



**HAL**  
open science

# Investigation of the interactions between the numerical and the modeling errors in the Homogenized Dirichlet Projection Method

Nicolas Moës, John Tinsley Oden, Tarek Zohdi

► **To cite this version:**

Nicolas Moës, John Tinsley Oden, Tarek Zohdi. Investigation of the interactions between the numerical and the modeling errors in the Homogenized Dirichlet Projection Method. *Computer Methods in Applied Mechanics and Engineering*, 1998, 159 (1-2), pp.79-101. 10.1016/S0045-7825(98)80104-7 . hal-01007001

**HAL Id: hal-01007001**

**<https://hal.science/hal-01007001>**

Submitted on 11 Jan 2023

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

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



Distributed under a Creative Commons Attribution - NonCommercial 4.0 International License

# Investigation of the interactions between the numerical and the modeling errors in the Homogenized Dirichlet Projection Method

Nicolas Moës \*, J. Tinsley Oden, Tarek I. Zohdi

*The Texas Institute for Computational and Applied Mathematics, The University of Texas at Austin, Taylor Hall 2400,  
Austin, TX 78712, USA*

The Homogenized Dirichlet Projection Method, HDPM, was developed in [1,2] as a systematic technique for analyzing highly-heterogeneous elastic structures. The method can provide an analysis of structures composed of composite materials with very complex microstructure at a fraction of the cost of solving the full fine-scale model. In the present investigation, the HDPM is revisited to take into account the unavoidable numerical errors produced in finite element approximations of the associated boundary value problems.

The total error of the HDPM, which now takes into account both for the modeling and numerical errors, is split into several terms, each accounting for a parameter in the method. The parameters are: the choice of the homogenized material property, the partition into subdomains, the coarse finite element mesh used to solve the homogenized problem and the fine meshes used to solve the subdomain problems. Numerical experiments are carried out on 1-D problems for which the exact solutions are easily calculated. The experiments reveal that the influence of the coarse and fine meshes are very different.

When free of numerical error the HDPM is based on four main results. We rewrite these in the framework of the error in the constitutive law. This leads to a clear mechanical interpretation of the results. Moreover, it allows us to extend the results to nonlinear constitutive laws and to find new properties of the HDPM. Some of the theoretical results are validated for specific cases involving known numerical error.

Finally, explicit a posteriori upper bounds are derived for the total and numerical error. A simple adaptive strategy is presented for choosing the fine mesh and the subdomain partition. The strategy is tested on a 1-D model problem.

## 1. Introduction

The Homogenized Dirichlet Projection Method (HDPM) [1,2] is an efficient method for solving large heterogeneous elastic problems. Basically, the method consists first in solving a homogenized problem with uniform material properties. Then, by computing an explicit error bound, the distance between the homogenized and exact solution is estimated. On subdomains where this distance is too large, a local analysis is performed using the homogenized solution as boundary conditions on the subdomains. This produces a perturbation of the homogenized solution which is closer to the exact solution. We refer as HDPM solution, the perturbation of the homogenized solution. Thus, the method introduces two types of problems: a regularized one with a uniform material property and a local analysis on subdomains

---

\*Corresponding author.

with the real fine-scale material properties. In practice, these two problems cannot be solved exactly, therefore numerical errors are introduced.

This paper is the first experimental investigation on the interactions between the numerical and modeling errors in the HDPM. Some theoretical results were already presented in [2]. We consider that the homogenized solution is obtained through the classical finite element method using a coarse mesh ( $H$ -mesh), and that the local solutions, over the subdomains, are obtained with fine meshes ( $h$ -meshes). In this investigation, we neglect the error introduced by the use of iterative solvers and numerical integration of the stiffness matrices. In order to isolate the numerical influence, we split the total error into several terms. Numerical experiments are then carried out on 1-D problems, for which the exact solutions are easily calculated.

The HDPM is based on four results [1,2]:

- an explicit upper bound of the distance between the homogenized solution and the exact fine-scale solution in the energy norm;
- the fact that the HDPM solution is closer in the energy norm to the exact solution than the original homogenized solution;
- a local bound on the difference between the HDPM and the homogenized solutions;
- an upper bound on the difference between the HDPM and exact solution in the energy norm.

These four results hold when the method is free of numerical errors. An important question then is to determine what can be said in the case of numerical errors. This question is also addressed in the present paper.

The exact and homogenized solutions differ because they do not satisfy the same constitutive law. Thus, the quality of the homogenized solution may be estimated by the way it satisfies the exact constitutive law. This fact makes it possible to rewrite all the main results of the HDPM using the error-in-the-constitutive-law concept [3], thereby leading to a clear mechanical interpretation of these results. Moreover, written in this framework, the HDPM may be extended to problems involving nonlinear constitutive laws.

Since the modeling error may be expressed as an error in the constitutive law, we are naturally led to estimate the numerical error in the same way [3]. Computable a posteriori upper bounds are obtained for the total and numerical errors. The estimates are then used in a simple adaptive strategy to adapt the  $h$ -mesh size and the subdomain size.

The paper is organized as follows. In Section 2, an outline of the HDPM is given, which takes into account the numerical approximations. A decomposition of the total error is also presented. In Section 3, numerical experiments are carried out on a 1-D problem to analyze the influence of the numerical errors. In Section 4, the main results of the HDPM are restated in the error-in-the-constitutive-law framework. They are then studied in the presence of numerical errors in Section 5. Finally, a posteriori error estimation and a simple adaptive strategy for choosing the mesh for the local problems and the subdomain partition are proposed and tested on 1-D problems in Section 6.

## 2. HDPM with finite element approximation

### 2.1. The reference problem

We consider a material body composed of a linearly-elastic material in static equilibrium under the action of given body forces  $\mathbf{f}_g$  and surface tractions  $\mathbf{t}_g$ . The domain,  $\Omega$ , occupied by the material body is considered regular: a simply-connected domain with Lipschitz boundary  $\partial\Omega$ . The boundary  $\partial\Omega$  consists of a portion  $\Gamma_u$  where displacements  $\mathbf{u}_g$  are prescribed and a portion  $\Gamma_t$  where tractions  $\mathbf{t}_g$  are prescribed,

$$\partial\Omega = \overline{\Gamma_u \cup \Gamma_t}, \quad \Gamma_u \cap \Gamma_t = \emptyset.$$

The problem to be solved on  $\Omega$  is to find a triple  $(\mathbf{u}, \boldsymbol{\epsilon}, \boldsymbol{\sigma})$  such that the kinematic constraints (1), the equilibrium equation (2) and the constitutive law (3) hold:

$$\boldsymbol{\epsilon} = \frac{1}{2}(\nabla \mathbf{u} + (\nabla \mathbf{u})^t) \quad \text{on } \Omega \quad \text{and} \quad \mathbf{u} = \mathbf{u}_g \quad \text{on } \Gamma_u; \quad (1)$$

$$\nabla \cdot \boldsymbol{\sigma} = -\mathbf{f}_g \quad \text{on } \Omega \quad \text{and} \quad \boldsymbol{\sigma} \cdot \mathbf{n} = \mathbf{t}_g \quad \text{on } \Gamma_t; \quad (2)$$

$$\boldsymbol{\sigma} = \mathbf{E} \boldsymbol{\epsilon} \quad \text{on } \Omega. \quad (3)$$

Here  $\mathbf{u}$ ,  $\boldsymbol{\epsilon}$  and  $\boldsymbol{\sigma}$  are the displacement, strain and stress, respectively and  $\mathbf{n}$  is the outward normal to the boundary. The elasticity tensor  $\mathbf{E}$  is a function of the position  $\mathbf{x}$  i.e.  $\mathbf{E} = \mathbf{E}(\mathbf{x})$ , and describes the microstructure of the material.

We remark that by eliminating  $\boldsymbol{\epsilon}$  and  $\boldsymbol{\sigma}$  in (1)–(3), and by taking into account the symmetries of  $\mathbf{E}$ , we have the classical elasticity problem for the displacement field,

$$-\nabla \cdot (\mathbf{E} \nabla \mathbf{u}) = \mathbf{f}_g \quad \text{on } \Omega, \quad \mathbf{u} = \mathbf{u}_g \quad \text{on } \Gamma_u, \quad \mathbf{E} \nabla \mathbf{u} \cdot \mathbf{n} = \mathbf{t}_g \quad \text{on } \Gamma_t. \quad (4)$$

A weak formulation of (4) is as follows:

$$\text{Find } \mathbf{u} \in V \text{ such that } \mathcal{B}(\mathbf{u}, \mathbf{v}) = \mathcal{F}(\mathbf{v}) \quad \forall \mathbf{v} \in V^o. \quad (5)$$

where

$$V = \{\mathbf{v} \in \mathbf{H}^1(\Omega) : \mathbf{v} = \mathbf{u}_g \text{ on } \Gamma_u\}, \quad V^o = \{\mathbf{v} \in \mathbf{H}^1(\Omega) : \mathbf{v} = 0 \text{ on } \Gamma_u\}, \quad (6)$$

$$\mathcal{B}(\mathbf{u}, \mathbf{v}) = \int_{\Omega} (\mathbf{E} \nabla \mathbf{u}) : \nabla \mathbf{v} \, d\mathbf{x}, \quad \mathcal{F}(\mathbf{v}) = \int_{\Omega} \mathbf{f}_g \cdot \mathbf{v} \, d\mathbf{x} + \int_{\Gamma_t} \mathbf{t}_g \cdot \mathbf{v} \, d\mathbf{s}. \quad (7)$$

## 2.2. The Homogenized Dirichlet Projection Method

Taking into account the unavoidable numerical errors, the HDPM can be summarized by four steps. The presentation which follows is valid for any space dimension but is illustrated for the 1-D problem, Fig. 1.

- Step 1: The so-called fine-scale problem, that is characterized by a highly heterogeneous material, Fig. 1, is replaced by an homogenized problem, Fig. 2, written:

Find a triple  $(\mathbf{u}^0, \boldsymbol{\epsilon}^0, \boldsymbol{\sigma}^0)$  such that

$$\boldsymbol{\epsilon}^0 = \frac{1}{2}(\nabla \mathbf{u}^0 + (\nabla \mathbf{u}^0)^t) \quad \text{on } \Omega \quad \text{and} \quad \mathbf{u}^0 = \mathbf{u}_g \quad \text{on } \Gamma_u; \quad (8)$$

$$\nabla \cdot \boldsymbol{\sigma}^0 = -\mathbf{f}_g \quad \text{on } \Omega \quad \text{and} \quad \boldsymbol{\sigma}^0 \cdot \mathbf{n} = \mathbf{t}_g \quad \text{on } \Gamma_t; \quad (9)$$

$$\boldsymbol{\sigma}^0 = \mathbf{E}^0 \boldsymbol{\epsilon}^0 \quad \text{on } \Omega. \quad (10)$$

$\mathbf{E}^0$  is the homogenized elasticity tensor, which, for convenience, is assumed to be constant over the domain. We note that a weak form of (8)–(10) analogous to (5) can be obtained with  $\mathcal{B}(\mathbf{u}, \mathbf{v})$  replaced by  $\mathcal{B}_0(\mathbf{u}^0, \mathbf{v}) = \int_{\Omega} (\mathbf{E}^0 \nabla \mathbf{u}^0) : \nabla \mathbf{v} \, d\mathbf{x}$ .

- Step 2: In general, the homogenized problem cannot be solved exactly, so numerical approximations are introduced. For example, using the finite element method on a mesh parameterized by a mesh size  $H$  (Fig. 3), the finite element solution is denoted by  $(\mathbf{u}^{0,H}, \boldsymbol{\epsilon}^{0,H}, \boldsymbol{\sigma}^{0,H})$  and satisfies

$$\boldsymbol{\epsilon}^{0,H} = \frac{1}{2}(\nabla \mathbf{u}^{0,H} + (\nabla \mathbf{u}^{0,H})^t) \quad \text{on } \Omega \quad \text{and} \quad \mathbf{u}^{0,H} = \mathbf{u}_g \quad \text{on } \Gamma_u;$$



Fig. 1. The reference problem: geometry, loading and material property. The exact solution is denoted by  $(\mathbf{u}, \boldsymbol{\epsilon}, \boldsymbol{\sigma})$ .

Fig. 2. The homogenized problem whose exact solution is denoted by  $(\mathbf{u}^0, \boldsymbol{\epsilon}^0, \boldsymbol{\sigma}^0)$ .

$$\int_{\Omega} \boldsymbol{\sigma}^{0,H} : \nabla \mathbf{v} \, d\mathbf{x} - \int_{\Omega} \mathbf{f}_g \cdot \mathbf{v} \, d\mathbf{x} - \int_{\Gamma_t} \mathbf{t}_g \cdot \mathbf{v} \, d\mathbf{s} = 0 \quad \forall \mathbf{v} \in V^{0,H}; \quad (11)$$

$$\boldsymbol{\sigma}^{0,H} = \mathbf{E}^0 \boldsymbol{\epsilon}^{0,H} \quad \text{on } \Omega \quad (12)$$

where  $V^{0,H} \subset V^0$  is an admissible finite element displacement space parameterized by the mesh size  $H$ . We assume that  $V^{0,H}$  belongs to a family of subspaces constructed so that  $\mathbf{u}^{0,H} \rightarrow \mathbf{u}^0 \in V$  as  $H \rightarrow 0$ . Also note that the finite element strain and stress are only well defined on the interior of each element. In other words, there are jumps in the tractions.

- Step 3: Once the homogenized problem is solved, the domain is partitioned into  $N$  subdomains  $\Omega_k$ ,  $k = 1, 2, \dots, N$ ,

$$\bigcup_{k=1}^N \overline{\Omega_k} = \overline{\Omega}, \quad \Omega_i \cap \Omega_j = \emptyset, \quad i \neq j.$$

Fig. 4 shows a uniform partition of the domain in subdomains of size  $\Delta$ . The subdomains are considered regular, simply connected, with Lipschitz boundary. On each subdomain it is possible to estimate the quality of the homogenized solution [1]. On subdomains where the quality of the homogenized solution is poor compared with the exact solution, a local analysis is performed. The local problem on the  $k^{\text{th}}$  subdomain,  $\Omega_k$ , reads:

Find the triple  $(\tilde{\mathbf{u}}_k^{0,H}, \tilde{\boldsymbol{\epsilon}}_k^{0,H}, \tilde{\boldsymbol{\sigma}}_k^{0,H})$  such that

$$\tilde{\boldsymbol{\epsilon}}_k^{0,H} = \frac{1}{2} (\nabla \tilde{\mathbf{u}}_k^{0,H} + (\nabla \tilde{\mathbf{u}}_k^{0,H})^t) \quad \text{on } \Omega_k \quad \text{and} \quad \tilde{\mathbf{u}}_k^{0,H} = \mathbf{u}^{0,H} \quad \text{on } \partial\Omega_k \setminus \partial\Omega, \quad (13)$$

$$\tilde{\mathbf{u}}_k^{0,H} = \mathbf{u}_g \quad \text{on } \partial\Omega_k \cap \Gamma_u; \quad (14)$$

$$\nabla \cdot \tilde{\boldsymbol{\sigma}}_k^{0,H} = -\mathbf{f}_g \quad \text{on } \Omega_k \quad \text{and} \quad \tilde{\boldsymbol{\sigma}}_k^{0,H} \cdot \mathbf{n} = \mathbf{t}_g \quad \text{on } \partial\Omega_k \cap \Gamma_t; \quad (15)$$

$$\tilde{\boldsymbol{\sigma}}_k^{0,H} = \mathbf{E} \tilde{\boldsymbol{\epsilon}}_k^{0,H} \quad \text{on } \Omega_k \quad (16)$$

where  $\partial\Omega_k$  denotes the boundary of  $\Omega_k$ . The homogenized displacements,  $\mathbf{u}^{0,H}$ , are imposed on the boundary of the subdomain, except on the boundary,  $\partial\Omega$ , where the actual boundary conditions are satisfied. Note also that the local analysis is solved with the actual microstructure described by  $\mathbf{E}$ .

For the sake of simplicity, we assume that the local analysis is performed on every subdomain, but this is seldom the case in application of the HDPM. The final solution  $(\tilde{\mathbf{u}}^{0,H}, \tilde{\boldsymbol{\epsilon}}^{0,H}, \tilde{\boldsymbol{\sigma}}^{0,H})$  over the domain is then constructed in the following manner

$$\begin{aligned} \tilde{\mathbf{u}}^{0,H} &= \mathbf{u}^{0,H} + \sum_{k=1}^N \mathcal{E}_k (\tilde{\mathbf{u}}_k^{0,H} - \mathbf{u}^{0,H}); \\ \tilde{\boldsymbol{\epsilon}}^{0,H} &= \boldsymbol{\epsilon}^{0,H} + \sum_{k=1}^N \mathcal{E}_k (\tilde{\boldsymbol{\epsilon}}_k^{0,H} - \boldsymbol{\epsilon}^{0,H}); \\ \tilde{\boldsymbol{\sigma}}^{0,H} &= \boldsymbol{\sigma}^{0,H} + \sum_{k=1}^N \mathcal{E}_k (\tilde{\boldsymbol{\sigma}}_k^{0,H} - \boldsymbol{\sigma}^{0,H}). \end{aligned}$$

where  $\mathcal{E}_k$  is a scalar function defined on  $\Omega$  with a value of unity on  $\Omega_k$  and zero elsewhere.

Let us define  $\Gamma_{\text{int}}$  as the union of the boundaries of all the subdomain modulo the boundary of the domain

$$\Gamma_{\text{int}} = \left( \bigcup_{k=1}^N \partial\Omega_k \right) \setminus \partial\Omega.$$

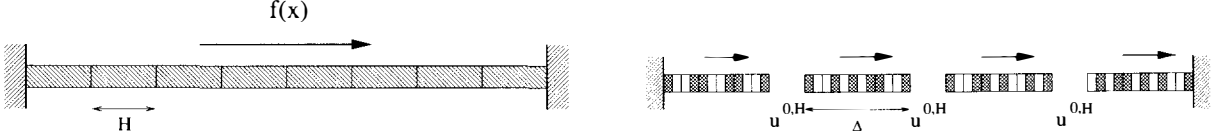


Fig. 3. The finite element problem associated with the homogenized problem. The mesh size is  $H$  and the finite element solution is  $(\mathbf{u}^{0,H}, \boldsymbol{\epsilon}^{0,H}, \boldsymbol{\sigma}^{0,H})$ .

Fig. 4. The subdomain problems. The subdomain size is  $\Delta$ . The finite element displacements  $\mathbf{u}^{0,H}$  are applied at the boundaries of the subdomains, in the interior of the domain. The exact solution of the subdomain problems is denoted by  $(\bar{\mathbf{u}}^{0,H}, \bar{\boldsymbol{\epsilon}}^{0,H}, \bar{\boldsymbol{\sigma}}^{0,H})$ .

By construction, the displacement field  $\bar{\mathbf{u}}^{0,H}$  is continuous across  $\Gamma_{\text{int}}$ . The normal stresses  $\bar{\boldsymbol{\sigma}}^{0,H} \mathbf{n}$  are not in general continuous across  $\Gamma_{\text{int}}$ . Let  $\Omega_1$  and  $\Omega_2$  be two subdomains having a common boundary. The traction jump  $[\bar{\boldsymbol{\sigma}}^{0,H} \cdot \mathbf{n}]$  on this boundary is defined by

$$[\bar{\boldsymbol{\sigma}}^{0,H} \cdot \mathbf{n}] = \bar{\boldsymbol{\sigma}}_1^{0,H} \cdot \mathbf{n}_1 + \bar{\boldsymbol{\sigma}}_2^{0,H} \cdot \mathbf{n}_2$$

where  $\mathbf{n}_i$  is the outward normal to the boundary for the subdomain  $i$  and  $\bar{\boldsymbol{\sigma}}_i^{0,H}$  is the stress state in the subdomain  $i$  on the boundary ( $i = 1, 2$ ).

Finally, if we assume that no numerical error was introduced in Step 2, the solution of the Step 3 just described will be denoted  $(\bar{\mathbf{u}}_k^0, \bar{\boldsymbol{\epsilon}}_k^0, \bar{\boldsymbol{\sigma}}_k^0)$  for each subdomain and  $(\bar{\mathbf{u}}^0, \bar{\boldsymbol{\epsilon}}^0, \bar{\boldsymbol{\sigma}}^0)$  on the whole domain. The solution  $(\bar{\mathbf{u}}_k^0, \bar{\boldsymbol{\epsilon}}_k^0, \bar{\boldsymbol{\sigma}}_k^0)$  satisfies problem (13)–(16) with all the  $H$  superscripts removed.

- Step 4: In fact, as in the case of the homogenized problem, the subdomain problems cannot be solved exactly, introducing further error components. On each subdomain,  $\Omega_k$ , the exact solution,  $(\bar{\mathbf{u}}_k^{0,H}, \bar{\boldsymbol{\epsilon}}_k^{0,H}, \bar{\boldsymbol{\sigma}}_k^{0,H})$ , of the problem (13)–(16) is approximated by a finite element solution parameterized by a mesh size  $h$  (Fig. 5). This approximate solution is denoted by  $(\tilde{\mathbf{u}}_k^{0,H,h}, \tilde{\boldsymbol{\epsilon}}_k^{0,H,h}, \tilde{\boldsymbol{\sigma}}_k^{0,H,h})$  and satisfies

$$\tilde{\boldsymbol{\epsilon}}_k^{0,H,h} = \frac{1}{2} (\nabla \tilde{\mathbf{u}}_k^{0,H,h} + (\nabla \tilde{\mathbf{u}}_k^{0,H,h})^t) \quad \text{on } \Omega_k, \quad \tilde{\mathbf{u}}_k^{0,H,h} = \mathbf{u}^{0,H} \quad \text{on } \partial\Omega_k \setminus \partial\Omega; \quad (17)$$

$$\tilde{\mathbf{u}}_k^{0,H,h} = \mathbf{u}_g \quad \text{on } \partial\Omega_k \cap \Gamma_u; \quad (18)$$

$$\int_{\Omega_k} \tilde{\boldsymbol{\sigma}}_k^{0,H,h} : \nabla \mathbf{v} \, d\mathbf{x} - \int_{\Omega_k} \mathbf{f}_g \cdot \mathbf{v} \, d\mathbf{x} - \int_{\Gamma_t \cap \partial\Omega_k} \mathbf{t}_g \cdot \mathbf{v} \, d\mathbf{s} = 0 \quad \forall \mathbf{v} \in V_k^{0,h}; \quad (19)$$

$$\tilde{\boldsymbol{\sigma}}_k^{0,H,h} = \mathbf{E} \tilde{\boldsymbol{\epsilon}}_k^{0,H,h} \quad \text{on } \Omega_k \quad (20)$$

where

$$V_k^{0,h} = \{\mathbf{v} \in V_k^h : \mathbf{v} = 0 \text{ on } \partial\Omega_k \setminus \Gamma_t\}.$$

$V_k^h$  is a finite element displacement space parameterized by the mesh size  $h$  on subdomain  $\Omega_k$ .

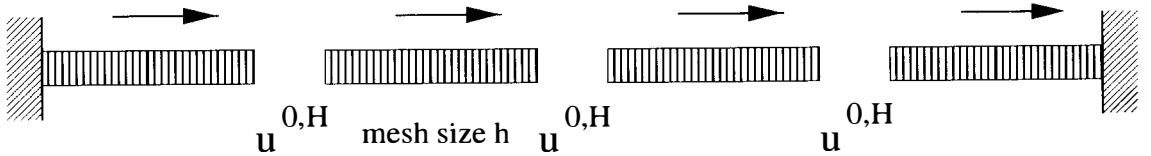


Fig. 5. The finite element problems associated to the subdomain problems. The mesh size is  $h$  and the finite element solution is denoted by  $(\tilde{\mathbf{u}}^{0,H,h}, \tilde{\boldsymbol{\epsilon}}^{0,H,h}, \tilde{\boldsymbol{\sigma}}^{0,H,h})$ .

The final solution of the HDPM is denoted by  $(\tilde{\mathbf{u}}^{0,H,h}, \tilde{\boldsymbol{\epsilon}}^{0,H,h}, \tilde{\boldsymbol{\sigma}}^{0,H,h})$  and is defined in the following manner:

$$\begin{aligned}\tilde{\mathbf{u}}^{0,H,h} &= \mathbf{u}^{0,H} + \sum_{k=1}^N \mathcal{E}_k(\tilde{\mathbf{u}}_k^{0,H,h} - \mathbf{u}^{0,H}); \\ \tilde{\boldsymbol{\epsilon}}^{0,H,h} &= \boldsymbol{\epsilon}^{0,H} + \sum_{k=1}^N \mathcal{E}_k(\tilde{\boldsymbol{\epsilon}}_k^{0,H,h} - \boldsymbol{\epsilon}^{0,H}); \\ \tilde{\boldsymbol{\sigma}}^{0,H,h} &= \boldsymbol{\sigma}^{0,H} + \sum_{k=1}^N \mathcal{E}_k(\tilde{\boldsymbol{\sigma}}_k^{0,H,h} - \boldsymbol{\sigma}^{0,H}).\end{aligned}$$

As before, we assume that  $V_k^h \subset V_k$  are members of appropriate families of subspace constructed so that  $\tilde{\mathbf{u}}_k^{0,H,h} \rightarrow \tilde{\mathbf{u}}_k^{0,H}$  in  $V_k$  as  $h \rightarrow 0$ , where

$$V_k = \{\mathbf{v} \in \mathbf{H}^1(\Omega_k) : \mathbf{v} = \mathbf{u}^{0,H} \text{ on } \partial\Omega_k \setminus \partial\Omega \text{ and } \mathbf{v} = \mathbf{u}_g \text{ on } \partial\Omega_k \cap \Gamma_u\}.$$

### 2.3. Decomposition of the error

Let us summarize the various solutions and their finite element approximations:

- $(\mathbf{u}, \boldsymbol{\epsilon} = \boldsymbol{\epsilon}(\mathbf{u}), \boldsymbol{\sigma} = \mathbf{E}\boldsymbol{\epsilon})$  = the exact solution, also called fine-scale solution;
  - $(\mathbf{u}^0, \boldsymbol{\epsilon}^0 = \boldsymbol{\epsilon}(\mathbf{u}^0), \boldsymbol{\sigma}^0 = \mathbf{E}^0\boldsymbol{\epsilon}^0)$  = the homogenized solution;
  - $(\mathbf{u}^{0,H}, \boldsymbol{\epsilon}^{0,H} = \boldsymbol{\epsilon}(\mathbf{u}^{0,H}), \boldsymbol{\sigma}^{0,H} = \mathbf{E}^0\boldsymbol{\epsilon}^{0,H})$  = the coarse-mesh finite element approximation of the homogenized solution;
  - $(\tilde{\mathbf{u}}^0, \tilde{\boldsymbol{\epsilon}}^0 = \boldsymbol{\epsilon}(\tilde{\mathbf{u}}^0), \tilde{\boldsymbol{\sigma}}^0 = \mathbf{E}\tilde{\boldsymbol{\epsilon}}^0)$  = the perturbed solution defined on  $N$  subdomains on which the exact  $\mathbf{u}^0$  is prescribed as boundary data, on the interior;
  - $(\tilde{\mathbf{u}}^{0,H}, \tilde{\boldsymbol{\epsilon}}^{0,H} = \boldsymbol{\epsilon}(\tilde{\mathbf{u}}^{0,H}), \tilde{\boldsymbol{\sigma}}^{0,H} = \mathbf{E}\tilde{\boldsymbol{\epsilon}}^{0,H})$  = the perturbed solution defined on the  $N$  subdomains on which the approximate homogenized solution  $\mathbf{u}^{0,H}$  is prescribed as boundary data, on the interior;
  - $(\tilde{\mathbf{u}}^{0,H,h}, \tilde{\boldsymbol{\epsilon}}^{0,H,h} = \boldsymbol{\epsilon}(\tilde{\mathbf{u}}^{0,H,h}), \tilde{\boldsymbol{\sigma}}^{0,H,h} = \mathbf{E}\tilde{\boldsymbol{\epsilon}}^{0,H,h})$  = the fine-mesh approximation of  $(\tilde{\mathbf{u}}^{0,H}, \tilde{\boldsymbol{\epsilon}}^{0,H}, \tilde{\boldsymbol{\sigma}}^{0,H})$ .
- $(\tilde{\mathbf{u}}^{0,H,h}, \tilde{\boldsymbol{\epsilon}}^{0,H,h}, \tilde{\boldsymbol{\sigma}}^{0,H,h})$  is referred as the HDPM solution. The error in the energy norm associated to this solution reads:

$$e = \|\mathbf{u} - \tilde{\mathbf{u}}^{0,H,h}\|_{E(\Omega)}, \quad \text{where } \|\mathbf{v}\|_{E(\Omega)}^2 = \int_{\Omega} \mathbf{E}\nabla\mathbf{v} : \nabla\mathbf{v} \, \mathbf{d}\mathbf{x}. \quad (21)$$

This error depends on four factors. Two are modeling-related: the homogenized material tensor  $\mathbf{E}^0$  and the partition in subdomains; and two are numerically-related: the  $H$ -mesh to solve the homogenized problem and the  $h$ -meshes to solve the subdomain problems.

We first separate the error into a modeling and a numerical part:

**PROPERTY 1.** If  $\tilde{\mathbf{u}}^{0,H,h} = \mathbf{u}^{0,H}$  on  $\Gamma_{\text{int}}$  then

$$e^2 = \underbrace{\|\mathbf{u} - \tilde{\mathbf{u}}^0\|_{E(\Omega)}^2}_{\text{modeling}} + \underbrace{\|\tilde{\mathbf{u}}^0 - \tilde{\mathbf{u}}^{0,H,h}\|_{E(\Omega)}^2}_{\text{numerical}} + 2 \underbrace{\int_{\Gamma_{\text{int}}} [\tilde{\boldsymbol{\sigma}}^0 \cdot \mathbf{n}] \cdot (\mathbf{u}^{0,H} - \mathbf{u}^0) \, \mathbf{d}\mathbf{s}}_{\text{coupling}}. \quad (22)$$

**PROOF.** Owing to the definition of the norm  $\|\cdot\|_{E(\Omega)}$  in (21), we have

$$e^2 = \|\mathbf{u} - \tilde{\mathbf{u}}^0\|_{E(\Omega)}^2 + \|\tilde{\mathbf{u}}^0 - \tilde{\mathbf{u}}^{0,H,h}\|_{E(\Omega)}^2 + 2 \int_{\Omega} \mathbf{E}(\nabla\tilde{\mathbf{u}}^0 - \nabla\mathbf{u}) : (\nabla\tilde{\mathbf{u}}^{0,H,h} - \nabla\tilde{\mathbf{u}}^0) \, \mathbf{d}\mathbf{x}.$$

The third term of the right hand side can be transformed as follows:

$$\int_{\Omega} \mathbf{E}(\nabla\tilde{\mathbf{u}}^0 - \nabla\mathbf{u}) : (\nabla\tilde{\mathbf{u}}^{0,H,h} - \nabla\tilde{\mathbf{u}}^0) \, \mathbf{d}\mathbf{x} = \int_{\Omega} (\tilde{\boldsymbol{\sigma}}^0 - \boldsymbol{\sigma}) : (\nabla\tilde{\mathbf{u}}^{0,H,h} - \nabla\tilde{\mathbf{u}}^0) \, \mathbf{d}\mathbf{x}$$

$$\begin{aligned}
&= - \underbrace{\int_{\Omega} (\nabla \cdot \tilde{\boldsymbol{\sigma}}^0 - \nabla \cdot \boldsymbol{\sigma}) \cdot (\tilde{\mathbf{u}}^{0,H,h} - \tilde{\mathbf{u}}^0) \, \mathbf{d}\mathbf{x}}_0 \\
&\quad + \int_{\Gamma_{\text{int}}} [(\tilde{\boldsymbol{\sigma}}^0 - \boldsymbol{\sigma}) \cdot \mathbf{n}] \cdot (\tilde{\mathbf{u}}^{0,H,h} - \tilde{\mathbf{u}}^0) \, \mathbf{d}\mathbf{s} \\
&= \int_{\Gamma_{\text{int}}} [\tilde{\boldsymbol{\sigma}}^0 \cdot \mathbf{n}] \cdot (\mathbf{u}^{0,H} - \mathbf{u}^0) \, \mathbf{d}\mathbf{s}.
\end{aligned}$$

We have used the fact that  $\mathbf{u}^0 = \tilde{\mathbf{u}}^0$  on  $\Gamma_{\text{int}}$ .  $\square$

The coupling term can be positive or negative. It represents the work done by the jump in traction  $[(\tilde{\boldsymbol{\sigma}}^0 \cdot \mathbf{n})]$  moving through the difference in the numerical and exact homogenized displacement. Thus, it depends both on the modeling and numerical errors. On the contrary, if we isolate the numerical error coming from the local analysis, we get the following decomposition:

*PROPERTY 2.* If  $\tilde{\mathbf{u}}^{0,H,h} = \tilde{\mathbf{u}}^{0,H}$  on  $\Gamma_{\text{int}}$ , then

$$e^2 = \underbrace{\|\mathbf{u} - \tilde{\mathbf{u}}^{0,H}\|_{E(\Omega)}^2}_{\text{mod.+num. (H-mesh)}} + \underbrace{\|\tilde{\mathbf{u}}^{0,H} - \tilde{\mathbf{u}}^{0,H,h}\|_{E(\Omega)}^2}_{\text{num. (h-mesh)}}. \quad (23)$$

*PROOF.* We have

$$e^2 = \|\mathbf{u} - \tilde{\mathbf{u}}^{0,H}\|_{E(\Omega)}^2 + \|\tilde{\mathbf{u}}^{0,H} - \tilde{\mathbf{u}}^{0,H,h}\|_{E(\Omega)}^2 + 2 \int_{\Omega} \mathbf{E}(\nabla \tilde{\mathbf{u}}^{0,H} - \nabla \mathbf{u}) : (\nabla \tilde{\mathbf{u}}^{0,H,h} - \nabla \tilde{\mathbf{u}}^{0,H}) \, \mathbf{d}\mathbf{x}$$

and

$$\begin{aligned}
\int_{\Omega} \mathbf{E}(\nabla \tilde{\mathbf{u}}^{0,H} - \nabla \mathbf{u}) : (\nabla \tilde{\mathbf{u}}^{0,H,h} - \nabla \tilde{\mathbf{u}}^{0,H}) \, \mathbf{d}\mathbf{x} &= \int_{\Omega} (\tilde{\boldsymbol{\sigma}}^{0,H} - \boldsymbol{\sigma}) : (\nabla \tilde{\mathbf{u}}^{0,H,h} - \nabla \tilde{\mathbf{u}}^{0,H}) \, \mathbf{d}\mathbf{x} \\
&= - \underbrace{\int_{\Omega} (\nabla \cdot \tilde{\boldsymbol{\sigma}}^{0,H} - \nabla \cdot \boldsymbol{\sigma}) \cdot (\tilde{\mathbf{u}}^{0,H,h} - \tilde{\mathbf{u}}^{0,H}) \, \mathbf{d}\mathbf{x}}_0 \\
&\quad + \int_{\Gamma_{\text{int}}} [(\tilde{\boldsymbol{\sigma}}^{0,H} - \boldsymbol{\sigma}) \cdot \mathbf{n}] \cdot (\tilde{\mathbf{u}}^{0,H,h} - \tilde{\mathbf{u}}^{0,H}) \, \mathbf{d}\mathbf{s} \\
&= 0. \quad \square
\end{aligned}$$

Let us now decompose the modeling and numerical errors appearing in Property 1. The modeling error can be decomposed into the error due to the homogenization process minus that gained by doing the local analysis (Property 3). The numerical error can be decomposed into the error due to the  $H$ -mesh approximation and the error due to the  $h$ -mesh approximation (Property 4).

*PROPERTY 3.*

$$\underbrace{\|\mathbf{u} - \tilde{\mathbf{u}}^0\|_{E(\Omega)}^2}_{\text{modeling error}} = \underbrace{\|\mathbf{u} - \mathbf{u}^0\|_{E(\Omega)}^2}_{\text{homogenization}} - \underbrace{\|\mathbf{u}^0 - \tilde{\mathbf{u}}^0\|_{E(\Omega)}^2}_{\text{sudom. analysis}}. \quad (24)$$

*PROOF.* We have

$$\|\mathbf{u} - \tilde{\mathbf{u}}^0\|_{E(\Omega)}^2 = \|\mathbf{u} - \mathbf{u}^0\|_{E(\Omega)}^2 - \|\mathbf{u}^0 - \tilde{\mathbf{u}}^0\|_{E(\Omega)}^2 + 2 \int_{\Omega} \mathbf{E}(\nabla \tilde{\mathbf{u}}^0 - \nabla \mathbf{u}) : (\nabla \tilde{\mathbf{u}}^0 - \nabla \mathbf{u}^0) \, \mathbf{d}\mathbf{x}$$

and

$$\int_{\Omega} \mathbf{E}(\nabla \tilde{\mathbf{u}}^0 - \nabla \mathbf{u}) : (\nabla \tilde{\mathbf{u}}^0 - \nabla \mathbf{u}^0) \, \mathbf{d}\mathbf{x} = \int_{\Omega} (\tilde{\boldsymbol{\sigma}}^0 - \boldsymbol{\sigma}) : (\nabla \tilde{\mathbf{u}}^0 - \nabla \mathbf{u}^0) \, \mathbf{d}\mathbf{x}$$



$$\begin{aligned}
&= \int_{\Gamma_{\text{int}}} [(\tilde{\boldsymbol{\sigma}}^0 - \boldsymbol{\sigma}) \cdot \mathbf{n}] \cdot (\tilde{\mathbf{u}}^0 - \mathbf{u}^0) \, ds \\
&= 0
\end{aligned}$$

owing to the fact that  $\mathbf{u}^0 = \tilde{\mathbf{u}}^0$  on  $\Gamma_{\text{int}}$ .  $\square$

*PROPERTY 4.* If  $\tilde{\mathbf{u}}^{0,H,h} = \tilde{\mathbf{u}}^{0,H}$  on  $\Gamma_{\text{int}}$  then

$$\underbrace{\|\tilde{\mathbf{u}}^0 - \tilde{\mathbf{u}}^{0,H,h}\|_{E(\Omega)}^2}_{\text{numerical}} = \underbrace{\|\tilde{\mathbf{u}}^0 - \tilde{\mathbf{u}}^{0,H}\|_{E(\Omega)}^2}_{\text{num. (H-mesh)}} + \underbrace{\|\tilde{\mathbf{u}}^{0,H} - \tilde{\mathbf{u}}^{0,H,h}\|_{E(\Omega)}^2}_{\text{num. (h-mesh)}}. \quad (25)$$

*PROOF.* We have

$$\begin{aligned}
\|\tilde{\mathbf{u}}^0 - \tilde{\mathbf{u}}^{0,H,h}\|_{E(\Omega)}^2 &= \|\tilde{\mathbf{u}}^0 - \tilde{\mathbf{u}}^{0,H}\|_{E(\Omega)}^2 + \|\tilde{\mathbf{u}}^{0,H} - \tilde{\mathbf{u}}^{0,H,h}\|_{E(\Omega)}^2 \\
&\quad + 2 \int_{\Omega} \mathbf{E}(\nabla \tilde{\mathbf{u}}^{0,H} - \nabla \tilde{\mathbf{u}}^0) : (\nabla \tilde{\mathbf{u}}^{0,H,h} - \nabla \tilde{\mathbf{u}}^{0,H}) \, dx
\end{aligned}$$

and

$$\begin{aligned}
\int_{\Omega} \mathbf{E}(\nabla \tilde{\mathbf{u}}^{0,H} - \nabla \tilde{\mathbf{u}}^0) : (\nabla \tilde{\mathbf{u}}^{0,H,h} - \nabla \tilde{\mathbf{u}}^{0,H}) \, dx &= \int_{\Omega} (\tilde{\boldsymbol{\sigma}}^{0,H} - \tilde{\boldsymbol{\sigma}}^0) : (\nabla \tilde{\mathbf{u}}^{0,H,h} - \nabla \tilde{\mathbf{u}}^{0,H}) \, dx \\
&= \int_{\Gamma_{\text{int}}} [(\tilde{\boldsymbol{\sigma}}^{0,H} - \tilde{\boldsymbol{\sigma}}^0) \cdot \mathbf{n}] \cdot (\tilde{\mathbf{u}}^{0,H,h} - \tilde{\mathbf{u}}^{0,H}) \, ds \\
&= 0.
\end{aligned}$$

Gathering Properties 1, 3 and 4, we finally get the following decomposition for the error:

*PROPERTY 5.* If  $\tilde{\mathbf{u}}^{0,H,h} = \tilde{\mathbf{u}}^{0,H} = \mathbf{u}^{0,H}$  on  $\Gamma_{\text{int}}$ , then

$$\begin{aligned}
e^2 &= \|\mathbf{u} - \mathbf{u}^0\|_{E(\Omega)}^2 - \|\mathbf{u}^0 - \tilde{\mathbf{u}}^0\|_{E(\Omega)}^2 + \|\tilde{\mathbf{u}}^0 - \tilde{\mathbf{u}}^{0,H}\|_{E(\Omega)}^2 + \|\tilde{\mathbf{u}}^{0,H} - \tilde{\mathbf{u}}^{0,H,h}\|_{E(\Omega)}^2 \\
&\quad + 2 \int_{\Gamma_{\text{int}}} [\tilde{\boldsymbol{\sigma}}^0 \cdot \mathbf{n}] \cdot (\mathbf{u}^{0,H} - \mathbf{u}^0) \, ds. \quad \square
\end{aligned} \quad (26)$$

Let us establish precisely the meaning of each of these terms:

- The first term,  $\|\mathbf{u} - \mathbf{u}^0\|_{E(\Omega)}$ , is the error introduced by replacing  $\mathbf{E}$  by  $\mathbf{E}^0$  in the original problem. It depends on  $\mathbf{E}^0$  and it is zero if  $\mathbf{E}^0$  and  $\mathbf{E}$  coincide.
- The second term,  $\|\mathbf{u}^0 - \tilde{\mathbf{u}}^0\|_{E(\Omega)}$ , is what we gain (minus sign in (26)) in solving the subdomain problems. It depends on  $\mathbf{E}^0$  and the subdomain partition, symbolically denoted by  $\Delta$ . As the subdomain size tends to zero,  $\Delta \rightarrow 0$  symbolically, this term tends to zero since  $\tilde{\mathbf{u}}^0 \rightarrow \mathbf{u}^0$  and it is zero if  $\mathbf{E}^0$  and  $\mathbf{E}$  coincide. Finally note that the two first terms in the right hand side of (26) form a quantity that is always greater or equal to zero, by Property 3.
- The third term,  $\|\tilde{\mathbf{u}}^0 - \tilde{\mathbf{u}}^{0,H}\|_{E(\Omega)}$ , takes into account the numerical errors introduced in solving the homogenized problems. It depends on  $\mathbf{E}^0$ ,  $\Delta$  and the  $H$ -mesh size. It tends to zero as  $H \rightarrow 0$  since  $\mathbf{u}^{0,H} \rightarrow \mathbf{u}^0$ . But is *not* zero in general if  $\mathbf{E}^0$  and  $\mathbf{E}$  coincide. We have

$$\|\tilde{\mathbf{u}}^0 - \tilde{\mathbf{u}}^{0,H}\|_{E(\Omega)}^2 = \int_{\Gamma_{\text{int}}} [(\tilde{\boldsymbol{\sigma}}^{0,H} - \tilde{\boldsymbol{\sigma}}^0) \cdot \mathbf{n}] \cdot (\mathbf{u}^{0,H} - \mathbf{u}^0) \, ds.$$

Thus, the third term depends on the quality of the finite element solution  $\mathbf{u}^{0,H}$  only on the boundary of the subdomains,  $\Gamma_{\text{int}}$ .

- The fourth term,  $\|\tilde{\mathbf{u}}^{0,H} - \tilde{\mathbf{u}}^{0,H,h}\|_{E(\Omega)}$ , takes into account the numerical errors introduced by solving approximately the subdomain problems. It tends to zero as  $h$  tends to zero since  $\tilde{\mathbf{u}}^{0,H,h} \rightarrow \tilde{\mathbf{u}}^{0,H}$  and it is not zero, in general, if  $\mathbf{E}^0$  and  $\mathbf{E}$  coincide.

- The last term in (26) is the coupling term whose sign can be positive or negative. It tends to zero as  $H \rightarrow 0$  and it is zero if  $\mathbf{E}^0$  and  $\mathbf{E}$  coincide. Thus, this term exhibits a coupling between the numerical and modeling errors.

### 3. Influence of the numerical errors

We study here the influence of the  $h$ - and  $H$ -mesh on the error in connection with a model problem. The influence of the subdomain size is also studied since it has many features similar to the influence of the  $h$ -mesh. The study is carried out on the 1-D problem described in Fig. 1. The problem to be solved on  $(0, 1)$  is to find a triple  $(u, \epsilon, \sigma)$  such that

$$\epsilon(x) = \frac{du(x)}{dx}, \quad x \in (0, 1) \quad \text{and} \quad u(0) = 0, \quad u(1) = 0; \quad (27)$$

$$\frac{d\sigma(x)}{dx} = -f(x), \quad f(x) = 1, \quad x \in (0, 1); \quad (28)$$

$$\sigma(x) = E\epsilon(x), \quad x \in (0, 1). \quad (29)$$

The Young modulus  $E = E(x)$  corresponds to a two-phase material with equal volumetric distribution of each phase. More precisely, we consider that the rod is composed of ‘particles’ of size  $d = 1/P$ ,  $P$  being the total number of particles. These particles have a Young modulus of  $E_1$  or  $E_2$ . The partition of the phases on the rod is generated randomly.

#### 3.1. Influence of the $h$ -mesh

The  $h$ -mesh size influences only the fourth term in the decomposition of the error (26). We introduce the notations

$$\epsilon_{\text{num},h} = \frac{\|\tilde{\mathbf{u}}^{0,H} - \bar{\mathbf{u}}^{0,H,h}\|_{E(\Omega)}}{\|\mathbf{u}\|_{E(\Omega)}}, \quad \epsilon = \frac{\|\mathbf{u} - \tilde{\mathbf{u}}^{0,H,h}\|_{E(\Omega)}}{\|\mathbf{u}\|_{E(\Omega)}}$$

and we first consider the case of one subdomain. So,  $\epsilon = \epsilon_{\text{num},h} = \epsilon_{\text{num},h}(h)$ . We thus study the behavior of the finite element method for a heterogeneous material. Piecewise linear finite elements are considered.

##### 3.1.1. Non-matching mesh

Fig. 6 shows the behavior of  $\epsilon_{\text{num},h}$  for uniform meshes non matching the particles boundaries (the number of particles is a power of 10 and the number of elements a power of 2). We observe that the error is almost constant during a range of mesh sizes and then starts to decrease with a rate of convergence  $O(\sqrt{h})$ . The ratio of the size of the element and the size of the particles is about three or four when the error starts to decrease.

The almost constant value of the error before the decrease does not depend on the size of particles nor on the size of the elements. It only depends on  $\tau$ , the ‘mismatch ratio’  $E_1/E_2$ . Fig. 7 shows two

Table 1

Evolution of the error (%) with the number of elements for several values of the mismatch  $\tau$ . The number of particles is  $10^4$

$\tau$	$\frac{ 1-\tau }{(1+\tau)} * 10^2$	Number of elements				
		1	4	16	64	256
100	98.0	100	98.2	98.0	98.0	98.0
10	81.8	100	83.1	81.9	81.7	81.4
5	66.7	100	69.1	66.8	66.6	66.2
2	33.3	100	41.2	33.8	33.3	33.0

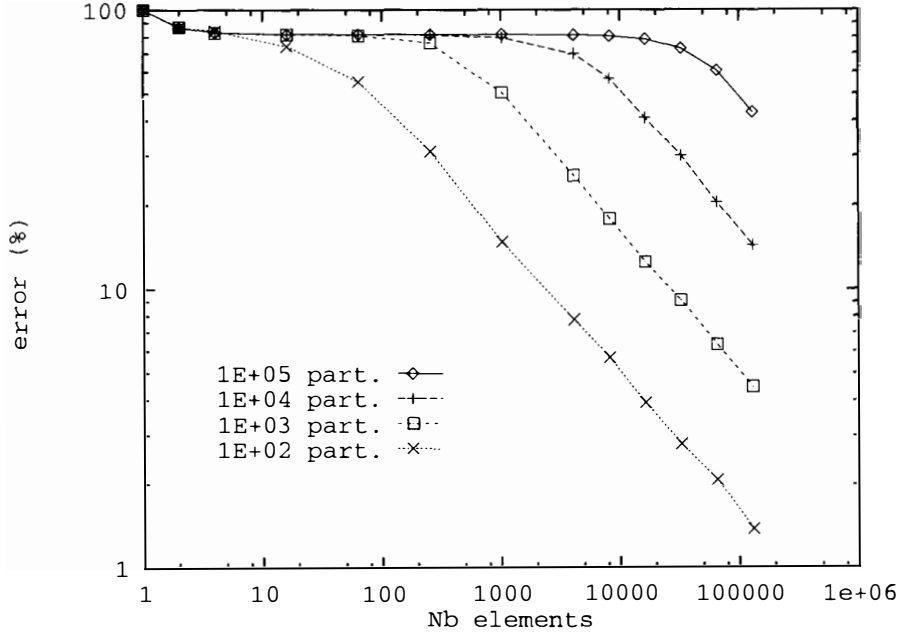


Fig. 6. Evolution of the numerical error  $\epsilon_{\text{num},h}$  (%) with the number of elements in the  $h$ -mesh for several numbers of particles. The mismatch is  $\tau = 10$ .

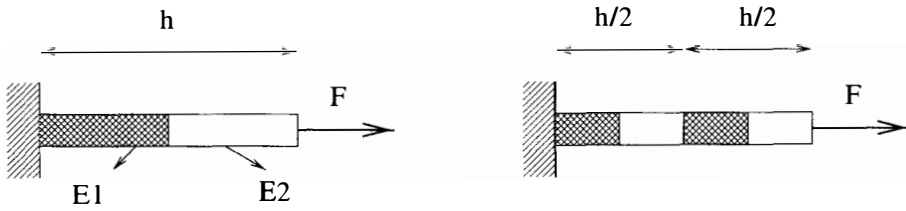


Fig. 7. Two finite element problems having the same numerical error.

finite element problems for which the numerical errors are identically the same. The relative error is  $C(\tau) = |1 - \tau|/(1 + \tau)$ . Although our 1-D problem is different from the one described in Fig. 7, the value  $C(\tau)$  is very close to the actual value for our problem (Table 1). Note that when  $\tau \rightarrow +\infty$ , the error stays at a level of 100% before a significant decrease.

The error exhibits three main trends (Fig. 8). First,  $O(h)$  convergence is observed until an error  $|1 - \tau|/(1 + \tau)$  is reached (when  $\tau = 1$  we observe an order of convergence of unity all the way as  $h \rightarrow 0$ ). Then, the error decreases very slowly. Finally, when  $h$  is around three times the size of the particle we observe a rate of convergence of  $O(\sqrt{h})$ . This rate is in agreement with the classical convergence results of the finite element method.

The following empirical formula gives an idea of the number of elements needed per particle, for this problem, in order to achieve a numerical error  $\epsilon_0$

$$\frac{d}{h} \simeq \frac{1}{r_{\text{crit}}} \left( \frac{|1 - \tau|}{(1 + \tau)} \right)^2 \frac{1}{\epsilon_0^2}.$$

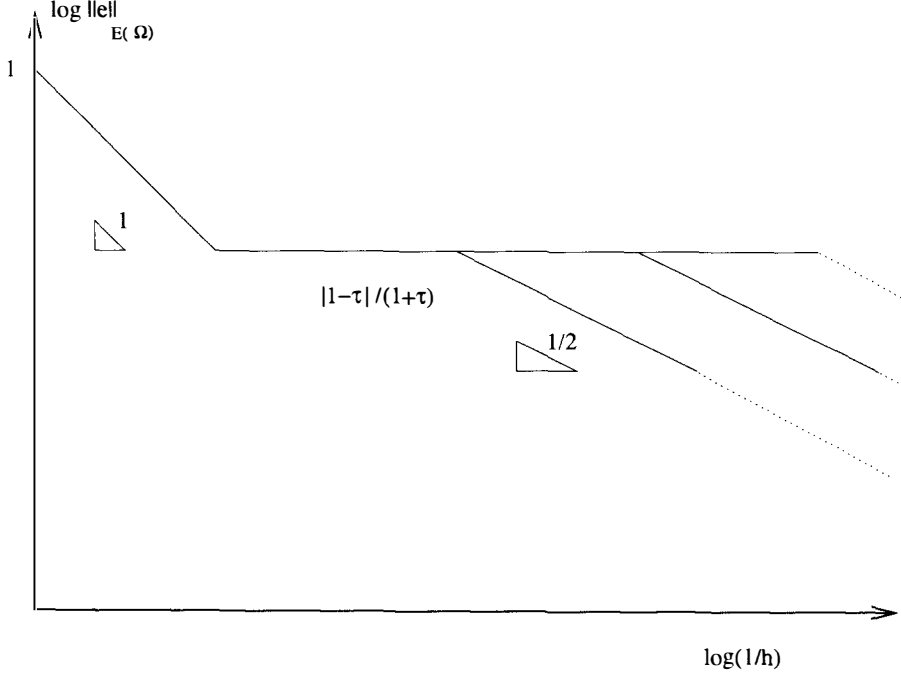


Fig. 8. The three main behavior phases of the numerical error  $\epsilon_{\text{num},h}$  with non-matching meshes.

$d$  is the size of the particle and  $r_{\text{crit}}$  is the critical ratio  $h/d$  needed to obtain a significant decrease of the error ( $r_{\text{crit}} \sim 3$ ). With  $\tau = 10$ , about 5 elements per particle are needed for 20% error, about 20 for 10% and about 85 for 5%. Due to the poor convergence rates the numerical effort needed is very large.

Concerning the influence of the subdomain size  $\Delta$ , we observe that  $\epsilon_{\text{num},h}$  increases slightly for a given  $h$ -mesh when  $\Delta$  is decreased, if  $\mathbf{E}^0 = \langle \mathbf{E}^{-1} \rangle^{-1}$ . The opposite behavior is observed if  $\mathbf{E}^0 = \langle \mathbf{E} \rangle$ . However, in general, whatever the value of  $\Delta$  and the choice of  $\mathbf{E}^0$ , the same general behavior as depicted in Fig. 9 may be observed. The only difference is that the limiting value  $|1 - \tau|/(1 + \tau)$  is no longer valid.

### 3.1.2. Matching mesh

We now consider the case where the mesh does take into account the boundary of the particles, i.e. a node is placed at each boundary between two particles. Fig. 9 shows the evolution of the error with the number of elements compared with the previous case where the mesh was not taking into account the boundaries of the particles. The error is now dramatically smaller and does not depend on the number of particles. The error is simply given by the mesh size  $h$ .

### 3.2. Influence of the $H$ -mesh

We denote this influence by  $\delta_H$ . It gathers the third and last term of the decomposition (26)

$$\delta_H = \epsilon_{\text{num},H}^2 + \gamma_H \quad (30)$$

where

$$\epsilon_{\text{num},H} = \frac{\|\tilde{\mathbf{u}}^0 - \tilde{\mathbf{u}}^{0,H}\|_{E(\Omega)}}{\|\mathbf{u}\|_{E(\Omega)}}, \quad \gamma_H = \frac{2}{\|\mathbf{u}\|_{E(\Omega)}^2} \int_{\Gamma_{\text{int}}} [\tilde{\boldsymbol{\sigma}}^0 \cdot \mathbf{n}] \cdot (\mathbf{u}^{0,H} - \mathbf{u}^0) \, ds. \quad (31)$$

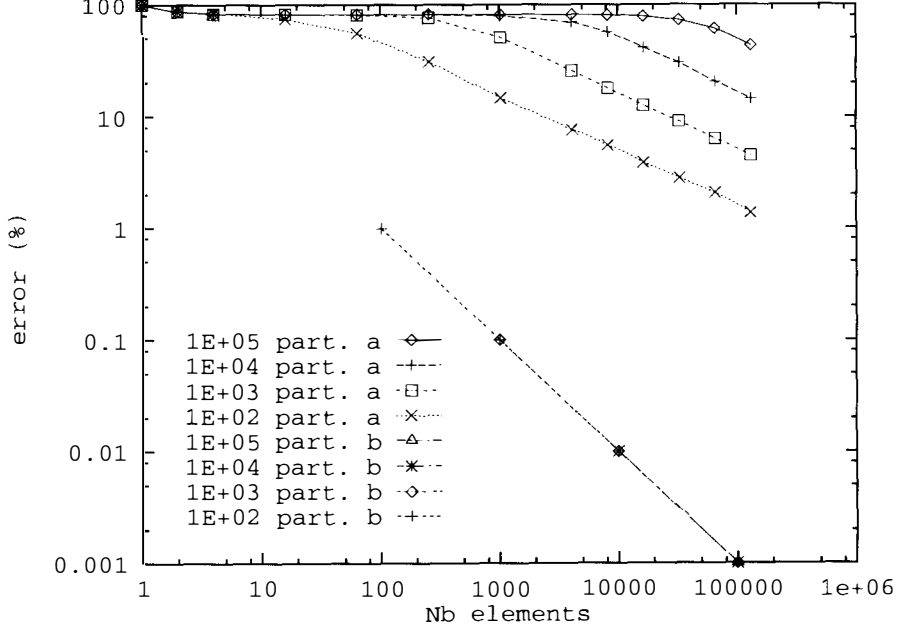


Fig. 9. Evolution of the numerical error  $\epsilon_{\text{num},h}$  (%) with the number of elements in the  $h$ -mesh for several numbers of particles. The mismatch is  $\tau = 10$ . The a and b cases correspond to meshes not matching and matching the boundary of the particles, respectively.

If the assumptions of the Property 5 holds, one can easily show that

$$\delta_H = \frac{1}{\|\mathbf{u}\|_{E(\Omega)}^2} \int_{\Gamma_{\text{int}}} [(\tilde{\boldsymbol{\sigma}}^{0,H} + \tilde{\boldsymbol{\sigma}}^0) \cdot \mathbf{n}] \cdot (\mathbf{u}^{0,H} - \mathbf{u}^0) \, ds. \quad (32)$$

For the 1-D problem under consideration, the finite element solution  $\mathbf{u}^{0,H}$  is such that  $\mathbf{u}^{0,H} = \mathbf{u}^0$  at the nodes. So  $\delta_H$  is zero since the subdomains boundaries are chosen so as to coincide with nodes of the  $H$ -mesh. However, the classical numerical error in the energy norm  $\|\mathbf{u}^0 - \mathbf{u}^{0,H}\|_{E^0(\Omega)}$  associated to  $\mathbf{u}^{0,H}$  is not zero. This remark highlights the fact that  $\delta_H$  depends only on the quality of  $\mathbf{u}^{0,H}$  on the boundary of the subdomains and not over the whole domain.

In order to purposely avoid exact finite element values at node, we modify the problem by introducing a non constant section  $S(x)$ . The problem is now

$$\frac{d}{dx}(S(x)E(x)\frac{du(x)}{dx}) = -f(x), \quad u(0) = 0, \quad u(1) = 0 \quad \text{with } S(x) = e^x \text{ and } f(x) = e^{-x}$$

and the homogenized problem is

$$\frac{d}{dx}(S(x)E^0 \frac{du^0(x)}{dx}) = -f(x), \quad u^0(0) = 0, \quad u^0(1) = 0.$$

Tables 2 and 3 show the behavior of  $\epsilon_{\text{num},H}$  and  $\gamma_H$  as the  $H$ -mesh is refined for two different choices of  $\mathbf{E}^0$ .

We observe that:

- The sign of  $\gamma_H$  can be positive or negative depending on the choice of  $\mathbf{E}^0$ ;
- $\epsilon_{\text{num},H}$  is very stable with respect to the number of subdomains.  $\gamma_H$  is very stable in the case  $\mathbf{E}^0 = \langle \mathbf{E} \rangle$  but not in the case  $\mathbf{E}^0 = \langle \mathbf{E}^{-1} \rangle^{-1}$ ;

Table 2

Influence of the number of elements in the  $H$ -mesh on  $\epsilon_{\text{num},H}$  and  $\gamma_H$  for a variable number of subdomains ( $N$ ).  $\mathbf{E}^0 = \langle \mathbf{E}^{-1} \rangle^{-1}$  and  $\tau = 10$

		Number of elements in the $H$ -mesh (1000 particles)							
		$N$	4	8	16	32	64	128	256
$\epsilon_{\text{num},H}$	4		1.00(-2)	2.53(-3)	6.34(-4)	1.59(-4)	3.97(-5)	9.91(-6)	2.48(-6)
	16				6.66(-4)	1.67(-4)	4.16(-5)	1.04(-5)	2.60(-6)
	64						4.28(-5)	1.07(-5)	2.67(-6)
$\gamma_H$	4		-2.70(-4)	-6.79(-5)	-1.70(-5)	-4.25(-6)	-1.06(-6)	-2.66(-7)	-6.64(-8)
	16				-4.19(-5)	-1.05(-5)	-2.62(-6)	-6.55(-7)	-1.64(-7)
	64						-6.69(-6)	-1.67(-6)	-4.18(-7)

Table 3

Influence of the number of elements in the  $H$ -mesh on  $\epsilon_{\text{num},H}$  and  $\gamma_H$  for a variable number of subdomains ( $N$ ).  $\mathbf{E}^\bullet = \langle \mathbf{E} \rangle$  and  $\tau = 10$

		Number of elements in the $H$ -mesh (1000 particles)							
		$N$	4	8	16	32	64	128	256
$\epsilon_{\text{num},H}$	4		3.33(-3)	8.37(-4)	2.10(-4)	5.24(-5)	1.31(-5)	3.28(-6)	8.19(-7)
	16				2.20(-4)	5.50(-5)	1.38(-5)	3.44(-6)	8.60(-7)
	64						1.41(-5)	3.53(-6)	8.84(-7)
$\gamma_H$	4		4.26(-3)	1.07(-3)	2.68(-4)	6.70(-5)	1.67(-5)	4.19(-6)	1.05(-6)
	16				2.88(-4)	7.19(-5)	1.80(-5)	4.49(-6)	1.12(-6)
	64						1.77(-5)	4.42(-6)	1.11(-6)

- $\epsilon_{\text{num},H}$  and  $\gamma_H$  are of the order of  $O(H^2)$  as  $H \rightarrow 0$ . This comes from the fact that we have super-convergence at the nodes:  $(\mathbf{u}^{0,H} - \mathbf{u}^0) \sim O(H^2)$ .
- Since  $\epsilon_{\text{num},H}$  and  $\gamma_H$  converge at the same rate, the influence of  $\epsilon_{\text{num},H}$  is quickly negligible in front of the influence of  $\gamma_H$  in  $\delta_H$ , see (30).  
Finally, to compare the influence of the  $H$ -mesh and of the  $h$ -mesh on the error, we must compare  $\delta_H$  to  $\epsilon_{\text{num},h}^2$  and observe that:
  - $\epsilon_{\text{num},h}$  is always positive;  $\delta_H$  is not necessarily positive or negative;
  - $\epsilon_{\text{num},h}$  depends on the quality of  $\tilde{\mathbf{u}}^{0,H,h}$  on each subdomain;  $\delta_H$  depends on the quality of  $\tilde{\mathbf{u}}^{0,H}$  only on the boundary of the subdomains;
  - The rate of convergence of  $\delta_H$  as  $H \rightarrow 0$  depends on the regularity of the homogenized solution  $\mathbf{u}^0$ . On the contrary, the rate of convergence of  $\epsilon_{\text{num},h}$  as  $h \rightarrow 0$  depends on the regularity of the local solution  $\tilde{\mathbf{u}}^{0,H}$  and the type of mesh (matching or non matching mesh).

### 3.3. Influence of the subdomain size

Let us introduce the notations

$$\epsilon_{\text{mod},o} = \frac{\|\mathbf{u} - \mathbf{u}^0\|_{E(\Omega)}}{\|\mathbf{u}\|_{E(\Omega)}}, \quad \epsilon_{\text{mod},\Delta} = \frac{\|\mathbf{u}^0 - \tilde{\mathbf{u}}^0\|_{E(\Omega)}}{\|\mathbf{u}\|_{E(\Omega)}}. \quad (33)$$

Fig. 10 shows the behavior of  $\epsilon_{\text{mod},\Delta}$  as the size of the subdomains,  $\Delta$ , tends to zero. One can observe, first, a very slow decrease of  $\epsilon_{\text{mod},\Delta}$ . Then, when the size of the subdomains is around three or four times the size of the particles,  $\epsilon_{\text{mod},\Delta}$  decreases as the square root of the size of the subdomains. Fig. 10 is very similar to Fig. 6. The critical subdomain size to obtain a significant decrease of  $\epsilon_{\text{mod},\Delta}$  is the same as the critical  $h$  size to obtain a significant decrease of the numerical error  $\epsilon_{\text{num},h}$ : this critical length is three to four times the size of the particles. Furthermore, the asymptotic rates of convergence are the same.

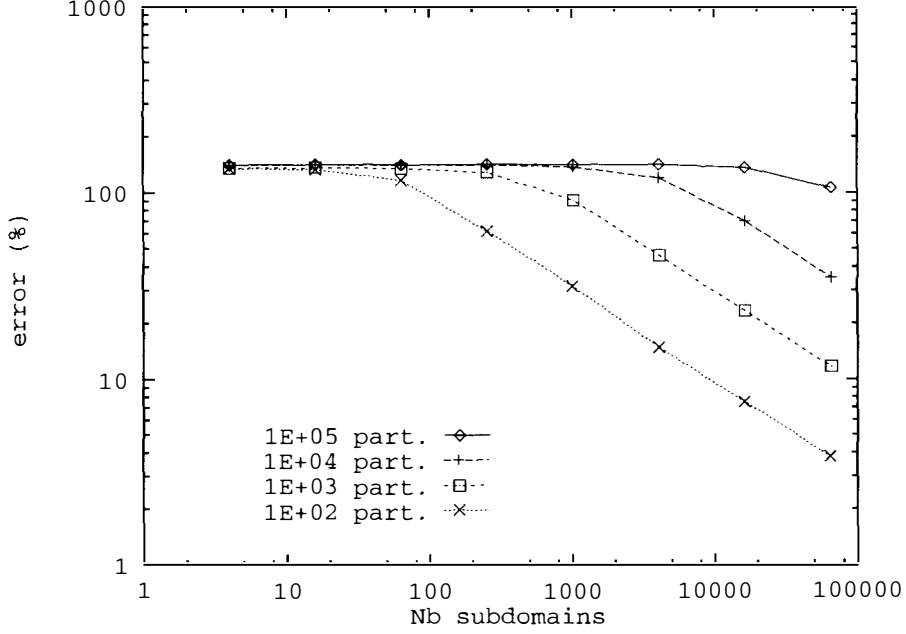


Fig. 10. Evolution of  $\epsilon_{\text{mod},\Delta}$  (%) with the number of subdomains, for several numbers of particles. The mismatch is  $\tau = 10$  and  $E^\bullet = (E^{-1})^{-1}$ .

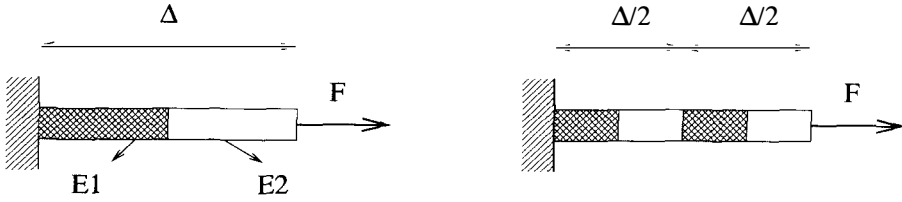


Fig. 11. These two cases have the same  $\epsilon_{\text{mod},\Delta}$  error.

There is however, a conceptual difference between  $\epsilon_{\text{mod},\Delta}$  and  $\epsilon_{\text{num},h}$ :  $\epsilon_{\text{mod},\Delta}$  is what we gain by solving exactly the fine scale problems and  $\epsilon_{\text{num},h}$  is the numerical error introduced by solving the local problem. Thus, the size of interest for the subdomain is *at least* three or four times the size of the particles and the size of interest for  $h$  is *at most* three or four times the size of the particles. The critical length, three to four times the size of particles, separates the long length range (homogenized problem) and the short length range (local analysis).

Fig. 11 shows two cases for which  $\epsilon_{\text{mod},\Delta}$  is the same:  $\epsilon_{\text{mod},\Delta} = 0.5|1 - \tau|/\sqrt{\tau}$ . With  $\tau = 10$ , we have an error of 142.3%. This value is very close to the constant error value in Fig. 10. The relationship between Figs. 11 and 7 is noteworthy.

#### 4. The error-in-the-constitutive-law framework

The homogenized solution and the exact solution both fulfill the kinematic constraints (1) and the equilibrium equation (2). They differ because they do not satisfy the same constitutive law. A way to

measure the distance between the two solutions is to compute the distance in the energy norm, between the exact and homogenized displacement:  $e = \|\mathbf{u} - \mathbf{u}^0\|_{E(\Omega)}$ . Another way is to measure the way the couple  $(\boldsymbol{\sigma}, \boldsymbol{\epsilon}^0)$  satisfies the fine-scale constitutive law.

The error on the constitutive law concept is very general and relies on a strong mechanical interpretation. It has been applied to many types of problem. Let us cite the design of error estimator for finite element computation [3] and the adjustment of finite element models using vibration tests [4]. For our purposes, this concept can be summarized as follows. Suppose we have at hand a stress–strain couple  $(\hat{\boldsymbol{\sigma}}, \hat{\boldsymbol{\epsilon}})$  satisfying the kinematic constraints (1) and the equilibrium equation (2). This couple will be the exact solution of the problem if and only if it satisfies the actual constitutive law:

$$\hat{\boldsymbol{\sigma}} = \mathbf{E}\hat{\boldsymbol{\epsilon}} \quad \text{on } \Omega. \quad (34)$$

Thus, the quality of the couple  $(\hat{\boldsymbol{\sigma}}, \hat{\boldsymbol{\epsilon}})$  can be measured by the way it satisfies (34). Let us introduce the quantity  $\eta$  defined by

$$\eta(\hat{\boldsymbol{\sigma}}, \hat{\boldsymbol{\epsilon}}) = \varphi^*(\hat{\boldsymbol{\sigma}}) + \varphi(\hat{\boldsymbol{\epsilon}}) - \hat{\boldsymbol{\sigma}} : \hat{\boldsymbol{\epsilon}}, \quad \text{where } \varphi^*(\hat{\boldsymbol{\sigma}}) = \frac{1}{2} \hat{\boldsymbol{\sigma}} : \mathbf{E}^{-1} \hat{\boldsymbol{\sigma}} \quad \text{and} \quad \varphi(\hat{\boldsymbol{\epsilon}}) = \frac{1}{2} \hat{\boldsymbol{\epsilon}} : \mathbf{E} \hat{\boldsymbol{\epsilon}}.$$

$\varphi(\hat{\boldsymbol{\epsilon}})$  is the strain energy and  $\varphi^*(\hat{\boldsymbol{\sigma}})$  is the complementary energy, its dual via the classical Legendre–Fenchel transformation,

$$\varphi^*(\hat{\boldsymbol{\sigma}}) = \sup_{\boldsymbol{\epsilon}'} (\hat{\boldsymbol{\sigma}} : \boldsymbol{\epsilon}' - \varphi(\boldsymbol{\epsilon}')).$$

It follows that  $\eta$  has the following two classical properties:

$$\eta(\hat{\boldsymbol{\sigma}}, \hat{\boldsymbol{\epsilon}}) \geq 0 \quad \forall (\hat{\boldsymbol{\sigma}}, \hat{\boldsymbol{\epsilon}}); \quad (35)$$

$$\eta(\hat{\boldsymbol{\sigma}}, \hat{\boldsymbol{\epsilon}}) = 0 \Leftrightarrow \hat{\boldsymbol{\sigma}} = \mathbf{E}\hat{\boldsymbol{\epsilon}}. \quad (36)$$

Thus, a measure of the absolute error associated with the couple  $(\hat{\boldsymbol{\sigma}}, \hat{\boldsymbol{\epsilon}})$  may be defined by:

$$Y(\hat{\boldsymbol{\sigma}}, \hat{\boldsymbol{\epsilon}}) = \left\{ 2 \int_{\Omega} \eta(\hat{\boldsymbol{\sigma}}, \hat{\boldsymbol{\epsilon}}) \, d\mathbf{x} \right\}^{1/2} \quad (37)$$

We shall refer to this error-in-the-constitutive-law approach as the ECL framework.

It is interesting to note that

$$\|\mathbf{u} - \mathbf{u}^0\|_{E(\Omega)} = Y(\boldsymbol{\sigma}, \boldsymbol{\epsilon}^0) \quad \text{and} \quad \|\mathbf{u}\|_{E(\Omega)} = Y(\boldsymbol{\sigma}, 0). \quad (38)$$

In other words, the distance between  $\mathbf{u}$  and  $\mathbf{u}^0$  in the energy norm, is equivalent to the way the couple  $(\boldsymbol{\sigma}, \boldsymbol{\epsilon}^0)$  satisfies the fine-scale constitutive law as measured by the functional  $Y(\cdot, \cdot)$ .

We shall now rewrite and demonstrate the main results of the HDPM in the ECL framework. This will lead to a mechanical interpretation of the results. We will also establish a new result.

- The first result of the HDPM is that it is possible to compute an explicit a posteriori upper bound for the homogenization error [1]. In particular, the exact solution need not be known to compute the bound:

$$\|\mathbf{u} - \mathbf{u}^0\|_{E(\Omega)} \leq \zeta \quad \text{where} \quad \zeta = \|\mathcal{I}_0 \nabla \mathbf{u}^0\|_{E(\Omega)}, \quad \mathcal{I}_0 = \mathbf{I} - \mathbf{E}^{-1} \mathbf{E}^0. \quad (39)$$

After some manipulations, this result can be proved to be equivalent to

$$Y(\boldsymbol{\sigma}, \boldsymbol{\epsilon}^0) \leq Y(\boldsymbol{\sigma}^0, \boldsymbol{\epsilon}^0). \quad (40)$$

In other words, an upper bound is obtained by replacing in the error in the constitutive law the exact stress field by the homogenized field.  $Y(\boldsymbol{\sigma}^0, \boldsymbol{\epsilon}^0)$  measures the way the homogenized solution satisfies the actual behavior. The following Lemma is the key to demonstrate all the main results of the HDPM using the ECL framework:

**LEMMA 1.** Let  $(\boldsymbol{\sigma}_1, \boldsymbol{\epsilon}_1)$  and  $(\boldsymbol{\sigma}_2, \boldsymbol{\epsilon}_2)$  be two stress–strain couples satisfying the following orthogonality condition:



$$\int_{\Omega} (\boldsymbol{\sigma}_1 - \boldsymbol{\sigma}_2) : (\boldsymbol{\epsilon}_1 - \boldsymbol{\epsilon}_2) \, d\mathbf{x} = 0. \quad (41)$$

Then, we have

$$Y^2(\boldsymbol{\sigma}_1, \boldsymbol{\epsilon}_1) + Y^2(\boldsymbol{\sigma}_2, \boldsymbol{\epsilon}_2) = Y^2(\boldsymbol{\sigma}_1, \boldsymbol{\epsilon}_2) + Y^2(\boldsymbol{\sigma}_2, \boldsymbol{\epsilon}_1). \quad \square$$

The pairs  $(\boldsymbol{\sigma}, \boldsymbol{\epsilon})$  and  $(\boldsymbol{\sigma}^0, \boldsymbol{\epsilon}^0)$  satisfy the orthogonality condition. Moreover,  $Y^2(\boldsymbol{\sigma}, \boldsymbol{\epsilon}) = 0$ . So, by the Lemma 1, we have

$$Y^2(\boldsymbol{\sigma}^0, \boldsymbol{\epsilon}^0) = Y^2(\boldsymbol{\sigma}, \boldsymbol{\epsilon}^0) + Y^2(\boldsymbol{\sigma}^0, \boldsymbol{\epsilon}) \quad (42)$$

which proves the upper-bound.  $\square$

- The second principle of the HDPM is that the modeling error is reduced by carrying out the subdomain analysis [1]. Suppose that the subdomain analysis is performed on the subdomain  $\Omega_k$ . We have

$$\|\mathbf{u} - \tilde{\mathbf{u}}^0\|_{E(\Omega_k)} \leq \|\mathbf{u} - \mathbf{u}^0\|_{E(\Omega_k)}.$$

This result may be rewritten

$$Y_k(\boldsymbol{\sigma}, \tilde{\boldsymbol{\epsilon}}^0) \leq Y_k(\boldsymbol{\sigma}, \boldsymbol{\epsilon}^0) \quad (43)$$

where

$$Y_k(\hat{\boldsymbol{\sigma}}, \hat{\boldsymbol{\epsilon}}) = \left\{ 2 \int_{\Omega_k} \eta(\hat{\boldsymbol{\sigma}}, \hat{\boldsymbol{\epsilon}}) \, d\mathbf{x} \right\}^{1/2}.$$

*PROOF OF (43).* The pairs  $(\boldsymbol{\sigma}, \boldsymbol{\epsilon}^0)$  and  $(\tilde{\boldsymbol{\sigma}}^0, \tilde{\boldsymbol{\epsilon}}^0)$  satisfy the orthogonality condition on  $\Omega_k$ . Then by Lemma 1, since  $Y_k(\tilde{\boldsymbol{\sigma}}^0, \tilde{\boldsymbol{\epsilon}}^0) = 0$ , we have

$$Y_k^2(\boldsymbol{\sigma}, \boldsymbol{\epsilon}^0) = Y_k^2(\boldsymbol{\sigma}, \tilde{\boldsymbol{\epsilon}}^0) + Y_k^2(\tilde{\boldsymbol{\sigma}}^0, \boldsymbol{\epsilon}^0), \quad (44)$$

proving the result.  $\square$

- The  $\zeta$  quantity defined in the first result is not, in general, an upper bound locally, i.e. the following inequality is not satisfied in general over a given subdomain,  $\Omega_k$ :

$$\|\mathbf{u} - \mathbf{u}^0\|_{E(\Omega_k)} \leq \zeta_k \quad \text{where} \quad \zeta_k = \|\mathcal{I}_0 \nabla \mathbf{u}^0\|_{E(\Omega_k)}. \quad (45)$$

However, the following inequality, which is the third result of the HDPM [1], holds:

$$\|\tilde{\mathbf{u}}^0 - \mathbf{u}^0\|_{E(\Omega_k)} \leq \zeta_k. \quad (46)$$

In other words, it is possible to determine locally where the local solution process will produce a significant change in the homogenized solution.

In terms of the ECL, we have

$$Y_k(\tilde{\boldsymbol{\sigma}}^0, \boldsymbol{\epsilon}^0) \leq Y_k(\boldsymbol{\sigma}^0, \boldsymbol{\epsilon}^0). \quad (47)$$

*PROOF OF (47).* Applying the Lemma 1 with the pairs  $(\tilde{\boldsymbol{\sigma}}^0, \tilde{\boldsymbol{\epsilon}}^0)$  and  $(\boldsymbol{\sigma}^0, \boldsymbol{\epsilon}^0)$ , we have

$$Y_k^2(\boldsymbol{\sigma}^0, \boldsymbol{\epsilon}^0) = Y_k^2(\tilde{\boldsymbol{\sigma}}^0, \boldsymbol{\epsilon}^0) + Y_k^2(\boldsymbol{\sigma}^0, \tilde{\boldsymbol{\epsilon}}^0). \quad \square$$

- The fourth principle is that an upper bound also exists for the modeling error obtained after the subdomain analysis [2]:

$$\|\mathbf{u} - \tilde{\mathbf{u}}^0\|_{E(\Omega)} \leq \psi \quad \text{where} \quad \psi^2 = 2(\mathcal{J}(\tilde{\mathbf{u}}^0) - \mathcal{J}(\mathbf{u}^0)) + \zeta^2. \quad (48)$$

$\mathcal{J}(\cdot)$  is the potential energy associated to the displacement field  $\cdot$  :

$$\mathcal{J}(\cdot) = \int_{\Omega} \varphi(\boldsymbol{\epsilon}(\cdot)) \, d\mathbf{x} - \int_{\Omega} \mathbf{f}_g \cdot (\cdot) \, d\mathbf{x} - \int_{\Gamma_t} \mathbf{t}_g \cdot (\cdot) \, d\mathbf{s}.$$

The result (48) may be rewritten

$$Y(\boldsymbol{\sigma}, \tilde{\boldsymbol{\epsilon}}^0) \leq Y(\boldsymbol{\sigma}^0, \tilde{\boldsymbol{\epsilon}}^0). \quad (49)$$

In other words, as for the first result, an upper bound is obtained by replacing the exact stress by the homogenized one.

*PROOF OF (49).* Applying the Lemma 1 with the pairs  $(\boldsymbol{\sigma}, \boldsymbol{\epsilon})$  and  $(\boldsymbol{\sigma}^0, \tilde{\boldsymbol{\epsilon}}^0)$ , we have

$$Y^2(\boldsymbol{\sigma}^0, \tilde{\boldsymbol{\epsilon}}^0) = Y^2(\boldsymbol{\sigma}, \tilde{\boldsymbol{\epsilon}}^0) + Y^2(\boldsymbol{\sigma}^0, \boldsymbol{\epsilon}). \quad \square \quad (50)$$

The ECL framework has allowed us to find a new mechanical interpretation of the HDPM results: the bounds are obtained by using the homogenized stress field in the error in the constitutive law expression. In the first and last results, it is the global exact stress field  $\boldsymbol{\sigma}$  that is replaced by the homogenized one (see (40) and (49)) and in the third result it is the local stress field  $\hat{\boldsymbol{\sigma}}^0$  that is replaced by the homogenized stress (see (47)).

We also note that all the HDPM results now extend to the nonlinear case. In the nonlinear case, all the proofs remain valid, only the expression of the potential  $\varphi$  and  $\varphi^*$  change. Let us inquire the meaning of the  $Y$  functional in the nonlinear case. Assuming that  $\hat{\boldsymbol{\sigma}}$  satisfies the equilibrium equation (2) and  $\hat{\mathbf{u}}$  satisfies the kinematic constraints (1), we can write

$$\frac{1}{2}Y^2(\hat{\boldsymbol{\sigma}}, \boldsymbol{\epsilon}(\hat{\mathbf{u}})) = \mathcal{J}(\hat{\mathbf{u}}) - \Pi(\hat{\boldsymbol{\sigma}}) \quad (51)$$

where  $\mathcal{J}(\hat{\mathbf{u}})$  is the potential energy and  $\Pi(\hat{\boldsymbol{\sigma}})$  the complementary potential energy:

$$\mathcal{J}(\hat{\mathbf{u}}) = \int_{\Omega} \varphi(\boldsymbol{\epsilon}(\hat{\mathbf{u}})) \, d\mathbf{x} - \int_{\Omega} \mathbf{f}_g \cdot (\hat{\mathbf{u}}) \, d\mathbf{x} - \int_{\Gamma_t} \mathbf{t}_g \cdot (\hat{\mathbf{u}}) \, d\mathbf{s}, \quad (52)$$

$$\Pi(\hat{\boldsymbol{\sigma}}) = - \int_{\Omega} \varphi^*(\hat{\boldsymbol{\sigma}}) \, d\mathbf{x} + \int_{\Gamma_u} (\hat{\boldsymbol{\sigma}} \cdot \mathbf{n}) \cdot \mathbf{u}_g \, d\mathbf{s}. \quad (53)$$

Thus  $Y$  is linked to the difference between the potential energy and the complementary potential energy. It is always positive and zero if and only if the couple  $(\hat{\boldsymbol{\sigma}}, \hat{\boldsymbol{\epsilon}} = \boldsymbol{\epsilon}(\hat{\mathbf{u}}))$  satisfies the behavior described by the relation:

$$\varphi^*(\hat{\boldsymbol{\sigma}}) + \varphi(\hat{\boldsymbol{\epsilon}}) - \hat{\boldsymbol{\sigma}} : \hat{\boldsymbol{\epsilon}} = 0. \quad (54)$$

The main results of the HDPM are stated in terms of explicit upper-bounds. The quality of these bounds, is defined by the following effectivity indices, all greater than or equal to one:

$$\theta_0 = \frac{\zeta}{\|\mathbf{u} - \mathbf{u}^0\|_{E(\Omega)}}, \quad \tilde{\theta}_0 = \frac{\psi}{\|\mathbf{u} - \tilde{\mathbf{u}}^0\|_{E(\Omega)}}, \quad \theta_{s,k} = \frac{\zeta_k}{\|\tilde{\mathbf{u}}^0 - \mathbf{u}^0\|_{E(\Omega_k)}},$$

or equivalently, in terms of the error in the constitutive law:

$$\theta_0 = \frac{Y(\boldsymbol{\sigma}^0, \boldsymbol{\epsilon}^0)}{Y(\boldsymbol{\sigma}, \boldsymbol{\epsilon}^0)}, \quad \tilde{\theta}_0 = \frac{Y(\boldsymbol{\sigma}^0, \tilde{\boldsymbol{\epsilon}}^0)}{Y(\boldsymbol{\sigma}, \tilde{\boldsymbol{\epsilon}}^0)}, \quad \theta_{s,k} = \frac{Y_k(\boldsymbol{\sigma}^0, \boldsymbol{\epsilon}^0)}{Y_k(\tilde{\boldsymbol{\sigma}}^0, \boldsymbol{\epsilon}^0)}.$$

We may also define the global sensitivity effectivity index,  $\theta_s$ :

$$\theta_s = \frac{\zeta}{\|\tilde{\mathbf{u}}^0 - \mathbf{u}^0\|_{E(\Omega)}} = \frac{Y(\boldsymbol{\sigma}^0, \boldsymbol{\epsilon}^0)}{Y(\tilde{\boldsymbol{\sigma}}^0, \boldsymbol{\epsilon}^0)}.$$

Of course,  $\theta_s$  is bounded from below by the best effectivity over the subdomain and from above by the worst effectivity index:

$$\min_{k=1, \dots, N} \theta_{s,k} \leq \theta_s \leq \max_{k=1, \dots, N} \theta_{s,k}.$$

The following property establishes that the effectivity index on the  $\psi$  estimate,  $\tilde{\theta}_0$ , is always greater than the effectivity index on the  $\zeta$  estimate,  $\theta_0$ . The global sensitivity index,  $\theta_s$ , is also always greater than  $\theta_0$ . Finally,  $\theta_0$ ,  $\tilde{\theta}_0$  and  $\theta_s$  are linked in a relation:

**PROPERTY 6.**

- $\theta_0 \leq \tilde{\theta}_0$  and  $\theta_0 = 1 \Leftrightarrow \tilde{\theta}_0 = 1$ ;
- $\theta_0 \leq \theta_s$ ;
- $(\theta_0^2 - \theta_s^2)(\theta_s^2 - \theta_0^2) = \theta_0^2(\theta_0^2 - 1)$ .

*PROOF.* Applying the Lemma 1 to the pairs  $(\boldsymbol{\sigma}^0, \boldsymbol{\epsilon}^0)$  and  $(\tilde{\boldsymbol{\sigma}}^0, \tilde{\boldsymbol{\epsilon}}^0)$ , we obtain since  $Y(\tilde{\boldsymbol{\sigma}}^0, \tilde{\boldsymbol{\epsilon}}^0) = 0$

$$Y^2(\boldsymbol{\sigma}^0, \boldsymbol{\epsilon}^0) = Y^2(\boldsymbol{\sigma}^0, \tilde{\boldsymbol{\epsilon}}^0) + Y^2(\tilde{\boldsymbol{\sigma}}^0, \boldsymbol{\epsilon}^0). \quad (55)$$

Summing relation (44) over all the subdomains yields

$$Y^2(\boldsymbol{\sigma}, \boldsymbol{\epsilon}^0) = Y^2(\boldsymbol{\sigma}, \tilde{\boldsymbol{\epsilon}}^0) + Y^2(\tilde{\boldsymbol{\sigma}}^0, \boldsymbol{\epsilon}^0). \quad (56)$$

Dividing relation (55) by relation (56), we obtain

$$\theta_0^2 = \frac{Y^2(\boldsymbol{\sigma}^0, \boldsymbol{\epsilon}^0)}{Y^2(\boldsymbol{\sigma}, \boldsymbol{\epsilon}^0)} = \frac{Y^2(\boldsymbol{\sigma}^0, \tilde{\boldsymbol{\epsilon}}^0) + Y^2(\tilde{\boldsymbol{\sigma}}^0, \boldsymbol{\epsilon}^0)}{Y^2(\boldsymbol{\sigma}, \tilde{\boldsymbol{\epsilon}}^0) + Y^2(\tilde{\boldsymbol{\sigma}}^0, \boldsymbol{\epsilon}^0)} \leq \frac{Y^2(\boldsymbol{\sigma}^0, \tilde{\boldsymbol{\epsilon}}^0)}{Y^2(\boldsymbol{\sigma}, \tilde{\boldsymbol{\epsilon}}^0)} = \tilde{\theta}_0^2 \quad (57)$$

and clearly  $\tilde{\theta}_0 = 1 \Leftrightarrow \tilde{\theta}_0 = 1$ . The second relation is obvious since

$$Y(\tilde{\boldsymbol{\sigma}}^0, \boldsymbol{\epsilon}^0) \leq Y(\boldsymbol{\sigma}, \boldsymbol{\epsilon}^0)$$

and the third one is obtained by eliminating the quantity  $Y(\tilde{\boldsymbol{\sigma}}^0, \boldsymbol{\epsilon}^0)$  between the two relations:

$$\theta_s = \frac{Y(\boldsymbol{\sigma}^0, \boldsymbol{\epsilon}^0)}{Y(\tilde{\boldsymbol{\sigma}}^0, \boldsymbol{\epsilon}^0)} \quad \text{and} \quad \tilde{\theta}_0^2 = \frac{Y^2(\boldsymbol{\sigma}^0, \boldsymbol{\epsilon}^0) - Y^2(\tilde{\boldsymbol{\sigma}}^0, \boldsymbol{\epsilon}^0)}{Y^2(\boldsymbol{\sigma}, \boldsymbol{\epsilon}^0) - Y^2(\tilde{\boldsymbol{\sigma}}^0, \boldsymbol{\epsilon}^0)}. \quad \square$$

Table 4 shows the influence of the number of subdomains on the modeling error

$$\epsilon = \frac{\|\mathbf{u} - \tilde{\mathbf{u}}^0\|_{E(\Omega)}}{\|\mathbf{u}\|_{E(\Omega)}} = \frac{Y(\boldsymbol{\sigma}, \tilde{\boldsymbol{\epsilon}}^0)}{Y(\boldsymbol{\sigma}, 0)}$$

and on the effectivity index  $\tilde{\theta}_0$ . We see that  $\tilde{\theta}_0$  tends monotonically to  $\theta_0$  as the subdomain size is reduced.

**Table 4**

Evolution of the modeling error  $\epsilon$  and the effectivity index  $\tilde{\theta}_0$  with respect to the number of subdomains,  $N$ , for three different values of  $\tau$ .  $\mathbf{E}^0 = \langle \mathbf{E}^{-1} \rangle^{-1}$ . When  $N = +\infty$ ,  $\epsilon = \frac{\|\mathbf{u} - \mathbf{u}^0\|_{E(\Omega)}}{\|\mathbf{u}\|_{E(\Omega)}} = \frac{Y(\boldsymbol{\sigma}, \boldsymbol{\epsilon}^0)}{Y(\boldsymbol{\sigma}, 0)}$  and  $\tilde{\theta}_0 = \theta_0$

$N$	$\tau = 10$		$\tau = 50$		$\tau = 100$	
	$\epsilon$	$\tilde{\theta}_0$	$\epsilon$	$\tilde{\theta}_0$	$\epsilon$	$\tilde{\theta}_0$
1	0.	$+\infty$	0.	$+\infty$	0.	$+\infty$
2	2.702(-4)	44.70039184	3.114(-4)	45.66017533	3.169(-4)	45.79441773
4	1.739(-2)	1.21748697	2.046(-2)	1.21758072	2.088(-2)	1.21759302
16	3.505(-2)	1.05768583	4.138(-2)	1.05734892	4.225(-2)	1.05730230
64	7.672(-2)	1.01231204	9.128(-2)	1.01205426	9.329(-2)	1.01201730
256	1.499(-1)	1.00323777	1.809(-1)	1.00308147	1.854(-1)	1.00305648
1024	3.232(-1)	1.00069796	4.959(-1)	1.00041075	6.074(-1)	1.00028513
4096	7.648(-1)	1.00012465	1.692	1.00003527	2.356	1.00001896
$+\infty$	1.421	1.00003609	3.460	1.00000844	4.944	1.00000431

## 5. Numerical results for the HDPM

The results of the HDPM method rely on the assumption that the homogenized and the local solutions are free of numerical errors. In practice, this is not the case. Thus, we seek for the meaning of the results, mainly the two first results, when numerical errors occur. We keep the error-in-the-constitutive-law framework allowing us to deal with both linear or nonlinear constitutive laws.

It is interesting to note that the inequality expressed in (40) still holds if the exact fields are replaced by the approximate fields. Defining  $(\boldsymbol{\sigma}^H, \boldsymbol{\epsilon}^H)$  as the finite element solution obtained with the real material and the  $H$ -mesh, we have

*PROPERTY 7.*

$$Y(\boldsymbol{\sigma}^H, \boldsymbol{\epsilon}^{0,H}) \leq Y(\boldsymbol{\sigma}^{0,H}, \boldsymbol{\epsilon}^{0,H}).$$

*PROOF.* The proof is obtained by the Lemma 1 using the pairs  $(\boldsymbol{\sigma}^H, \boldsymbol{\epsilon}^H)$  and  $(\boldsymbol{\sigma}^{0,H}, \boldsymbol{\epsilon}^{0,H})$ . The orthogonality condition (41) is indeed satisfied at the finite element level.  $\square$

Defining the numerical effectivity index,  $\theta_0^H$ , as

$$\theta_0^H = \frac{Y(\boldsymbol{\sigma}^{0,H}, \boldsymbol{\epsilon}^{0,H})}{Y(\boldsymbol{\sigma}^H, \boldsymbol{\epsilon}^H)},$$

we thus have  $\theta_0^H \geq 1$  like  $\theta_0 \geq 1$ . Numerical experiments carried out in [1] for 1-D and 3-D problems confirm the property, although in this study an iterative solver and approximate spatial integration were used in the 3-D case. It was also noticed in [1] that  $\theta_0^H$  is very stable with respect to the mesh size  $H$ . Finally, note that, as the  $\zeta$  estimate  $(Y(\boldsymbol{\sigma}^0, \boldsymbol{\epsilon}^0))$ ,  $Y(\boldsymbol{\sigma}^{0,H}, \boldsymbol{\epsilon}^{0,H})$  is zero if  $E$  and  $E^0$  coincide, even though numerical errors occur.

By performing the local analysis, we know from the second result of the HDPM that we obtain a perturbation of the homogenized solution,  $\tilde{\mathbf{u}}^0$ , that is closer to the exact solution than  $\mathbf{u}^0$ . This might not hold for the numerical approximation of  $\tilde{\mathbf{u}}^0$  because it is not an exact solution of the subdomain problems. Fortunately, we have the following result:

*PROPERTY 8.* Let  $\mathbf{u}_k^{0,H}$  be the restriction of  $\mathbf{u}^{0,H}$  to  $\Omega_k$  and  $\boldsymbol{\epsilon}_k^{0,H} = \boldsymbol{\epsilon}(\mathbf{u}_k^{0,H})$ . If  $(\mathbf{u}_k^{0,H} - \tilde{\mathbf{u}}_k^{0,H,h}) \in V_k^h$  then

$$Y_k(\boldsymbol{\sigma}, \tilde{\boldsymbol{\epsilon}}_k^{0,H,h}) \leq Y_k(\boldsymbol{\sigma}, \boldsymbol{\epsilon}_k^{0,H}).$$

*PROOF.* We apply Lemma 1 with the pairs  $(\tilde{\boldsymbol{\sigma}}_k^{0,H,h}, \tilde{\boldsymbol{\epsilon}}_k^{0,H,h})$  and  $(\boldsymbol{\sigma}, \boldsymbol{\epsilon}_k^{0,H})$ . Since  $(\mathbf{u}_k^{0,H} - \tilde{\mathbf{u}}_k^{0,H,h}) \in V_k^h$  the orthogonality condition holds.  $\square$

Finally, we stress that two quantities equivalent when free of numerical errors may become quite different in the case of numerical error. For instance, as seen in the previous section, the following two expressions of  $\psi$  are equivalent:

$$\begin{aligned} \psi &= \left( 2(\mathcal{J}(\tilde{\mathbf{u}}^0) - \mathcal{J}(\mathbf{u}^0)) + \|\mathcal{I}_0 \nabla \mathbf{u}^0\|_{E(\Omega)}^2 \right)^{1/2}, \\ \psi &= Y(\boldsymbol{\sigma}^0, \tilde{\boldsymbol{\epsilon}}^0). \end{aligned}$$

If we inject the approximate solution, we get

$$\begin{aligned} \psi^{\text{num1}} &= \left( 2(\mathcal{J}(\tilde{\mathbf{u}}^{0,H,h}) - \mathcal{J}(\mathbf{u}^{0,H})) + \|\mathcal{I}_0 \nabla \mathbf{u}^{0,H}\|_{E(\Omega)}^2 \right)^{1/2}, \\ \psi^{\text{num2}} &= Y(\boldsymbol{\sigma}^{0,H}, \tilde{\boldsymbol{\epsilon}}^{0,H,h}). \end{aligned}$$

$\psi^{\text{num1}}$  and  $\psi^{\text{num2}}$  no longer coincide since the first one may be the square root of a negative value, whereas the second one is always well defined.

## 6. A posteriori error estimation and adaptive strategy

From (23), we recall that

$$\|\mathbf{u} - \tilde{\mathbf{u}}^{0,H,h}\|_{E(\Omega)}^2 = \|\mathbf{u} - \tilde{\mathbf{u}}^{0,H}\|_{E(\Omega)}^2 + \|\tilde{\mathbf{u}}^{0,H} - \tilde{\mathbf{u}}^{0,H,h}\|_{E(\Omega)}^2$$

which can be rewritten

$$Y^2(\boldsymbol{\sigma}, \tilde{\boldsymbol{\epsilon}}^{0,H,h}) = Y^2(\boldsymbol{\sigma}, \tilde{\boldsymbol{\epsilon}}^{0,H}) + Y^2(\hat{\boldsymbol{\sigma}}^{0,H}, \tilde{\boldsymbol{\epsilon}}^{0,H,h}). \quad (58)$$

The Properties 9, 10 and 11 below give upper-bounds for the three terms in (58). The spaces  $\mathcal{U}$  and  $\mathcal{S}$  describe the regularity imposed to the displacement and stress field, respectively.

*PROPERTY 9.*

$$Y(\boldsymbol{\sigma}, \tilde{\boldsymbol{\epsilon}}^{0,H,h}) \leq Y(\hat{\boldsymbol{\sigma}}, \tilde{\boldsymbol{\epsilon}}^{0,H,h}) \quad \forall \hat{\boldsymbol{\sigma}} \in \mathcal{S}_{\text{ad}}$$

where

$$\mathcal{S}_{\text{ad}} = \left\{ \boldsymbol{\sigma} \in \mathbb{L}_{\text{sym}}^2(\Omega) : \int_{\Omega} \boldsymbol{\sigma} : \nabla \mathbf{v} \, d\mathbf{x} - \mathcal{F}(\mathbf{v}) = 0 \quad \forall \mathbf{v} \in V \right\}, \quad (59)$$

$$\mathbb{L}_{\text{sym}}^2 = \left\{ \boldsymbol{\tau} = \{\tau_{ij}\} \in (L^2(\Omega))^{n \times n}, \boldsymbol{\tau} = \boldsymbol{\tau}^t \right\}. \quad (60)$$

*PROOF.* With the pairs  $(\boldsymbol{\sigma}, \boldsymbol{\epsilon})$  and  $(\hat{\boldsymbol{\sigma}}, \tilde{\boldsymbol{\epsilon}}^{0,H,h})$ , Lemma 1 gives

$$Y^2(\hat{\boldsymbol{\sigma}}, \tilde{\boldsymbol{\epsilon}}^{0,H,h}) = Y^2(\boldsymbol{\sigma}, \tilde{\boldsymbol{\epsilon}}^{0,H,h}) + Y^2(\hat{\boldsymbol{\sigma}}, \boldsymbol{\epsilon}). \quad \square \quad (61)$$

*PROPERTY 10.*

$$Y(\boldsymbol{\sigma}, \tilde{\boldsymbol{\epsilon}}^{0,H}) \leq Y(\hat{\boldsymbol{\sigma}}, \tilde{\boldsymbol{\epsilon}}^{0,H}) \quad \forall \hat{\boldsymbol{\sigma}} \in \mathcal{S}_{\text{ad}}.$$

*PROOF.* With the pairs  $(\boldsymbol{\sigma}, \boldsymbol{\epsilon})$  and  $(\hat{\boldsymbol{\sigma}}, \tilde{\boldsymbol{\epsilon}}^{0,H})$ , we get from Lemma 1

$$Y^2(\hat{\boldsymbol{\sigma}}, \tilde{\boldsymbol{\epsilon}}^{0,H}) = Y^2(\boldsymbol{\sigma}, \tilde{\boldsymbol{\epsilon}}^{0,H}) + Y^2(\hat{\boldsymbol{\sigma}}, \boldsymbol{\epsilon}). \quad \square \quad (62)$$

*PROPERTY 11.*

$$Y(\tilde{\boldsymbol{\sigma}}^{0,H}, \tilde{\boldsymbol{\epsilon}}^{0,H,h}) \leq Y(\hat{\tilde{\boldsymbol{\sigma}}}^{0,H,h}, \tilde{\boldsymbol{\epsilon}}^{0,H,h}) \quad \forall \hat{\tilde{\boldsymbol{\sigma}}}^{0,H,h} \in \tilde{\mathcal{S}}_{\text{ad}}$$

where

$$\begin{aligned} \tilde{\mathcal{S}}_{\text{ad}} = \prod_{k=1}^N \{ & \boldsymbol{\tau} \in \mathbb{L}_{\text{sym}}^2(\Omega_k) : \int_{\Omega_k} \boldsymbol{\tau} : \nabla \mathbf{v} \, d\mathbf{x} = \int_{\Omega_k} \mathbf{f}_g \cdot \mathbf{v} \, d\mathbf{x} + \int_{\Gamma_k \cap \partial\Omega_k} \mathbf{t}_g \cdot \mathbf{v} \, d\mathbf{s}, \\ & \forall \mathbf{v} \in \mathbf{H}^1(\Omega_k), \mathbf{v} = 0 \text{ on } \Gamma_{\text{int}} \cup \Gamma_{\text{u}} \} \end{aligned} \quad (63)$$

*PROOF.* With the pairs  $(\tilde{\boldsymbol{\sigma}}^{0,H}, \tilde{\boldsymbol{\epsilon}}^{0,H})$  and  $(\hat{\tilde{\boldsymbol{\sigma}}}^{0,H,h}, \tilde{\boldsymbol{\epsilon}}^{0,H,h})$ , we have by Lemma 1

$$Y^2(\hat{\tilde{\boldsymbol{\sigma}}}^{0,H,h}, \tilde{\boldsymbol{\epsilon}}^{0,H,h}) = Y^2(\tilde{\boldsymbol{\sigma}}^{0,H}, \tilde{\boldsymbol{\epsilon}}^{0,H,h}) + Y^2(\hat{\tilde{\boldsymbol{\sigma}}}^{0,H,h}, \tilde{\boldsymbol{\epsilon}}^{0,H}). \quad \square \quad (64)$$

Note that this last Property is also valid on each subdomain. Finally, it is worth mentioning that for linear elasticity, the relations (61), (62) and (64) may also be obtained using the hyper-circle theorem [5].

Let us define the three following effectivity indices, all greater than one:

$$\theta = \frac{\hat{e}}{e}, \quad \theta_{\text{mod}} = \frac{\hat{e}_{\text{mod}}}{e_{\text{mod}}}, \quad \theta_{\text{num}} = \frac{\hat{e}_{\text{num}}}{e_{\text{num}}}$$

where

$$\hat{e} = Y(\hat{\boldsymbol{\sigma}}, \tilde{\boldsymbol{\epsilon}}^{0,H,h}), \quad \hat{e}_{\text{mod}} = Y(\hat{\boldsymbol{\sigma}}, \tilde{\boldsymbol{\epsilon}}^{0,H}), \quad \hat{e}_{\text{num}} = Y(\hat{\boldsymbol{\sigma}}, \tilde{\boldsymbol{\epsilon}}^{0,H,h}), \quad (65)$$

$$e = Y(\boldsymbol{\sigma}, \tilde{\boldsymbol{\epsilon}}^{0,H,h}), \quad e_{\text{mod}} = Y(\boldsymbol{\sigma}, \tilde{\boldsymbol{\epsilon}}^{0,H}), \quad e_{\text{num}} = Y(\tilde{\boldsymbol{\sigma}}^{0,H}, \tilde{\boldsymbol{\epsilon}}^{0,H,h}). \quad (66)$$

We have the following property:

*PROPERTY 12.* If  $\hat{e}$  and  $\hat{e}_{\text{mod}}$  are computed with the same stress field  $\hat{\boldsymbol{\sigma}} \in \mathcal{S}_{\text{ad}}$ , we have:  $\theta \leq \theta_{\text{mod}}$  and  $\theta = 1. \Leftrightarrow \theta_{\text{mod}} = 1$ .

*PROOF.* With the pairs  $(\hat{\boldsymbol{\sigma}}, \tilde{\boldsymbol{\epsilon}}^{0,H,h})$  and  $(\tilde{\boldsymbol{\sigma}}^{0,H}, \tilde{\boldsymbol{\epsilon}}^{0,H,h})$ , Lemma 1 gives

$$Y^2(\hat{\boldsymbol{\sigma}}, \tilde{\boldsymbol{\epsilon}}^{0,H,h}) = Y^2(\hat{\boldsymbol{\sigma}}, \tilde{\boldsymbol{\epsilon}}^{0,H}) + Y^2(\tilde{\boldsymbol{\sigma}}^{0,H}, \tilde{\boldsymbol{\epsilon}}^{0,H,h}). \quad (67)$$

Dividing (67) by (58), we get

$$\theta^2 = \frac{Y^2(\hat{\boldsymbol{\sigma}}, \tilde{\boldsymbol{\epsilon}}^{0,H,h})}{Y^2(\boldsymbol{\sigma}, \tilde{\boldsymbol{\epsilon}}^{0,H,h})} = \frac{Y^2(\hat{\boldsymbol{\sigma}}, \tilde{\boldsymbol{\epsilon}}^{0,H}) + Y^2(\tilde{\boldsymbol{\sigma}}^{0,H}, \tilde{\boldsymbol{\epsilon}}^{0,H,h})}{Y^2(\boldsymbol{\sigma}, \tilde{\boldsymbol{\epsilon}}^{0,H}) + Y^2(\tilde{\boldsymbol{\sigma}}^{0,H}, \tilde{\boldsymbol{\epsilon}}^{0,H,h})} \leq \frac{Y^2(\hat{\boldsymbol{\sigma}}, \tilde{\boldsymbol{\epsilon}}^{0,H})}{Y^2(\boldsymbol{\sigma}, \tilde{\boldsymbol{\epsilon}}^{0,H})} = \theta_{\text{mod}}^2 \quad (68)$$

and clearly  $\theta_{\text{mod}} = 1. \Leftrightarrow \theta = 1. \quad \square$

From the finite element stresses  $\boldsymbol{\sigma}^{0,H}$ , it is possible to build explicitly a stress field  $\hat{\boldsymbol{\sigma}}^{0,H}$  belonging to  $\mathcal{S}_{\text{ad}}$  [3]. Similarly, it is possible to build from the finite element stresses  $\tilde{\boldsymbol{\sigma}}^{0,H,h}$  stresses  $\hat{\boldsymbol{\sigma}}^{0,H,h}$  belonging to  $\tilde{\mathcal{S}}_{\text{ad}}$ . Thus, the upper bound  $\hat{e}$  and  $\hat{e}_{\text{num}}$  can be practically computed since  $\tilde{\boldsymbol{\epsilon}}^{0,H,h}$  is known. On the contrary,  $\hat{e}_{\text{mod}}$  cannot be computed since  $\tilde{\boldsymbol{\epsilon}}^{0,H}$  is unknown. However, when  $\hat{e}^2 \geq \hat{e}_{\text{num}}^2$ , we can evaluate  $e_{\text{mod}}$  by  $\tilde{e}_{\text{mod}}$ :

$$\tilde{e}_{\text{mod}}^2 = \hat{e}^2 - \hat{e}_{\text{num}}^2$$

and define the effectivity index

$$\tilde{\theta}_{\text{mod}} = \tilde{e}_{\text{mod}}/e_{\text{mod}}.$$

This latter effectivity is not necessarily greater than one.

Table 5 gives the results for the estimated errors and the effectivity indices for the model problem. The relative errors  $\hat{\epsilon}$ ,  $\hat{\epsilon}_{\text{num}}$  and  $\tilde{\epsilon}_{\text{mod}}$  are defined by

$$\hat{\epsilon} = \hat{e}/\|\mathbf{u}\|_{E(\Omega)}, \quad \hat{\epsilon}_{\text{num}} = \hat{e}_{\text{num}}/\|\mathbf{u}\|_{E(\Omega)}, \quad \tilde{\epsilon}_{\text{mod}} = \tilde{e}_{\text{mod}}/\|\mathbf{u}\|_{E(\Omega)}.$$

One can see that the effectivity indices,  $\theta$  and  $\theta_{\text{num}}$ , are close to one, especially for  $\theta$ . For the last two meshes, the modeling error can be evaluated by the difference between the estimated total error and the numerical errors and the modeling effectivity index is very good.

### 6.1. A simple adaptive strategy

An effective adaptive strategy for the HDPM should be able to (1) select the best suited homogenized material property  $\mathbf{E}^0$ , (2) partition  $\Omega$  into subdomains, (3) produce the  $H$ -mesh and (4) the  $h$ -mesh in

Table 5

Estimated errors (%) and effectivity indices with a growing number of elements in the  $h$ -mesh. There are 1000 particles, 256 subdomains and the mismatch is  $\tau = 10$

Nb elts	$\hat{\epsilon}$	$\theta$	$\hat{\epsilon}_{\text{num}}$	$\theta_{\text{num}}$	$\tilde{\epsilon}_{\text{mod}}$	$\tilde{\theta}_{\text{mod}}$
256	137.3	1.0001	213.3	1.66		
1024	92.07	1.0001	111.5	1.41		
4096	57.39	1.0004	34.65	1.07	45.76	0.970
16384	49.75	1.0005	15.57	1.02	47.27	0.999

Table 6

An example of simple adaptive strategy. 1000 particles,  $\mathbf{E}^\bullet = \langle \mathbf{E}^{-1} \rangle^{-1}$  and  $\tau = 10$ .

$N$	$\hat{\epsilon}$	$\tilde{\epsilon}_{\text{mod}}$	$\hat{\epsilon}_{\text{num}}$
256	49.77	47.27	15.57
128	32.21	29.16	13.66
64	21.17	16.94	12.70

order to minimize the cost of the computation for a prescribed total accuracy. The problem of choosing  $\mathbf{E}^0$  and the partition is already a difficult task by itself. Therefore, we make the following assumptions:

- The ‘best’ homogenized material property is known;
- No numerical errors are introduced by the  $H$ -mesh;
- We consider a uniform subdomain partition and uniform  $h$ -meshes;
- The subdomain analysis is carried out on each subdomain.

We try to adapt the subdomain size and the  $h$ -mesh size to reach a prescribed total accuracy  $\epsilon_0\%$ . Instead of minimizing the cost of the computation, we impose a given sharing of the modeling and numerical errors. This sharing is described by the parameter  $\alpha$  ( $0 < \alpha < 1$ ):  $\epsilon_{\text{mod}} = (\alpha)^{1/2}\epsilon_0$  and  $\epsilon_{\text{num}} = (1 - \alpha)^{1/2}\epsilon_0$ .

We propose the following simple strategy:

- Step 1: Start with a subdomain size which is at least three or four times the size of a particle;
- Step 2: With this size of subdomain, use standard  $h$ -adaptive finite element method to optimize the mesh to reach the prescribed numerical accuracy;
- Step 3: Iteratively reduce the number of subdomains, keeping the  $h$ -mesh fixed, until the prescribed modeling accuracy is reached.

As an example, we consider our 1-D problem with 1000 particles and  $\tau = 10$ . The homogenized modulus  $\mathbf{E}^0$  is taken as  $\mathbf{E}^0 = \langle \mathbf{E}^{-1} \rangle^{-1}$  and no errors are introduced by the  $H$ -mesh since the exact displacements values are obtained at the nodes (cf Section 3.2). We target a 20% error divided into 14.14% for modeling error and 14.14% for the numerical error ( $\alpha = 0.5$ ). We start with 256 subdomains and we optimize the mesh. Note that for our 1-D model problem, if we work with meshes matching the particles boundaries the numerical error is quickly very small. In order to get substantial numerical error, we choose to work with non matching meshes (our 1-D problem is academic and it is highly advised for 2- and 3-D problem to use meshes matching the boundaries of the heterogeneities). With 16,384 elements, we get  $\hat{\epsilon}_{\text{num}} = 15.57\%$ . The modeling error is  $\tilde{\epsilon}_{\text{mod}} = 47.27\%$ . The number of subdomains,  $N$ , is decreased to 128 and then to 64. The final error is 21.17%. Table 6 gives a summary of the adaptive process.

## 7. Conclusions

In the Homogenized Dirichlet Projection Method, numerical errors occur when solving the homogenized problem ( $H$ -mesh) and when performing the local analysis ( $h$ -mesh). The influences of the  $h$ - and  $H$ -mesh on the error are completely different. The influence of the  $H$ -mesh is expressed in terms of the difference between the numerical and exact homogenized displacement on the boundary of the subdomains. As  $H \rightarrow 0$ , the total error decreases or increases depending on the choice made for the homogenized material property. Conversely, the influence of the  $h$ -mesh is expressed as the distance in the energy norm between the exact and numerical solution in displacement of the local subdomain problems. If the nodes of the  $h$ -mesh and the particle boundaries are not matching, the rate of convergence of the numerical error is very poor,  $O(\sqrt{h})$  for piecewise linear elements. Moreover, this rate of convergence is achieved only if the mesh size is *smaller* than a critical length being three to four times the size of the particles. On the contrary, if the  $h$ -mesh matches the particle boundaries, the rate is  $O(h)$  for piecewise linear elements and the critical mesh size is no longer present.

Concerning the modeling error, a critical length also appears, the same as for the numerical error: the subdomain size should be *bigger* than three to four times the size of the particles to allow an efficient decomposition into subdomains, for our 1-D model problem.

Using the ECL concept, we are able to give a mechanical interpretation of the main results of the HDPM and to extend them to nonlinear constitutive laws.

Finally, computable upper bounds for the total and numerical error are obtained and the effectivity indices obtained for our 1-D model problem are close to one. A simple adaptive strategy is also proposed to choose the size of the subdomain and the  $h$ -mesh size.

## Acknowledgment

The authors gratefully acknowledge the support of this work by the U.S. Office of Naval Research under Grant N00014-95-1-0401 and the National Science Foundation under grant ECS-9422707.

## References

- [1] Tarek I. Zohdi, J.T. Oden, and Gregory J. Rodin, Hierarchical modeling of heterogeneous bodies, TICAM report 96-21, Texas Institute for Computational and Applied Mathematics, Austin, 1996. Also to appear in *Comput. Methods Appl. Mech. Engrg.*
- [2] J.T. Oden and Tarek I. Zohdi, Analysis of elastic structures composed of highly heterogeneous materials, TICAM report 96-56, Texas Institute for Computational and Applied Mathematics, Austin, 1996. Also to appear in *Comput. Methods Appl. Mech. Engrg.*
- [3] P. Ladevèze and D. Leguillon, Error estimate procedure in the finite element method and application, *SIAM J. Numer. Anal.* 20(3) (1983) 485–509.
- [4] P. Ladevèze and M. Reynier, A localisation method of stiffness errors for the adjustment of F.E. models, In Special issue, 12th ASME Mechanical Vibration and Noise Conference, Montreal, 1989.
- [5] W. Prager and J.L. Synge, Approximation in elasticity based on the concept of functions space, *Quart. Appl. Math.* 5 (1947) 261–269.