteed bounds of the homogenized matrix.

In section 3.4, the connection between discretized primal-dual formulations is studied. It provides a computational simplification: the solution of the dual formulation, necessary for the lower bound, can be avoided; moreover, for the odd number of discretization points, it can be fully omitted.

The upper-lower bounds of the homogenized matrix are calculated in Section 3.5. Although, their exact values, dependent on arbitrary non-regular material coefficients, are able to exactly evaluate only for some special cases. Hence, we propose their approximations; moreover, a careful computation still guarantees the upper-lower bounds structure.

Numerical examples in last Section 4 validate the theory and demonstrate the upper-lower bounds structure of homogenized matrices.

2. Preliminaries to homogenization and to guaranteed bounds. In this section, we summarize the well established theory of homogenization of periodic media and add the part about upper and lower bounds of the homogenized material coefficients. In the sequel, letter d denotes the dimension of the model problem, assuming $d = 2, 3$; Greek letters α, β are reserved to indices relating dimension, thus ranging $1, \dots, d$ — further, it is for simplicity omitted.

Sets \mathbb{C}^d and \mathbb{R}^d are the spaces of complex and real vectors with canonical basis $\{\epsilon_\alpha\}$ and are equipped with Lebesgue measure $d\mathbf{x}$. We denote by $|\Omega|_d$ the d -dimensional Lebesgue measure of a measurable set $\Omega \subset \mathbb{R}^d$. Standard Euklidean norm $\|\cdot\|_2$ on \mathbb{C}^d is induced by scalar product $(\mathbf{u}, \mathbf{v})_{\mathbb{C}^d} = \sum_\alpha u_\alpha \overline{v_\alpha}$ for $\mathbf{u}, \mathbf{v} \in \mathbb{C}^d$.

Set $\mathbb{R}_{\text{spd}}^{d \times d}$ denotes the space of symmetric positive definite matrices of size $d \times d$ with norm $\|\mathbf{A}\|_2 = \max_{\mathbf{x} \in \mathbb{R}^d, \|\mathbf{x}\|=1} \|\mathbf{A}\mathbf{x}\|_2$ that equals to the largest eigenvalue.

Function $f : \mathbb{R}^d \mapsto \mathbb{R}$ is \mathbf{Y} -periodic (with period $\mathbf{Y} \in \mathbb{R}^d$) if $f(\mathbf{x} + \mathbf{Y} \odot \mathbf{k}) = f(\mathbf{x})$ for arbitrary $\mathbf{x} \in \mathbb{R}^d$, $\mathbf{k} \in \mathbb{Z}^d$, where operator \odot denotes the element-wise multiplication. Then, \mathbf{Y} -periodic functions suffice to define only on a periodic unit cell (PUC), the set defined as $\mathcal{Y} := (-Y_\alpha, Y_\alpha)_{\alpha=1}^d \subset \mathbb{R}^d$. We will identify two integrable functions which are equal almost everywhere. The mean value of function $\mathbf{v} \in L_{\text{per}}^2(\mathcal{Y}; \mathbb{R}^d)$ over periodic unit cell \mathcal{Y} is denoted as $\langle \mathbf{v} \rangle := \frac{1}{|\mathcal{Y}|_d} \int_{\mathcal{Y}} \mathbf{v}(\mathbf{x}) d\mathbf{x} \in \mathbb{R}^d$.

We define space $C_{\text{per}}(\mathcal{Y}; \mathbb{X})$ of continuous \mathbf{Y} -periodic functions $\mathbb{R}^d \mapsto \mathbb{X}$, where \mathbb{X} is some finite dimensional vector space, e.g. $\mathbb{C}, \mathbb{R}, \mathbb{C}^d, \mathbb{R}^d$. Vector valued functions, for $\mathbb{X} = \mathbb{C}^d$ or $\mathbb{X} = \mathbb{R}^d$, are denoted with small bold letters \mathbf{v} and have components v_α .

Spaces $L_{\text{per}}^2(\mathcal{Y}; \mathbb{X})$ or $L_{\text{per}}^\infty(\mathcal{Y}; \mathbb{R}_{\text{spd}}^{d \times d})$ are composed of functions $\mathbf{v} : \mathbb{R}^d \mapsto \mathbb{X}$ or $\mathbf{A} : \mathbb{R}^d \mapsto \mathbb{R}_{\text{spd}}^{d \times d}$ having \mathbf{Y} -periodic, measurable components v_α or $A_{\alpha\beta}$ and having finite norm, i.e. $\|\mathbf{v}\|_{L_{\text{per}}^2(\mathcal{Y}; \mathbb{X})} < \infty$ or $\|\mathbf{A}\|_{L_{\text{per}}^\infty(\mathcal{Y}; \mathbb{R}_{\text{spd}}^{d \times d})} < \infty$. The first norm is deduced from scalar product $(\mathbf{u}, \mathbf{v})_{L_{\text{per}}^2(\mathcal{Y}; \mathbb{X})} = \frac{1}{|\mathcal{Y}|_d} \int_{\mathcal{Y}} (\mathbf{u}(\mathbf{x}), \mathbf{v}(\mathbf{x}))_{\mathbb{X}} d\mathbf{x}$ while the second norm

is defined as $\|\mathbf{A}\|_{L^\infty(\mathcal{Y}; \mathbb{R}^{d \times d})} = \text{ess sup}_{\mathbf{x} \in \mathcal{Y}} \|\mathbf{A}(\mathbf{x})\|_2$. If there is no ambiguity, both the norms and the scalar products are denoted with subscript L^2_{per} or L^∞_{per} rather than $L^2_{\text{per}}(\mathcal{Y}; \mathbb{X})$ or $L^\infty_{\text{per}}(\mathcal{Y}; \mathbb{R}^{d \times d})$.

Next, we define the spaces of Helmholtz decomposition of $L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d) = \mathcal{U} \oplus^\perp \mathcal{E} \oplus^\perp \mathcal{J}$, i.e. the spaces of constant, curl-free with zero mean, and divergence free with zero mean fields

$$\mathcal{U} = \{\mathbf{v} \in L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d) : \mathbf{v}(\mathbf{x}) = \text{const}\}, \quad (2.1a)$$

$$\mathcal{E} = \{\mathbf{v} \in L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d) : \nabla \times \mathbf{v} = \mathbf{0}, \langle \mathbf{v} \rangle = \mathbf{0}\}, \quad (2.1b)$$

$$\mathcal{J} = \{\mathbf{v} \in L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d) : \nabla \cdot \mathbf{v} = 0, \langle \mathbf{v} \rangle = \mathbf{0}\}, \quad (2.1c)$$

where differential operator $\nabla = (\frac{\partial}{\partial x_\alpha})_{\alpha=1}^d$ is meant in the distributional sense. For dimension $d \neq 3$, the curl-free condition in (2.1b) means $(\nabla \times \mathbf{v})_{\alpha\beta} := \frac{\partial v_\alpha}{\partial x_\beta} - \frac{\partial v_\beta}{\partial x_\alpha} = 0$. Since space \mathcal{U} consists of constant functions, we identify spaces \mathcal{U} and \mathbb{R}^d ; it validates the operation such as $\mathbf{E} + \mathbf{v} \in L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)$ for $\mathbf{E} \in \mathbb{R}^d$, $\mathbf{v} \in L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)$ and $\mathbf{A}\mathbf{J} \in \mathbb{R}^d$ for $\mathbf{A} \in \mathbb{R}^{d \times d}_{\text{spd}}$ and $\mathbf{J} \in \mathcal{U}$.

2.1. Homogenization in primal and dual formulations. The theory is demonstrated for a scalar problem modeling: diffusion, stationary heat transfer, or electric conductivity — our choice. The well established theory of periodic homogenization for linear elliptic partial differential equations can be found in [1, 13, 4]. Although, we focus only on the numerical solution of microscopic fields and the consequence evaluation of the homogenized matrix and its bounds.

NOTATION 2.1. *Here and in the sequel, $\mathbf{A} \in L^\infty(\mathcal{Y}, \mathbb{R}^{d \times d}_{\text{spd}})$ denotes symmetric and uniformly elliptic¹ material coefficients of electric conductivity, $\mathbf{e} \in \mathcal{E}$ and $\mathbf{j} \in \mathcal{J}$ perturbation of electric field and electric current resp., and $\mathbf{E}, \mathbf{J} \in \mathcal{U}$ their macroscopic counterparts. Then their summation $(\mathbf{E} + \mathbf{e})$ and $(\mathbf{J} + \mathbf{j})$ represent microscopic fields.*

REMARK 2.2. *The variables from previous notations, being additionally sufficiently smooth, satisfy the differential equations*

$$\mathbf{J} + \mathbf{j} = \mathbf{A}(\mathbf{E} + \mathbf{e}) \quad \nabla \cdot \mathbf{j} = 0 \quad \nabla \times \mathbf{e} = 0 \quad (2.2)$$

that clarify the definition of subspaces \mathcal{U}, \mathcal{E} , and \mathcal{J} ; the addition of periodic boundary conditions and prescription of macroscopic loads $\mathbf{E}, \mathbf{J} \in \mathbb{R}^d$ sets the homogenization problem.

DEFINITION 2.3 (Homogenization problem). *The primal and the dual homogenization problem states: find homogenized matrices $\mathbf{A}_{\text{eff}}, \mathbf{A}_{\text{eff},D} \in \mathbb{R}^{d \times d}$ such that for arbitrary fixed macroscopic quantities $\mathbf{E}, \mathbf{J} \in \mathbb{R}^d$, the following relations hold*

$$(\mathbf{A}_{\text{eff}} \mathbf{E}, \mathbf{E})_{\mathbb{R}^d} = \inf_{\mathbf{e} \in \mathcal{E}} (\mathbf{A}(\mathbf{E} + \mathbf{e}), \mathbf{E} + \mathbf{e})_{L^2_{\text{per}}} \quad (2.3a)$$

$$(\mathbf{A}_{\text{eff},D}^{-1} \mathbf{J}, \mathbf{J})_{\mathbb{R}^d} = \inf_{\mathbf{j} \in \mathcal{J}} (\mathbf{A}^{-1}(\mathbf{J} + \mathbf{j}), \mathbf{J} + \mathbf{j})_{L^2_{\text{per}}} \quad (2.3b)$$

REMARK 2.4. *The homogenized matrices are symmetric positive definite and they equal to one another $\mathbf{A}_{\text{eff}} = \mathbf{A}_{\text{eff},D}$ — it follows from the perturbation duality*

¹There exists positive constant $c_A > 0$ such that inequality $c_A \|u\|_2^2 \leq (\mathbf{A}(\mathbf{x})u, u)_{\mathbb{R}^d}$ holds for almost all $x \in \mathcal{Y}$ and all nonzero $u \in \mathbb{R}^d$.

theorem, see e.g. [7, 27, 6] and compare it to Lem. 3.24 and 3.25 for a discrete setting. The symmetry comes from upcoming Eq. (2.6) and symmetry of \mathbf{A} . The positive definiteness is the consequence of the uniform ellipticity of \mathbf{A} and simultaneously the consequence of Voight $(\mathbf{A}_{\text{eff}}\mathbf{E}, \mathbf{E})_{\mathbb{R}^d} \leq (\mathbf{A}\mathbf{E}, \mathbf{E})_{L^2_{\text{per}}}$ and Reuss $(\mathbf{A}_{\text{eff,D}}^{-1}\mathbf{J}, \mathbf{J})_{\mathbb{R}^d} \leq (\mathbf{A}^{-1}\mathbf{J}, \mathbf{J})_{L^2_{\text{per}}}$ bounds according to [29, 25]; the bounds come from (2.3) for $\mathbf{e} = \mathbf{j} \equiv 0$.

The minimizers of the variational formulations (2.3a) and (2.3b) can also be found as a solution of weak formulations with existence and uniqueness provided by Lax-Milgram lemma. Thanks to linearity, the minimizers can be found only for unitary macroscopic loads.

DEFINITION 2.5 (auxiliary problems). We say that $\tilde{\mathbf{e}}^{(\alpha)} \in \mathcal{E}$ and $\tilde{\mathbf{j}}^{(\alpha)} \in \mathcal{J}$ are unitary minimizers if

$$(\mathbf{A}\tilde{\mathbf{e}}^{(\alpha)}, \mathbf{v})_{L^2_{\text{per}}} = -(\mathbf{A}\boldsymbol{\epsilon}_\alpha, \mathbf{v})_{L^2_{\text{per}}}, \quad \forall \mathbf{v} \in \mathcal{E}, \quad (2.4)$$

$$(\mathbf{A}^{-1}\tilde{\mathbf{j}}^{(\alpha)}, \mathbf{v})_{L^2_{\text{per}}} = -(\mathbf{A}^{-1}\boldsymbol{\epsilon}_\alpha, \mathbf{v})_{L^2_{\text{per}}}, \quad \forall \mathbf{v} \in \mathcal{J}. \quad (2.5)$$

Next, we define unitary microscopic fields $\mathbf{e}^{(\alpha)} := \boldsymbol{\epsilon}_\alpha + \tilde{\mathbf{e}}^{(\alpha)}$ and $\mathbf{j}^{(\alpha)} := \boldsymbol{\epsilon}_\alpha + \tilde{\mathbf{j}}^{(\alpha)}$.

REMARK 2.6 (Consequences of the linearity). Minimizers $\tilde{\mathbf{e}}^{(\mathbf{E})} \in \mathcal{E}$ and $\tilde{\mathbf{j}}^{(\mathbf{J})} \in \mathcal{J}$ of the homogenization problems for macroscopic fields $\mathbf{E}, \mathbf{J} \in \mathbb{R}^d$ can be obtained, due to linear structure, from unitary minimizers

$$\tilde{\mathbf{e}}^{(\mathbf{E})} = \sum_{\alpha} E_{\alpha} \tilde{\mathbf{e}}^{(\alpha)}, \quad \tilde{\mathbf{j}}^{(\mathbf{J})} = \sum_{\alpha} J_{\alpha} \tilde{\mathbf{j}}^{(\alpha)}.$$

Alike, the components of homogenized material coefficients states

$$(\mathbf{A}_{\text{eff}})_{\alpha\beta} = (\mathbf{A}\mathbf{e}^{(\beta)}, \mathbf{e}^{(\alpha)})_{L^2_{\text{per}}}, \quad (\mathbf{A}_{\text{eff}}^{-1})_{\alpha\beta} = (\mathbf{A}^{-1}\mathbf{j}^{(\beta)}, \mathbf{j}^{(\alpha)})_{L^2_{\text{per}}}. \quad (2.6)$$

REMARK 2.7. The dual unitary microscopic fields $\mathbf{j}^{(\beta)}$ can be expressed as a linear combination of primal ones $\mathbf{e}^{(\alpha)}$, hence $\mathbf{j}^{(\beta)} = \mathbf{A} \sum_{\alpha=1}^d E_{\alpha} \mathbf{e}^{(\alpha)}$ where $\mathbf{E} = \mathbf{A}_{\text{eff}}^{-1} \boldsymbol{\epsilon}_{\beta}$. Similarly to Rem. 2.4, it comes from the perturbation duality theorem, see e.g. [7, 27, 6] and compare it to Lem. 3.24 and 3.25 for a discrete setting.

2.2. Upper-lower bounds of homogenized matrix. The upper-lower bounds obtained from a posteriori error estimates were introduced by Dvořák [5, 6] for a scalar problem and later independently by Więckowski [32] for linear elasticity; this section provides a summary of results in [6]. In what follows, relation $\mathbf{C} \preceq \mathbf{D}$ between matrices $\mathbf{C}, \mathbf{D} \in \mathbb{R}_{\text{spd}}^{d \times d}$ stands for ordering in the sense of quadratic forms; it is equivalent to $\mathbf{E} \cdot \mathbf{C}\mathbf{E} \leq \mathbf{E} \cdot \mathbf{D}\mathbf{E}$ for all $\mathbf{E} \in \mathbb{R}^d$. In this section, we will work with some conforming approximations of unitary minimizers $\tilde{\mathbf{e}}^{(\alpha)}$ and $\tilde{\mathbf{j}}^{(\alpha)}$, namely with $\tilde{\mathbf{e}}_{\mathbf{N}}^{(\alpha)} = (\mathbf{e}_{\mathbf{N}}^{(\alpha)} - \boldsymbol{\epsilon}_\alpha) \in \mathcal{E}$ and $\tilde{\mathbf{j}}_{\mathbf{N}}^{(\alpha)} = (\mathbf{j}_{\mathbf{N}}^{(\alpha)} - \boldsymbol{\epsilon}_\alpha) \in \mathcal{J}$; parameter \mathbf{N} represent inverse of discretization size of FEM or the number of discretization points in our case of FFT-based FEM, for detail see 3.3.

DEFINITION 2.8. We say that matrices $\bar{\mathbf{A}}_{\text{eff}, \mathbf{N}}, \underline{\mathbf{A}}_{\text{eff}, \mathbf{N}} \in \mathbb{R}^{d \times d}$ defined as

$$(\bar{\mathbf{A}}_{\text{eff}, \mathbf{N}})_{\alpha\beta} = (\mathbf{A}\mathbf{e}_{\mathbf{N}}^{(\beta)}, \mathbf{e}_{\mathbf{N}}^{(\alpha)})_{L^2_{\text{per}}}, \quad (\underline{\mathbf{A}}_{\text{eff}, \mathbf{N}}^{-1})_{\alpha\beta} = (\mathbf{A}^{-1}\mathbf{j}_{\mathbf{N}}^{(\beta)}, \mathbf{j}_{\mathbf{N}}^{(\alpha)})_{L^2_{\text{per}}}$$

are upper and lower bounds of homogenized properties \mathbf{A}_{eff} .

Correctness of the definition, the fact that they are truly the upper-lower bounds, is stated in Lem. 2.10 that is based on a following statement.

LEMMA 2.9. *Let $\mathbf{C}, \mathbf{D} \in \mathbb{R}_{\text{spd}}^{d \times d}$ be such that $\mathbf{C} \preceq \mathbf{D}$. Then $\mathbf{D}^{-1} \preceq \mathbf{C}^{-1}$.*

Proof. From assumption $\mathbf{C} \preceq \mathbf{D}$, we claim that $\mathbf{R}^T \mathbf{C} \mathbf{R} \preceq \mathbf{R}^T \mathbf{D} \mathbf{R}$ for any regular matrix $\mathbf{R} \in \mathbb{R}^{d \times d}$; it comes from $(\mathbf{R}\mathbf{E}) \cdot \mathbf{C} (\mathbf{R}\mathbf{E}) \leq (\mathbf{R}\mathbf{E}) \cdot \mathbf{D} (\mathbf{R}\mathbf{E})$ holding for arbitrary $\mathbf{E} \in \mathbb{R}^d$ and the invertibility of matrix \mathbf{R} , namely the property $\mathbf{R}(\mathbb{R}^d) = \mathbb{R}^d$.

Since \mathbf{D} is positive definite, matrix $\mathbf{D}^{-\frac{1}{2}}$ exists. Consequently, the multiplication of inequality $\mathbf{C} \preceq \mathbf{D}$ by matrix $\mathbf{D}^{-\frac{1}{2}}$ produces inequality $\mathbf{D}^{-\frac{1}{2}} \mathbf{C} \mathbf{D}^{-\frac{1}{2}} \preceq \mathbf{I}$. Matrix $\mathbf{D}^{-\frac{1}{2}} \mathbf{C} \mathbf{D}^{-\frac{1}{2}}$ possesses all eigenvalues real (from symmetry) and smaller or equal to 1. Thus its inverse $\mathbf{D}^{\frac{1}{2}} \mathbf{C}^{-1} \mathbf{D}^{\frac{1}{2}}$ have eigenvalues larger or equal to 1. The next multiplication of inequality $\mathbf{I} \preceq \mathbf{D}^{\frac{1}{2}} \mathbf{C}^{-1} \mathbf{D}^{\frac{1}{2}}$ by matrix $\mathbf{D}^{-\frac{1}{2}}$ reveals the proof. \square

LEMMA 2.10. *The matrices from previous definition 2.8 are symmetric positive definite and satisfy the upper-lower bound structure*

$$\mathbf{A}_{\text{eff}} \preceq \overline{\mathbf{A}}_{\text{eff}, \mathbf{N}}, \quad \mathbf{A}_{\text{eff}}^{-1} \preceq \underline{\mathbf{A}}_{\text{eff}, \mathbf{N}}^{-1}, \quad \underline{\mathbf{A}}_{\text{eff}, \mathbf{N}} \preceq \mathbf{A}_{\text{eff}} \preceq \overline{\mathbf{A}}_{\text{eff}, \mathbf{N}}. \quad (2.7)$$

Proof. The first two inequalities follow from the minimization problems (2.3) that are evaluated approximately, for a particular choice of microscopic fields $\tilde{\mathbf{e}}_{\mathbf{N}}^{(\alpha)}$, $\tilde{\mathbf{j}}_{\mathbf{N}}^{(\alpha)}$ rather than minimizers $\tilde{\mathbf{e}}^{(\alpha)}$, $\tilde{\mathbf{e}}^{(\alpha)}$. The last inequality is a consequence of previous Lem. 2.9.

The symmetry of the homogenized matrices comes from their definition 2.8 and the symmetry of material coefficients \mathbf{A} . The positive definiteness is a consequence of first two inequalities in (2.7) with the positive definiteness of homogenized matrix $\mathbf{A}_{\text{eff}} \in \mathbb{R}_{\text{spd}}^{d \times d}$. \square

DEFINITION 2.11 (Approximate homogenized matrix with guaranteed error). *The mean of the upper-lower bounds from Def. 2.8, $\mathbf{A}_{\text{eff}, \mathbf{N}} = \frac{1}{2}(\overline{\mathbf{A}}_{\text{eff}, \mathbf{N}} + \underline{\mathbf{A}}_{\text{eff}, \mathbf{N}})$, is called approximate homogenized matrix with guaranteed error $\mathbf{D}_{\mathbf{N}} = \frac{1}{2}(\overline{\mathbf{A}}_{\text{eff}, \mathbf{N}} - \underline{\mathbf{A}}_{\text{eff}, \mathbf{N}})$.*

LEMMA 2.12 (Element-wise upper-lower bounds). *The upper-lower bounds from Def. 2.8 imply element-wise bounds, explicitly for diagonal components*

$$\underline{\mathbf{A}}_{\text{eff}, \mathbf{N}, \alpha\alpha} \leq \mathbf{A}_{\text{eff}, \alpha\alpha} \leq \overline{\mathbf{A}}_{\text{eff}, \mathbf{N}, \alpha\alpha}, \quad (2.8)$$

and for non-diagonal components, i.e. for $\alpha \neq \beta$

$$\underline{\mathbf{A}}_{\text{eff}, \mathbf{N}, \alpha\beta} - D_{\text{eff}, \mathbf{N}, \alpha\alpha} - D_{\text{eff}, \mathbf{N}, \beta\beta} \leq \mathbf{A}_{\text{eff}, \alpha\beta} \leq \overline{\mathbf{A}}_{\text{eff}, \mathbf{N}, \alpha\beta} + D_{\text{eff}, \mathbf{N}, \alpha\alpha} + D_{\text{eff}, \mathbf{N}, \beta\beta}.$$

Proof. The proof for the diagonal terms (2.8) comes from the inequality (2.7) tested with $\boldsymbol{\epsilon}_\alpha$. The estimates for the non-diagonal terms come from equality

$$2\mathbf{A}_{\text{eff}, \alpha\beta} = (\boldsymbol{\epsilon}_\alpha + \boldsymbol{\epsilon}_\beta) \cdot \mathbf{A}_{\text{eff}} (\boldsymbol{\epsilon}_\alpha + \boldsymbol{\epsilon}_\beta) - \mathbf{A}_{\text{eff}, \alpha\alpha} - \mathbf{A}_{\text{eff}, \beta\beta},$$

the first inequality in (2.7) tested with $\boldsymbol{\epsilon}_\alpha + \boldsymbol{\epsilon}_\beta$, i.e. $(\boldsymbol{\epsilon}_\alpha + \boldsymbol{\epsilon}_\beta) \cdot \mathbf{A}_{\text{eff}} (\boldsymbol{\epsilon}_\alpha + \boldsymbol{\epsilon}_\beta) \leq (\boldsymbol{\epsilon}_\alpha + \boldsymbol{\epsilon}_\beta) \cdot \overline{\mathbf{A}}_{\text{eff}, \mathbf{N}} (\boldsymbol{\epsilon}_\alpha + \boldsymbol{\epsilon}_\beta)$, and the inequalities for diagonal components (2.8). An analogy yields the lower bound. \square

LEMMA 2.13 (Properties of trace operator). *Let $\mathbf{B}, \mathbf{C}, \mathbf{D} \in \mathbb{R}_{\text{spd}}^{d \times d}$. Then*

$$\text{tr}(\mathbf{B} + \mathbf{C}) = \text{tr} \mathbf{B} + \text{tr} \mathbf{C}, \quad \text{tr}(\mathbf{BCD}) \leq (\text{tr} \mathbf{B})(\text{tr} \mathbf{C})(\text{tr} \mathbf{D}).$$

Moreover, let $\mathbf{C} \preceq \mathbf{D}$, then $0 \leq \text{tr}(\mathbf{D} - \mathbf{C}) \leq (\text{tr} \mathbf{D})^2 (\mathbf{C}^{-1} - \mathbf{D}^{-1})$.

Proof. The proof can be found in [6, Lemmas 4.3 and 4.4]. \square

LEMMA 2.14 (The rate of convergence of homogenized properties). *The trace of guaranteed error \mathbf{D}_N , from Def. 2.11, satisfies following inequality*

$$\mathrm{tr} \mathbf{D}_N \leq \|\mathbf{A}\|_{L^\infty_{\mathrm{per}}} \sum_{\alpha} \|\tilde{\mathbf{e}}^{(\alpha)} - \tilde{\mathbf{e}}_N^{(\alpha)}\|_{L^2_{\mathrm{per}}}^2 + (\mathrm{tr} \mathbf{A}_{\mathrm{eff}})^2 \|\mathbf{A}^{-1}\|_{L^\infty_{\mathrm{per}}} \sum_{\alpha} \|\tilde{\mathbf{j}}^{(\alpha)} - \tilde{\mathbf{j}}_N^{(\alpha)}\|_{L^2_{\mathrm{per}}}^2.$$

Proof. First, we prove two statements

$$0 \leq \bar{\mathbf{A}}_{\mathrm{eff},N,\alpha\alpha} - \mathbf{A}_{\mathrm{eff},\alpha\alpha} \leq \|\mathbf{A}\|_{L^\infty} \|\tilde{\mathbf{e}}^{(\alpha)} - \tilde{\mathbf{e}}_N^{(\alpha)}\|_{L^2_{\mathrm{per}}}^2, \quad (2.9a)$$

$$0 \leq \underline{\mathbf{A}}_{\mathrm{eff},N,\alpha\alpha}^{-1} - \mathbf{A}_{\mathrm{eff},\alpha\alpha}^{-1} \leq \|\mathbf{A}^{-1}\|_{L^\infty} \|\tilde{\mathbf{j}}^{(\alpha)} - \tilde{\mathbf{j}}_N^{(\alpha)}\|_{L^2_{\mathrm{per}}}^2. \quad (2.9b)$$

Both inequalities follow from Lem. 2.10 and Hölder inequality; we demonstrate the calculation for the first one

$$\begin{aligned} 0 \leq \bar{\mathbf{A}}_{\mathrm{eff},N,\alpha\alpha} - \mathbf{A}_{\mathrm{eff},\alpha\alpha} &= |(\mathbf{A}\mathbf{e}^{(\alpha)}, \mathbf{e}^{(\alpha)})_{L^2_{\mathrm{per}}} - (\mathbf{A}\mathbf{e}_N^{(\alpha)}, \mathbf{e}_N^{(\alpha)})_{L^2_{\mathrm{per}}}| \\ &\leq |(\mathbf{A}(\tilde{\mathbf{e}}^{(\alpha)} - \tilde{\mathbf{e}}_N^{(\alpha)}), (\tilde{\mathbf{e}}^{(\alpha)} - \tilde{\mathbf{e}}_N^{(\alpha)}))_{L^2_{\mathrm{per}}}| \\ &\leq \|\mathbf{A}\|_{L^\infty} \|\tilde{\mathbf{e}}^{(\alpha)} - \tilde{\mathbf{e}}_N^{(\alpha)}\|_{L^2_{\mathrm{per}}}^2. \end{aligned}$$

Using the properties of the trace operator stated in Lem. 2.13 and previously proven statements (2.9), we finish the proof with direct calculation

$$\begin{aligned} \mathrm{tr} \mathbf{D}_N &= \mathrm{tr}(\bar{\mathbf{A}}_{\mathrm{eff},N} - \mathbf{A}_{\mathrm{eff}}) + \mathrm{tr}(\mathbf{A}_{\mathrm{eff}} - \underline{\mathbf{A}}_{\mathrm{eff},N}) \\ &\leq \mathrm{tr}(\bar{\mathbf{A}}_{\mathrm{eff},N} - \mathbf{A}_{\mathrm{eff}}) + (\mathrm{tr} \mathbf{A}_{\mathrm{eff}})^2 \mathrm{tr}(\underline{\mathbf{A}}_{\mathrm{eff},N}^{-1} - \mathbf{A}_{\mathrm{eff}}^{-1}) \\ &\leq \|\mathbf{A}\|_{L^\infty_{\mathrm{per}}} \sum_{\alpha} \|\tilde{\mathbf{e}}^{(\alpha)} - \tilde{\mathbf{e}}_N^{(\alpha)}\|_{L^2_{\mathrm{per}}}^2 + (\mathrm{tr} \mathbf{A}_{\mathrm{eff}})^2 \|\mathbf{A}^{-1}\|_{L^\infty_{\mathrm{per}}} \sum_{\alpha} \|\tilde{\mathbf{j}}^{(\alpha)} - \tilde{\mathbf{j}}_N^{(\alpha)}\|_{L^2_{\mathrm{per}}}^2. \end{aligned}$$

\square

REMARK 2.15. *The trace of \mathbf{D}_N will converges to zero for $\min_{\alpha} N_{\alpha} \rightarrow \infty$, if approximate minimizers $\tilde{\mathbf{e}}_N^{\alpha}$, $\tilde{\mathbf{j}}_N^{\alpha}$ converge to minimizers $\tilde{\mathbf{e}}^{\alpha}$, $\tilde{\mathbf{e}}^{\alpha}$ for $\min_{\alpha} N_{\alpha} \rightarrow \infty$. Moreover, the trace operator is a norm on the set of symmetric positive definite matrices showing the convergence of \mathbf{D}_N to zero for any matrix norm.*

3. Guaranteed bounds using FFT-based FEM. This section is the core of the work. It provides the theory for the arbitrary precise guaranteed bounds of the homogenized matrix calculated with local fields provided by the FFT-based FEM. We start in Sec. 3.1 with the definition of finite dimensional spaces, the spaces of trigonometric polynomials. Then we follow with fully discrete spaces that are the analogue to Helmholtz decomposition spaces. The FFT-based FEM is described in Sec. 3.3; additionally to [31], the theory for the non-odd number of discretization points is provided. Sec. 3.4 is dedicated to the connection between the primal and the dual formulations in the fully discrete setting. Finally, the calculation of the upper-lower bounds is explained in Sec. 3.5.

In the sequel, vector $\mathbf{N} \in \mathbb{N}^d$ is reserved for a number of discretization points, then scalar $|\mathbf{N}|_{\Pi} := \prod_{\alpha} N_{\alpha}$ denotes the number of degrees of freedom. If N_{α} is odd (even) for all α we talk about the odd (even) number of discretization points, otherwise about the non-odd ones. A reduced and a full index sets state for

$$\mathbb{Z}_{\mathbf{N}}^d = \left\{ \mathbf{k} \in \mathbb{Z}^d : |k_{\alpha}| < \frac{N_{\alpha}}{2} \right\}, \quad \underline{\mathbb{Z}}_{\mathbf{N}}^d = \left\{ \mathbf{k} \in \mathbb{Z}^d : -\frac{N_{\alpha}}{2} \leq k_{\alpha} < \frac{N_{\alpha}}{2} \right\}. \quad (3.1)$$

A multi-index notation is employed, in which \mathbb{R}^N represents $\mathbb{R}^{N_1 \times \dots \times N_d}$. Set $\mathbb{R}^{d \times N}$ represents the space of vectors \mathbf{v} with components v_α^n and $\mathbb{R}^{d \times d \times N \times N}$ the space of matrices \mathbf{A} with components $A_{\alpha\beta}^{nm}$ for α, β and $\mathbf{n}, \mathbf{m} \in \mathbb{Z}_N^d$. Next, vectors $\mathbf{v}^n \in \mathbb{R}^d$ for $\mathbf{n} \in \mathbb{Z}_N^d$ and $\mathbf{v}_\alpha \in \mathbb{R}^N$ for α represent subvectors of \mathbf{v} with components v_α^n . Analogically, submatrices $\mathbf{A}^{nm} \in \mathbb{R}^{d \times d}$ and $\mathbf{A}_{\alpha\beta} \in \mathbb{R}^{N \times N}$ can be defined. A scalar product on set $\mathbb{R}^{d \times N}$ is defined as $(\mathbf{u}, \mathbf{v})_{\mathbb{R}^{d \times N}} := \sum_\alpha \sum_{\mathbf{n} \in \mathbb{Z}_N^d} u_\alpha^n v_\alpha^n$ and matrix \mathbf{A} by vector \mathbf{v} multiplication as $(\mathbf{A}\mathbf{v})_\alpha^n := \sum_\beta \sum_{\mathbf{m} \in \mathbb{Z}_N^d} A_{\alpha\beta}^{nm} v_\beta^m$. Matrix \mathbf{A} is symmetric positive definite if relation $A_{\alpha\beta}^{mn} = A_{\beta\alpha}^{nm}$ holds for all components and inequality $(\mathbf{A}\mathbf{v}, \mathbf{v})_{\mathbb{R}^{d \times N}} > 0$ applies for arbitrary $\mathbf{v} \in \mathbb{R}^{d \times N}$. We use the serif font for vectors $\mathbf{v} \in \mathbb{R}^{d \times N}$ and matrices $\mathbf{A} \in \mathbb{R}^{d \times d \times N \times N}$ to distinguish from vectors $\mathbf{E} \in \mathbb{R}^d$ and matrices $\mathbf{A}_{\text{eff}} \in \mathbb{R}^{d \times d}$ and from vector valued functions $\mathbf{v} \in L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)$. In order to differentiate vectors and matrices for different number of discretization points N , we write them with subscript, i.e. \mathbf{v}_N and \mathbf{A}_N . Finally, operator \oplus^\perp denotes the direct sum of mutually orthogonal subspaces, e.g. $\mathbb{R}^d = \epsilon_1 \oplus^\perp \epsilon_2 \oplus^\perp \dots \oplus^\perp \epsilon_d$.

3.1. Trigonometric polynomials and their properties. In this section, we introduce the finite dimensional space of trigonometric polynomials and their properties. Definitions and lemmas adopted from [26] are amended for the non-odd number of discretization points N in order to satisfy the conformity of discretization — the requirement for the guaranteed upper-lower bounds.

NOTATION 3.1 (DFT). For $N \in \mathbb{N}^d$ we define, up to constant, unitary matrices $\mathbf{F}_N, \mathbf{F}_N^{-1} \in \mathbb{C}^{d \times d \times N \times N}$ of the Discrete Fourier transform (DFT) and its inverse (iDFT) as

$$\mathbf{F}_N = \frac{1}{|N|_\Pi} (\delta_{\alpha\beta} \omega_N^{-mn})_{\alpha,\beta=1,\dots,d}^{\mathbf{m},\mathbf{n} \in \mathbb{Z}_N^d} \quad \mathbf{F}_N^{-1} = (\delta_{\alpha\beta} \omega_N^{mn})_{\alpha,\beta=1,\dots,d}^{\mathbf{m},\mathbf{n} \in \mathbb{Z}_N^d}$$

where $\delta_{\alpha\beta}$ is Kronecker delta and $\omega_N^{mn} = \exp\left(2\pi i \sum_\alpha \frac{m_\alpha n_\alpha}{N_\alpha}\right)$ with $\mathbf{m}, \mathbf{n} \in \mathbb{Z}^d$ are their components.

DEFINITION 3.2 (nodal points, basis functions). Let $N \in \mathbb{N}^d$. We define nodal points of periodic unit cell $\mathbf{x}_N^n = \sum_\alpha \frac{2Y_\alpha n_\alpha}{N_\alpha} \epsilon_\alpha$ and Fourier and shape basis functions

$$\varphi_n(\mathbf{x}) = \exp\left(\pi i \sum_\alpha \frac{n_\alpha x_\alpha}{Y_\alpha}\right), \quad \varphi_{N,\mathbf{n}}(\mathbf{x}) = \frac{1}{|N|_\Pi} \sum_{\mathbf{m} \in \mathbb{Z}_N^d} \omega_N^{-mn} \varphi_m(\mathbf{x}), \quad \mathbf{n} \in \mathbb{Z}_N^d.$$

LEMMA 3.3 (Properties of φ_m and $\varphi_{N,\mathbf{m}}$). Let $\mathbf{m}, \mathbf{n} \in \mathbb{Z}_N^d$, then

$$(\varphi_m, \varphi_n)_{L^2_{\text{per}}} = \delta_{mn} \quad \varphi_n(\mathbf{x}_N^{\mathbf{m}}) = \omega_N^{mn} \quad (3.2a)$$

$$\varphi_{N,\mathbf{m}}(\mathbf{x}_N^{\mathbf{n}}) = \delta_{mn} \quad (\varphi_{N,\mathbf{m}}, \varphi_{N,\mathbf{n}})_{L^2_{\text{per}}} = \frac{\delta_{mn}}{|N|_\Pi} \quad (3.2b)$$

Proof. The proof, which can be also found in [31], is the consequence of direct calculation; the proof of last equality comes from orthogonality of vectors $(\omega_N^{mn})_{\mathbf{m} \in \mathbb{Z}_N^d}$ for $\mathbf{n} \in \mathbb{Z}_N^d$ in \mathbb{C}^N . \square

DEFINITION 3.4 (Trigonometric polynomials). For $\mathbf{N} \in \mathbb{N}^d$, we define the spaces of trigonometric polynomials $\mathcal{T}_{\mathbf{N}}$, $\tilde{\mathcal{T}}_{\mathbf{N}}$ and their vector valued versions $\mathcal{T}_{\mathbf{N}}^d, \tilde{\mathcal{T}}_{\mathbf{N}}^d$ as

$$\begin{aligned} \mathcal{T}_{\mathbf{N}} &= \left\{ \sum_{\mathbf{n} \in \mathbb{Z}_{\mathbf{N}}^d} \hat{v}^{\mathbf{n}} \varphi_{\mathbf{n}} : \hat{v}^{\mathbf{n}} \in \mathbb{C}, \hat{v}^{\mathbf{n}} = \overline{\hat{v}^{-\mathbf{n}}} \right\}, & \mathcal{T}_{\mathbf{N}}^d &= \{ \mathbf{v} : v_{\alpha} \in \mathcal{T}_{\mathbf{N}} \}. \\ \tilde{\mathcal{T}}_{\mathbf{N}} &= \left\{ \sum_{\mathbf{n} \in \mathbb{Z}_{\mathbf{N}}^d} v^{\mathbf{n}} \varphi_{\mathbf{N}, \mathbf{n}} : v^{\mathbf{n}} \in \mathbb{R} \right\}, & \tilde{\mathcal{T}}_{\mathbf{N}}^d &= \{ \mathbf{v} : v_{\alpha} \in \tilde{\mathcal{T}}_{\mathbf{N}} \}. \end{aligned}$$

DEFINITION 3.5 (Interpolation projection). Interpolation operator $\mathcal{Q}_{\mathbf{N}} : C_{\text{per}}(\mathcal{Y}; \mathbb{R}^d) \rightarrow L^2_{\text{per}}(\mathcal{Y}; \mathbb{C}^d)$ is defined as

$$\mathcal{Q}_{\mathbf{N}}[f] = \sum_{\mathbf{m} \in \mathbb{Z}_{\mathbf{N}}^d} f(\mathbf{x}_{\mathbf{N}}^{\mathbf{m}}) \varphi_{\mathbf{N}, \mathbf{m}}.$$

LEMMA 3.6. Interpolation operator $\mathcal{Q}_{\mathbf{N}}$ is a projection and its image is $\tilde{\mathcal{T}}_{\mathbf{N}}^d$.

Proof. It comes from the definition of operator $\mathcal{Q}_{\mathbf{N}}$ in Def. 3.5, space $\tilde{\mathcal{T}}_{\mathbf{N}}^d$ in Def. 3.4, and second property in Lem. (3.2b). \square

DEFINITION 3.7. Operator $\mathcal{I}_{\mathbf{N}} : \tilde{\mathcal{T}}_{\mathbf{N}}^d \rightarrow \mathbb{R}^{d \times N}$ stocks the values of the trigonometric polynomials at the nodal points to a vector $\mathcal{I}_{\mathbf{N}}[\mathbf{v}_{\mathbf{N}}] = (\mathbf{v}_{\mathbf{N}, \alpha}(\mathbf{x}_{\mathbf{N}}^{\mathbf{n}}))_{\alpha=1, \dots, d}^{\mathbf{n} \in \mathbb{Z}_{\mathbf{N}}^d}$.

LEMMA 3.8. Operator $\mathcal{I}_{\mathbf{N}}$ from the previous definition is an isomorphism.

Proof. The proof is the consequence of Def. 3.7 and of the second property in Eq. 3.2a. \square

REMARK 3.9 (Connection of representation). Trigonometric polynomial $\mathbf{v}_{\mathbf{N}} \in \tilde{\mathcal{T}}_{\mathbf{N}}^d$ can be uniquely expressed using both the Fourier coefficients and the function values at the nodal points

$$\mathbf{v}_{\mathbf{N}} = \sum_{\mathbf{m} \in \mathbb{Z}_{\mathbf{N}}^d} \mathbf{v}_{\mathbf{N}}(\mathbf{x}_{\mathbf{N}}^{\mathbf{m}}) \varphi_{\mathbf{N}, \mathbf{m}} = \sum_{\mathbf{n} \in \mathbb{Z}_{\mathbf{N}}^d} \hat{\mathbf{v}}_{\mathbf{N}}(\mathbf{n}) \varphi_{\mathbf{n}}. \quad (3.3)$$

with a connection through the DFT as $\hat{\mathbf{v}}_{\mathbf{N}} = \mathbf{F}_{\mathbf{N}} \mathbf{v}_{\mathbf{N}}$, where $\mathbf{v}_{\mathbf{N}} := \mathcal{I}_{\mathbf{N}}[\mathbf{v}_{\mathbf{N}}]$ and the vector of Fourier coefficients $\hat{\mathbf{v}}_{\mathbf{N}}$ has components $(\hat{\mathbf{v}}_{\mathbf{N}})_{\alpha}^{\mathbf{m}} := \hat{v}_{\mathbf{N}, \alpha}(\mathbf{m})$. Thus, space $\tilde{\mathcal{T}}_{\mathbf{N}}^d$ can be alternatively characterized with the Fourier coefficients as

$$\tilde{\mathcal{T}}_{\mathbf{N}}^d = \left\{ \sum_{\mathbf{m} \in \mathbb{Z}_{\mathbf{N}}^d} \hat{\mathbf{v}}_{\mathbf{N}}^{\mathbf{m}} \varphi_{\mathbf{m}} : \hat{\mathbf{v}}_{\mathbf{N}} \in \mathbf{F}_{\mathbf{N}}(\mathbb{R}^{d \times N}) \right\}.$$

REMARK 3.10. The trigonometric polynomials are real valued if the Fourier coefficients obey conjugate symmetry $\hat{\mathbf{v}}(\mathbf{n}) = \overline{\hat{\mathbf{v}}(-\mathbf{n})}$ for all $\mathbf{n} \in \mathbb{Z}^d$. In Def. 3.4, it is valid only for trigonometric polynomials $\mathcal{T}_{\mathbf{N}} \subset L^2_{\text{per}}(\mathcal{Y})$ and $\mathcal{T}_{\mathbf{N}}^d \subset L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)$.

A peculiar situation occurs for space $\tilde{\mathcal{T}}_{\mathbf{N}}^d$. If \mathbf{N} is odd, both spaces coincide $\mathcal{T}_{\mathbf{N}}^d = \tilde{\mathcal{T}}_{\mathbf{N}}^d$ as the index sets do $\mathbb{Z}_{\mathbf{N}}^d = \underline{\mathbb{Z}}_{\mathbf{N}}^d$; generally, the inclusion $\mathcal{T}_{\mathbf{N}}^d \subseteq \tilde{\mathcal{T}}_{\mathbf{N}}^d$ holds. Unfortunately, space $\tilde{\mathcal{T}}_{\mathbf{N}}^d$ fails to be real valued $\tilde{\mathcal{T}}_{\mathbf{N}}^d \not\subseteq L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)$ because the Fourier coefficients with frequencies $\mathbf{n} \in \underline{\mathbb{Z}}_{\mathbf{N}}^d \setminus \mathbb{Z}_{\mathbf{N}}^d$ miss the opposite counterpart with frequencies $-\mathbf{n}$.

3.2. Fully discrete spaces. In this section, we define fully discrete spaces — the spaces storing the values of the trigonometric polynomials at the nodal points. We show their connection to constant, curl-free, and divergence-free spaces with zero mean introduced in Eq. (2.1).

DEFINITION 3.11 (Fully discrete projections). *Let $\hat{\mathbf{F}}^{(i)}(\mathbf{n})$ for $i = 0, 1, 2$ and $\mathbf{n} \in \mathbb{Z}^d$ be the Fourier coefficients of integral kernels from Def. A.1. We define block diagonal matrices $\hat{\mathbf{G}}^{(0)}, \hat{\mathbf{G}}_0^{(i)}, \hat{\mathbf{G}}_I^{(i)} \in \mathbb{R}^{d \times d \times N \times N}$ for $i = 1, 2$ as*

$$\begin{aligned} (\hat{\mathbf{G}}^{(0)})_{\alpha\beta}^{mn} &= \hat{\mathbf{F}}_{\alpha\beta}^{(0)}(\mathbf{m})\delta_{mn} \\ (\hat{\mathbf{G}}_0^{(i)})_{\alpha\beta}^{mn} &= \begin{cases} \hat{\mathbf{F}}_{\alpha\beta}^{(i)}(\mathbf{m})\delta_{mn}, & \text{for } \mathbf{m} \in \mathbb{Z}_N \\ 0, & \text{for } \mathbf{m} \in \mathbb{Z}^d \setminus \mathbb{Z}_N \end{cases} \\ (\hat{\mathbf{G}}_I^{(i)})_{\alpha\beta}^{mn} &= \begin{cases} \hat{\mathbf{F}}_{\alpha\beta}^{(i)}(\mathbf{m})\delta_{mn}, & \text{for } \mathbf{m} \in \mathbb{Z}_N \\ \delta_{\alpha\beta}\delta_{mn}, & \text{for } \mathbf{m} \in \mathbb{Z}^d \setminus \mathbb{Z}_N \end{cases} \end{aligned}$$

where $\mathbf{m}, \mathbf{n} \in \mathbb{Z}_N^d$. Next, we define the matrices without hat as similarity transformation with matrix of DFT \mathbf{F} , i.e. $\mathbf{G}^{(0)} = \mathbf{F}^{-1}\hat{\mathbf{G}}^{(0)}\mathbf{F}$, $\mathbf{G}_0^{(i)} = \mathbf{F}^{-1}\hat{\mathbf{G}}_0^{(i)}\mathbf{F}$ and $\mathbf{G}_I^{(i)} = \mathbf{F}^{-1}\hat{\mathbf{G}}_I^{(i)}\mathbf{F}$.

LEMMA 3.12. *Matrices $\mathbf{G}^{(0)}, \mathbf{G}_0^{(i)}, \mathbf{G}_I^{(i)}$ for $i = 1, 2$, defined in Def. 3.11, are orthogonal projections.*

Proof. From the matrices in Def. 3.11, we deduce

$$\mathbf{I} = \mathbf{G}^{(0)} + \mathbf{G}_0^{(1)} + \mathbf{G}_I^{(2)}, \quad \mathbf{I} = \mathbf{G}^{(0)} + \mathbf{G}_I^{(1)} + \mathbf{G}_0^{(2)}.$$

with the help of the Fourier coefficients of integral kernels in Def. A.1. Moreover, a direct calculation shows that both of the triples $\hat{\mathbf{G}}^{(0)}, \hat{\mathbf{G}}_0^{(1)}, \hat{\mathbf{G}}_I^{(2)}$ and $\hat{\mathbf{G}}^{(0)}, \hat{\mathbf{G}}_I^{(1)}, \hat{\mathbf{G}}_0^{(2)}$ are mutually orthogonal projections, see also [17]. \square

DEFINITION 3.13 (Finite dimensional subspaces). *With the previously defined projections, we introduce the subspaces of space $\mathbb{R}^{d \times N}$*

$$\begin{aligned} \mathbb{U}_N &= \mathbf{G}^{(0)}\mathbb{R}^{d \times N} & \mathbb{E}_N &= \mathbf{G}_0^{(1)}\mathbb{R}^{d \times N} & \mathbb{J}_N &= \mathbf{G}_0^{(2)}\mathbb{R}^{d \times N} \\ & & \bar{\mathbb{E}}_N &= \mathbf{G}_I^{(1)}\mathbb{R}^{d \times N} & \bar{\mathbb{J}}_N &= \mathbf{G}_I^{(2)}\mathbb{R}^{d \times N} \end{aligned}$$

and their trigonometric relatives

$$\begin{aligned} \mathcal{U}_N &= \mathcal{I}_N^{-1}[\mathbb{U}_N] & \mathcal{E}_N &= \mathcal{I}_N^{-1}[\mathbb{E}_N] & \mathcal{J}_N &= \mathcal{I}_N^{-1}[\mathbb{J}_N] \\ \bar{\mathcal{E}}_N &= \mathcal{I}_N^{-1}[\bar{\mathbb{E}}_N] & \bar{\mathcal{J}}_N &= \mathcal{I}_N^{-1}[\bar{\mathbb{J}}_N] \end{aligned}$$

In Fig. 3.1, we introduce the relation diagram of subspaces based on the following

$$\begin{array}{ccccc} \mathcal{I}_N^{-1}[\mathbb{R}^{d \times N}] = \bar{\mathcal{T}}_N^d & \begin{array}{c} \text{only for odd } N \\ \subsetneq \end{array} & L_{\text{per}}^2(\mathcal{Y}; \mathbb{R}^d) & \begin{array}{c} \text{only for odd } N \\ \supsetneq \end{array} & \bar{\mathcal{T}}_N^d = \mathcal{I}_N^{-1}[\mathbb{R}^{d \times N}] \\ \parallel & & \parallel & & \parallel \\ \mathcal{I}_N^{-1}[\mathbb{U}_N] = \mathcal{U}_N & = & \mathcal{U} & = & \mathcal{U}_N = \mathcal{I}_N^{-1}[\mathbb{U}_N] \\ \oplus^\perp & & \oplus^\perp & & \oplus^\perp \\ \mathcal{I}_N^{-1}[\mathbb{E}_N] = \mathcal{E}_N & \begin{array}{c} \subsetneq \\ \supsetneq \end{array} & \mathcal{E} & \begin{array}{c} \text{only for odd } N \\ \supsetneq \end{array} & \bar{\mathcal{E}}_N = \mathcal{I}_N^{-1}[\bar{\mathbb{E}}_N] \\ \oplus^\perp & & \oplus^\perp & & \oplus^\perp \\ \mathcal{I}_N^{-1}[\bar{\mathbb{J}}_N] = \bar{\mathcal{J}}_N & \begin{array}{c} \text{only for odd } N \\ \subsetneq \end{array} & \mathcal{J} & \supsetneq & \mathcal{J}_N = \mathcal{I}_N^{-1}[\mathbb{J}_N] \end{array}$$

FIG. 3.1. The scheme of subspaces

two lemmas. The first one is a base for the fully discrete formulation of GAwNI in Sec. 3.3 and especially for a treatment with dual formulation Sec. 3.4.

LEMMA 3.14. *For the subspaces from Def. 3.13, the following three conditions hold:*

(i) *Space $\mathbb{R}^{d \times N}$ can be split into three mutually orthogonal subspaces*

$$\mathbb{R}^{d \times N} = \mathbb{U}_N \oplus^\perp \mathbb{E}_N \oplus^\perp \mathbb{J}_N, \quad \mathbb{R}^{d \times N} = \mathbb{U}_N \oplus^\perp \bar{\mathbb{E}}_N \oplus^\perp \mathbb{J}_N. \quad (3.4)$$

(ii) *The subspaces with tilde enlarge the original one, i.e. $\mathbb{E}_N \subseteq \bar{\mathbb{E}}_N$, $\mathbb{J}_N \subseteq \bar{\mathbb{J}}_N$.*

(iii) *If $N \in \mathbb{N}^d$ is odd (for all elements) it simplifies to $\mathbb{E}_N = \bar{\mathbb{E}}_N$ and $\mathbb{J}_N = \bar{\mathbb{J}}_N$.*

Proof. The proof is based on Lem. 3.12. Since both properties, constitution of identity and orthogonality, hold, all subspaces from the lemma are the subsets of set $\mathbb{R}^{d \times N}$ rather than set $\mathbb{C}^{d \times N}$; decomposition (3.4) thus holds. Finally, if N is odd, the index sets $\underline{\mathbb{Z}}_N^d$ and \mathbb{Z}_N^d coincide and the same stand for projections, $\mathbf{G}_0^{(1)} = \mathbf{G}_I^{(1)}$ and $\mathbf{G}_0^{(2)} = \mathbf{G}_I^{(2)}$, see Def. 3.11. \square

LEMMA 3.15. *The scheme stated in Fig. 3.1 holds true.*

Proof. First, the middle column is the Helmholtz decomposition, in our special case provided in Lem. A.2.

Next, we utilize previous Lem. 3.14 and isomorphism of operator \mathcal{I}_N in Lem. 3.8 which prove the rest of the columns.

The equality between spaces in the first two and the last two columns comes from Def. 3.13, hence the relations between subspaces in second, third, and fourth column have to be established. Relation $\tilde{\mathcal{F}}_N^d \subsetneq L_{\text{per}}^2(\mathcal{Y}; \mathbb{R}^{d \times N})$ holding only for odd N is discussed in Rem. 3.10. Equality $\mathcal{U}_N = \mathcal{U}$ is trivial as it contains only constant fields — the Fourier representation is exact.

Relation $\mathcal{E}_N \subseteq \mathcal{E}$ comes as a consequence of projection $\mathbf{G}_0^{(1)}$ that is deduced from the kernel $\hat{\mathbf{T}}^{(1)}$ of continuous projection defined in Appendix. A; defective frequencies $\mathbf{n} \in \underline{\mathbb{Z}}_N^d \setminus \mathbb{Z}_N^d$ are erased, see Rem. 3.10. The last inclusion $\mathcal{J}_N \subseteq \mathcal{J}$ is an analogue. \square

REMARK 3.16. *The previous proof yields the alternative characterization of the conforming subspaces: $\mathcal{E}_N = \mathcal{E} \cap \mathcal{F}_N^d$ and $\mathcal{J}_N = \mathcal{J} \cap \mathcal{F}_N^d$.*

3.3. FFT-based Finite Element Method. This section provides the overview of the FFT-based FEM. The theory for odd number of discretization points N is described in [31] including convergence results. Here, we extend the situation for non-odd number of discretization points N in a way to provide conforming approximations.

DEFINITION 3.17 (GAwNI). *Let material coefficients \mathbf{A} have continuous coefficients. Galerkin approximations with numerical integration (GAwNI) of the primal and the dual homogenization problems, Def. 2.3, state as: find discrete homogenized matrices $\mathbf{A}_{\text{eff},N}^{\text{FFTH}}, \mathbf{A}_{\text{eff},D,N}^{\text{FFTH}} \in \mathbb{R}^{d \times d}$ satisfying following relations for arbitrary macroscopic loads $\mathbf{E}, \mathbf{J} \in \mathbb{R}^d$*

$$(\mathbf{A}_{\text{eff},N}^{\text{FFTH}} \mathbf{E}, \mathbf{E})_{\mathbb{R}^d} = \inf_{\mathbf{e}_N \in \mathcal{E}_N} (\mathcal{Q}_N[\mathbf{A}(\mathbf{E} + \mathbf{e}_N)], \mathbf{E} + \mathbf{e}_N)_{L_{\text{per}}^2}, \quad (3.5a)$$

$$((\mathbf{A}_{\text{eff},D,N}^{\text{FFTH}})^{-1} \mathbf{J}, \mathbf{J})_{\mathbb{R}^d} = \inf_{\mathbf{j}_N \in \mathcal{J}_N} (\mathcal{Q}_N[\mathbf{A}^{-1}(\mathbf{J} + \mathbf{j}_N)], \mathbf{J} + \mathbf{j}_N)_{L_{\text{per}}^2}. \quad (3.5b)$$

REMARK 3.18. *The scalar products on the right-hand side in (3.5) are real valued and hence the discrete homogenized matrices are. Although, functions $\mathcal{Q}_N[\mathbf{A}^{-1}(\mathbf{E} +$*

\mathbf{e}_N), $\mathcal{Q}_N[\mathbf{A}^{-1}(\mathbf{J} + \mathbf{j}_N)] \in \tilde{\mathcal{T}}_N^d$ generally fail to be only real valued, see Rem. 3.10, the defective Fourier coefficients with frequencies $\mathbf{n} \in \mathbb{Z}_N^d \setminus \mathbb{Z}_N^d$ are eliminated by the space of test functions $\tilde{\mathcal{T}}_N^d$.

DEFINITION 3.19 (Fully discrete formulations of GAwNI). Find $\tilde{\mathbf{A}}_{\text{eff},N}^{\text{FFTH}}, \tilde{\mathbf{A}}_{\text{eff},D,N}^{\text{FFTH}} \in \mathbb{R}^{d \times d}$ satisfying following relations for arbitrary macroscopic loads $\mathbf{E}, \mathbf{J} \in \mathbb{R}^d$

$$(\tilde{\mathbf{A}}_{\text{eff},N}^{\text{FFTH}} \mathbf{E}, \mathbf{E})_{\mathbb{R}^d} = \frac{1}{|\mathbb{N}|_{\Pi}} \inf_{\mathbf{e}_N \in \mathbb{E}_N} (\mathbf{A}_N(\mathbf{E}_N + \mathbf{e}_N), \mathbf{E}_N + \mathbf{e}_N)_{\mathbb{R}^d \times \mathbb{N}} \quad (3.6a)$$

$$((\tilde{\mathbf{A}}_{\text{eff},D,N}^{\text{FFTH}})^{-1} \mathbf{J}, \mathbf{J})_{\mathbb{R}^d} = \frac{1}{|\mathbb{N}|_{\Pi}} \inf_{\mathbf{j}_N \in \mathbb{J}_N} (\mathbf{A}_N^{-1}(\mathbf{J}_N + \mathbf{j}_N), \mathbf{J}_N + \mathbf{j}_N)_{\mathbb{R}^d \times \mathbb{N}} \quad (3.6b)$$

where $\mathbf{E}_N = \mathcal{I}_N[\mathbf{E}] \in \mathbb{U}_N$, $\mathbf{J}_N = \mathcal{I}_N[\mathbf{J}] \in \mathbb{J}_N$, and $\mathbf{A}_N \in \mathbb{R}^{d \times d \times N \times N}$ with components $\mathbf{A}_{N,\alpha\beta}^{mn} = \mathbf{A}_{\alpha\beta}(\mathbf{x}_N^m) \delta_{mn}$ for α, β and $\mathbf{m}, \mathbf{n} \in \mathbb{Z}_N^d$.

REMARK 3.20. Minimizers $\tilde{\mathbf{e}}_N^{(\alpha)}, \tilde{\mathbf{j}}_N^{(\alpha)}$, and $\tilde{\mathbf{e}}_N^{(\alpha)}, \tilde{\mathbf{j}}_N^{(\alpha)}$ of both formulations corresponding to unitary macroscopic fields $\tilde{\mathbf{e}}_N$ are called discrete unitary minimizers. The existence and uniqueness is provided in following lemma.

LEMMA 3.21. The homogenized matrices of both previous formulations, Def. 3.17 and 3.19, coincide $\mathbf{A}_{\text{eff},N}^{\text{FFTH}} = \tilde{\mathbf{A}}_{\text{eff},N}^{\text{FFTH}}, \mathbf{A}_{\text{eff},D,N}^{\text{FFTH}} = \tilde{\mathbf{A}}_{\text{eff},D,N}^{\text{FFTH}}$. Moreover, discrete minimizers $\tilde{\mathbf{e}}_N^{(\alpha)}, \tilde{\mathbf{j}}_N^{(\alpha)}$ and $\tilde{\mathbf{e}}_N^{(\alpha)}, \tilde{\mathbf{j}}_N^{(\alpha)}$ of both formulations exist, are unique, and are connected to each other $\mathcal{I}_N[\tilde{\mathbf{e}}_N] = \tilde{\mathbf{e}}_N, \mathcal{I}_N[\tilde{\mathbf{j}}_N] = \tilde{\mathbf{j}}_N$.

Proof. Although, the proof is a generalization of that in [31], it follows the same ideas: the linearity of scalar product with the definition of interpolation operator \mathcal{Q}_N and the property of space basis functions $\varphi_{N,\mathbf{k}}$, and the second property in (3.2b). For completeness, we provide calculation

$$\begin{aligned} & (\mathcal{Q}_N[\mathbf{A}(\mathbf{E} + \tilde{\mathbf{e}}_N)], \mathbf{E} + \tilde{\mathbf{e}}_N)_{L_{\text{per}}^2} = \\ & = \left(\sum_{\mathbf{m} \in \mathbb{Z}_N^d} \mathbf{A}(\mathbf{x}_N^{\mathbf{m}}) [\mathbf{E} + \tilde{\mathbf{e}}_N(\mathbf{x}_N^{\mathbf{m}})] \varphi_{N,\mathbf{m}}, \sum_{\mathbf{n} \in \mathbb{Z}_N^d} [\mathbf{E} + \tilde{\mathbf{e}}_N(\mathbf{x}_N^{\mathbf{n}})] \varphi_{N,\mathbf{n}} \right)_{L_{\text{per}}^2} \\ & = \sum_{\mathbf{m} \in \mathbb{Z}_N^d} \sum_{\mathbf{n} \in \mathbb{Z}_N^d} \mathbf{A}(\mathbf{x}_N^{\mathbf{m}}) [\mathbf{E} + \tilde{\mathbf{e}}_N(\mathbf{x}_N^{\mathbf{m}})] [\mathbf{E} + \tilde{\mathbf{e}}_N(\mathbf{x}_N^{\mathbf{n}})] (\varphi_{N,\mathbf{m}}, \varphi_{N,\mathbf{n}})_{L_{\text{per}}^2} \\ & = \sum_{\mathbf{m} \in \mathbb{Z}_N^d} \sum_{\mathbf{n} \in \mathbb{Z}_N^d} \mathbf{A}_N^{\mathbf{m}} [\mathbf{E}_N^{\mathbf{m}} + \tilde{\mathbf{e}}_N^{\mathbf{m}}] [\mathbf{E}_N^{\mathbf{n}} + \tilde{\mathbf{e}}_N^{\mathbf{n}}] \frac{\delta_{mn}}{|\mathbb{N}|_{\Pi}} = \frac{1}{|\mathbb{N}|_{\Pi}} (\mathbf{A}_N[\mathbf{E}_N + \tilde{\mathbf{e}}_N], \mathbf{E}_N + \tilde{\mathbf{e}}_N)_{\mathbb{R}^d \times \mathbb{N}} \end{aligned}$$

The existence and the uniqueness of minimizers are provided due to symmetric and positive definite matrix \mathbf{A}_N . \square

REMARK 3.22 (Solution of GAwNI by Conjugate gradients). According to [33], the minimizers of fully discrete formulations are equivalent to weak formulations, similar to those in Def. 2.5. Moreover, the solutions can be found by the means of Conjugate gradients applied to linear systems $\mathbf{C}\mathbf{x} = \mathbf{b}$ and $\mathbf{C}'\mathbf{x}' = \mathbf{b}'$ defined for particular α as

$$\underbrace{\mathbf{G}_0^{(1)} \mathbf{A}_N}_{\mathbf{C}} \underbrace{\mathbf{e}_N^{(\alpha)}}_{\mathbf{x}} = \underbrace{-\mathbf{G}_0^{(1)} \mathbf{A}_N \mathbf{E}_N^{(\alpha)}}_{\mathbf{b}} \quad \underbrace{\mathbf{G}_0^{(2)} \mathbf{A}_N^{-1}}_{\mathbf{C}'} \underbrace{\mathbf{j}_N^{(\alpha)}}_{\mathbf{x}'} = \underbrace{-\mathbf{G}_0^{(2)} \mathbf{A}_N^{-1} \mathbf{J}_N^{(\alpha)}}_{\mathbf{b}'}$$

for initial approximations $\mathbf{x}_0, \mathbf{x}'_0$ that have to belong to the appropriate subspaces, $\mathbb{E}_N, \mathbb{J}_N$ resp. Easily, Conjugate gradients minimizes the quadratic forms in the fully discrete formulations 3.6.

LEMMA 3.23 (Convergence of discrete minimizers). *Let material coefficients \mathbf{A} be continuous and sufficiently regular in order to minimizers $\mathbf{e}^{(\alpha)}$, $\mathbf{j}^{(\alpha)}$, Def. 2.5, be sufficiently regular — particularly having all weak partial derivatives up to order μ in the space $L^2_{\text{per}}(\mathcal{Y}; \mathbb{R}^d)$. Then the sequence of discrete minimizers $\mathbf{e}_N^{(\alpha)}$, $\mathbf{j}_N^{(\alpha)}$, Def. 3.17, converge to the minimizers, i.e.*

$$\|\mathbf{e} - \mathbf{e}_N\|_{L^2_{\text{per}}} \leq C \left(\min_{\alpha} N_{\alpha} \right)^{-\mu} \rightarrow 0 \quad \text{for } \min_{\alpha} N_{\alpha} \rightarrow \infty$$

where constant C is independent of N ; it depends only on the material coefficients, its positive definiteness, norm, and regularity.

3.4. Connection of primal and dual formulations. This section is dedicated to the connection between the primal and the dual formulations in the fully discrete setting (3.6). General theory for dual problems can be found for example in [7]. We start with the statement of general lemma.

LEMMA 3.24 (Perturbation duality theorem — page 54 in [7]). *Let \mathbb{V} and \mathbb{X} be Hilbert spaces and $\Phi : \mathbb{V} \times \mathbb{X} \rightarrow \mathbb{R}$ be a continuous functional, convex, coercive for the first variable ($\lim_{\mathbf{v} \in \mathbb{V}, \|\mathbf{v}\| \rightarrow \infty} \Phi(\mathbf{v}, 0) = \infty$), and satisfying following property: there exist $\mathbf{v}_0 \in \mathbb{V}$ such that $\Phi(\mathbf{v}_0, \cdot)$ is finite and continuous around $0 \in \mathbb{V}$. Then the primal and the dual problems*

$$\min_{\mathbf{v} \in \mathbb{V}} \Phi(\mathbf{v}, 0) \qquad \max_{\mathbf{x}^* \in \mathbb{X}} -\Phi^*(0; \mathbf{x}^*) \qquad (3.7)$$

have the solutions and the extremal values are equal to one another. Here Φ^* is the usual Fenchel conjugate function defined on $\mathbb{V} \times \mathbb{X}$ as

$$\Phi^*(\mathbf{v}^*; \mathbf{x}^*) := \sup_{\mathbf{v} \in \mathbb{V}, \mathbf{x} \in \mathbb{X}} [(\mathbf{v}^*, \mathbf{v})_{\mathbb{V}} + (\mathbf{x}^*, \mathbf{x})_{\mathbb{X}} - \Phi(\mathbf{v}, \mathbf{x})]. \qquad (3.8)$$

The following lemma is an application of previous lemma on the homogenization problem in fully discrete setting; it is sufficiently general to apply for arbitrary number of discretization points. Both the lemma and the proof are analogy to Proposition 2.2 and Corollary 2.3 in [6] that are stated for continuous formulations.

LEMMA 3.25 (Transformation to dual formulation). *Let $\mathbf{A}_N \in \mathbb{R}^{d \times d \times N \times N}$ be symmetric positive definite matrix, \mathbb{V}_N be a proper nontrivial subspace of $\mathbb{U}_N^{\perp} = \{\mathbf{v} \in \mathbb{R}^{d \times N} : \mathbf{v} \cdot \mathbf{u} = 0 \text{ for all } \mathbf{u} \in \mathbb{U}_N\}$, and the following primal problem be set: find $\tilde{\mathbf{A}}_{\text{eff}} \in \mathbb{R}^d$ satisfying the following relation for arbitrary macroscopic load $\mathbf{E} \in \mathbb{R}^d$*

$$(\tilde{\mathbf{A}}_{\text{eff}} \mathbf{E}, \mathbf{E})_{\mathbb{R}^d} = \frac{1}{|\mathbf{N}|_{\Pi}} \inf_{\mathbf{e}_N \in \mathbb{V}_N} (\mathbf{A}_N(\mathbf{E}_N + \mathbf{e}_N), \mathbf{E}_N + \mathbf{e}_N)_{\mathbb{R}^{d \times N}}$$

where $\mathbf{E}_N := \mathcal{I}_N[\mathbf{E}] \in \mathbb{U}_N$. Then the problem is equivalent to the dual problem: find $\tilde{\mathbf{A}}_{\text{eff}} \in \mathbb{R}^d$ satisfying the following relation for arbitrary macroscopic load $\mathbf{J} \in \mathbb{R}^d$

$$(\tilde{\mathbf{A}}_{\text{eff}}^{-1} \mathbf{J}, \mathbf{J})_{\mathbb{R}^d} = \frac{1}{|\mathbf{N}|_{\Pi}} \inf_{\mathbf{j}_N \in \mathbb{W}_N} (\mathbf{A}_N^{-1}(\mathbf{J}_N + \mathbf{j}_N), \mathbf{J}_N + \mathbf{j}_N)_{\mathbb{R}^{d \times N}}$$

where $\mathbf{J}_N = \mathcal{I}_N[\mathbf{J}] \in \mathbb{U}_N$ and $\mathbb{W}_N = (\mathbb{U}_N \oplus \mathbb{V}_N)^{\perp}$. Moreover, between macroscopic fields \mathbf{E} , \mathbf{J} and between minimizers $\tilde{\mathbf{e}}_N^{(\mathbf{E})}$, $\tilde{\mathbf{j}}_N^{(\mathbf{J})}$ of the primal and the dual formulations, the following relations hold

$$\mathbf{J} = \tilde{\mathbf{A}}_{\text{eff}} \mathbf{E}, \qquad \mathbf{J}_N + \tilde{\mathbf{j}}_N^{(\mathbf{J})} = \mathbf{A}_N[\mathbf{E}_N + \tilde{\mathbf{e}}_N^{(\mathbf{E})}]. \qquad (3.9)$$

Proof. First, we define a function $\Phi : \mathbb{V}_N \times \mathbb{R}^{d \times N} \rightarrow \mathbb{R}$ as

$$\Phi(\mathbf{e}_N, \mathbf{x}) := \frac{1}{2} (\mathbf{A}_N [\mathbf{E}_N + \mathbf{e}_N + \mathbf{x}], [\mathbf{E}_N + \mathbf{e}_N + \mathbf{x}])_{\mathbb{R}^{d \times N}}.$$

Since matrix \mathbf{A}_N is symmetric positive definite and the underlying spaces are finite dimensional, the assumptions of previous Lem. 3.25 are satisfied. Hence the primal formulation reformulated with Φ to

$$(\mathbf{A}_{\text{eff}} \mathbf{E}, \mathbf{E})_{\mathbb{R}^d} = \frac{2}{|\mathbf{N}|_{\Pi}} \inf_{\mathbf{e}_N \in \mathbb{V}_N} \Phi(\mathbf{e}_N, \mathbf{0}),$$

is equivalent to the dual formulation

$$(\mathbf{A}_{\text{eff}} \mathbf{E}, \mathbf{E})_{\mathbb{R}^d} = \frac{2}{|\mathbf{N}|_{\Pi}} \sup_{\mathbf{x}^* \in \mathbb{R}^{d \times N}} -\Phi^*(\mathbf{0}, \mathbf{x}^*)$$

where $\mathbf{0} \in \mathbb{R}^{d \times N}$ is a vector with all components equal to zero and $\Phi^* : \mathbb{V}_N \times \mathbb{R}^{d \times N} \rightarrow \mathbb{R}$ is the Fenchel conjugate function, see Eq. (3.8).

Using substitution $\mathbf{x}' = \mathbf{E}_N + \mathbf{e}_N + \mathbf{x}$ where \mathbf{x}' covers the whole space $\mathbb{R}^{d \times N}$, we deduce

$$\begin{aligned} (\mathbf{A}_{\text{eff}} \mathbf{E}, \mathbf{E})_{\mathbb{R}^d} &= \frac{2}{|\mathbf{N}|_{\Pi}} \sup_{\mathbf{x}^* \in \mathbb{R}^{d \times N}} -\Phi^*(\mathbf{0}, \mathbf{x}^*) \\ &= \frac{2}{|\mathbf{N}|_{\Pi}} \sup_{\mathbf{x}^* \in \mathbb{R}^{d \times N}} \left[- \sup_{\mathbf{e}_N \in \mathbb{V}_N, \mathbf{x}' \in \mathbb{R}^{d \times N}} \left((\mathbf{x}^*, \mathbf{x}' - \mathbf{E}_N - \mathbf{e}_N)_{\mathbb{R}^{d \times N}} - \frac{1}{2} (\mathbf{A}_N \mathbf{x}', \mathbf{x}')_{\mathbb{R}^{d \times N}} \right) \right] \\ &= \frac{2}{|\mathbf{N}|_{\Pi}} \sup_{\mathbf{x}^* \in \mathbb{R}^{d \times N}} \left[(\mathbf{x}^*, \mathbf{E}_N)_{\mathbb{R}^{d \times N}} + \inf_{\mathbf{e}_N \in \mathbb{V}_N} (\mathbf{x}^*, \mathbf{e}_N)_{\mathbb{R}^{d \times N}} \right. \\ &\quad \left. - \sup_{\mathbf{x}' \in \mathbb{R}^{d \times N}} \left((\mathbf{x}^*, \mathbf{x}')_{\mathbb{R}^{d \times N}} - \frac{1}{2} (\mathbf{A}_N \mathbf{x}', \mathbf{x}')_{\mathbb{R}^{d \times N}} \right) \right] \end{aligned}$$

We focus on the supreme in the last equation where the equilibrium point satisfies $\mathbf{A}_N \mathbf{x}' = \mathbf{x}^*$. Since \mathbf{A}_N is symmetric positive definite, we have $\mathbf{x}' = \mathbf{A}_N^{-1} \mathbf{x}^*$. Therefore, the inner supremum is simplified to

$$\sup_{\mathbf{x}' \in \mathbb{R}^{d \times N}} \left((\mathbf{x}^*, \mathbf{x}')_{\mathbb{R}^{d \times N}} - \frac{1}{2} (\mathbf{A}_N \mathbf{x}', \mathbf{x}')_{\mathbb{R}^{d \times N}} \right) = \frac{1}{2} (\mathbf{A}_N^{-1} \mathbf{x}^*, \mathbf{x}^*)_{\mathbb{R}^{d \times N}}.$$

The inner infimum equals to minus the indicator function of V_N^\perp , explicitly

$$\inf_{\mathbf{e}_N \in \mathbb{V}_N} (\mathbf{x}^*, \mathbf{e}_N)_{\mathbb{R}^{d \times N}} = \begin{cases} 0, & \text{for } \mathbf{x}^* \in V_N^\perp \\ -\infty, & \text{otherwise.} \end{cases}$$

Hence, we can omit the inf term while restrict the supremum to space V_N^\perp . It leads to

$$(\mathbf{A}_{\text{eff}} \mathbf{E}, \mathbf{E})_{\mathbb{R}^d} = \frac{2}{|\mathbf{N}|_{\Pi}} \sup_{\mathbf{x}^* \in V_N^\perp} \left[(\mathbf{x}^*, \mathbf{E}_N)_{\mathbb{R}^{d \times N}} - \frac{1}{2} (\mathbf{A}_N^{-1} \mathbf{x}^*, \mathbf{x}^*)_{\mathbb{R}^{d \times N}} \right].$$

Since the following inclusion holds $\mathbb{U}_N \subsetneq V_N^\perp$, space V_N^\perp can be decomposed $V_N^\perp = \mathbb{U}_N \oplus^\perp \mathbb{W}_N$ where $\mathbb{W}_N = \{\mathbf{w} \in V_N^\perp : \mathbf{w} \cdot \mathbf{u} = 0 \text{ for all } \mathbf{u} \in \mathbb{U}_N\}$. Hence, we state

$$(\mathbf{A}_{\text{eff}} \mathbf{E}, \mathbf{E})_{\mathbb{R}^d} = \frac{2}{|\mathbf{N}|_{\Pi}} \sup_{\mathbf{J}_N \in \mathbb{U}_N} \left[(\mathbf{J}_N, \mathbf{E}_N)_{\mathbb{R}^{d \times N}} - \inf_{\mathbf{j}_N \in \mathbb{W}_N} \frac{1}{2} (\mathbf{A}_N^{-1} (\mathbf{J}_N + \mathbf{j}_N), \mathbf{J}_N + \mathbf{j}_N)_{\mathbb{R}^{d \times N}} \right] \blacksquare$$

where $\mathbf{x}^* = \mathbf{J}_N + \mathbf{j}_N$ such that $\mathbf{J}_N \in \mathbb{U}_N$, $\mathbf{j}_N \in \mathbb{W}_N$. Then we define matrix $\mathbf{B}_{\text{eff}} \in \mathbb{R}^{d \times d}$ satisfying for arbitrary $\mathbf{J} \in \mathbb{R}^d$

$$(\mathbf{B}_{\text{eff}} \mathbf{J}, \mathbf{J})_{\mathbb{R}^d} = \frac{1}{|\mathbb{N}|_{\Pi}} \inf_{\mathbf{j}_N \in \mathbb{V}_N} (\mathbf{A}_N^{-1}(\mathbf{J}_N + \mathbf{j}_N), \mathbf{J}_N + \mathbf{j}_N)_{\mathbb{R}^{d \times N}}$$

where $\mathbf{J}_N = \mathcal{I}_N[\mathbf{J}]$.

Now, we show $\mathbf{B}_{\text{eff}} = \mathbf{A}_{\text{eff}}^{-1}$, the first identity in Eq. (3.9). With obvious identity

$$\frac{1}{|\mathbb{N}|_{\Pi}} (\mathbf{J}_N, \mathbf{E}_N)_{\mathbb{R}^{d \times N}} = (\mathbf{J}, \mathbf{E})_{\mathbb{R}^d},$$

the dual problem becomes

$$(\mathbf{A}_{\text{eff}} \mathbf{E}, \mathbf{E})_{\mathbb{R}^d} = \sup_{\mathbf{J} \in \mathbb{R}^d} [2(\mathbf{J}, \mathbf{E})_{\mathbb{R}^d} - (\mathbf{B}_{\text{eff}} \mathbf{J}, \mathbf{J})_{\mathbb{R}^d}]$$

and comply with the equilibrium state $\mathbf{B}_{\text{eff}} \mathbf{J} = \mathbf{E}$. The matrix \mathbf{B}_{eff} is symmetric positive definite as \mathbf{A}_N^{-1} is (see e.g. [4]), whence substitution of $\mathbf{J} = \mathbf{B}_{\text{eff}}^{-1} \mathbf{E}$ into the dual formulation leads to required identity.

The second identity in Eq. (3.9) follows from equations obtained during the proof, particularly $\mathbf{A}_N(\mathbf{E}_N + \mathbf{v} + \mathbf{x}) = \mathbf{x}^* = \mathbf{J}_N + \mathbf{j}_N$, and the fact that the primal formulation is obtained for $\mathbf{x} = 0$. \square

COROLLARY 3.26 (Special case of the odd number of discretization points N).

Let $N \in \mathbb{N}^d$ be odd, and the fully discrete formulations 3.6 be defined. Then:

- (i) Both the primal and the dual homogenized matrices coincide $\mathbf{A}_{\text{eff},N}^{\text{FFTH}} = \mathbf{A}_{\text{eff},D,N}^{\text{FFTH}}$.
- (ii) Primal and dual discrete minimizers $\tilde{\mathbf{e}}_N^{(\alpha)}$, $\tilde{\mathbf{j}}_N^{(\alpha)}$ are related

$$\boldsymbol{\epsilon}_\beta + \tilde{\mathbf{j}}_N^{(\beta)} = \mathbf{A}_N \sum_{\alpha} E_{\alpha}(\boldsymbol{\epsilon}_{\alpha} + \tilde{\mathbf{e}}_N^{(\alpha)}) \quad (3.10)$$

where $\mathbf{E} = (\mathbf{A}_{\text{eff},N}^{\text{FFTH}})^{-1} \boldsymbol{\epsilon}_\beta$.

Proof. The proof is the consequence of Lem. 3.25 for $\mathbb{V}_N = \mathbb{E}_N$, $\mathbb{W}_N = \mathbb{J}_N$, and decomposition $\mathbb{R}^{d \times N} = \mathbb{U}_N \oplus \mathbb{E}_N \oplus \mathbb{J}_N$ stated in Lem. 3.14. \square

COROLLARY 3.27 (General case of the arbitrary number of discretization points N). Let we have following homogenization problems: find $\tilde{\mathbf{A}}_{\text{eff},N}^{\text{FFTH}}$, $\tilde{\mathbf{A}}_{\text{eff},D,N}^{\text{FFTH}} \in \mathbb{R}^{d \times d}$ such that

$$((\tilde{\mathbf{A}}_{\text{eff},N}^{\text{FFTH}})^{-1} \mathbf{J}, \mathbf{J})_{\mathbb{R}^d} = \frac{1}{|\mathbb{N}|_{\Pi}} \inf_{\mathbf{j}_N \in \mathbb{J}_N} (\mathbf{A}_N^{-1}(\mathbf{J}_N + \mathbf{j}_N), \mathbf{J}_N + \mathbf{j}_N)_{\mathbb{R}^{d \times N}} \quad (3.11a)$$

$$(\tilde{\mathbf{A}}_{\text{eff},D,N}^{\text{FFTH}} \mathbf{E}, \mathbf{E})_{\mathbb{R}^d} = \frac{1}{|\mathbb{N}|_{\Pi}} \inf_{\mathbf{e}_N \in \mathbb{E}_N} (\mathbf{A}_N(\mathbf{E}_N + \mathbf{e}_N), \mathbf{E}_N + \mathbf{e}_N)_{\mathbb{R}^d} \quad (3.11b)$$

and let $\tilde{\mathbf{e}}_N^{(\alpha)}$ and $\tilde{\mathbf{j}}_N^{(\alpha)}$ be their approximate unitary minimizers. Then the following holds:

- (i) The fully discrete primal and dual formulations, Eq. 3.6a and 3.6b in Def. 3.19 are equivalent to those here in Eq. (3.11a) and (3.11b) resp. in the sense the homogenized matrices coincide $\mathbf{A}_{\text{eff},N}^{\text{FFTH}} = \tilde{\mathbf{A}}_{\text{eff},N}^{\text{FFTH}}$ and $\mathbf{A}_{\text{eff},D,N}^{\text{FFTH}} = \tilde{\mathbf{A}}_{\text{eff},D,N}^{\text{FFTH}}$.

(ii) The discrete unitary minimizers $\tilde{\mathbf{e}}_N^{(\beta)}$, $\tilde{\mathbf{j}}_N^{(\alpha)}$ of Eq. 3.6 can be expressed as a linear combination of these minimizers $\bar{\mathbf{e}}_N^{(\alpha)}$, $\bar{\mathbf{j}}_N^{(\alpha)}$ of Eq. (3.11) as

$$\boldsymbol{\epsilon}_\beta + \tilde{\mathbf{e}}_N^{(\beta)} = \mathbf{A}_N^{-1} \sum_{\alpha} J_{\alpha}(\boldsymbol{\epsilon}_\alpha + \bar{\mathbf{j}}_N^{(\alpha)}) \in \mathbb{E}_N \quad (3.12a)$$

$$\boldsymbol{\epsilon}_\beta + \tilde{\mathbf{j}}_N^{(\beta)} = \mathbf{A}_N \sum_{\alpha} E_{\alpha}(\boldsymbol{\epsilon}_\alpha + \bar{\mathbf{e}}_N^{(\alpha)}) \in \mathbb{J}_N \quad (3.12b)$$

where the macroscopic quantities are set to $\mathbf{E} := (\mathbf{A}_{\text{eff,D},N}^{\text{FFTH}})^{-1} \boldsymbol{\epsilon}_\beta$ and $\mathbf{J} := \mathbf{A}_{\text{eff},N}^{\text{FFTH}} \boldsymbol{\epsilon}_\beta$.

(iii) The primal and the dual homogenized matrices can be compared

$$\mathbf{A}_{\text{eff,D},N}^{\text{FFTH}} \preceq \mathbf{A}_{\text{eff},N}^{\text{FFTH}}. \quad (3.13)$$

Proof. The proof is mainly the consequence of Lem. 3.25 and 3.14: for the primal fully discrete formulation with $\mathbb{V}_N = \mathbb{E}_N$, $\mathbb{W}_N = \bar{\mathbb{J}}_N$ and for the dual one with $\mathbb{V}_N = \mathbb{J}_N$, $\mathbb{W}_N = \bar{\mathbb{E}}_N$ — see Fig. 3.1 with the scheme of the subspaces. The inequality in 3.13 comes from relation $\mathbb{E}_N \subseteq \bar{\mathbb{E}}_N$, see Lem. 3.14, and a relation

$$\frac{1}{|\mathbb{N}|_{\Pi}} \inf_{\mathbf{e}_N \in \bar{\mathbb{E}}_N} (\mathbf{A}_N(\mathbf{E}_N + \mathbf{e}_N), \mathbf{E}_N + \mathbf{e}_N)_{\mathbb{R}^{d \times N}} \leq \frac{1}{|\mathbb{N}|_{\Pi}} \inf_{\mathbf{e}_N \in \mathbb{E}_N} (\mathbf{A}_N(\mathbf{E}_N + \mathbf{e}_N), \mathbf{E}_N + \mathbf{e}_N)_{\mathbb{R}^{d \times N}}$$

holding for arbitrary fixed $\mathbf{E}_N \in \mathbb{U}_N$. \square

REMARK 3.28. The effective matrices $\mathbf{A}_{\text{eff},N}^{\text{FFTH}}$ and $\mathbf{A}_{\text{eff,D},N}^{\text{FFTH}}$ can be compared with relation \preceq , in the sense of quadratic norms, to none of the matrices \mathbf{A}_{eff} , $\underline{\mathbf{A}}_{\text{eff},N}$, and $\bar{\mathbf{A}}_{\text{eff},N}$ as is numerically shown in Sec. 4.1.

REMARK 3.29. In engineering literature relating FFT-based homogenization, e.g. [19, 20, 16], the criterion for the numerical convergence of the primal approximate minimizers \mathbf{e}_N is based on an equilibrium condition of the dual fields $\mathbf{A}_N \mathbf{e}_N$ controlling to be divergence-free. However, this criterion is reasonable only for the odd-number of discretization points, cf. Eq. (3.10), as observed in [20]; they also offer a remedy that exactly corresponds to the formulation in Eq. (3.11a) — the dual fields (3.12a) are then, if a convergence is reached, in appropriate subspace \mathcal{J}_N , the space of divergence-free fields.

3.5. Calculation of upper-lower bounds. The calculation of the upper-lower bounds of the homogenized matrix consists of the integral evaluation of type $(\mathbf{A} \mathbf{e}_N^{(\alpha)}, \mathbf{e}_N^{(\beta)})_{L_{\text{per}}^2}$ occurring in Def. 2.8. Generally, the integral cannot be evaluated in a closed form because of non-specific material coefficients. The idea is to adjust material coefficients to calculate the integrals accurately and efficiently and simultaneously keep the upper-lower bounds structure.

For an easier orientation among various homogenized matrices, we refer to their scheme in Fig. 3.2. The matrices \mathbf{A}_{eff} , $\underline{\mathbf{A}}_{\text{eff},N}$, $\bar{\mathbf{A}}_{\text{eff},N}$, $\mathbf{A}_{\text{eff},N}$, and \mathbf{D}_N introduced in Def. 2.3, 2.8, 2.11 are in no relation to matrices $\mathbf{A}_{\text{eff},N}^{\text{FFTH}}$, $\mathbf{A}_{\text{eff,D},N}^{\text{FFTH}}$ from Def. 3.17, see Rem. 3.28.

In this section, we introduce approximations of upper-lower bounds $\tilde{\mathbf{A}}_{\text{eff,D},N}^{\text{lin},M}$, $\tilde{\mathbf{A}}_{\text{eff},N}^{\text{lin},M}$ based on piecewise bilinear material coefficients and homogenized matrices $\underline{\underline{\mathbf{A}}}_{\text{eff},N}^{\text{lin},M}$, $\underline{\mathbf{A}}_{\text{eff},N}^{\text{lin},M}$ based on piecewise constant material coefficients defined in a way to guaranty bounds. All four introduced matrices can be calculated efficiently by FFT algorithm, see Lem. 3.31 and Lem. 3.32.

$$\begin{array}{ccccccccccc}
& & & \tilde{\mathbf{A}}_{\text{eff},D,N}^{\text{lin},M} & & \tilde{\mathbf{A}}_{\text{eff},N}^{\text{lin},M} & & & & & \\
& & & \Downarrow & & \Downarrow & & & & & \\
0 & \preceq & \underline{\mathbf{A}}_{\text{eff},N}^{\text{con},M} & \preceq & \underline{\mathbf{A}}_{\text{eff},N} & \preceq & \mathbf{A}_{\text{eff}} & \preceq & \overline{\mathbf{A}}_{\text{eff},N} & \preceq & \overline{\overline{\mathbf{A}}}_{\text{eff},N}^{\text{con},M} \\
& & & \parallel & & \parallel & & & & & \\
& & & \mathbf{A}_{\text{eff},N} - \mathbf{D}_N & \preceq & \mathbf{A}_{\text{eff},N} & \preceq & \mathbf{A}_{\text{eff},N} + \mathbf{D}_N & & & \\
& & & \Downarrow & & \Downarrow & & & & & \\
& & & \mathbf{A}_{\text{eff},D,N}^{\text{FFTH}} & \underset{\text{equality if } N \text{ is odd}}{\preceq} & \mathbf{A}_{\text{eff},N}^{\text{FFTH}} & & & & &
\end{array}$$

FIG. 3.2. The overview of homogenized material bounds

LEMMA 3.30 (Sufficient condition for the upper-lower bounds). *Let $\mathbf{A} \in L_{\text{per}}^{\infty}(\mathcal{Y}; \mathbb{R}_{\text{spd}}^{d \times d})$ be material coefficients and $\overline{\mathbf{A}}, \underline{\mathbf{A}} \in L_{\text{per}}^{\infty}(\mathcal{Y}; \mathbb{R}^{d \times d})$ its upper and lower approximations satisfying*

$$\underline{\mathbf{A}}(\mathbf{x}) \preceq \mathbf{A}(\mathbf{x}) \preceq \overline{\mathbf{A}}(\mathbf{x}), \quad \text{for almost all } \mathbf{x} \in \mathcal{Y}. \quad (3.14)$$

Let $\tilde{\mathbf{e}}_N^{(\alpha)} \in \mathcal{E}_N$ and $\tilde{\mathbf{j}}_N^{(\alpha)} \in \mathcal{J}_N$ be unitary minimizers for material coefficients \mathbf{A} , cf. Def. 2.5. Then matrices $\overline{\overline{\mathbf{A}}}_{\text{eff}}, \underline{\underline{\mathbf{A}}}_{\text{eff}} \in \mathbb{R}^{d \times d}$, defined as

$$(\overline{\overline{\mathbf{A}}}_{\text{eff},N})_{\alpha\beta} = (\overline{\mathbf{A}}(\boldsymbol{\epsilon}_{\beta} + \tilde{\mathbf{e}}_N^{(\beta)}), \boldsymbol{\epsilon}_{\alpha} + \tilde{\mathbf{e}}_N^{(\alpha)})_{L_{\text{per}}^2}, \quad (3.15a)$$

$$(\underline{\underline{\mathbf{A}}}_{\text{eff},N}^{-1})_{\alpha\beta} = (\underline{\mathbf{A}}^{-1}(\boldsymbol{\epsilon}_{\beta} + \tilde{\mathbf{j}}_N^{(\beta)}), \boldsymbol{\epsilon}_{\alpha} + \tilde{\mathbf{j}}_N^{(\alpha)})_{L_{\text{per}}^2}, \quad (3.15b)$$

comply with the upper-lower bound structure, i.e.

$$\underline{\underline{\mathbf{A}}}_{\text{eff},N} \preceq \underline{\mathbf{A}}_{\text{eff},N} \preceq \mathbf{A}_{\text{eff}} \preceq \overline{\mathbf{A}}_{\text{eff},N} \preceq \overline{\overline{\mathbf{A}}}_{\text{eff},N}.$$

Proof. The inner inequalities $\underline{\mathbf{A}}_{\text{eff},N} \preceq \mathbf{A}_{\text{eff}} \preceq \overline{\mathbf{A}}_{\text{eff},N}$ are already proven in Lem. 2.10; the rest easily arise from assumed inequality (3.14) that is kept under integration. \square

Next lemma provides a way for the calculation of the homogenized matrices by the FFT routine, Lem. 3.32. It requires the material coefficients to be expressed as a linear combination of some basis functions concentrated on the set of nodal points.

LEMMA 3.31 (Calculation of homogenized matrices). *Let $\mathbf{u}_N, \mathbf{v}_N \in \mathcal{T}_N^d$ be trigonometric polynomials and $\mathbf{A}_M \in L_{\text{per}}^{\infty}(\mathcal{Y}; \mathbb{R}^{d \times d})$ for $M \in \mathbb{N}^d$ be function explicitly expressed as*

$$\mathbf{A}_M(\mathbf{x}) = \sum_{\mathbf{n} \in \mathbb{Z}_M^d} \psi(\mathbf{x} + \mathbf{x}_M^{\mathbf{n}}) \mathbf{A}_M^{\mathbf{n}}, \quad \mathbf{x} \in \mathcal{Y}$$

where $\psi \in L_{\text{per}}^{\infty}(\mathcal{Y}; \mathbb{R})$ is some basis function and $\mathbf{A}_M \in \mathbb{R}^{d \times d \times M}$. Then the integrals of the type occurring in Eq. (3.15) can be calculated as

$$(\mathbf{A}_M \mathbf{u}_N, \mathbf{v}_N)_{L_{\text{per}}^2} = \frac{1}{|\mathcal{Y}|_d} \sum_{\alpha, \beta} \sum_{\mathbf{m} \in \mathbb{Z}_{2N}^d} w(\mathbf{m}) \hat{u}_{N, \alpha, \beta}^{\mathbf{m}} \hat{A}_{\alpha\beta}^{\mathbf{m}} \quad (3.16)$$

where integration weight $w(\mathbf{m})$ is defined as $w(\mathbf{m}) := \int_{\mathcal{Y}} \psi(\mathbf{x}) \varphi_{\mathbf{m}}(\mathbf{x})$ and factors

$\widehat{u}_{N,\beta,\alpha}^m, \widehat{A}_{\alpha\beta}^m$ are defined as

$$\widehat{u}_{N,\beta,\alpha}^m = \frac{1}{2|\mathbb{N}|_{\Pi}} \sum_{\mathbf{k} \in \mathbb{Z}_{2N}^d} \mathbf{u}_{N,\beta}(\mathbf{x}_{2N}^{\mathbf{k}}) \mathbf{v}_{N,\alpha}(\mathbf{x}_{2N}^{\mathbf{k}}) \omega_{2N}^{-m\mathbf{k}} \quad (3.17a)$$

$$\widehat{A}_{\alpha\beta}^m = \sum_{\mathbf{n} \in \mathbb{Z}_M^d} A_{\alpha\beta}^{\mathbf{n}} \omega_M^{-m\mathbf{n}} \quad (3.17b)$$

Proof. First, we note that $\int_{\mathcal{Y}} \psi(\mathbf{x} + \mathbf{x}_M^{\mathbf{n}}) \varphi_m(\mathbf{x}) d\mathbf{x} = \int_{\mathcal{Y}} \psi(\mathbf{x}) \varphi_m(\mathbf{x} - \mathbf{x}_M^{\mathbf{n}}) d\mathbf{x} = w(\mathbf{m}) \omega_M^{-m\mathbf{n}}$ where $\mathbf{m}, \mathbf{n} \in \mathbb{Z}^d$. Since $\mathbf{u}_{N,\beta}, \mathbf{v}_{N,\alpha} \in \mathcal{T}_N$, for their multiplication holds $\mathbf{u}_{N,\beta} \mathbf{v}_{N,\alpha} \in \mathcal{T}_{2N}$. Thus, it can be expressed as interpolation through $2N$ nodal points

$$\mathbf{u}_{N,\beta} \mathbf{v}_{N,\alpha} = \sum_{\mathbf{m} \in \mathbb{Z}_{2N}^d} \widehat{u}_{N,\beta,\alpha}^m \varphi_m \quad (3.18)$$

Then, the direct calculation finish the proof

$$\begin{aligned} (\mathbf{A}_M \mathbf{u}_N, \mathbf{v}_N)_{L^2_{\text{per}}} &= \frac{1}{|\mathcal{Y}|_d} \sum_{\alpha,\beta} \int_{\mathcal{Y}} A_{M,\alpha\beta}(\mathbf{x}) \mathbf{u}_{N,\beta}(\mathbf{x}) \mathbf{v}_{N,\alpha}(\mathbf{x}) d\mathbf{x} \\ &= \frac{1}{|\mathcal{Y}|_d} \sum_{\alpha,\beta} \sum_{\mathbf{n} \in \mathbb{Z}_M^d} \sum_{\mathbf{m} \in \mathbb{Z}_{2N}^d} A_{\alpha\beta}^{\mathbf{n}} \widehat{u}_{N,\beta,\alpha}^m \int_{\mathcal{Y}} \psi(\mathbf{x} + \mathbf{x}_M^{\mathbf{n}}) \varphi_m(\mathbf{x}) d\mathbf{x} \\ &= \frac{1}{|\mathcal{Y}|_d} \sum_{\alpha,\beta} \sum_{\mathbf{m} \in \mathbb{Z}_{2N}^d} \widehat{u}_{N,\beta,\alpha}^m w(\mathbf{m}) \sum_{\mathbf{n} \in \mathbb{Z}_M^d} A_{\alpha\beta}^{\mathbf{n}} \omega_M^{-m\mathbf{n}}. \end{aligned}$$

□

LEMMA 3.32 (Homogenized matrices by FFT algorithm). *Let the assumptions from the previous lemma be satisfied and, in addition, let $M \in \mathbb{N}^d$ be such that $M_{\alpha} \geq 2N_{\alpha}$ then the formula in Eq. (3.16) can be calculated using the FFT algorithm of size $2N$ and M resp.*

Proof. In Eq. (3.17a), the function values of $\mathbf{u}_{N,\beta}, \mathbf{v}_{N,\alpha}$ at nodal points $\mathbf{x}_{2N}^{\mathbf{k}}$ can be calculated with inverse DFT of size $2N$ — the Fourier coefficients with frequencies $\mathbb{Z}_{2N}^d - \mathbb{Z}_N^d$ are completed with nils. Then Eq. (3.17) is, up to a constant, the DFT on space \mathbb{R}^{2N} and \mathbb{R}^M resp. — compare it to DFT matrix \mathbf{F}_N^d in Def. 3.1 acting on space $\mathbb{R}^{d \times N}$. The Fourier coefficients of $\widehat{u}_{N,\beta,\alpha}^m$ equal to zero for $\mathbf{m} \in \mathbb{Z}^d \setminus \mathbb{Z}_{2N}^d$, it reveals requirement $M_{\alpha} \geq 2N_{\alpha}$. □

Up to now, we have shown that the homogenized properties can be calculated effectively using the FFT algorithm if the material coefficients are expressed as the linear combination of basis functions concentrated on the set of nodal points. Now, we show, in Def. 3.33, some examples of basis functions, to be utilized in Lem. 3.31, that can be used to calculate the homogenized matrices; the choice depends on the possibility to analytically express the integral weights, cf. Lem. 3.34.

DEFINITION 3.33 (Constant and bilinear basis functions). *Let $M \in \mathbb{R}^d$ be a parameter such that $M_{\alpha} > 1$. We say that functions $\chi_M, \text{tri}_M \in L^{\infty}_{\text{per}}(\mathcal{Y}; \mathbb{R})$, defined on \mathcal{Y} as*

$$\chi_M(\mathbf{x}) = \begin{cases} 1, & |x_{\alpha}| < \frac{Y_{\alpha}}{M_{\alpha}} \text{ for all } \alpha \\ 0, & \text{otherwise} \end{cases}, \quad \text{tri}_M(\mathbf{x}) = \prod_{\alpha} \max\{1 - |\frac{x_{\alpha} M_{\alpha}}{2Y_{\alpha}}|, 0\},$$

are, one by one, a constant and a bilinear basis functions.

LEMMA 3.34 (Weights of numerical integration). *Let χ_M and tri_M be the basis functions from Def. 3.33. Then for $\mathbf{m} \in \mathbb{Z}^d$ we state*

$$\begin{aligned} w_M^0(\mathbf{m}) &:= \int_{\mathcal{Y}} \chi_M(\mathbf{x}) \varphi_{\mathbf{m}}(\mathbf{x}) d\mathbf{x} = \prod_{\alpha} \frac{2Y_{\alpha}}{M_{\alpha}} \text{sinc}\left(\frac{m_{\alpha}}{M_{\alpha}}\right) \\ w_M^1(\mathbf{m}) &:= \int_{\mathcal{Y}} \text{tri}_M(\mathbf{x}) \varphi_{\mathbf{m}}(\mathbf{x}) d\mathbf{x} = \prod_{\alpha} \frac{2Y_{\alpha}}{M_{\alpha}} \text{sinc}^2\left(\frac{m_{\alpha}}{M_{\alpha}}\right) \end{aligned}$$

where $\text{sinc}(x) := \begin{cases} 1, & x = 0 \\ \frac{\sin(\pi x)}{\pi x}, & x \neq 0 \end{cases}$.

Proof. For $\mathbf{m} \in \mathbb{Z}^d$ we calculate

$$\begin{aligned} \int_{\mathcal{Y}} \chi_M \varphi_{\mathbf{m}}(\mathbf{x}) d\mathbf{x} &= \prod_{\alpha} \int_{|x_{\alpha}| < \frac{Y_{\alpha}}{M_{\alpha}}} \exp(i\pi \frac{x_{\alpha} m_{\alpha}}{Y_{\alpha}}) d\mathbf{x}_{\alpha} = \prod_{\alpha} \left[\frac{Y_{\alpha}}{i\pi m_{\alpha}} \exp(i\pi \frac{x_{\alpha} m_{\alpha}}{Y_{\alpha}}) \right]_{-\frac{Y_{\alpha}}{M_{\alpha}}}^{\frac{Y_{\alpha}}{M_{\alpha}}} \\ &= \prod_{\alpha} \left[\frac{2Y_{\alpha}}{M_{\alpha}} \frac{\sin(\frac{\pi m_{\alpha}}{M_{\alpha}})}{\frac{\pi m_{\alpha}}{M_{\alpha}}} \right] = \prod_{\alpha} \frac{2Y_{\alpha}}{M_{\alpha}} \text{sinc}\left(\frac{m_{\alpha}}{M_{\alpha}}\right) \end{aligned}$$

Integral weights $w_M^1(\mathbf{m})$ for bilinear basis functions tri_M are calculated accordingly.

□

REMARK 3.35. *Since we realize that the integral weights can be calculated for the basis functions shifted by $\mathbf{h} \in \mathbb{R}^d$, e.g.*

$$\int_{\mathcal{Y}} \chi_N(\mathbf{x} + \mathbf{h}) \varphi_{\mathbf{m}}(\mathbf{x}) d\mathbf{x} = w_N^0(\mathbf{m}) \varphi_{\mathbf{m}}(-\mathbf{h}),$$

we can calculate the upper and lower bounds exactly if the material coefficients are expressed as some linear combination of the shifted basis functions from Def. 3.33. We note that the shifts by particular \mathbf{h} enable to calculate the bounds using the FFT algorithm, cf. Lem. 3.31 and Lem. 3.32.

Next, the basis functions stated in Def. 3.33 are not the only suitable ones, for example the circle basis function defined for $\mathbf{x} \in \mathcal{Y}$ as

$$\text{circ}_r(\mathbf{x}) = \begin{cases} 1, & \|\mathbf{x}\|_2 \leq r \\ 0, & \text{otherwise} \end{cases}$$

produces the integral weight with the Bessel function of the first kind B_1 , i.e. for $\xi_{\alpha}(\mathbf{m}) = \frac{m_{\alpha}}{Y_{\alpha}}$

$$w_N^{\text{circ}}(\mathbf{m}) := \int_{\mathcal{Y}} \text{circ}_r(\mathbf{x}) \varphi_{\mathbf{m}}(\mathbf{x}) d\mathbf{x} = \begin{cases} \pi r^2, & \mathbf{m} = \mathbf{0} \\ r^2 \frac{B_1(2\pi r \|\xi(\mathbf{m})\|_2)}{r \|\xi(\mathbf{m})\|_2}, & \text{otherwise} \end{cases}$$

Now, we will define approximations of the guaranteed bounds with the basis functions stated in Def. 3.33; a piecewise constant approximation is defined in a way to still provide the guaranteed bounds while a piecewise bilinear approximation only approximate these bounds, see Lem. 3.30.

DEFINITION 3.36 (Constant and bilinear approximation of material coefficients). Let $\mathbf{A} \in L_{\text{per}}^{\infty}(\mathcal{Y}; \mathbb{R}_{\text{spd}}^{d \times d})$ be material coefficients and $\mathbf{M} \in \mathbb{N}$ a parameter. Then we define functions $\overline{\mathbf{A}}^{\text{con},M}, \underline{\mathbf{B}}^{\text{con},M}, \mathbf{A}^{\text{lin},M}, \mathbf{B}^{\text{lin},M} \in L_{\text{per}}^{\infty}(\mathcal{Y}; \mathbb{R}^{d \times d})$ for $\mathbf{x} \in \mathcal{Y}$ as

$$\begin{aligned}\overline{\mathbf{A}}^{\text{con},M}(\mathbf{x}) &= \sum_{\mathbf{n} \in \underline{\mathbb{Z}}_M^d} \chi_M(\mathbf{x} + \mathbf{x}_M^n) \overline{\mathbf{A}}_M^{\text{con},M,n}, \\ \underline{\mathbf{B}}^{\text{con},M}(\mathbf{x}) &= \sum_{\mathbf{n} \in \underline{\mathbb{Z}}_M^d} \chi_M(\mathbf{x} + \mathbf{x}_M^n) \underline{\mathbf{B}}^{\text{con},M,n}, \\ \mathbf{A}^{\text{lin},M}(\mathbf{x}) &= \sum_{\mathbf{n} \in \underline{\mathbb{Z}}_M^d} \text{tri}_M(\mathbf{x} + \mathbf{x}_M^n) \mathbf{A}(\mathbf{x}_M^n), \\ \mathbf{B}^{\text{lin},M}(\mathbf{x}) &= \sum_{\mathbf{n} \in \underline{\mathbb{Z}}_M^d} \text{tri}_M(\mathbf{x} + \mathbf{x}_M^n) \mathbf{A}^{-1}(\mathbf{x}_M^n),\end{aligned}$$

where matrices $\overline{\mathbf{A}}_M^{\text{con},M,n}, \underline{\mathbf{B}}^{\text{con},M,n} \in \mathbb{R}^{d \times d}$ are defined as

$$\overline{\mathbf{A}}_M^{\text{con},M,n} = \|\mathbf{A}(\mathbf{x} + \mathbf{x}_M^n)\|_{L^{\infty}(\Omega_M; \mathbb{R}_{\text{spd}}^{d \times d})} \mathbf{I}, \quad \underline{\mathbf{B}}^{\text{con},M,n} = \|\mathbf{A}^{-1}(\mathbf{x} + \mathbf{x}_M^n)\|_{L^{\infty}(\Omega_M; \mathbb{R}_{\text{spd}}^{d \times d})} \mathbf{I},$$

noting that L_{per}^{∞} -norm is restricted on $\Omega_M = \prod_{\alpha} \left(-\frac{Y_{\alpha}}{M_{\alpha}}, \frac{Y_{\alpha}}{M_{\alpha}}\right)$. The matrices are called, one by one, upper constant, lower constant, upper bilinear, and lower bilinear approximation of material coefficients.

LEMMA 3.37 (Constant approximation of material coefficients). The constant approximations of the material coefficients, Def. 3.36, satisfy the sufficient condition in Lem. 3.30 for guaranteeing bounds.

Proof. It is necessary to show

$$\mathbf{A}(\mathbf{x}) \preceq \overline{\mathbf{A}}^{\text{con},M}(\mathbf{x}), \quad \mathbf{A}^{-1}(\mathbf{x}) \preceq \underline{\mathbf{B}}^{\text{con},M}(\mathbf{x}), \quad (\underline{\mathbf{B}}^{\text{con},M})^{-1}(\mathbf{x}) \preceq \mathbf{A}(\mathbf{x}) \preceq \overline{\mathbf{A}}^{\text{con},M}(\mathbf{x}),$$

however, it is a direct consequence of the definition of the approximated material coefficients. \square

DEFINITION 3.38 (Bounds and its approximation of homogenized material coefficients). Let $\tilde{\mathbf{e}}_N^{(\alpha)} \in \mathcal{E}_N$ and $\tilde{\mathbf{j}}_N^{(\alpha)} \in \mathcal{J}_N$ be approximations of unitary minimizers, e.g. from Def. 2.5, and $\overline{\mathbf{A}}^{\text{con},M}, \underline{\mathbf{B}}^{\text{con},M}, \mathbf{A}^{\text{lin},M}, \mathbf{B}^{\text{lin},M}$ for $\mathbf{M} \in \mathbb{N}^d$ be approximations of material coefficients from Def. 3.36. Then we define the bounds $\overline{\underline{\mathbf{A}}}_{\text{eff},N}^{\text{con},M}, \underline{\underline{\mathbf{A}}}_{\text{eff},N}^{\text{con},M}, \tilde{\underline{\mathbf{A}}}_{\text{eff},N}^{\text{lin},M}, \tilde{\underline{\mathbf{A}}}_{\text{eff},D,N}^{\text{lin},M} \in \mathbb{R}^{d \times d}$ of the homogenized matrix as

$$\overline{\underline{\mathbf{A}}}_{\text{eff},N,\alpha\beta}^{\text{con},M} = (\overline{\mathbf{A}}^{\text{con},M} \mathbf{e}_N^{(\beta)}, \mathbf{e}_N^{(\alpha)})_{L_{\text{per}}^2}, \quad (\underline{\underline{\mathbf{A}}}_{\text{eff},N}^{\text{con},M})_{\alpha\beta}^{-1} = (\underline{\mathbf{B}}^{\text{con},M} \mathbf{j}_N^{(\beta)}, \mathbf{j}_N^{(\alpha)})_{L_{\text{per}}^2}, \quad (3.19a)$$

$$\tilde{\underline{\mathbf{A}}}_{\text{eff},N,\alpha\beta}^{\text{lin},M} = (\mathbf{A}^{\text{lin},M} \mathbf{e}_N^{(\beta)}, \mathbf{e}_N^{(\alpha)})_{L_{\text{per}}^2}, \quad (\tilde{\underline{\mathbf{A}}}_{\text{eff},D,N}^{\text{lin},M})_{\alpha\beta}^{-1} = (\mathbf{B}^{\text{lin},M} \mathbf{j}_N^{(\beta)}, \mathbf{j}_N^{(\alpha)})_{L_{\text{per}}^2}. \quad (3.19b)$$

THEOREM 3.39 (Guaranteed bounds for piecewise constant material coefficients). The bounds $\underline{\underline{\mathbf{A}}}_{\text{eff},N}^{\text{con},M}, \overline{\underline{\mathbf{A}}}_{\text{eff},N}^{\text{con},M} \in \mathbb{R}^{d \times d}$ from the previous definition are the guaranteed bounds, i.e.

$$\underline{\underline{\mathbf{A}}}_{\text{eff},N}^{\text{con},M} \preceq \mathbf{A}_{\text{eff}} \preceq \overline{\underline{\mathbf{A}}}_{\text{eff},N}^{\text{con},M}, \quad (3.20)$$

and can be calculated with the FFT algorithm.

Proof. Relation (3.20) is a corollary of Lem. 3.37. The possibility to calculate it using the FFT algorithm is a consequence of formula (3.16) in Lem. 3.31 and discussion in Lem. 3.32. \square

REMARK 3.40 (Homogenized bounds). *The bounds obtained with the bilinear approximations of the material coefficients are only approximation of the bounds from Def. 2.8, explicitly $\tilde{\mathbf{A}}_{\text{eff},\mathbf{N}}^{\text{lin},M} \approx \overline{\mathbf{A}}_{\text{eff},\mathbf{N}}$, $\tilde{\mathbf{A}}_{\text{eff},\mathbf{D},\mathbf{N}}^{\text{lin},M} \approx \underline{\mathbf{A}}_{\text{eff},\mathbf{N}}$. Nevertheless, it is still possible to calculate it using the FFT algorithm, cf. Lem. 3.32.*

4. Numerical experiments. This section is dedicated to numerical experiments. We discuss the practical aspects of computing the homogenized matrices, and then verify the theoretical results: the upper-lower bounds structure of homogenized matrices in Section 4.1, and the rate of convergence, particularly convergence of approximate solutions to continuous one and the convergence of the homogenized matrices, Sec. 4.2.

Finally, the numerical experiments are compared to the p-version of the Finite Element Method provided in [6], see Sec. 4.3.

REMARK 4.1. *The calculations are provided for a 2-dimensional problem with (2, 2)-periodic material coefficients defined on periodic unit cell $\mathcal{Y} = (-1, 1) \times (-1, 1) \subset \mathbb{R}^2$ as*

$$\mathbf{A}(\mathbf{x}) = \mathbf{I}[1 + \rho f(\mathbf{x})], \quad \mathbf{x} \in \mathcal{Y},$$

where $\mathbf{I} \in \mathbb{R}_{\text{spd}}^{d \times d}$ is the identity matrix, ρ corresponds to the phase ratio and is taken as 10 except Sec. 4.3 where it is taken as $(10^{-3} - 1)$, and $f : \mathcal{Y} \rightarrow \mathbb{R}$ is a scalar nonnegative function defined explicitly for particular experiment — it controls the shape of inclusions and the regularity of material coefficients.

Moreover, discrete minimizers $\mathbf{e}_{\mathbf{N}}^{(\alpha)} \in \mathcal{E}_{\mathbf{N}}$ and $\mathbf{j}_{\mathbf{N}}^{(\alpha)} \in \mathcal{J}_{\mathbf{N}}$ from Def. 3.17 are obtained for both odd and even number of discretization points, namely $N = (n, n)$ where either $n \in \{5 \cdot 3^\alpha | \alpha \in \mathbb{N}_0, 0 \leq \alpha \leq 6\}$ or $n \in \{2^\alpha | \alpha \in \mathbb{N}, 2 \leq \alpha \leq 10\}$; both the sets of numbers are geometric series carrying ratios between successive terms λ equal to either 3 or 2.

LOGARITHM 4.2 (Calculation of guaranteed bounds of homogenized material properties). *Let $\mathbf{A}_{\mathbf{M}} \in L^\infty_{\text{per}}(\mathcal{Y}; \mathbb{R}^{d \times d})$ be material coefficients. The algorithm for calculation of the homogenized matrices is composed of several steps:*

- (i) *Set the number of discretization points \mathbf{N} and assemble matrices $\mathbf{A}_{\mathbf{N}}$, $\mathbf{A}_{\mathbf{N}}^{-1}$, $\hat{\mathbf{G}}_{\mathbf{0}}^{(1)}$, $\hat{\mathbf{G}}_{\mathbf{0}}^{(2)} \in \mathbb{R}^{d \times d \times N \times N}$ defined in Def. 3.19 and 3.11. Since they are block diagonal, only the diagonals are stored and the matrix by vector multiplication is provided as element-wise multiplication.*
- (ii) *For α , find discrete minimizers $\tilde{\mathbf{e}}_{\mathbf{N}}^{(\alpha)} \in \mathbb{E}_{\mathbf{N}}$, $\tilde{\mathbf{j}}_{\mathbf{N}}^{(\alpha)} \in \mathbb{J}_{\mathbf{N}}$ as the solutions of linear systems described in Rem. 3.22 for unitary macroscopic loads $\boldsymbol{\epsilon}_\alpha$. The convergence criterion is based on the norm of residuum $\|\mathbf{r}\|_{\text{CG}} \leq \varepsilon$ where $\|\mathbf{r}\|_{\text{CG}}^2 := |\mathbf{N}|_{\Pi}^{-1}(\mathbf{r}, \mathbf{r})_{\mathbb{R}^{d \times \mathbf{N}}}$; the yielding value is set as $\varepsilon = 10^{-10}$ and initial approximate vectors of CG are set as zeros.*
- (iii) *Calculate, if possible as stated in Rem. 3.35, the exact upper and lower bounds $\overline{\mathbf{A}}_{\text{eff},\mathbf{N}}$, $\underline{\mathbf{A}}_{\text{eff},\mathbf{N}}$, see Def. 2.8. Otherwise, evaluate the approximations of the upper-lower bounds $\overline{\mathbf{A}}_{\text{eff},\mathbf{N}}^{\text{lin},M}$, $\underline{\mathbf{A}}_{\text{eff},\mathbf{N}}^{\text{lin},M}$, $\tilde{\mathbf{A}}_{\text{eff},\mathbf{N}}^{\text{con},M}$, $\tilde{\mathbf{A}}_{\text{eff},\mathbf{D},\mathbf{N}}^{\text{con},M}$, Def. 3.38, according to Lem. 3.31 for some sufficiently large $M \in \mathbb{R}^d$, cf. Lem. 3.32.*

REMARK 4.3 (Convergence criterion). *The norm for residuum $\|\mathbf{r}\|_{\text{CG}}$, due to Plancherel's theorem, equals to $\|\mathcal{I}_{\mathbf{N}}^{-1}[\mathbf{r}]\|_{L^2_{\text{per}}}$. The yielding value is set as small as possible to diminish an error caused by an inaccuracy in the solution of the linear systems.*

REMARK 4.4 (Avoiding the solution of dual formulation). *If \mathbf{N} is odd, the dual discrete minimizers $\tilde{\mathbf{j}}_{\mathbf{N}}^{(\alpha)}$ can be obtained from Eq. 3.10 based on the assumption that the original minimizers $\tilde{\mathbf{e}}_{\mathbf{N}}^{(\alpha)}$ are the exact solutions of the corresponding linear systems, see Rem. 3.22. In reality, the linear systems are solved only approximately, thus it fails the dual minimizers to be from appropriate subspace $\tilde{\mathbf{j}}_{\mathbf{N}}^{(\alpha)} \notin \mathbb{J}_{\mathbf{N}}$. It can be saved with projection operator $\mathbf{G}_{\mathbf{0}}^{(2)}$ and, in case $\mathbf{A}_{\mathbf{N}}$ is badly conditioned, by providing couple of iterations of the dual formulation (3.6b).*

REMARK 4.5 (Interpolation operator $\mathcal{Q}_{\mathbf{N}}$ for non-continuous functions). *The discrete formulations, Def. 3.17, require the interpolation operator $\mathcal{Q}_{\mathbf{N}}$ to be well defined; it takes the function values and thus, originally, the operator is well defined, for example, on continuous functions.*

We assume the piecewise constant material coefficients and define the discrete formulations as it states in Def. 3.19. Generally, it still provides the upper-lower bounds, however, it can fail to converge. Alternatively, some regularization of material coefficients, see [31], can be performed to obtain convergence, nevertheless, generally arbitrary slow.

4.1. The behavior of the bounds of the homogenized matrix. In this section, we validate the properties of the upper-lower bounds of the homogenized matrices for two phase materials with three types of inclusions characterized by scalar functions f , see Rem. 4.1; explicitly, they are defined as

- square (S): $f(\mathbf{x}) = \begin{cases} 1, & \|\mathbf{x}\|_{\infty} < \frac{3}{5} \\ 0, & \text{otherwise} \end{cases}$,
- square (S1): $f(\mathbf{x}) = \begin{cases} 1, & \|\mathbf{x}\|_{\infty} < \frac{3}{4} \\ 0, & \text{otherwise} \end{cases}$,
- square (S2): $f(\mathbf{x}) = \begin{cases} 1, & \|\mathbf{x}\|_{\infty} \leq \frac{3}{4} \\ 0, & \text{otherwise} \end{cases}$.

Fig. 4.1 depicts the interface between phases and the nodal points sets, $\{\mathbf{x}_{\mathbf{N}}^n \in \mathcal{Y} : \mathbf{n} \in \mathbb{Z}_{\mathbf{N}}^d\}$, for particular \mathbf{N} . The squares (S1) and (S2) differ only at the interface having the 2-dimensional Lebesgue's measure equal to zero, thus insignificant for the upper and lower bounds, see Def. 2.8.

The square (S) receive the interface exactly between the nodal points — it approximate well an inclusion contrary to the squares (S1) and (S2) representing the extreme cases of an inclusion approximation. The nodal points lie exactly on the boundary causing the difference in the discrete formulation and consequently in discrete minimizers $\mathbf{e}_{\mathbf{N}}^{(\alpha)}, \mathbf{j}_{\mathbf{N}}^{(\alpha)}$ and upper-lower bounds $\overline{\mathbf{A}}_{\text{eff},\mathbf{N}}, \underline{\mathbf{A}}_{\text{eff},\mathbf{N}}$.

In Fig. 4.1, we demonstrate the properties of the homogenized matrices for their particular diagonal component. Inequality $\underline{\mathbf{A}}_{\text{eff},\mathbf{N}} \preceq \overline{\mathbf{A}}_{\text{eff},\mathbf{N}}$ stated in Lem. 2.10 is satisfied and the error, difference between them, is approaching zero supporting Lem. 2.14.

For odd \mathbf{N} in Fig. 4.1(a), the primal and the dual approximate homogenized matrices coincide $\mathbf{A}_{\text{eff},\mathbf{N}}^{\text{FFTH}} = \mathbf{A}_{\text{eff},\mathbf{D},\mathbf{N}}^{\text{FFTH}}$ as stated in Theorem 3.26. Moreover, it approximates properly the real homogenized coefficients \mathbf{A}_{eff} even for small \mathbf{N} compared to

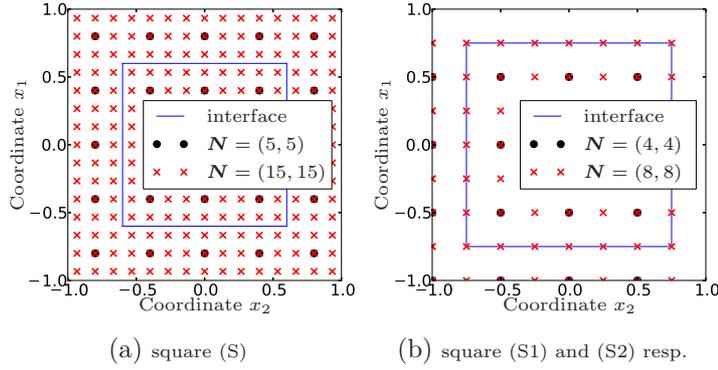


FIG. 4.1. Periodic unit cell with nodal points and interfaces between phases

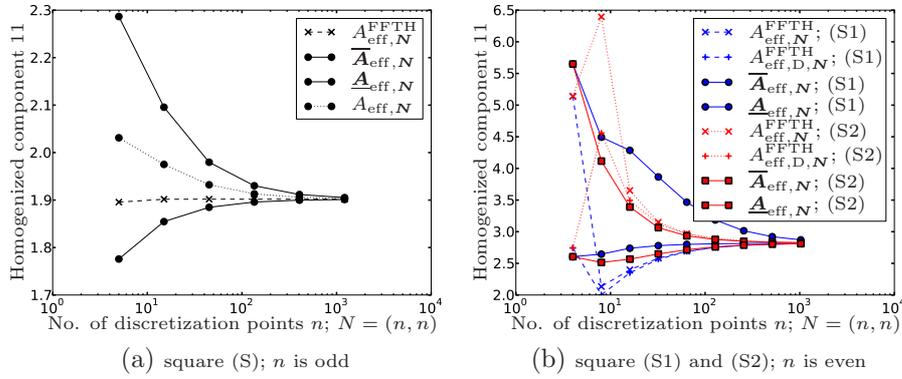


FIG. 4.2. The upper-lower bounds of the homogenized matrix for both the odd and the even number of discretization points

the mean value $\mathbf{A}_{\text{eff},N} = \frac{1}{2}(\underline{\mathbf{A}}_{\text{eff},N} + \overline{\mathbf{A}}_{\text{eff},N})$ that overestimates.

For even N in Fig. 4.1(b), the approximate homogenized matrices satisfy a sharp inequality $\mathbf{A}_{\text{eff},D,N}^{\text{FFTH}} < \mathbf{A}_{\text{eff},N}^{\text{FFTH}}$ for both squares (S1) and (S2) confirming (3.13) in Theorem 3.27. Both the homogenized matrices $\mathbf{A}_{\text{eff},N}^{\text{FFTH}}$, $\mathbf{A}_{\text{eff},D,N}^{\text{FFTH}}$ either overestimates or underestimates even over the upper or the lower bounds. Exception is a case $N = (4, 4)$ when the material coefficients coincide at the nodal points for both squares (S1) and (S2).

However, the primal and the dual homogenized matrices, in Eq. (3.6a) and (3.11b), differ substantially for small N alike subspaces \mathcal{E}_N and $\bar{\mathcal{E}}_N$ do. Hence, the mean of the upper-lower bounds $\mathbf{A}_{\text{eff},N}$ provides the more accurate result. Albeit, from the design perspective, either the lower or the upper bound is chosen as the most reliable homogenized property depending on a design demand: resistance or conductivity.

4.2. Rate of convergence. In this section, we focus on the rate of convergence (RoC) depending primarily on the regularity of material coefficient. Thus, according to Rem. 4.1, we define material coefficients through scalar functions f as

- circle (C): $f(\mathbf{x}) = \begin{cases} 1, & \|\mathbf{x}\|_2 < \frac{3}{5} \\ 0, & \text{otherwise} \end{cases}$

- cone (E): $f(\mathbf{x}) = \begin{cases} 1 - \|\mathbf{x}\|_2, & \|\mathbf{x}\|_2 < 1 \\ 0, & \text{otherwise} \end{cases}$,
- hummock (H): $f(\mathbf{x}) = \begin{cases} 1 - 2\|\mathbf{x}\|_2^2, & \|\mathbf{x}\|_2 < \frac{1}{2} \\ 2(1 - \|\mathbf{x}\|_2)^2, & \frac{1}{2} \leq \|\mathbf{x}\|_2 < 1, \\ 0, & \text{otherwise} \end{cases}$,
- standard mollifier (M): $f(\mathbf{x}) = \begin{cases} \exp(1 - \frac{1}{1-\|\mathbf{x}\|_2}), & \|\mathbf{x}\|_2 < 1 \\ 0, & \text{otherwise} \end{cases}$;

the cut through the periodic unit cell, for $x_2 \equiv 0$, can be observed in Fig. 4.3. The regularity of the material coefficients are based on the regularity of scalar functions f ; circle (C) is piecewise constant, cone (E) has piecewise constant the first derivative, hummock (H) has piecewise constant the second derivative while standard mollifier (M) is infinitely smooth, having continuous all derivatives.

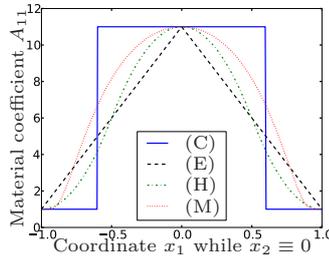


FIG. 4.3. The material coefficients at PUC for $x_2 \equiv 0$ showing their smoothness

First, we verify the RoC of $\|e^{(\alpha)} - e_{\mathbf{N}}^{(\alpha)}\|_{L^2_{\text{per}}} \leq C(\frac{n}{2})^{-\mu}$ from Lem. 3.23 for $\alpha = 1$. Unfortunately, the exact solution $e^{(1)}$ is possible to obtain only in a very special cases, hence, we approximate it by a solution calculated with a very fine grid, i.e. $e^{(1)} \approx e_{\bar{\mathbf{N}}}^{(1)}$ where we have chosen $\bar{\mathbf{N}} = (2005, 2005)$. It enables to calculate the approximations of order μ , particularly

$$\mu_{\mathbf{N}} = \log \left(\frac{\|e_{\bar{\mathbf{N}}}^{(1)} - e_{\mathbf{N}}^{(1)}\|_{L^2_{\text{per}}}}{\|e_{\bar{\mathbf{N}}}^{(1)} - e_{\lambda \bar{\mathbf{N}}}^{(1)}\|_{L^2_{\text{per}}}} \right) (\log \lambda)^{-1}, \quad (4.1)$$

where the factors λ equals to 2 or 3, cf. Rem. 4.1. Norms $\|\cdot\|_{L^2_{\text{per}}}$ in (4.1) can be calculated exactly using Plancherel's theorem as all functions are the trigonometric polynomials.

TABLE 4.1
The rate of convergence $\mu_{\mathbf{N}}$ for the odd number of discretization points

RoC	$f \setminus n$	5	15	45	135	405	limit
$\mu_{\mathbf{N}}$	(C)	0.568	0.579	0.511	0.544	0.565	0.5
$\mu_{\mathbf{N}}$	(E)	1.083	1.410	1.522	1.527	1.557	1.5
$\mu_{\mathbf{N}}$	(H)	2.504	2.554	2.476	2.478	2.499	2.5
$\mu_{\mathbf{N}}$	(M)	1.491	3.009	6.777	10.266	6.091	—

The rate of convergence can also be studied from the upper-lower bounds of the homogenized matrix. According to Lem. 2.14, it converges with rate 2μ that represents the smallest value of following rates $\|e^{(\alpha)} - e_{\mathbf{N}}^{(\alpha)}\|_{L^2_{\text{per}}}^2$ and $\|J^{(\alpha)} - J_{\mathbf{N}}^{(\alpha)}\|_{L^2_{\text{per}}}^2$ for all α .

TABLE 4.2
The rate of convergence μ_N for the even number of discretization points

RoC	$f \setminus n$	4	8	16	32	64	128	256	512	limit
μ_N	(C)	-0.139	1.111	0.315	0.484	0.553	0.578	0.467	0.603	0.5
μ_N	(E)	1.067	1.210	1.293	1.407	1.483	1.515	1.522	1.525	1.5
μ_N	(H)	2.326	2.964	2.652	2.468	2.443	2.452	2.469	2.486	2.5
μ_N	(M)	1.475	1.710	2.714	4.812	7.246	8.895	12.656	4.858	—

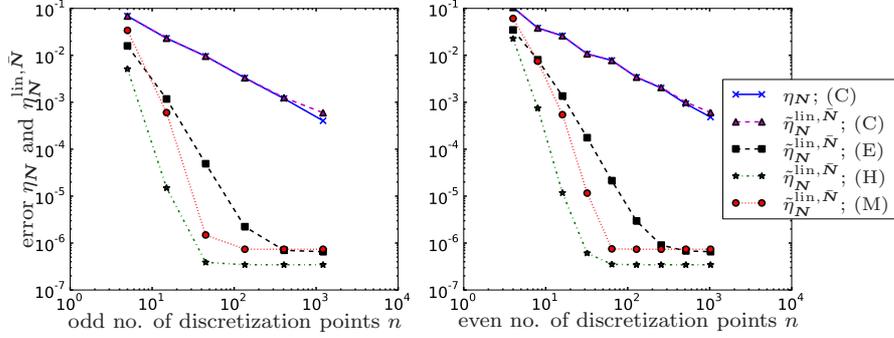


FIG. 4.4. Progress in error η_N for increasing in the number of discretization points $N = (n, n)$

Hence, we define

$$\bar{\mu}_N = \frac{1}{2} \log \left(\frac{\text{tr } D_N}{\text{tr } D_{\lambda N}} \right) (\log \lambda)^{-1}, \quad (4.2)$$

where D_N is the error from Def. 2.11 and λ is taken as 2 or 3.

Nevertheless, it can be calculated only for special material coefficients \mathbf{A} , e.g. for problem (C), cf. Rem. 3.35. Thus, we define the alternative

$$\tilde{\mu}_N^{\text{lin}, \bar{N}} = \frac{1}{2} \log \left(\frac{\text{tr } \tilde{D}_N^{\text{lin}, \bar{N}}}{\text{tr } \tilde{D}_{\lambda N}^{\text{lin}, \bar{N}}} \right) (\log \lambda)^{-1}, \quad (4.3)$$

where $\tilde{D}_N^{\text{lin}, \bar{N}} = \frac{1}{2} \left| \tilde{\mathbf{A}}_{\text{eff}, N}^{\text{lin}, \bar{N}} - \tilde{\mathbf{A}}_{\text{eff}, D, N}^{\text{lin}, \bar{N}} \right|$ and $\bar{N} = (2005, 2005)$.

Analogically, we define the normalized errors

$$\eta_N := \frac{\text{tr } D_N}{\text{tr } \mathbf{A}_{\text{eff}, M}} \quad (4.4a)$$

$$\tilde{\eta}_N^{\text{lin}, \bar{N}} := \frac{\text{tr } \tilde{D}_N^{\text{lin}, \bar{N}}}{\text{tr } \mathbf{A}_{\text{eff}, M}} \quad (4.4b)$$

where $M = (1215, 1215)$ or $M = (1024, 1024)$ depending on N being odd or even, and \bar{N} is taken as previously.

Now, we will discuss and compare the results. First, there is no significant observation between even and odd number of discretization points N in the rates of convergence, compare Tab. 4.1, 4.3 versus Tab. 4.2, 4.4, and see also almost straight lines of the normalized errors in Fig. 4.4.

The rates of convergence suffer from an inaccuracy for large N ; the rates of $\|e^{(1)} - e_N^{(1)}\|_{L^2_{\text{per}}}$ stated in Tab. 4.1 and 4.2 are depreciated due to approximation of

TABLE 4.3

The rate of convergence (RoC) from the guaranteed bounds for odd $\mathbf{N} = (n, n)$

RoC	$f \setminus n$	5	15	45	135	405	theory
$\bar{\mu}_{\mathbf{N}}$	(C)	0.495	0.402	0.484	0.460	0.496	0.5
$\tilde{\mu}_{\mathbf{N}}^{\text{lin}, \tilde{\mathbf{N}}}$	(C)	0.495	0.401	0.483	0.441	0.338	0.5
$\tilde{\mu}_{\mathbf{N}}^{\text{lin}, \tilde{\mathbf{N}}}$	(E)	1.185	1.446	1.406	0.526	0.032	1.5
$\tilde{\mu}_{\mathbf{N}}^{\text{lin}, \tilde{\mathbf{N}}}$	(H)	2.655	1.662	0.054	0.000	0.000	2.5
$\tilde{\mu}_{\mathbf{N}}^{\text{lin}, \tilde{\mathbf{N}}}$	(M)	1.836	2.732	0.319	0.000	0.000	???

TABLE 4.4

The rate of convergence (RoC) from the guaranteed bounds for even $\mathbf{N} = (n, n)$

RoC	$f \setminus n$	4	8	16	32	64	128	256	512	limit
$\bar{\mu}_{\mathbf{N}}$	(C)	0.727	0.276	0.641	0.230	0.596	0.379	0.561	0.469	0.5
$\tilde{\mu}_{\mathbf{N}}^{\text{lin}, \tilde{\mathbf{N}}}$	(C)	0.727	0.276	0.641	0.231	0.593	0.363	0.524	0.360	0.5
$\tilde{\mu}_{\mathbf{N}}^{\text{lin}, \tilde{\mathbf{N}}}$	(E)	1.055	1.284	1.475	1.520	1.427	0.853	0.207	0.026	1.5
$\tilde{\mu}_{\mathbf{N}}^{\text{lin}, \tilde{\mathbf{N}}}$	(H)	2.456	2.999	2.133	0.396	0.015	0.000	0.000	0.000	2.5
$\tilde{\mu}_{\mathbf{N}}^{\text{lin}, \tilde{\mathbf{N}}}$	(M)	1.511	1.892	2.774	1.975	0.013	0.000	-0.000	0.000	???

minimizers $\mathbf{e}^{(1)} \approx \mathbf{e}_{\mathbf{N}}^{(1)}$ and the rates of the upper-lower bounds, see Tab. 4.3 and 4.3, due to approximations $\bar{\mathbf{A}}_{\text{eff}, \mathbf{N}} \approx \tilde{\mathbf{A}}_{\text{eff}, \mathbf{N}}^{\text{lin}, \tilde{\mathbf{N}}}$ and $\underline{\mathbf{A}}_{\text{eff}, \mathbf{N}} \approx \tilde{\mathbf{A}}_{\text{eff}, \mathbf{N}}^{\text{lin}, \tilde{\mathbf{N}}}$. The exception is the rate $\bar{\mu}_{\mathbf{N}}$ in (4.2), calculated for circle (C), that can be compared with its approximation $\tilde{\mu}_{\mathbf{N}}^{\text{lin}, \tilde{\mathbf{N}}}$ in (4.3), see Tab. 4.3 and 4.4. Both the rates coincide for two digits up to $\mathbf{N} = (45, 45)$ for the odd case and up to $\mathbf{N} = (64, 64)$ for the even case. It can also be compared in terms of errors (4.4) shown in Fig. 4.4 — the corresponding lines differ significantly only for $\mathbf{N} = (1215, 1215)$ and $\mathbf{N} = (1215, 1215)$ resp.

Further, the errors in Fig. 4.4 — observed for problems (E), (H), and (M) — reach a limit state about the order 10^{-6} . It is primarily caused by an inaccuracy in approximation $\mathbf{e}^{(1)} \approx \mathbf{e}_{\mathbf{N}}^{(1)}$. It also depreciate the rates in Tab. 4.3, 4.4. However, the rates give evidence in the parts where errors in Fig. 4.4 produces the straight lines.

Concluding, the rates of convergence, calculated as in Eq. (4.1), (4.2), and (4.3), are comparable and depend on the regularity of material coefficients. Particularly, problem (C) reaches the rate $\frac{1}{2}$, while (E) $\frac{3}{2}$, and the last (H) and (M) even higher — about $\frac{5}{2}$ — however, this value is highly influenced by the inaccuracy in its determination.

4.3. Comparison with p-version of FEM. In this section, we compare the upper-lower bounds calculated with FFT-based FEM to p-version of FEM by [6]. The comparison is made for material coefficients defined according to Rem. 4.1 through scalar functions:

- circle (C): $f(\mathbf{x}) = \begin{cases} 1, & \|\mathbf{x}\|_2 < \pi^{-\frac{1}{2}} \\ 0, & \text{otherwise} \end{cases}$
- square (S): $f(\mathbf{x}) = \begin{cases} 1, & \|\mathbf{x}\|_\infty \leq \frac{1}{2} \\ 0, & \text{otherwise} \end{cases}$
- rectangle (R): $f(\mathbf{x}) = \begin{cases} 1, & |x_1| \leq \frac{\sqrt{2}}{4} \wedge |x_2| \leq \frac{\sqrt{2}}{2} \\ 0, & \text{otherwise} \end{cases}$

Contrary to the previous examples, the conductivity of circle, square, and rectangle inclusion is set to 10^{-3} that intends to model void. The number of degrees of freedom

for the p-version of FEM reaches 959 for circle and 511 for square and rectangle while the number of discretization points for FFT-based FEM was taken as $\mathbf{N} = (1215, 1215)$, hence substantially larger.

However, p-version of FEM is still significantly better in the terms of errors η , see Tab. 4.5. It is mainly caused by better approximation properties of the p-version of FEM, the shapes of inclusions are well suited for this method. Although, p-version of FEM only approximates circle inclusion contrary to the FFT-based method. Moreover, the FFT-based FEM is influenced by variational crime caused by the numerical integration in Def. 3.19.

TABLE 4.5

The comparison of the homogenized matrices between the p-version of FEM and the FFT-based FEM

method problem\property	p-FEM $A_{\text{eff},11}^{\text{FEM}}$	p-FEM $A_{\text{eff},22}^{\text{FEM}}$	p-FEM η^{FEM}	FFTH $A_{\text{eff},\mathbf{N},11}$	FFTH $A_{\text{eff},\mathbf{N},22}$	FFTH $\eta_{\mathbf{N}}$
(C)	0.600	0.600	1.166582e-04	0.588	0.588	2.676e-02
(S)	0.578	0.578	3.611941e-03	0.476	0.476	2.149e-01
(R)	0.425	0.671	6.403170e-03	0.346	0.569	1.951e-01

Nevertheless, the p-version of FEM with divergence-free subspaces is mostly suitable to 2-dimensional problems. It is an opportunity for the FFT-based FEM that manage higher dimensional problems without any additional effort, especially for data provided as voxel images.

Acknowledgments. The authors are thankful to Jaroslav Haslinger for pointing out the works of Jan Dvořák, [5, 6].

Appendix A. Continuous projections on solenoidal and curl-free spaces.

DEFINITION A.1. For $i \in \{0, 1, 2\}$, we define operators $\mathcal{G}^{(i)}[\cdot] : L_{\text{per}}^2(\mathcal{Y}; \mathbb{R}^d) \rightarrow L_{\text{per}}^2(\mathcal{Y}; \mathbb{R}^d)$ as convolution

$$\mathcal{G}^{(i)}[\mathbf{v}](\mathbf{x}) := \int_{\mathcal{Y}} \Gamma^{(i)}(\mathbf{x} - \mathbf{y}) \mathbf{v}(\mathbf{y}) d\mathbf{y} = \sum_{\mathbf{n} \in \mathbb{Z}^d} \hat{\Gamma}^{(i)}(\mathbf{n}) \hat{\mathbf{v}}(\mathbf{n}) \varphi_{\mathbf{n}}(\mathbf{x})$$

where $\hat{\mathbf{v}}_{\alpha} := (v_{\alpha}, \varphi_{\mathbf{n}})_{L_{\text{per}}^2}$ denotes the Fourier coefficients for $\varphi_{\mathbf{n}}(\mathbf{x}) = \exp(i\pi \sum_{\alpha=1}^d \frac{k_{\alpha} x_{\alpha}}{Y_{\alpha}})$.

Integral kernels $\Gamma^{(i)}$ are easily expressed in the Fourier space; the matrices $\hat{\Gamma}^{(i)}(\mathbf{n}) \in \mathbb{R}^{d \times d}$ of the Fourier coefficients reads

$$\hat{\Gamma}^{(0)}(\mathbf{n}) = \begin{cases} \mathbf{I} \\ \mathbf{0} \end{cases} \quad \hat{\Gamma}^{(1)}(\mathbf{n}) = \begin{cases} \mathbf{0} \\ \frac{\boldsymbol{\xi}(\mathbf{n}) \otimes \boldsymbol{\xi}(\mathbf{n})}{\boldsymbol{\xi}(\mathbf{n}) \cdot \boldsymbol{\xi}(\mathbf{n})} \end{cases} \quad \hat{\Gamma}^{(2)}(\mathbf{n}) = \begin{cases} \mathbf{0}, & \text{for } \mathbf{n} = \mathbf{0} \\ \mathbf{I} - \frac{\boldsymbol{\xi}(\mathbf{n}) \otimes \boldsymbol{\xi}(\mathbf{n})}{\boldsymbol{\xi}(\mathbf{n}) \cdot \boldsymbol{\xi}(\mathbf{n})}, & \text{for } \mathbf{n} \in \mathbb{Z}^d \setminus \{\mathbf{0}\} \end{cases}$$

where $\mathbf{I} \in \mathbb{R}^{d \times d}$ is the identity matrix, $\mathbf{0}$ denotes either a vector or a matrix with zero components, and $\xi_{\alpha}(\mathbf{n}) = \frac{n_{\alpha}}{Y_{\alpha}}$ for period \mathbf{Y} being consistent with periodic unit cell $\mathcal{Y} = \prod_{\alpha} (-Y_{\alpha}, Y_{\alpha}) \subset \mathbb{R}^d$.

LEMMA A.2. Operators $\mathcal{G}^{(i)}$ from previous definition are mutually orthogonal projections on, step-by-step, \mathcal{U}, \mathcal{E} , and \mathcal{J} — the subspaces of $L_{\text{per}}^2(\mathcal{Y}; \mathbb{R}^d)$ defined in Eq. (2.1).

Proof. In [31], we show in detail that $\mathcal{G}^{(1)}$ is a projection onto \mathcal{E} , the other cases are analogical. It is based on mutual orthogonality of $\hat{\Gamma}^{(i)}(\mathbf{n})$ for particular $\mathbf{n} \in \mathbb{Z}^d$ that can be found for example in [17]. \square

REFERENCES

- [1] A. BENSOUSSAN, G. PAPANICOLAOU, AND J. LIONS, *Asymptotic analysis for periodic structures*, vol. 5, North Holland, 1978.
- [2] S. BRISARD AND L. DORMIEUX, *FFT-based methods for the mechanics of composites: A general variational framework*, Computational Materials Science, 49 (2010), pp. 663–671.
- [3] ———, *Combining Galerkin approximation techniques with the principle of Hashin and Shtrikman to derive a new FFT-based numerical method for the homogenization of composites*, Computer Methods in Applied Mechanics and Engineering, (2012).
- [4] D. CIORANESCU AND P. DONATO, *An Introduction to Homogenization*, Oxford Lecture Series in Mathematics and Its Applications, Oxford University Press, 1999.
- [5] J. DVOŘÁK, *Optimization of composite materials*, Master's thesis, The Charles University in Prague, June 1993.
- [6] ———, *A reliable numerical method for computing homogenized coefficients*, tech. rep., available at CiteSeerX <http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.45.1190.1995>.
- [7] I. EKELAND AND R. TEMAM, *Convex Analysis and Variational Problems*, SIAM, 1976.
- [8] 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 (1999), pp. 41–47.
- [9] Z. HASHIN AND S. SHTRIKMAN, *On some variational principles in anisotropic and nonhomogeneous elasticity*, Journal of the Mechanics and Physics of Solids, 10 (1962), pp. 335–342.
- [10] ———, *A variational approach to the theory of the effective magnetic permeability of multiphase materials*, Journal of Applied Physics, 33 (1962), pp. 3125–3131.
- [11] ———, *A variational approach to the theory of the elastic behaviour of multiphase materials*, Journal of the Mechanics and Physics of Solids, 11 (1963), pp. 127–140.
- [12] H. HOANG-DUC AND G. BONNET, *Effective properties of viscoelastic heterogeneous periodic media: an approximate solution accounting for the distribution of heterogeneities.*, Mechanics of Materials, (2012).
- [13] V. JIKOV, S. KOZLOV, AND O. OLEINIK, *Homogenization of Differential Operators and Integral Functionals*, Springer-Verlag, 1994.
- [14] S. LEE, R. LEBENSOHN, AND A. ROLLETT, *Modeling the viscoplastic micromechanical response of two-phase materials using Fast Fourier Transforms*, International Journal of Plasticity, 27 (2011), pp. 707–727.
- [15] J. LI, X. TIAN, AND R. ABDELMOULA, *A damage model for crack prediction in brittle and quasi-brittle materials solved by the FFT method*, International journal of fracture, (2012), pp. 1–12.
- [16] J. C. MICHEL, H. MOULINEC, AND P. SUQUET, *A computational method based on augmented Lagrangians and fast Fourier transforms for composites with high contrast*, CMES-Computer Modeling in Engineering & Sciences, 1 (2000), pp. 79–88.
- [17] G. W. MILTON, *The Theory of Composites*, vol. 6 of Cambridge Monographs on Applied and Computational Mathematics, Cambridge University Press, Cambridge, UK, 2002.
- [18] V. MONCHIET AND G. BONNET, *A polarization-based FFT iterative scheme for computing the effective properties of elastic composites with arbitrary contrast*, International Journal for Numerical Methods in Engineering, 89 (2012), pp. 1419–1436.
- [19] H. MOULINEC AND P. SUQUET, *A fast numerical method for computing the linear and nonlinear mechanical properties of composites*, Comptes rendus de l'Académie des sciences. Série II, Mécanique, physique, chimie, astronomie, 318 (1994), pp. 1417–1423.
- [20] ———, *A numerical method for computing the overall response of nonlinear composites with complex microstructure*, Computer Methods in Applied Mechanics and Engineering, 157 (1997), pp. 69–94.
- [21] J. NOVÁK, A. KUČEROVÁ, AND J. ZEMAN, *Microstructural enrichment functions based on stochastic wang tilings*, arXiv preprint arXiv:1110.4183, (2011).
- [22] J. NĚMEČEK, V. KRÁLÍK, AND J. VONDŘEJC, *Micromechanical analysis of heterogeneous structural materials*, Cement and Concrete Composites, (2012).
- [23] J. NĚMEČEK, V. KRÁLÍK, AND J. VONDŘEJC, *A two-scale micromechanical model for aluminium foam based on results from nanoindentation*, (2012). Submitted.
- [24] J. ODEN, T. BELYTSCHKO, I. BABUŠKA, AND T. HUGHES, *Research directions in computational mechanics*, Computer Methods in Applied Mechanics and Engineering, 192 (2003), pp. 913–922.
- [25] A. REUSS AND Z. ANGNEW, *A calculation of the bulk modulus of polycrystalline materials*, Math Meth, 9 (1929), p. 55.

- [26] J. SARANEN AND G. VAINIKKO, *Periodic Integral and Pseudodifferential Equations with Numerical Approximation*, Springer Monographs Mathematics, 2000.
- [27] P. SUQUET, *A dual method in homogenization: application to elastic media*, J. Mech. Theor. Appl, 79 (1982), p. 98.
- [28] P. SUQUET, H. MOULINEC, O. CASTELNAU, M. MONTAGNAT, N. LAHELLEC, F. GRENNERAT, P. DUVAL, AND R. BRENNER, *Multi-scale modeling of the mechanical behavior of polycrystalline ice under transient creep*, Procedia IUTAM, 3 (2012), pp. 64–78.
- [29] W. VOIGT, *Lehrbuch der kristallphysik*, vol. 34, BG Teubner, 1910.
- [30] J. VONDŘEJC, J. ZEMAN, AND I. MAREK, *Analysis of a Fast Fourier Transform based method for modeling of heterogeneous materials*, Lecture Notes in Computer Science, 7116 (2012), pp. 512–522.
- [31] ———, *FFT-based finite element method for homogenization*, (2013). In preparation.
- [32] Z. WIEÇKOWSKI, *Dual finite element methods in mechanics of composite materials*, Journal of Theoretical and Applied Mechanics, 2 (1995), pp. 233–252.
- [33] 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 (2010), pp. 8065–8071.