

**Original citation:**

Langwallner, Bernhard, Ortner, Christoph and Süli, Endre. (2010) Existence and convergence results for the Galerkin approximation of an electronic density functional. *Mathematical Models and Methods in Applied Sciences*, Vol.20 (No.12). pp. 2237-2265.

**Permanent WRAP url:**

<http://wrap.warwick.ac.uk/43806>

**Copyright and reuse:**

The Warwick Research Archive Portal (WRAP) makes this work of researchers of the University of Warwick available open access under the following conditions. Copyright © and all moral rights to the version of the paper presented here belong to the individual author(s) and/or other copyright owners. To the extent reasonable and practicable the material made available in WRAP has been checked for eligibility before being made available.

Copies of full items can be used for personal research or study, educational, or not-for-profit purposes without prior permission or charge. Provided that the authors, title and full bibliographic details are credited, a hyperlink and/or URL is given for the original metadata page and the content is not changed in any way.

**Publisher's statement:**

Electronic version of an article published as (2010) Existence and convergence results for the Galerkin approximation of an electronic density functional. *Mathematical Models and Methods in Applied Sciences*, Vol.20 (No.12). pp. 2237-2265

<http://dx.doi.org/10.1142/S021820251000491X> © World Scientific Publishing Company,  
<http://www.worldscientific.com/worldscinet/m3as>

**note on versions:**

The version presented here may differ from the published version or, version of record, if you wish to cite this item you are advised to consult the publisher's version. Please see the 'permanent WRAP url' above for details on accessing the published version and note that access may require a subscription.

For more information, please contact the WRAP Team at: [publications@warwick.ac.uk](mailto:publications@warwick.ac.uk)

warwick**publications**wrap  
  
highlight your research

<http://wrap.warwick.ac.uk/>

Mathematical Models and Methods in Applied Sciences  
© World Scientific Publishing Company

## Existence and Convergence Results for the Galerkin Approximation of an Electronic Density Functional

BERNHARD LANGWALLNER

*University of Oxford, Mathematical Institute,  
24-29 St Giles', Oxford OX1 3LB, UK.  
langwallner@maths.ox.ac.uk*

CHRISTOPH ORTNER

*University of Oxford, Mathematical Institute,  
24-29 St Giles', Oxford OX1 3LB, UK.  
ortner@maths.ox.ac.uk*

ENDRE SÜLI

*University of Oxford, Mathematical Institute,  
24-29 St Giles', Oxford OX1 3LB, UK.  
suli@maths.ox.ac.uk*

Received (Day Month Year)  
Revised (Day Month Year)  
Communicated by (xxxxxxxxxx)

We formulate and analyze a model for the study of finite clusters of atoms or localized defects in infinite crystals based on orbital-free density functional theory. We show that the resulting constrained optimization problem has a minimizer and we provide a careful analysis of the solubility of the associated Euler–Lagrange equation. Based on these results, and using tools from saddle-point theory and nonlinear analysis, we then show that a Galerkin discretization has a solution that converges to the correct limit.

*Keywords:* Nonlinear eigenvalue problem, Galerkin discretization, Thomas–Fermi type functionals

AMS Subject Classification: 65N60, 65N25, 65Z05

### 1. Introduction

Quantum mechanics is the accepted microscopic theory of atoms, molecules, and solids. Quantum mechanical simulations have therefore become an indispensable tool in physics, chemistry, and materials science. In particular, they are considered *ab initio*, that is, no empirical input is necessary. Computations based on the Schrödinger equation (the governing equation of quantum mechanics) require the solution of a linear partial differential equation in a high-dimensional configuration space. This fact dramatically limits the size of the problems that can be tackled.<sup>22</sup>

2 *Bernhard Langwallner, Christoph Ortner, and Endre Süli*

Possibly the most widely employed remedy is known as *Density Functional Theory* (DFT).<sup>21,23</sup> Its theoretical justification rests on the fact that all physical ground state properties of a quantum mechanical system can directly be computed from the ground state electron density rather than the electronic wave function. The former is defined on  $\mathbb{R}^3$  whereas the latter is defined on  $\mathbb{R}^{3N}$ ,  $N$  being the number of electrons.

The most demanding task in DFT calculations is the computation of the kinetic energy. It involves repeated numerical solution of single electron eigenvalue problems and subsequent orthogonalization of the resulting wave functions. In order to avoid this huge cost, approximate forms have also been developed for the kinetic energy. The resulting models are commonly referred to as orbital-free density functional theory.<sup>13,26,30</sup> Early models of that class include Thomas–Fermi type functionals; see for example Finnis.<sup>13</sup> Orbital-free DFT has been observed to work well for systems with weakly varying electron density, most prominently aluminium, see Wang, Govind & Carter.<sup>33,34</sup>

A simple example of a density functional is given by the Thomas–Fermi–Dirac–von Weizsäcker (TFDW) functional. Let  $\rho$  be the electron density,  $R = \{R_i, i = 1, \dots, N_{\text{nuc}}\}$  the set of nucleus positions and  $Z_i \in \mathbb{N}, i = 1, \dots, N_{\text{nuc}}$ , the charges of the nuclei. Then the energy is written as

$$E(\rho, R) = T_s(\rho) + E_{\text{xc}}(\rho) + E_{\text{H}}(\rho) + E_{\text{ext}}(\rho, R) + E_{\text{zz}}(R). \quad (1.1)$$

Here,  $T_s$  is the kinetic energy

$$T_s(\rho) = \frac{3}{10}(3\pi^2)^{2/3} \int_{\mathbb{R}^3} \rho^{5/3}(x) \, dx + \frac{\lambda}{8} \int_{\mathbb{R}^3} \frac{|\nabla \rho(x)|^2}{\rho(x)} \, dx,$$

where  $\lambda > 0$  is a parameter. The exchange correlation energy  $E_{\text{xc}}$  in the so-called local density approximation<sup>9,31</sup> has the form

$$E_{\text{xc}}(\rho) = \int_{\mathbb{R}^3} (-C_x \rho(x)^{4/3} + \varepsilon_c(\rho(x))) \rho(x) \, dx,$$

where  $\varepsilon_c$  is a phenomenological correction term. We note that this term renders  $E(\cdot, R)$  non-convex.

The remaining terms are of electrostatic nature. The Hartree term

$$E_{\text{H}}(\rho) = \frac{1}{2} \int_{\mathbb{R}^3} \int_{\mathbb{R}^3} \frac{\rho(x)\rho(y)}{|x-y|} \, dx \, dy$$

incorporates interaction among electrons. The energy of electrons in the electrostatic field generated by the nuclei is accounted for by

$$E_{\text{ext}}(\rho, R) = \int_{\mathbb{R}^3} \rho(x) V_{\text{ext}}(x) \, dx, \quad \text{where} \quad V_{\text{ext}}(x) = \sum_i \frac{-Z_i}{|x - R_i|}.$$

Finally, the electrostatic repulsion energy of the nuclei is

$$E_{\text{zz}}(R) = \frac{1}{2} \sum_{i=1}^{N_{\text{nuc}}} \sum_{\substack{j=1 \\ j \neq i}}^{N_{\text{nuc}}} \frac{Z_i Z_j}{|R_i - R_j|}.$$

For numerical or modeling reasons the nuclear point charges  $Z_i$  at  $R_i$  may be replaced by a smooth charge density  $\rho_n$  (see Gavini et al.<sup>18</sup>). For example,  $\rho_n$  may be written as a sum of compactly supported smooth functions centered at the positions  $R_i$ ,  $i = 1, \dots, N_{\text{nuc}}$ , (the dependence of  $\rho_n$  on  $R$  will be suppressed)

$$\rho_n(x) = \sum_{i=1}^{N_{\text{nuc}}} Z_i \tilde{\rho}_0(x - R_i),$$

where  $\tilde{\rho}_0 \in C_0^\infty(\mathbb{R}^3)$ ,  $\tilde{\rho}_0 \geq 0$ , and  $\int \tilde{\rho}_0(x) dx = 1$ . Then, the repulsion energy takes the form

$$E_{\text{zz}}(R) = \frac{1}{2} \int_{\mathbb{R}^3} \int_{\mathbb{R}^3} \frac{\rho_n(x)\rho_n(y)}{|x-y|} dx dy.$$

This expression includes the potential energy of every nucleus in its own field. This is, however, only a constant contribution to the overall energy and may easily be subtracted. The potential for the nucleus-electron interaction is given by

$$V_{\text{ext}}(x) = - \int_{\mathbb{R}^3} \frac{\rho_n(y)}{|x-y|} dy,$$

which allows for a symmetric expression for the sum of all electrostatic terms

$$E_{\text{H}}(\rho) + E_{\text{ext}}(\rho, R) + E_{\text{zz}}(R) = \int_{\mathbb{R}^3} \int_{\mathbb{R}^3} \frac{(\rho(x) - \rho_n(x))(\rho(y) - \rho_n(y))}{|x-y|} dy dx.$$

The nonlocal nature of this term represents a numerical challenge for all density functional calculations. The given density functional has to be minimized subject to the constraints

$$\int_{\mathbb{R}^3} \rho(x) dx = N, \quad \text{and} \quad \rho \geq 0.$$

A survey of theoretical results in connection with Thomas–Fermi type models can be found in an article by Lieb.<sup>25</sup> The author considers different models of increasing complexity in  $\mathbb{R}^3$  and assesses their mathematical structure and physical validity as well as their relations to quantum mechanics.

Several numerical approximations of the functional (1.1) or related models have been proposed in the literature. For Galerkin discretizations, the crucial question is the choice of basis. The most popular basis sets are plane waves (i.e. Fourier modes<sup>7,35</sup>) and finite elements.<sup>18</sup> Plane waves can only be applied to periodic systems, which means, for example, that no defects can be simulated (usually, periodic arrays of defects are considered instead). On the other hand, the implementation of the Coulomb interaction kernel can be done very efficiently. Finite elements are not inherently periodic and allow for adaptivity in space, which is particularly useful for additional coarse-graining approximations.<sup>17</sup> Calculating the convolution in the electrostatic terms remains a challenge. Finite difference approximations were suggested in References 14, 15, 28, 29.

Little work can be found in the literature on the numerical analysis of these approximations. The only convergence result in a finite element context we are aware

4 *Bernhard Langwallner, Christoph Ortner, and Endre Süli*

of is a proof of  $\Gamma$ -convergence in Gavini et al.<sup>18</sup>, which contains no information about convergence rates. The latter, however, are crucial for understanding the efficiency of the numerical method. In Cancès, Chakir & Maday<sup>6,7</sup> the authors perform a rigorous numerical analysis of a planewave discretization of a density functional. That work also includes convergence rates.

In our present and ongoing work<sup>24</sup> we aim to fill this gap and provide a complete convergence theory for Galerkin discretizations of orbital-free electronic density functionals. In Section 2.1 we first formulate a model in a bounded domain that allows the simulation of finite clusters or isolated defects in an infinite medium. We prove existence of minimizers and give a careful analysis of first- and second-order optimality conditions in the remainder of Section 2. Our reformulation of the optimality system is particularly suitable for Galerkin finite element discretizations. The main result in Section 2.3 connects the stability of the minimization problem to the stability of this system. This result allows us, in Section 3, to prove the existence and convergence of Galerkin discretizations of the optimality system.

The present work is of, predominantly, theoretical/analytical nature. The important practical issues of optimal convergence rates and numerical integration will be addressed in a forthcoming article.<sup>24</sup> All of the results discussed herein carry over to the formulation *with* numerical integration.

## 2. Existence and Analysis of Minimizers

In this section we will suggest a mathematical model based on the functional described in the introduction and study the resulting minimization problem as well as the associated optimality conditions. First we discuss a few ideas that simplify the TFDW functional with regard to subsequent numerical approximation.

The constraint  $\rho \geq 0$  can be enforced by setting  $\rho = u^2$ . This substitution also has the advantage that the term involving  $\nabla\rho$  becomes easier to evaluate:

$$E(u, R) = \frac{\lambda}{2} \int_{\mathbb{R}^3} |\nabla u|^2 dx + C_{\text{TF}} \int_{\mathbb{R}^3} |u|^{10/3} dx - C_x \int_{\mathbb{R}^3} |u|^{8/3} dx \quad (2.1)$$

$$+ \int_{\mathbb{R}^3} \varepsilon_c(u^2)u^2 dx - \frac{1}{2} \int_{\mathbb{R}^3} \int_{\mathbb{R}^3} \frac{(u^2(x) - \rho_n(x))(u^2(y) - \rho_n(y))}{|x - y|} dx dy.$$

We also need to address the evaluation of the electrostatic term. The double integral can in principle be computed explicitly using the Fourier transform.<sup>14,16,35</sup> In Gavini et al.<sup>17,18</sup> the authors suggest a different approach that makes use of the special structure, respectively the physics, represented by the term. Note that the integral kernel in the last term in (2.1) is the Green's function of the Poisson equation in  $\mathbb{R}^3$ . Therefore, the electrostatic potential

$$\phi(x) := \int_{\mathbb{R}^3} \frac{(u^2(y) - \rho_n(y))}{|x - y|} dy$$

is simply the solution of the equation

$$-\frac{1}{4\pi} \Delta\phi = u^2 - \rho_n,$$

subject to homogeneous Dirichlet boundary condition at infinity. From this, it can be deduced that

$$\frac{1}{2} \int_{\mathbb{R}^3} \int_{\mathbb{R}^3} \frac{(u^2(x) - \rho_n(x))(u^2(y) - \rho_n(y))}{|x - y|} dx dy = \frac{1}{2} \int_{\mathbb{R}^3} (u^2 - \rho_n) \phi dx.$$

Formally, the right-hand side can also be written as

$$\frac{1}{2} \int_{\mathbb{R}^3} (u^2 - \rho_n) \phi dx = - \inf_{\varphi} \left[ \int_{\mathbb{R}^3} \frac{1}{8\pi} |\nabla \varphi|^2 - (u^2 - \rho_n) \varphi dx \right]. \quad (2.2)$$

This equality is referred to as the direct Coulomb formula in Gavini et al.<sup>18</sup>

### 2.1. Artificial boundary conditions

We now introduce a model to simulate finite clusters of atoms or localized phenomena in infinite crystals (e.g. vacancies, interstitials or dislocation loops). For both computational and analytical reasons, we would like to formulate the problem in a bounded open domain  $\Omega \subset \mathbb{R}^3$  rather than in the whole of  $\mathbb{R}^3$ . To this end, we will assume that we know the square root density  $u$  as well as the electrostatic potential  $\phi$  in  $\mathbb{R}^3 \setminus \Omega$ , which induces *artificial Dirichlet boundary conditions* for  $u$  and  $\phi$  on  $\partial\Omega$ . We discuss potential pitfalls of this approach in Remark 2.1 below.

There are two specific examples that we have in mind. If a finite cluster of atoms is studied, then we set  $u = \phi = 0$  in  $\mathbb{R}^3 \setminus \Omega$ . If we study a localized defect in an infinite crystal (e.g., a vacancy or a dislocation) then we let  $u$  and  $\phi$  on  $\partial\Omega$  be the square root density and electrostatic potential of a perfect crystal. (Minimizers of Thomas–Fermi type functionals for perfect crystals have been studied in Blanc, Le Bris & Lions<sup>2</sup> and Catto, Le Bris & Lions.<sup>8</sup>)

To make this concrete, we assume that we are given functions  $u_{\text{ex}}, \phi_{\text{ex}} \in H_{\text{loc}}^2(\mathbb{R}^3)$  and define the admissible set

$$A_u = \left\{ u \in u_{\text{ex}} + H_0^1(\Omega) : \|u\|_{L^2}^2 = N \right\}, \quad (2.3)$$

and the energy functional (suppressing the nuclei positions, which are held fixed)

$$E(u) = T(u) + X(u) + \Phi(u), \quad (2.4)$$

where

$$\begin{aligned} T(u) &= \frac{\lambda}{2} \|\nabla u\|_{L^2(\Omega)}^2, \\ X(u) &= \int_{\Omega} F(u) dx, \quad \text{and} \end{aligned} \quad (2.5)$$

$$\Phi(u) = - \inf_{\phi \in \phi_{\text{ex}} + H_0^1(\Omega)} \tilde{\Phi}(u, \phi) = - \inf_{\phi \in \phi_{\text{ex}} + H_0^1(\Omega)} \left[ \int_{\Omega} \frac{1}{8\pi} |\nabla \phi|^2 - (u^2 - \rho_n) \phi dx \right].$$

We have split the energy functional, in a way that is convenient for the analysis, into a quadratic, a nonlinear local, and a nonlocal part. In the original model the function  $F$  is given by

$$F(u) = C_{\text{TF}} |u|^{10/3} - C_x |u|^{8/3} + \varepsilon_c (u^2) u^2 \quad (2.6)$$

6 *Bernhard Langwallner, Christoph Ortner, and Endre Süli*

and combines a portion of the kinetic energy and the exchange-correlation energy. However, our results are independent of the precise form of  $F$  and only require a certain degree of smoothness and growth conditions, which we make precise below.

It is also worth remarking that in the case of homogeneous boundary conditions on  $u$  and  $\phi$  ( $u_{\text{ex}} = \phi_{\text{ex}} = 0$ ), the energy  $\Phi$  reduces to  $\Phi(u) = \frac{1}{2} \int_{\Omega} (u^2 - \rho_n) \phi \, dx$ , where  $\phi$  is the weak solution of the equation  $-\Delta \phi = 4\pi(u^2 - \rho_n)$  with homogeneous Dirichlet boundary condition. In this case  $\Phi$  is easily recognizable as the potential energy of the charge density  $(u^2 - \rho_n)$  in its own electrostatic field.

For future reference, we also define the constraint functional  $c : H^1(\Omega) \rightarrow \mathbb{R}$ ,

$$c(u) = \frac{1}{2} \left( \|u\|_{L^2}^2 - N \right).$$

Our goal is to solve the minimization problem

$$\min_{u \in A_u} E(u). \tag{2.7}$$

**Remark 2.1.** Our justification for this approach is that, if  $\bar{u}$  is a minimizer of the original problem, with associated electrostatic potential  $\bar{\phi}$ , then, setting  $u_{\text{ex}} = \bar{u}$  and  $\phi_{\text{ex}} = \bar{\phi}$ ,  $\bar{u}|_{\Omega}$  is a solution to (2.7) with associated electrostatic potential  $\bar{\phi}|_{\Omega}$ . In that sense, the problem with artificial boundary conditions is *consistent* with the original problem posed in  $\mathbb{R}^3$ .

Since we do not normally know the exact electron density and electrostatic potential outside of  $\Omega$ , we are essentially forced to make our ‘best guess’ for the problem at hand. This creates an error in the system, which cannot be controlled by adjusting the discretization. One usually hopes that, by choosing domains of increasing size, this error will shrink to zero, however, it is far from straightforward to establish this rigorously for any system other than a (near-)perfect crystal.

For simplicity, let us discuss the case of a finite cluster where we set  $u_{\text{ex}} = \phi_{\text{ex}} = 0$  in  $\mathbb{R}^3 \setminus \Omega$ . Hence, the question arises, if  $(\bar{u}, \bar{\phi})$  are the *exact* square root density and electrostatic potential of the system, how fast  $|\bar{u}(x)|$  and  $|\bar{\phi}(x)|$  decay as  $|x| \rightarrow \infty$ .

Whereas quantum mechanics suggests that the decay of  $u$  is exponential, this is less clear for  $\phi$ . Note, in particular, that variations of  $u$ , well inside  $\Omega$ , can in general create comparatively large variations of  $\phi$  in all of  $\mathbb{R}^3$ . For rather special configurations of the nuclei, we may hope that lower order multipole moments vanish, assuring a sufficiently fast decay of  $\phi$  (e.g., if there is a point symmetry).

In the present work we will simply assume that these deviations decay sufficiently fast and concentrate on aspects of numerical approximation theory. We are planning to study the issues pointed out in this remark in future work. In particular, we consider the present work a first step towards an analysis of a combined TFDW / Quasicontinuum model in the spirit of Garcia-Cervera, Lu & E<sup>16</sup> and Gavini et al.<sup>17</sup> For this type of analysis it will also be necessary to study deviations of  $u$  and  $\phi$  from the perfect crystal case introduced by localized defects.  $\square$

## 2.2. Existence of a minimizer

We assume from now on that the homogeneous Dirichlet problem is  $H^2$ -regular, that is, if  $f \in L^2(\Omega)$  and if  $v \in H_0^1(\Omega)$  is the solution of

$$(\nabla v, \nabla w) = (f, w) \quad \forall w \in H_0^1(\Omega), \quad (2.8)$$

then  $v \in H^2(\Omega)$ , and there exists a constant  $C_{\text{reg}}$ , independent of  $f$ , such that

$$\|v\|_{H^2} \leq C_{\text{reg}} \|f\|_{L^2}.$$

This is the case, for example, if  $\Omega$  is convex<sup>20</sup> or if  $\Omega$  is  $C^2$ -regular.<sup>12</sup> For notational convenience we define the solution operator of the Poisson equation with homogeneous Dirichlet boundary condition:  $(-\Delta_0)^{-1} : L^2(\Omega) \rightarrow H^2(\Omega) \cap H_0^1(\Omega)$ ,  $f \mapsto v$ , where  $v$  solves (2.8).

Throughout the analysis, the positions  $R_i$  of the nuclei are fixed. We assume that  $F \in C^2(\mathbb{R})$  and that it satisfies the growth condition

$$a_1 \leq F(t) \leq c_2 t^q + a_2 \quad \forall t \in \mathbb{R}, \quad (2.9)$$

with constants  $a_1, a_2 \in \mathbb{R}$ ,  $c_2 \geq 0$ , and  $3 \leq q < 6$ . Moreover, we assume that  $F''$  is locally Hölder continuous; more precisely, there exists a positive constant  $C$ , such that, for all  $t_1, t_2 \in \mathbb{R}$ ,

$$|F''(t_1) - F''(t_2)| \leq C(1 + |t_1|^{q-3} + |t_2|^{q-3})(|t_1 - t_2|^\alpha + |t_1 - t_2|), \quad (2.10)$$

where  $0 < \alpha \leq 1$ . We note that these conditions are satisfied if  $F$  is defined by (2.6).

The functional  $T : H^1(\Omega) \rightarrow \mathbb{R}$  is strongly continuous, twice continuously Fréchet differentiable and weakly lower semicontinuous.

It is also straightforward to verify that, under condition (2.10), the functional  $X : H^1(\Omega) \rightarrow \mathbb{R}$  defined by (2.5) is twice continuously Fréchet differentiable with

$$\langle X'(u), h \rangle = \int_{\Omega} F'(u) h \, dx, \quad \text{and} \quad \langle X''(u)h_1, h_2 \rangle = \int_{\Omega} F''(u)h_1 h_2 \, dx,$$

for  $h, h_1, h_2 \in H^1(\Omega)$ . Furthermore,  $X$  and  $X'$  are locally Lipschitz continuous, and  $X''$  is locally Hölder continuous, with

$$\|X''(u) - X''(v)\| \leq C(1 + \|u\|_{H^1}^{q-3} + \|v\|_{H^1}^{q-3})(\|u - v\|_{H^1}^\alpha + \|u - v\|_{H^1}), \quad (2.11)$$

for all  $u, v \in H^1(\Omega)$ . To prove this we begin as follows: for  $u, v, h_1, h_2 \in H^1(\Omega)$  we have that

$$\begin{aligned} |\langle (X''(u) - X''(v))h_1, h_2 \rangle| &\leq \int_{\Omega} |F''(u) - F''(v)| |h_1| |h_2| \, dx & (2.12) \\ &\leq C \int_{\Omega} |h_1| |h_2| (1 + |u|^{q-3} + |v|^{q-3})(|u - v|^\alpha + |u - v|) \, dx \\ &\leq C \|h_1\|_{L^4} \|h_2\|_{L^4} (1 + \| |u|^{q-3} \|_{L^4} + \| |v|^{q-3} \|_{L^4}) \\ &\quad \cdot (\| |u - v|^\alpha \|_{L^4} + \|u - v\|_{L^4}), \end{aligned}$$



8 *Bernhard Langwallner, Christoph Ortner, and Endre Süli*

where we have used a generalized version of Hölder's inequality. Now we know that  $\| |u|^\gamma \|_{L^k} \leq C \|u\|_{L^k}^\gamma$  for all  $0 \leq \gamma < 1$  and all  $k > 1$ :

$$\| |u|^\gamma \|_{L^k}^k = \int_{\Omega} 1 \cdot |u|^{\gamma k} dx \leq C \| |u|^{\gamma k} \|_{L^{1/\gamma}} = C \|u\|_{L^k}^{\gamma k}.$$

Applying this inequality to (2.12) and using the embedding  $H^1(\Omega)$  in  $L^4(\Omega)$  we deduce (2.11).

Since  $X$  is strongly continuous in  $L^q(\Omega)$  and since  $H^1(\Omega)$  is compactly embedded in  $L^p(\Omega)$ , for  $1 \leq p < 6$ , see Adams,<sup>1</sup> it follows also that  $X$  is sequentially weakly continuous on  $H^1(\Omega)$ .

Next we consider the functional  $\Phi$  defined in (2.5).

**Lemma 2.1.** *The functional  $\Phi : L^4(\Omega) \rightarrow \mathbb{R}$  is well-defined and continuous;  $\Phi$  is twice continuously Fréchet differentiable and the derivatives are, for  $u \in L^4(\Omega)$ , given by*

$$\begin{aligned} \langle \Phi'(u), h \rangle &= 2 \int_{\Omega} u \phi h dx \quad \text{for } h \in L^4(\Omega) \text{ and,} \\ \langle \Phi''(u)h_1, h_2 \rangle &= 2 \int_{\Omega} \phi h_1 h_2 + 8\pi u h_1 (-\Delta_0)^{-1}(u h_2) dx \quad \text{for } h_1, h_2 \in L^4(\Omega), \end{aligned}$$

where  $\phi$  is the weak solution of  $-\Delta \phi = 4\pi(u^2 - \rho_n)$ ,  $\phi|_{\partial\Omega} = \phi_{\text{ex}}$ . Furthermore,  $\Phi$  is bounded below on the set  $A_u$  defined in (2.3), and the restriction  $\Phi|_{H^1(\Omega)}$  is sequentially weakly continuous in  $H^1(\Omega)$ .

**Proof.** For every  $u \in L^4(\Omega)$ , the functional  $\tilde{\Phi}(u, \cdot)$  from (2.5) is clearly continuous, convex, coercive and weakly lower semicontinuous on  $V_\phi = \{\psi \in H^1(\Omega) : \psi|_{\partial\Omega} = \phi_{\text{ex}}\}$ . Thus, there exists a unique minimizer  $\phi_u$ . This minimizer satisfies the equation

$$(\nabla \phi_u, \nabla \psi) = 4\pi(u^2 - \rho_n, \psi) \quad \forall \psi \in H_0^1(\Omega),$$

with boundary condition  $\phi_u|_{\partial\Omega} = \phi_{\text{ex}}$ . The auxiliary function

$$\xi = \phi_{\text{ex}} - (-\Delta_0)^{-1}(-\Delta)\phi_{\text{ex}} \in \phi_{\text{ex}} + H_0^1(\Omega)$$

satisfies  $(\nabla \xi, \nabla \psi) = 0$  for all  $\psi \in H_0^1(\Omega)$ . From this, it follows that

$$\phi_u = 4\pi(-\Delta_0)^{-1}(u^2 - \rho_n) + \xi$$

and, moreover, after straightforward algebraic manipulations, that

$$\begin{aligned} \Phi(u) &= -\tilde{\Phi}(u, \phi_u) = 2\pi \int_{\Omega} (u^2 - \rho_n)(-\Delta_0)^{-1}(u^2 - \rho_n) dx \\ &\quad + \int_{\Omega} (u^2 - \rho_n)\xi dx - \frac{1}{8\pi} \int_{\Omega} |\nabla \xi|^2 dx. \end{aligned} \tag{2.13}$$

Differentiating with respect to  $u$  yields the expressions for  $\Phi'$  and  $\Phi''$  as given above.

Finally, we show that  $\Phi|_{A_u}$  is bounded below. Clearly, the first term on the right-hand side in (2.13) is non-negative, and the last term is a constant depending only on  $\phi_{\text{ex}}$ . Since we assumed that  $\phi_{\text{ex}} \in H_{\text{loc}}^2(\mathbb{R}^3)$  and that the Poisson problem is

$H^2$ -regular, it follows that  $\xi \in L^\infty(\Omega)$ . Therefore, the second term on the right-hand side of (2.13) can be bounded as follows:

$$\left| \int_{\Omega} (u^2 - \rho_n) \xi \, dx \right| \leq \|\xi\|_{L^\infty} (\|u\|_{L^2}^2 + \|\rho_n\|_{L^1}) = \|\xi\|_{L^\infty} (N + \|\rho_n\|_{L^1}).$$

Hence, we can deduce that  $\Phi(u) \geq C(N, \rho_n, \phi_{\text{ex}})$  for all  $u \in A_u$ .

The sequential weak continuity of  $\Phi|_{H^1(\Omega)}$  is a direct consequence of the compact embedding of  $H^1(\Omega)$  in  $L^4(\Omega)$  (see Theorem 6.3 in Adams<sup>1</sup>) and the strong continuity of  $\Phi$  on  $L^4(\Omega)$ .  $\square$

**Theorem 2.1.**  $E : A_u \rightarrow \mathbb{R}$  has at least one minimizer.

**Proof.** We apply the direct method of the calculus of variations.<sup>10</sup> First, we observe that  $E$  is coercive on  $A_u$ , which can be seen as follows:

$$\begin{aligned} E(u) &\geq \frac{\lambda}{2} \|\nabla u\|_{L^2}^2 + \int_{\Omega} a_1 \, dx + \Phi(u) \\ &\geq \frac{\lambda}{2} \|\nabla u\|_{L^2}^2 + a_1 |\Omega| - C(N, \rho_n, \phi_{\text{ex}}) \\ &\geq C_1 \|u\|_{H^1}^2 - C_2 \|u_{\text{ex}}\|_{H^1}^2 - C_3(\Omega, F, N, \rho_n, \phi_{\text{ex}}). \end{aligned}$$

Here we have used the growth condition (2.9) on  $F$ , the lower bound on  $\Phi(u)$ , established in Lemma 2.1, and Poincaré's inequality for  $u - u_{\text{ex}}$ . Hence, minimizing sequences of  $E$  are bounded in  $H^1(\Omega)$ , and we can find a weakly convergent subsequence. Since  $E$  is weakly lower semicontinuous as a sum of weakly lower semicontinuous functions, it follows that the weak limit of the subsequence is a minimizer; see for example Dacorogna.<sup>10</sup>  $\square$

We mention at this point that, since  $u$  was defined to be the square root of  $\rho$ , the minimizer in the case of homogeneous boundary conditions cannot be unique: if  $u$  minimizes  $E$ , then so does  $-u$ . The question whether there can be more than two minimizers of the energy (2.4) is beyond the scope of this article.

### 2.3. The Euler–Lagrange Equations

So far, we have shown existence of solutions to the minimization problem (2.7). Next, we are interested in the characterization of such points, i.e., in optimality conditions. For example, an article by Maurer & Zowe<sup>27</sup> provides necessary and sufficient optimality conditions for a large class of optimization problems in Banach spaces.

Throughout this section, the derivatives  $E'(\bar{u})$  and  $c'(\bar{u})$  are understood as elements of  $H^{-1}(\Omega) := H_0^1(\Omega)^*$ , that is, they operate on functions from  $H_0^1(\Omega)$ . This means, for example, that  $\ker c'(\bar{u}) \subset H_0^1(\Omega)$ .

Theorem 3.1 in Maurer & Zowe<sup>27</sup> yields the first-order necessary optimality condition

$$\langle E'(\bar{u}), v \rangle = 0 \quad \forall v \in \ker c'(\bar{u}) \subset H_0^1(\Omega),$$

10 *Bernhard Langwallner, Christoph Ortner, and Endre Süli*

provided that  $0 \in \text{int} \{ \langle c'(\bar{u}), v \rangle : v \in \mathbf{H}_0^1(\Omega) \} \subseteq \mathbb{R}$ . Since  $\langle c'(\bar{u}), v \rangle = (\bar{u}, v)$  and  $\bar{u} \neq 0$ , this condition is always satisfied. Theorem 3.2 in Maurer & Zowe<sup>27</sup> yields the existence of a *Lagrange multiplier*  $\bar{\mu} \in \mathbb{R}$  such that

$$E'(\bar{u}) + \bar{\mu}c'(\bar{u}) = 0 \in \mathbf{H}^{-1}(\Omega), \quad \text{and} \quad c(\bar{u}) = 0. \quad (2.14)$$

Using the definitions and results of the previous section we now rewrite the first-order optimality system in a way that is more suitable for numerical approximation. Let  $\bar{u} \in A_u$  be a local minimizer with associated Lagrange multiplier  $\bar{\mu} \in \mathbb{R}$  and let  $\bar{\phi} \in \mathbf{H}^1(\Omega)$  be the associated electrostatic potential, then Lemma 2.1 implies that  $\bar{u}$ ,  $\bar{\phi}$  and  $\bar{\mu} \in \mathbb{R}$  solve the nonlinear system

$$\begin{aligned} \lambda(\nabla u, \nabla v) + (F'(u), v) + 2(\phi u, v) + \mu(u, v) &= 0 \quad \forall v \in \mathbf{H}_0^1(\Omega), \\ \frac{1}{4\pi}(\nabla \phi, \nabla \psi) - (u^2 - \rho_n, \psi) &= 0 \quad \forall \psi \in \mathbf{H}_0^1(\Omega), \\ \frac{\nu}{2} \left( \int_{\Omega} u^2 \, dx - N \right) &= 0 \quad \forall \nu \in \mathbb{R} \end{aligned} \quad (2.15)$$

with the boundary conditions

$$u|_{\partial\Omega} = u_{\text{ex}}, \quad \text{and} \quad \phi|_{\partial\Omega} = \phi_{\text{ex}}.$$

We will focus on solving this system instead of (2.14) or the minimization problem (2.7). It has to be pointed out that solving (2.15) (or even (2.14)) is not equivalent to solving the minimization problem since the functional is non-convex. However, we will focus on those solutions of (2.15) which correspond to local minimizers of (2.7), making use of the second-order optimality condition (2.19).

We define the function spaces

$$\begin{aligned} \mathcal{Y} &= \mathbf{H}^1(\Omega) \times \mathbf{H}^1(\Omega) \times \mathbb{R}, \quad \mathcal{Y}_D = (u_{\text{ex}} + \mathbf{H}_0^1(\Omega)) \times (\phi_{\text{ex}} + \mathbf{H}_0^1(\Omega)) \times \mathbb{R}, \\ \mathcal{Y}_0 &= \mathbf{H}_0^1(\Omega) \times \mathbf{H}_0^1(\Omega) \times \mathbb{R}, \quad \mathcal{Y}_0^* = \mathbf{H}^{-1}(\Omega) \times \mathbf{H}^{-1}(\Omega) \times \mathbb{R}, \end{aligned}$$

so that the system (2.15) defines an operator  $\mathcal{F} : \mathcal{Y}_D \rightarrow \mathcal{Y}_0^*$ , rewritten as

$$\langle \mathcal{F}(u, \phi, \mu), (v, \psi, \nu) \rangle = 0 \quad \forall (v, \psi, \nu) \in \mathcal{Y}_0. \quad (2.16)$$

For future reference, we split  $\mathcal{F}$  into

$$\mathcal{F}(u, \phi, \mu) = \begin{pmatrix} -\lambda\Delta u \\ -\frac{1}{4\pi}\Delta\phi \\ 0 \end{pmatrix} + \mathcal{G}(u, \phi, \mu),$$

where  $\mathcal{G}$  contains all terms without derivatives, and where the Laplacian  $\Delta$  is understood as a linear map from  $\mathbf{H}^1(\Omega)$  to  $\mathbf{H}^{-1}(\Omega)$  in the following way:  $\langle -\Delta u, v \rangle = (\nabla u, \nabla v)$  for all  $v \in \mathbf{H}_0^1(\Omega)$ .

It can be shown easily that  $\mathcal{F}$  is Fréchet differentiable with derivative  $\mathcal{F}'(u, \phi, \mu) : \mathcal{Y}_0 \rightarrow \mathcal{Y}_0^*$ ,

$$\mathcal{F}'(u, \phi, \mu) \cdot (v, \psi, \nu) = \begin{pmatrix} -\lambda\Delta v + (F''(u) + 2\phi + \mu)v + 2u\psi + \nu u \\ -\frac{1}{4\pi}\Delta\psi - 2u v \\ (u, v) \end{pmatrix} \quad \forall (v, \psi, \nu) \in \mathcal{Y}_0,$$

and that  $\mathcal{F}'$  is locally Hölder continuous: there is a continuous function  $L_{\mathcal{F}'} : \mathbb{R} \rightarrow \mathbb{R}$  and  $\alpha \in (0, 1)$  such that

$$\|\mathcal{F}'(y_1) - \mathcal{F}'(y_2)\| \leq L_{\mathcal{F}'}(\|y_1\|_{\mathcal{Y}} + \|y_2\|_{\mathcal{Y}}) (\|y_1 - y_2\|_{\mathcal{Y}}^\alpha + \|y_1 - y_2\|_{\mathcal{Y}}) \quad (2.17)$$

for all  $y_1, y_2 \in \mathcal{Y}$ . This follows immediately from (2.10). A direct consequence of the differentiability is that  $\mathcal{F}$  is locally Lipschitz continuous: there is a continuous function  $L_{\mathcal{F}} : \mathbb{R} \rightarrow \mathbb{R}$  such that

$$\|\mathcal{F}(y_1) - \mathcal{F}(y_2)\|_{\mathcal{Y}_0^*} \leq L_{\mathcal{F}}(\|y_1\|_{\mathcal{Y}} + \|y_2\|_{\mathcal{Y}}) \|y_1 - y_2\|_{\mathcal{Y}}, \quad (2.18)$$

for all  $y_1, y_2 \in \mathcal{Y}$ .

At this point, we make some observations concerning the regularity of solutions to (2.15). The functions  $\bar{u}$  and  $\bar{\phi}$  solve equations of the type

$$-\lambda \Delta \bar{u} = g_1, \quad \text{and} \quad -\frac{1}{4\pi} \Delta \bar{\phi} = g_2,$$

subject to the boundary conditions  $\bar{u}|_{\partial\Omega} = u_{\text{ex}}$ ,  $\bar{\phi}|_{\partial\Omega} = \phi_{\text{ex}}$ , respectively. Here, the functions  $g_1$  and  $g_2$  obviously depend on  $\bar{u}$ ,  $\bar{\phi}$  and  $\bar{\mu}$ . From embedding theorems and the growth properties of  $F'$  we know that  $g_1, g_2 \in L^2(\Omega)$ . Since we assumed that the Poisson problem (2.8) is  $H^2$ -regular, we can deduce that  $\bar{u}, \bar{\phi} \in H^2(\Omega)$ . The Sobolev embedding theorem<sup>1</sup> states that  $H^2(\Omega) \subset C^{0,\gamma}(\bar{\Omega})$  for all  $0 < \gamma \leq 1/2$ . Hence,  $\bar{u}, \bar{\phi} \in C^{0,\gamma}(\bar{\Omega})$  for every  $\gamma \leq 1/2$ , and in particular  $\bar{u}, \bar{\phi} \in L^\infty(\Omega)$ .

#### 2.4. Second-order optimality conditions

The function  $\mathcal{L}(u, \mu) = E(u) + \mu c(u)$  is called a *Lagrangian*. Since both  $E$  and  $c$  are twice continuously Fréchet differentiable, the same holds for  $\mathcal{L}$ .

From Theorem 3.3 in Maurer & Zowe<sup>27</sup> we deduce that, if  $\bar{u}$  is a solution of (2.7), and  $\bar{\mu}$  is its associated Lagrange multiplier, then the necessary second-order optimality condition

$$\langle \nabla_{uu} \mathcal{L}(\bar{u}, \bar{\mu}) v, v \rangle \geq 0 \quad \forall v \in \ker c'(\bar{u})$$

holds. Conversely, if  $(\bar{u}, \bar{\mu})$  satisfies (2.14) as well as the sufficient second-order optimality condition

$$\langle \nabla_{uu} \mathcal{L}(\bar{u}, \bar{\mu}) v, v \rangle \geq \gamma \|\nabla v\|_{L^2}^2 \quad \forall v \in \ker c'(\bar{u}), \quad (2.19)$$

for some constant  $\gamma > 0$ , then  $\bar{u}$  is an isolated local minimizer of  $E$  in  $A_u$ , see Maurer & Zowe.<sup>27</sup> We call a critical point  $\bar{u} \in A_u$  that satisfies (2.19) a *uniform minimizer* of (2.7).

Written out explicitly, (2.19) reads

$$\lambda(\nabla v, \nabla v) + ((F''(\bar{u}) + 2\bar{\phi} + \bar{\mu})v, v) + 16\pi(\bar{u}v, (-\Delta_0)^{-1}\bar{u}v) \geq \gamma \|\nabla v\|_{L^2}^2, \quad (2.20)$$

for all  $v \in \ker c'(\bar{u})$ , where  $\bar{\phi}$  is the electrostatic potential associated with  $\bar{u}$ .

The next step is to prove that, if  $\bar{u}$  is a uniform local minimizer with associated electrostatic potential  $\bar{\phi}$  and Lagrange multiplier  $\bar{\mu}$ , then  $\mathcal{F}'(\bar{u}, \bar{\phi}, \bar{\mu})$  is an

12 *Bernhard Langwallner, Christoph Ortner, and Endre Süli*

isomorphism. This will be an important tool for the convergence analysis in the next section.

**Proposition 2.1.** *Let  $\bar{y} = (\bar{u}, \bar{\phi}, \bar{\mu}) \in \mathcal{Y}_D$  such that (2.20) holds with  $\gamma > 0$ . Then,  $\mathcal{F}'(\bar{y}) : \mathcal{Y}_0 \rightarrow \mathcal{Y}_0^*$  is an isomorphism.*

**Proof.** We need to show that the equation

$$\mathcal{F}'(\bar{u}, \bar{\phi}, \bar{\mu}) \cdot (v, \psi, \nu) = (f, g, \kappa) \quad (2.21)$$

is uniquely solvable in  $\mathcal{Y}_0$  for every  $(f, g, \kappa) \in \mathcal{Y}_0^*$ . To this end we define two bilinear forms  $a_{\bar{y}} : \mathbf{H}_0^1(\Omega)^2 \times \mathbf{H}_0^1(\Omega)^2 \rightarrow \mathbb{R}$  and  $b_{\bar{y}} : \mathbf{H}_0^1(\Omega)^2 \times \mathbb{R} \rightarrow \mathbb{R}$ ,

$$\begin{aligned} a_{\bar{y}}((v, \psi), (w, \chi)) &= \lambda(\nabla v, \nabla w) + ((F''(\bar{u}) + 2\bar{\phi} + \bar{\mu})v, w) + (2\bar{u}\psi, w) \\ &\quad + \frac{1}{4\pi}(\nabla\psi, \nabla\chi) - 2(\bar{u}v, \chi), \\ b_{\bar{y}}((v, \psi), \eta) &= (\bar{u}, v)\eta. \end{aligned}$$

Then, equation (2.21) takes the form of a saddle-point problem,

$$\begin{aligned} a_{\bar{y}}((v, \psi), (w, \chi)) + b_{\bar{y}}((w, \chi), \nu) &= \langle f, w \rangle + \langle g, \chi \rangle \quad \forall w, \chi \in \mathbf{H}_0^1(\Omega), \\ b_{\bar{y}}((v, \psi), \eta) &= \eta\kappa \quad \forall \eta \in \mathbb{R}. \end{aligned} \quad (2.22)$$

The bilinear forms  $a_{\bar{y}}$  and  $b_{\bar{y}}$  are continuous on  $\mathbf{H}_0^1(\Omega)^2 \times \mathbf{H}_0^1(\Omega)^2$  and  $\mathbf{H}_0^1(\Omega)^2 \times \mathbb{R}$ , respectively. We define

$$\begin{aligned} \ker b_{\bar{y}} &:= \{(v, \psi) \in \mathbf{H}_0^1(\Omega)^2 : b_{\bar{y}}((v, \psi), \eta) = 0 \quad \forall \eta \in \mathbb{R}\} \\ &= \{(v, \psi) \in \mathbf{H}_0^1(\Omega)^2 : v \in \ker c'(\bar{u})\}. \end{aligned}$$

For a saddle-point problem such as (2.22) there are well-known sufficient conditions for solvability; see Theorem 1.1 in Brezzi & Fortin.<sup>4</sup> The bilinear form  $b_{\bar{y}}$  has to satisfy an inf-sup condition of the form

$$\inf_{\nu \in \mathbb{R}} \sup_{(v, \psi) \in \mathbf{H}_0^1(\Omega)^2} \frac{b_{\bar{y}}((v, \psi), \nu)}{|\nu| \|\nabla v, \nabla \psi\|_{L^2}} \geq \kappa_b > 0,$$

and the linear operator associated with  $a_{\bar{y}}$  has to be invertible on  $\ker b_{\bar{y}}$ .

*Step 1. Inf-sup condition for  $b_{\bar{y}}$ .* Since  $\|\bar{u}\|_{L^2} = N$ , we have  $\bar{u} \neq 0$ , and therefore  $b_{\bar{y}}$  obeys an inf-sup condition on  $\mathbf{H}_0^1(\Omega)^2 \times \mathbb{R}$ :

$$\begin{aligned} \inf_{\nu \in \mathbb{R}} \sup_{(v, \psi) \in \mathbf{H}_0^1(\Omega)^2} \frac{b_{\bar{y}}((v, \psi), \nu)}{|\nu| \|\nabla v, \nabla \psi\|_{L^2}} &= \inf_{\nu \in \mathbb{R}} \sup_{(v, \psi) \in \mathbf{H}_0^1(\Omega)^2} \frac{\nu \int_{\Omega} \bar{u}v \, dx}{|\nu| \|\nabla v, \nabla \psi\|_{L^2}} \\ &\geq \frac{(\bar{u}, (-\Delta_0)^{-1}\bar{u})}{(\bar{u}, (-\Delta_0)^{-1}\bar{u})^{1/2}} =: \kappa_b > 0, \end{aligned}$$

where for a given  $\nu \neq 0$  we have chosen  $v = \text{sign}(\nu)(-\Delta_0)^{-1}\bar{u}$ , and  $\psi = 0$ .

*Step 2. Invertibility of  $a_{\bar{y}}$  on  $\ker b_{\bar{y}}$ .* The proof of the unique solvability of the variational principle: find  $(v, \psi) \in \ker b_{\bar{y}}$  such that

$$a_{\bar{y}}((v, \psi), (w, \chi)) = \langle f, w \rangle + \langle g, \chi \rangle \quad \forall (w, \chi) \in \ker b_{\bar{y}}, \quad (2.23)$$

where  $f, g \in H^{-1}(\Omega)$ , requires more work and really relies on the assumption that  $\bar{u}$  is a uniform minimizer.

First we show that solutions are unique. For  $f = g = 0$  we want to prove that the only possible solution is  $(v, \psi) = (0, 0)$ . Looking at the definition of  $a_{\bar{y}}$  we see that  $g = 0$  leads to  $\psi = 8\pi(-\Delta_0)^{-1}\bar{u}v$ . Substituting this into (2.23) and testing with  $w = v$  and  $\chi = 0$  we obtain

$$\lambda(\nabla v, \nabla v) + ((F''(\bar{u}) + 2\bar{\phi} + \bar{\mu})v, v) + 16\pi(\bar{u}v, (-\Delta_0)^{-1}\bar{u}v) = 0. \quad (2.24)$$

Since  $\bar{u}$  is assumed to be a uniform minimizer we know from (2.20) that the bilinear form on the left-hand side of (2.24) is coercive on  $\{w \in H_0^1(\Omega) : (w, \bar{u}) = 0\}$ , so we get  $v = 0$  and hence also  $\psi = 0$ .

Next we prove that for every  $f, g \in H^{-1}(\Omega)$  there exists a solution. Let  $v \in \ker c'(\bar{u})$  be the unique solution of

$$\begin{aligned} \lambda(\nabla v, \nabla w) + ((F''(\bar{u}) + 2\bar{\phi} + \bar{\mu})v, w) + 16\pi(\bar{u}w, (-\Delta_0)^{-1}\bar{u}v) \\ = \langle f, w \rangle - (4\pi\bar{u}(-\Delta_0)^{-1}g, w) \quad \forall w \in \ker c'(\bar{u}). \end{aligned}$$

This equation is uniquely solvable by the Lax–Milgram theorem,<sup>3</sup> again since the bilinear form (as a function of  $v$  and  $w$ ) on the left-hand side is coercive. Let

$$\psi = 4\pi(-\Delta_0)^{-1}(2\bar{u}v + g).$$

Combining the last two equations shows that  $v, \psi$  indeed solve equation (2.23). Thus, we have shown unique solvability in  $\ker b_{\bar{y}}$  for every  $f, g \in H^{-1}(\Omega)$ .

From Theorem 1.1 in Brezzi & Fortin<sup>4</sup> we can now conclude that  $\mathcal{F}'(\bar{y})$  is indeed an isomorphism from  $\mathcal{Y}_0$  to  $\mathcal{Y}_0^*$ .  $\square$

For future reference, we mention that the invertibility of  $a_{\bar{y}}$  on  $\ker b_{\bar{y}}$  shown in the proof is equivalent to the existence of a constant  $\kappa_a > 0$  such that

$$\begin{aligned} \inf_{(v, \psi) \in \ker b_{\bar{y}}} \sup_{(w, \chi) \in \ker b_{\bar{y}}} \frac{a_{\bar{y}}((v, \psi), (w, \chi))}{\|(\nabla v, \nabla \psi)\|_{L^2} \|(\nabla w, \nabla \chi)\|_{L^2}} \geq \kappa_a, \quad \text{and} \\ \inf_{(w, \chi) \in \ker b_{\bar{y}}} \sup_{(v, \psi) \in \ker b_{\bar{y}}} \frac{a_{\bar{y}}((v, \psi), (w, \chi))}{\|(\nabla v, \nabla \psi)\|_{L^2} \|(\nabla w, \nabla \chi)\|_{L^2}} \geq \kappa_a. \end{aligned} \quad (2.25)$$

For this equivalence see for example Proposition 1.2 in Brezzi & Fortin,<sup>4</sup> respectively the discussion following Remark 1.6 in Brezzi & Fortin.<sup>4</sup>

### 3. Galerkin Discretization

In this section, we propose a discretization of the minimization problem (2.7), which corresponds to a Galerkin discretization of the optimality system (2.15). We will show that, for sufficiently small values of the discretization parameter, the discretized problem has a solution and that as the discretization parameter tends to zero a sequence of numerical solutions converges to the continuous solution. Optimal convergence rates will be addressed in future work.

14 *Bernhard Langwallner, Christoph Ortner, and Endre Süli*

### 3.1. The Discretized Functional

Let  $(S_h)_{h \in (0,1]}$  be a family of finite-dimensional subspaces of  $H^1(\Omega)$  with the approximation property

$$\inf_{v \in S_h} \|\nabla(u - v)\|_{L^2} \leq Ch|u|_{H^2} \quad \forall u \in H^2(\Omega),$$

and define  $S_{h,0} = S_h \cap H_0^1(\Omega)$ . Let  $\mathcal{I}_h : H^2(\Omega) \rightarrow S_h$  be an approximation operator with

$$\|\phi - \mathcal{I}_h \phi\|_{H^1} \leq Ch|\phi|_{H^2} \quad \text{for all } \phi \in H^2(\Omega). \quad (3.1)$$

Moreover, let

$$u_{\text{ex},h} = \mathcal{I}_h u_{\text{ex}}|_{\bar{\Omega}}, \quad \text{and} \quad \phi_{\text{ex},h} = \mathcal{I}_h \phi_{\text{ex}}|_{\bar{\Omega}}.$$

We introduce an approximation of the energy functional (2.4) defined on  $S_h$  of the following form

$$E_h(u_h) = \frac{\lambda}{2} \int_{\Omega} |\nabla u_h|^2 dx + \int_{\Omega} F(u_h) dx + \Phi_h(u_h), \quad (3.2)$$

where

$$\Phi_h(u_h) = - \inf_{\phi_h \in \phi_{\text{ex},h} + S_{h,0}} \left[ \int_{\Omega} \frac{1}{8\pi} |\nabla \phi_h|^2 - \phi_h(u_h^2 - \rho_n) dx \right].$$

Let  $A_{u,h} := \{u_h \in u_{\text{ex},h} + S_{h,0} : \|u_h\|_{L^2}^2 = N\}$  be the set of discrete admissible functions. We consider the discretized minimization problem

$$\min_{u_h \in A_{u,h}} E_h(u_h). \quad (3.3)$$

As in the continuous case, we get the following optimality conditions: if  $\bar{u}_h$  is a (local) minimizer of  $E_h$  in  $A_{u,h}$ , then there exists a discrete electrostatic potential  $\bar{\phi}_h \in S_h$  and a Lagrange multiplier  $\bar{\mu}_h \in \mathbb{R}$  such that

$$\begin{aligned} \lambda(\nabla \bar{u}_h, \nabla v) + (F'(\bar{u}_h), v) + 2(\bar{\phi}_h \bar{u}_h, v) + \mu_h(\bar{u}_h, v) &= 0 \quad \forall v \in S_{h,0}, \\ \frac{1}{4\pi}(\nabla \bar{\phi}_h, \nabla \psi) - (\bar{u}_h^2 - \rho_n, \psi) &= 0 \quad \forall \psi \in S_{h,0}, \\ \frac{\nu}{2} \left( \int_{\Omega} \bar{u}_h^2 dx - N \right) &= 0 \quad \forall \nu \in \mathbb{R}, \end{aligned} \quad (3.4)$$

where  $\bar{u}_h$  and  $\bar{\phi}_h$  satisfy the boundary conditions

$$\bar{u}_h|_{\partial\Omega} = u_{\text{ex},h}, \quad \bar{\phi}_h|_{\partial\Omega} = \phi_{\text{ex},h}.$$

These discrete optimality conditions turn out to be the Galerkin discretization of the optimality system (2.15). We introduce the discrete function spaces

$$\begin{aligned} \mathcal{Y}_h &= S_h \times S_h \times \mathbb{R}, \quad \mathcal{Y}_{h,D} = (u_{\text{ex},h} + S_{h,0}) \times (\phi_{\text{ex},h} + S_{h,0}) \times \mathbb{R}, \\ \mathcal{Y}_{h,0} &= S_{h,0} \times S_{h,0} \times \mathbb{R}, \quad \mathcal{Y}_{h,0}^* = S_{h,0}^* \times S_{h,0}^* \times \mathbb{R}. \end{aligned}$$

In analogy to the continuous case we write the system (3.4) in the more compact form

$$\langle \mathcal{F}_h(\bar{u}_h, \bar{\phi}_h, \bar{\mu}_h), (v, \psi, \nu) \rangle = 0 \quad \forall (v, \psi, \nu) \in \mathcal{Y}_{h,0}, \quad (3.5)$$

where  $\mathcal{F}_h : \mathcal{Y}_{h,D} \rightarrow \mathcal{Y}_{h,0}^*$ . The operator  $\mathcal{F}_h$  will sometimes be split into two parts,

$$\mathcal{F}_h(u_h, \phi_h, \mu_h) = \begin{pmatrix} -\lambda \Delta_h u_h \\ -\frac{1}{4\pi} \Delta_h \phi_h \\ 0 \end{pmatrix} + \mathcal{G}(u_h, \phi_h, \mu_h),$$

where  $\mathcal{G}$  has the same form as in the continuous case but it is now restricted to  $\mathcal{Y}_h$ , and where the discrete Laplacian  $(-\Delta_h) : S_h \rightarrow S_{h,0}^*$  is defined by  $\langle -\Delta_h v_h, w_h \rangle = (\nabla v_h, \nabla w_h)$  for  $v_h \in S_h$  and all  $w_h \in S_{h,0}$ . The operator  $(-\Delta_{h,0})^{-1} : S_{h,0}^* \rightarrow S_{h,0}$  maps  $f$  to the solution  $\phi_h \in S_{h,0}$  of  $(\nabla \phi_h, \nabla v_h) = \langle f, v_h \rangle$  for all  $v_h \in S_{h,0}$ .

Differentiability of  $\mathcal{F}_h$  is easily shown. The derivative  $\mathcal{F}'_h$  is again Hölder continuous and takes the form

$$\mathcal{F}'_h(u_h, \phi_h, \mu_h) \cdot (v, \psi, \nu) = \begin{pmatrix} -\lambda \Delta_h v + (F''(u_h) + 2\phi_h + \mu_h)v + 2u_h \psi + \nu u_h \\ -\frac{1}{4\pi} \Delta_h \psi - 2u_h v \\ (u_h, v) \end{pmatrix},$$

for  $(v, \psi, \nu) \in \mathcal{Y}_h$ . Just as in the continuous case, this linear operator has saddle-point structure.

At this point we note that  $\mathcal{F}'_h$  may in fact be extended to the whole of  $\mathcal{Y}$ . Slightly abusing notation, we will write  $\mathcal{F}'_h(y) : \mathcal{Y}_{h,0} \rightarrow \mathcal{Y}_{h,0}^*$  for any  $y \in \mathcal{Y}$ , but we stress that this is still an operator between the discrete function spaces. The next result is the discrete counterpart of Proposition 2.1.

**Proposition 3.1.** *Let  $\bar{y} = (\bar{u}, \bar{\phi}, \bar{\mu})$  be a solution of (2.15) such that  $\bar{u}$  is a uniform minimizer of (2.7) that satisfies (2.19). Then, there exist  $h_0 \in (0, 1]$  and  $\delta > 0$  such that  $\mathcal{F}'_h(y) : \mathcal{Y}_{h,0} \rightarrow \mathcal{Y}_{h,0}^*$  is an isomorphism for every  $h \leq h_0$  and for any  $y \in B_\delta(\bar{y}) \subset \mathcal{Y}$ . Moreover, there is a constant  $M > 0$  such that*

$$\|\mathcal{F}'_h(y)^{-1}\| \leq M \quad \forall y \in B_\delta(\bar{y}) \quad \forall h \leq h_0. \quad (3.6)$$

**Proof.** We begin by showing invertibility of  $\mathcal{F}'_h(\bar{y}) : \mathcal{Y}_{h,0} \rightarrow \mathcal{Y}_{h,0}^*$ . To this end, we again interpret the problem in saddle-point form and prove an inf-sup inequality for  $b_{\bar{y}}$  and the invertibility of  $a_{\bar{y}}$  on  $S_{h,0}^2 \cap \ker b_{\bar{y}}$ .

*Step 1. Inf-sup condition for  $b_{\bar{y}}$ .* We have

$$\begin{aligned} \inf_{\nu \in \mathbb{R}} \sup_{(v, \psi) \in S_{h,0}^2} \frac{b_{\bar{y}}((v, \psi), \nu)}{|\nu| \|(\nabla v, \nabla \psi)\|_{L^2}} &= \inf_{\nu \in \mathbb{R}} \sup_{(v, \psi) \in S_{h,0}^2} \frac{\nu \int_{\Omega} \bar{u} v \, dx}{|\nu| \|(\nabla v, \nabla \psi)\|_{L^2}} \\ &\geq \frac{(\bar{u}, (-\Delta_{h,0})^{-1} \bar{u})}{\|\nabla(-\Delta_{h,0})^{-1} \bar{u}\|_{L^2}} =: \kappa_{b,h}, \end{aligned}$$



16 *Bernhard Langwallner, Christoph Ortner, and Endre Süli*

where, for given  $\nu \neq 0$ , we have chosen  $v = \text{sign}(\nu)(-\Delta_{h,0})^{-1}\bar{u}$  and  $\psi = 0$ . If  $h$  is sufficiently small then

$$\kappa_{b,h} = \frac{(\bar{u}, (-\Delta_{h,0})^{-1}\bar{u})}{\|\nabla(-\Delta_{h,0})^{-1}\bar{u}\|_{L^2}} \geq \frac{1}{2} (\bar{u}, (-\Delta_0)^{-1}\bar{u})^{1/2} = \kappa_b/2 > 0.$$

In particular, we deduce that the inf-sup constant  $\kappa_{b,h}$  is bounded away from zero if  $h$  is sufficiently small.

*Step 2.* Considering now  $a_{\bar{y}}$ , we have to prove that the system

$$a_{\bar{y}}((v, \psi), (w, \chi)) = \langle f, w \rangle + \langle g, \chi \rangle \quad \forall (w, \chi) \in \ker b_{\bar{y}} \cap S_{h,0}^2 \quad (3.7)$$

has a unique solution  $(v, \psi) \in \ker b_{\bar{y}} \cap S_{h,0}^2$  for every  $f, g \in H^{-1}(\Omega)$ . If the trial and test spaces were simply  $S_{h,0}^2$  instead of the constraint space  $\ker b_{\bar{y}} \cap S_{h,0}^2$ , then this would follow from a classical argument by Schatz.<sup>32</sup> The present case requires some modifications, which we study in detail in the Appendix. Lemmas 5.2 and 5.1 provide the regularity and approximation results in  $\ker b_{\bar{y}} \cap S_{h,0}^2$  necessary for an application of the Schatz argument, which is carried out in Theorem 5.1. Hence, we deduce that (3.7) is uniquely solvable, provided that  $h$  is small enough. Theorem 5.1 also implies the existence of an inf-sup constant  $\kappa_{a,h}$  for  $a_{\bar{y}}$  on  $K_h := S_{h,0}^2 \cap \ker b_{\bar{y}}$ , similarly as in (2.25), in the continuous case:

$$\inf_{(v,\psi) \in K_h} \sup_{(w,\chi) \in K_h} \frac{a_{\bar{y}}((v, \psi), (w, \chi))}{\|(\nabla v, \nabla \psi)\|_{L^2} \|(\nabla w, \nabla \chi)\|_{L^2}} \geq \kappa_{a,h} > 0, \quad \text{and}$$

$$\inf_{(w,\zeta) \in K_h} \sup_{(v,\chi) \in K_h} \frac{a_{\bar{y}}((v, \psi), (w, \chi))}{\|(\nabla v, \nabla \psi)\|_{L^2} \|(\nabla w, \nabla \chi)\|_{L^2}} \geq \kappa_{a,h} > 0.$$

Furthermore, Theorem 5.1 guarantees that  $\kappa_{a,h}$  is bounded away from zero as  $h \rightarrow 0$ .

*Step 3.* We have shown that for sufficiently small  $h$ ,  $h \leq h_0$  say,  $\mathcal{F}'_h(\bar{y})$  is an isomorphism. This also means that for every  $h \leq h_0$ ,  $\mathcal{F}'_h(\bar{y})$  satisfies the inf-sup conditions

$$\inf_{y_h \in \mathcal{Y}_h} \sup_{z_h \in \mathcal{Y}_h} \frac{\langle \mathcal{F}'_h(\bar{y})y_h, z_h \rangle}{\|y_h\|_{\mathcal{Y}} \|z_h\|_{\mathcal{Y}}} \geq \kappa_h \quad \text{and} \quad \inf_{z_h \in \mathcal{Y}_h} \sup_{y_h \in \mathcal{Y}_h} \frac{\langle \mathcal{F}'_h(\bar{y})y_h, z_h \rangle}{\|y_h\|_{\mathcal{Y}} \|z_h\|_{\mathcal{Y}}} \geq \kappa_h,$$

with  $\kappa_h > 0$ . Theorem 1.1 in Brezzi & Fortin<sup>4</sup> shows a way of bounding the inf-sup constant  $\kappa_h$  in terms of the inf-sup constants  $\kappa_{a,h}$  for  $a_{\bar{y}}$  and  $\kappa_{b,h}$  for  $b_{\bar{y}}$ . Since the latter are uniformly bounded away from zero as  $h \rightarrow 0$ , we deduce that the same holds for  $\kappa_h$ .

Since  $\mathcal{F}'_h$  satisfies the discrete equivalent of the Hölder condition (2.17), it follows that there exists a neighbourhood  $B_\delta(\bar{y}) \subset \mathcal{Y}$ , where  $\delta > 0$  is chosen sufficiently small, such that, for  $h \leq h_0$  and  $y \in B_\delta(\bar{y})$ ,  $\mathcal{F}'_h(y)$  is an isomorphism and such that the inf-sup constants for  $\mathcal{F}'_h(y)^{-1}$  are uniformly bounded away from zero. This implies (3.6).  $\square$

### 3.2. Existence and Convergence

The following convergence theorem constitutes the main result of this section. The proof uses ideas commonly used in the finite element literature on nonlinear problems; see for example Brezzi, Rappaz & Raviart<sup>5</sup> and Dobrowolski & Rannacher.<sup>11</sup>

**Theorem 3.1.** *Let  $\bar{u}$  be a minimizer of (2.7) that satisfies (2.19). Let  $\bar{\phi}$  and  $\bar{\mu}$  be, respectively, the associated electrostatic potential and Lagrange multiplier. Then, there exist  $h_0 \in (0, 1]$ ,  $\delta > 0$  such that the discretized problem (3.4) has a unique solution  $\bar{y}_h = (\bar{u}_h, \bar{\phi}_h, \bar{\mu}_h) \in \mathcal{Y}_{h,D}$  in the neighbourhood  $B_\delta(\bar{y}) \subset \mathcal{Y}$  for all  $h < h_0$ . Furthermore, there exists a constant  $C$  such that*

$$\|\bar{u} - \bar{u}_h\|_{H^1} + \|\bar{\phi} - \bar{\phi}_h\|_{H^1} + |\bar{\mu} - \bar{\mu}_h| \leq Ch.$$

**Proof.** The proof is divided into four steps.

*Step 1.* We show that for an approximation  $\Pi_h \bar{y} \in \mathcal{Y}_{h,D}$  of  $\bar{y} = (\bar{u}, \bar{\phi}, \bar{\mu})$ , we have  $\|\mathcal{F}_h(\Pi_h \bar{y})\|_{\mathcal{Y}_{h,0}^*} \leq C_1 h$  for sufficiently small  $h$  where  $C_1$  is independent of  $h$ .

Let  $u_h, \phi_h \in S_h$  be the Ritz projections of  $\bar{u}$  and  $\bar{\phi}$ , respectively, i.e., the solutions of the equations

$$(\nabla u_h, \nabla v) = (\nabla \bar{u}, \nabla v) \quad \forall v \in S_{h,0} \quad \text{and} \quad (\nabla \phi_h, \nabla v) = (\nabla \bar{\phi}, \nabla v) \quad \forall v \in S_{h,0},$$

with boundary conditions

$$u_h|_{\partial\Omega} = u_{\text{ex},h}, \quad \phi_h|_{\partial\Omega} = \phi_{\text{ex},h}.$$

In other words,  $\Delta_h u_h = \Delta \bar{u}|_{S_{h,0}}$  and  $\Delta_h \phi_h = \Delta \bar{\phi}|_{S_{h,0}}$ . We define  $\Pi_h \bar{y} = (u_h, \phi_h, \bar{\mu})$ . Convergence theory for the Poisson equation,  $\bar{u}, \bar{\phi} \in H^2(\Omega)$ , and the approximation property (3.1) of  $\mathcal{I}_h$  then lead to

$$\|\bar{y} - \Pi_h \bar{y}\|_{\mathcal{Y}} \leq Ch.$$

Using the fact that  $\mathcal{F}(\bar{y})|_{\mathcal{Y}_{h,0}} = 0$  and  $\mathcal{F}_h(\Pi_h \bar{y}) = \mathcal{F}(\Pi_h \bar{y})|_{\mathcal{Y}_{h,0}}$  we proceed as follows:

$$\|\mathcal{F}_h(\Pi_h \bar{y})\|_{\mathcal{Y}_{h,0}^*} = \|\mathcal{F}(\Pi_h \bar{y})|_{\mathcal{Y}_{h,0}} - \mathcal{F}(\bar{y})|_{\mathcal{Y}_{h,0}}\|_{\mathcal{Y}_{h,0}^*} \leq \|\mathcal{F}(\Pi_h \bar{y}) - \mathcal{F}(\bar{y})\|_{\mathcal{Y}_{h,0}^*}.$$

From the local Lipschitz continuity (2.18) of  $\mathcal{F}$  we deduce that

$$\|\mathcal{F}_h(\Pi_h \bar{y})\|_{\mathcal{Y}_{h,0}^*} \leq C \|\bar{y} - \Pi_h \bar{y}\|_{\mathcal{Y}} \leq C_1 h.$$

*Step 2.* In Proposition 3.1 we have shown that there is an open neighbourhood  $B_\delta(\bar{y})$  of  $\bar{y}$  in  $\mathcal{Y}$  and  $h_0 \in (0, 1]$  such that  $\mathcal{F}'_h(y) : \mathcal{Y}_{h,0} \rightarrow \mathcal{Y}_{h,0}^*$  is an isomorphism for all  $y \in B_\delta(\bar{y})$ ,  $h \leq h_0$  and  $\|\mathcal{F}'_h(y)^{-1}\| \leq M$ , uniformly for  $y \in B_\delta(\bar{y})$ . Moreover, we observe that  $\mathcal{F}'_h$  satisfies a Hölder continuity property similar to (2.17): there is  $L_{\bar{y},\delta}$  and  $\alpha \in (0, 1)$  such that

$$\|\mathcal{F}'_h(y_1) - \mathcal{F}'_h(y_2)\|_{\mathcal{Y}_{h,0}^*} \leq L_{\bar{y},\delta} (\|y_1 - y_2\|_{\mathcal{Y}}^\alpha + \|y_1 - y_2\|_{\mathcal{Y}}) \quad \forall y_1, y_2 \in B_\delta(\bar{y}).$$

*Step 3. Existence and uniqueness of a solution.* We want to show that there exists a locally unique solution of  $\mathcal{F}_h(y_h) = 0$ . The idea is to construct a contractive

18 *Bernhard Langwallner, Christoph Ortner, and Endre Süli*

mapping whose fixed point is the solution  $y_h$ . To this end, we rewrite this equation as

$$\mathcal{F}_h(y_h) - \mathcal{F}_h(y_0) = -\mathcal{F}_h(y_0),$$

and choose  $y_0 = \Pi_h \bar{y}$  so that the right-hand side is “small”. Linearization leads to

$$\mathcal{F}'_h(y_0)(y_h - y_0) = -\mathcal{F}_h(y_0) - \int_0^1 (\mathcal{F}'_h(y_0 + t(y_h - y_0)) - \mathcal{F}'_h(y_0)) dt (y_h - y_0).$$

We recall that  $\mathcal{F}'_h(y_0)$  is an isomorphism if  $h$  is sufficiently small. Let us assume in what follows that  $h$  is small enough such that  $\|\bar{y} - y_0\|_{\mathcal{Y}} \leq \delta/2$ . Then, for  $R < \delta/2$ , we define the map  $\mathcal{N} : B_R(y_0) \rightarrow \mathcal{Y}_{h,D}$  by

$$\mathcal{F}'_h(y_0)(\mathcal{N}(y) - y_0) = -\mathcal{F}_h(y_0) - \int_0^1 (\mathcal{F}'_h(y_0 + t(y - y_0)) - \mathcal{F}'_h(y_0)) dt (y - y_0).$$

We will show that  $\mathcal{N}$  is a contraction from  $B_R(y_0)$  into  $B_R(y_0)$  if  $R$  is chosen sufficiently small.

First, we prove that  $\mathcal{N}$  maps  $B_R(y_0)$  to  $B_R(y_0)$  for sufficiently small  $R$ . For each  $y \in B_R(y_0)$  we have, with  $\alpha \in (0, 1)$ , that

$$\begin{aligned} M^{-1} \|\mathcal{N}(y) - y_0\|_{\mathcal{Y}} &\leq \|\mathcal{F}_h(y_0)\|_{\mathcal{Y}_0^*} + R \int_0^1 \|\mathcal{F}'_h(y_0 + t(y - y_0)) - \mathcal{F}'_h(y_0)\| dt \\ &\leq C_2(h + RL_{\bar{y},\delta}(R + R^\alpha)), \end{aligned}$$

where we have used the stability property (3.6). To ensure that  $\mathcal{N}(y) \in B_R(y_0)$ , we need to bound  $C_2(h + RL_{\bar{y},\delta}(R + R^\alpha))$  by  $R/M$ . If  $R$  and  $h$  are sufficiently small, this obviously holds. It is also clear that  $R$  can be chosen independently of  $h$ .

Next, we show that  $\mathcal{N}$  is a contraction on  $B_R(y_0)$ . If  $\eta_1, \eta_2 \in B_R(\bar{y})$ , then

$$\begin{aligned} \mathcal{F}'_h(y_0)(\mathcal{N}(\eta_1) - \mathcal{N}(\eta_2)) &= \mathcal{F}_h(\eta_2) - \mathcal{F}_h(\eta_1) + \mathcal{F}'_h(y_0)(\eta_1 - \eta_2) \\ &= \int_0^1 [\mathcal{F}'_h(y_0) - \mathcal{F}'_h(\eta_1 + t(\eta_2 - \eta_1))] (\eta_1 - \eta_2) dt. \end{aligned}$$

Thus,  $\|\mathcal{N}(\eta_1) - \mathcal{N}(\eta_2)\|_{\mathcal{Y}}$  can be estimated as follows:

$$\begin{aligned} M^{-1} \|\mathcal{N}(\eta_1) - \mathcal{N}(\eta_2)\|_{\mathcal{Y}} &\leq \int_0^1 \|\mathcal{F}'_h(y_0) - \mathcal{F}'_h(\eta_1 + t(\eta_2 - \eta_1))\|_{\mathcal{Y}} dt \|\eta_1 - \eta_2\| \\ &\leq L(R + R^\alpha) \cdot \|\eta_1 - \eta_2\|_{\mathcal{Y}}. \end{aligned}$$

For sufficiently small  $R$  we obtain  $L(R + R^\alpha)M < 1$  and hence  $\mathcal{N}$  is a contraction on  $B_R(y_0)$ .

We can now use Banach’s Fixed Point Theorem<sup>36</sup> to obtain the existence and uniqueness of a fixed point  $\bar{y}_h$  of the map  $\mathcal{N} : B_R(y_0) \rightarrow B_R(y_0)$ . This fixed point  $\bar{y}_h$  is a solution of  $\mathcal{F}_h(y) = 0$ . For sufficiently small  $h$  this solution is in the neighborhood  $B_{2R}(\bar{y})$ :

$$\|\bar{y} - \bar{y}_h\|_{\mathcal{Y}} \leq \|\bar{y} - \Pi_h \bar{y}\|_{\mathcal{Y}} + \|\Pi_h \bar{y} - \bar{y}_h\|_{\mathcal{Y}} \leq Ch + R.$$

*Step 4.* Finally, convergence can be obtained by a minor modification of the above argument. If we let  $R = C_R h$  and  $C_R > MC_2$  we can repeat the previous steps and deduce  $\|\Pi_h \bar{y} - \bar{y}_h\|_{\mathcal{Y}} \leq C_R h$ . This shows

$$\|\bar{y} - \bar{y}_h\|_{\mathcal{Y}} \leq Ch + C_R h,$$

which concludes the proof.  $\square$

**Proposition 3.2.** *Under the same assumptions as in Theorem 3.1 and for sufficiently small  $h$ , the discrete solution  $\bar{u}_h \in A_{u,h}$  is a uniform minimizer of the discretized functional (3.2) over  $A_{u,h}$ .*

**Proof.** We define the two Lagrangians

$$\mathcal{L}(u, \mu) = E(u) + \mu c(u) \quad \text{and} \quad \mathcal{L}_h(u_h, \mu_h) = E_h(u) + \mu_h c(u_h).$$

Since  $\bar{u}$  is a uniform minimizer, we have  $\langle \nabla_{uu}^2 \mathcal{L}(\bar{u}, \bar{\mu}) v, v \rangle \geq \gamma \|\nabla v\|_{L^2}^2$  for all  $v \in \ker c'(\bar{u})$ .

Given  $v_h \in \ker c'(\bar{u}_h) \cap S_{h,0}$ , we have

$$\begin{aligned} \langle \nabla_{uu}^2 \mathcal{L}_h(\bar{u}_h, \bar{\mu}_h) v_h, v_h \rangle &= \langle \nabla_{uu}^2 \mathcal{L}(\bar{u}, \bar{\mu}) v, v \rangle + \langle (\nabla_{uu}^2 \mathcal{L}_h(\bar{u}_h, \bar{\mu}_h) - \nabla_{uu}^2 \mathcal{L}(\bar{u}, \bar{\mu})) v_h, v_h \rangle \\ &\quad + 2 \langle \nabla_{uu}^2 \mathcal{L}(\bar{u}, \bar{\mu}) v, (v_h - v) \rangle \\ &\quad + \langle \nabla_{uu}^2 \mathcal{L}(\bar{u}, \bar{\mu}) (v_h - v), (v_h - v) \rangle \end{aligned} \quad (3.8)$$

for arbitrary  $v \in \ker c'(\bar{u})$ .

Next, we prove that every  $v_h \in \ker c'(\bar{u}_h)$  can be approximated by  $v \in \ker c'(\bar{u})$  since  $\|\bar{u} - \bar{u}_h\|_{H^1} \leq C_1 h$ . Let  $\varphi = (-\Delta_0)^{-1} \bar{u} \in H_0^1(\Omega)$  and define

$$v = v_h - \frac{(\nabla \varphi, \nabla v_h)}{\|\nabla \varphi\|_{L^2}^2} \varphi.$$

It follows immediately that  $(v, \bar{u}) = 0$ , i.e.,  $v \in \ker c'(\bar{u})$ . A quick calculation using  $(\bar{u}_h, v_h) = 0$  leads to  $\|\nabla(v - v_h)\|_{L^2} \leq Ch \|\nabla v_h\|_{L^2}$ :

$$\begin{aligned} \|\nabla(v - v_h)\|_{L^2} &= \frac{|(\nabla \varphi, \nabla v_h)|}{\|\nabla \varphi\|_{L^2}} = \frac{|(\bar{u}, v_h)|}{\|\nabla \varphi\|_{L^2}} = \frac{|(\bar{u} - \bar{u}_h, v_h)|}{\|\nabla \varphi\|_{L^2}} \\ &\leq C \|\bar{u} - \bar{u}_h\|_{L^2} \frac{\|\nabla v_h\|_{L^2}}{\|\nabla \varphi\|_{L^2}} \leq Ch \|\nabla v_h\|_{L^2}, \end{aligned}$$

where  $\|\nabla \varphi\|_{L^2} > 0$  has been absorbed in the generic constant  $C$ . Here, we have used the Cauchy-Schwarz inequality and Poincaré's inequality  $\|v_h\|_{L^2} \leq C \|\nabla v_h\|_{L^2}$  for  $v_h \in S_{h,0}$ . Based on this result we can easily derive  $\|\nabla v\|_{L^2} \geq (1 - Ch) \|\nabla v_h\|_{L^2}$ .

With this choice of  $v$  we see that the first term on the right-hand side of (3.8) satisfies

$$\langle \nabla_{uu}^2 \mathcal{L}(\bar{u}, \bar{\mu}) v, v \rangle \geq \gamma(1 - Ch) \|\nabla v_h\|_{L^2}^2.$$

20 *Bernhard Langwallner, Christoph Ortner, and Endre Süli*

The third and fourth term on the right-hand side of (3.8) can be bounded by  $Ch \|\nabla v_h\|_{L^2}^2$ . For the remaining term, we get

$$\begin{aligned} \langle (\nabla_{uu}^2 \mathcal{L}_h(\bar{u}_h, \bar{\mu}_h) - \nabla_{uu}^2 \mathcal{L}(\bar{u}, \bar{\mu}))v_h, v_h \rangle &= \langle (X''(\bar{u}) - X''(\bar{u}_h))v_h, v_h \rangle \\ &\quad + \langle (\Phi''(\bar{u}) - \Phi''_h(\bar{u}_h))v_h, v_h \rangle + (\bar{\mu}_h - \bar{\mu}) \|v_h\|_{L^2}^2. \end{aligned}$$

The part involving the nonlinear local functional  $X$  can obviously be bounded by  $C(h + h^\alpha) \|\nabla v_h\|_{L^2}^2$ ; see (2.11). Using the expression for  $\Phi''$  presented in Lemma 2.1 and its discrete analogue, as well as the convergence of  $\bar{\phi}_h$  to  $\bar{\phi}$ , we see that also  $|\langle (\Phi''(\bar{u}) - \Phi''_h(\bar{u}_h))v_h, v_h \rangle| \leq Ch \|\nabla v_h\|_{L^2}^2$ . Finally, since  $|\bar{\mu}_h - \bar{\mu}| \leq Ch$  we get  $|\bar{\mu}_h - \bar{\mu}| \|v_h\|_{L^2}^2 \leq Ch \|\nabla v_h\|_{L^2}^2$  by Poincaré's inequality.

Summarising, we have established that

$$\langle \nabla_{uu}^2 \mathcal{L}_h(\bar{u}_h, \bar{\mu}_h)v_h, v_h \rangle \geq \gamma/2 \|\nabla v_h\|_{L^2}^2 \quad \forall v_h \in \ker c'(\bar{u}_h) \cap S_h,$$

for sufficiently small  $h$ . From Theorem 5.6 in Maurer & Zowe<sup>27</sup> we deduce that  $\bar{u}_h$  is a uniform minimizer of  $E_h$  subject to  $c(\bar{u}_h) = 0$ .  $\square$

#### 4. Numerical Example

We have implemented a two-dimensional version of the discretization described above using piecewise linear, respectively, cubic Lagrange finite elements. The two-dimensional functional is obtained by replacing  $8\pi$  with  $4\pi$  in the definition of  $\Phi$  (2.5) and the function  $F$  with  $F(u) = \pi u^4/2 - \frac{4\sqrt{2}}{3\pi^{5/2}}|u|^3$ , see Ghosh & Dhara.<sup>19</sup>

To solve the nonlinear system (3.4) we apply Newton's method. The justification for this choice emerges from the availability of good initial guesses for  $u$ . Physical insight suggests that the electron density is higher close to nuclei. In fact, close to nuclei we expect the electron density to be close to the case of an isolated atom. We therefore solve the spherically symmetric TFDW problem for a single atom first. The initial value for  $u$  is then the square root of the sum of spherically symmetric single atom electron densities centered around nuclei. An initial guess for  $\phi$  is obtained by solving  $-\Delta\phi = 4\pi(u^2 - \rho_n)$ . In all computations,  $\rho_n$  is the sum of Gauss functions with variance  $\sigma_0$  centered at nucleus positions. Given a sufficiently fine mesh, we observe that the Newton iteration enters the regime of fast local convergence immediately. We conclude that there is no need to apply globalization strategies.

Results of a computation involving a cluster of 39 atoms in two dimensions are documented in Figure 1. The configuration was obtained by slightly perturbing a hexagonal configuration. Parameters used in this computation were  $\Omega = (-30, 30)^2$ ,  $\lambda = 3.2$ ,  $\sigma_0 = 0.5$  and  $Z_i = 6$  for all  $i = 1, \dots, 39$ . The solution  $\bar{u}_h$  is shown on the left. Note that the solution is only shown on a subset of  $\Omega$  to highlight the interesting features. On the right-hand side relative errors of  $E$  and  $\mu$  are plotted against the number of degrees of freedom,  $N_{\text{DOF}}$ , for linear (P1), respectively, cubic (P3) finite elements on a series of uniform grids. We observe that both the energy  $E$  and the Lagrange multiplier  $\mu$  converge with order  $2p$  if  $p$  is the order of the

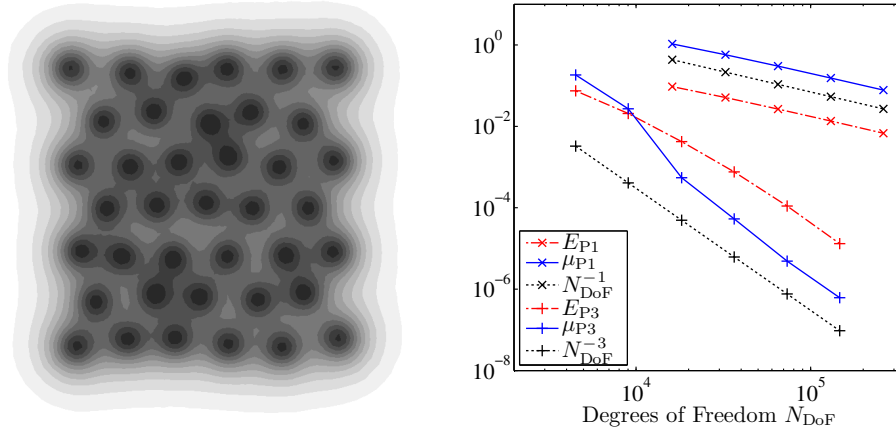


Fig. 1. Results of calculations for a cluster of 39 atoms in two dimensions:  $\bar{u}_h$  is shown on the left, relative errors of  $E$  and  $\mu$  are plotted on the right (along with lines indicating decay with  $N_{\text{DoF}}^{-1}$  and  $N_{\text{DoF}}^{-3}$ ). Since the results are in 2D, the discretization parameter satisfies  $h \sim N_{\text{DoF}}^{-1/2}$ . Hence, the plot suggests that the errors of  $E$  and  $\mu$  behave like  $\mathcal{O}(h^2)$  for P1 elements and  $\mathcal{O}(h^6)$  for P3 elements, respectively.

finite element space. This has been proven for plane wave and finite element discretizations of a slightly simpler class of nonlinear eigenvalue problems in Cancès, Chakir & Maday<sup>6</sup> (see Theorems 2 and 3). In Cancès, Chakir & Maday<sup>7</sup> (Theorem 3.1) the authors prove for a plane wave discretization of the Thomas–Fermi–von Weizsäcker functional that the convergence order of the Lagrange multiplier is the same as for the energy. A rigorous study of convergence rates for a finite element discretization with numerical integration will be provided in our forthcoming work Langwallner, Ortner & Süli.<sup>24</sup> The main difference between the work in Cancès, Chakir & Maday<sup>7</sup> and ours is that the plane wave discretization allows to treat the electrostatic interactions directly, whereas in the finite element discretization, we solve an enlarged system.

### Conclusions and Future Work

The Euler–Lagrange equations for the Thomas–Fermi–Dirac–von Weizsäcker functional can be rewritten as a non-monotone, semilinear elliptic system with a nonlinear constraint. In this paper, we have combined arguments used for linear saddle-point problems and linearization techniques based on the Inverse Function Theorem, to establish the existence and convergence of the sequence of solutions of a Galerkin discretization of this system. All results of the present work can be maintained if numerical integration of sufficiently high order is used. We will show this in a forthcoming article.<sup>24</sup>

In our current research we are developing a theory for optimal convergence rates

22 *Bernhard Langwallner, Christoph Ortner, and Endre Süli*

in situations where the Laplace operator has higher boundary regularity, and we are studying how the accuracy of the electronic density and electrostatic potential affects the accuracy of the nuclei when they are relaxed to their equilibrium positions. Another extension, which will require a careful study, is the effect of “cutting off” the electronic density and the electrostatic potential by imposing Dirichlet boundary values on  $\partial\Omega$ .

We view the present paper as the first in a series of papers, and as a preliminary step in our effort to develop a theory for coarse-graining the TFDW functional in the spirit of Garcia-Cervera, Lu & E<sup>16</sup> and Gavini et al.<sup>17</sup>

## 5. Appendix: An Indefinite Elliptic System with Constraint

In this appendix we generalize Schatz’ classical result on the Galerkin discretization of indefinite elliptic equations<sup>32</sup> to the constrained system, which we encountered in the analysis in Section 3.1.

We assume that  $\Omega \subset \mathbb{R}^d$ ,  $d \geq 1$ , is a bounded, open domain in which the Poisson problem is  $H^2(\Omega)$ -regular, that is, for any right-hand side  $f \in L^2(\Omega)$ , the solution  $u \in H_0^1(\Omega)$  of

$$(\nabla u, \nabla v) = (f, v) \quad \forall v \in H_0^1(\Omega)$$

belongs to  $H^2(\Omega) \cap H_0^1(\Omega)$  and

$$|u|_{H^2} \leq C_{\text{reg}} \|f\|_{L^2}. \quad (5.1)$$

With  $d = 3$  this is exactly the assumption we made in Section 2.2.

Let  $m \geq 1$  and let  $\mathbf{V} \subset H_0^1(\Omega; \mathbb{R}^m)$  be a subspace with co-dimension one. We assume that the linear constraint defining the subspace is given by a nonzero  $L^2$ -function  $\mathbf{g} \in L^2(\Omega; \mathbb{R}^m)$ , so that  $\mathbf{V}$  can be written as

$$\mathbf{V} = \{ \mathbf{v} \in H_0^1(\Omega; \mathbb{R}^m) : (\mathbf{v}, \mathbf{g}) = 0 \}.$$

Let  $a$  be a bilinear form on  $H_0^1(\Omega; \mathbb{R}^m) \times H_0^1(\Omega; \mathbb{R}^m)$  defined by

$$a(\mathbf{u}, \mathbf{v}) = (\nabla \mathbf{u}, \nabla \mathbf{v}) + (\mathbf{u}, \mathbf{M} \mathbf{v}) \quad \text{for } \mathbf{u}, \mathbf{v} \in H_0^1(\Omega; \mathbb{R}^m), \quad (5.2)$$

where  $\mathbf{M} \in L^\infty(\Omega; \mathbb{R}^{m \times m})$ .

We immediately see that  $a$  is bounded:

$$a(\mathbf{u}, \mathbf{v}) \leq C_a \|\nabla \mathbf{u}\|_{L^2} \|\nabla \mathbf{v}\|_{L^2} \quad \forall \mathbf{u}, \mathbf{v} \in \mathbf{V}, \quad (5.3)$$

and that there exists a constant  $K \geq 0$  such that  $a(\mathbf{u}, \mathbf{v}) + K(\mathbf{u}, \mathbf{v})$  is coercive,

$$a(\mathbf{u}, \mathbf{u}) + K(\mathbf{u}, \mathbf{u}) \geq \alpha \|\nabla \mathbf{u}\|_{L^2}^2 \quad \forall \mathbf{u}, \mathbf{v} \in \mathbf{V}. \quad (5.4)$$

We assume that for every  $\mathbf{f} \in L^2(\Omega; \mathbb{R}^m)$  there is a unique solution  $\mathbf{u} \in \mathbf{V}$  of the variational problem

$$a(\mathbf{u}, \mathbf{v}) = (\mathbf{f}, \mathbf{v}) \quad \forall \mathbf{v} \in \mathbf{V}. \quad (5.5)$$

The adjoint variational problem to (5.5) is given by

$$a(\mathbf{v}, \mathbf{u}) = (\mathbf{f}, \mathbf{v}) \quad \forall \mathbf{v} \in \mathbf{V}.$$

Note that the adjoint problem has the same form except that  $M$  is replaced by  $M^T$ . It follows by an argument involving the Fredholm alternative, and which carries over verbatim from the scalar case (see Theorem 6.2.4 in Evans<sup>12</sup>), that the adjoint problem also has a unique solution. This implies, in particular, the existence of  $\kappa > 0$  such that

$$\inf_{\mathbf{u} \in \mathbf{V}} \sup_{\mathbf{v} \in \mathbf{V}} \frac{a(\mathbf{u}, \mathbf{v})}{\|\mathbf{u}\|_{\mathbf{H}^1} \|\mathbf{v}\|_{\mathbf{H}^1}} \geq \kappa \quad \text{and} \quad \inf_{\mathbf{u} \in \mathbf{V}} \sup_{\mathbf{v} \in \mathbf{V}} \frac{a(\mathbf{v}, \mathbf{u})}{\|\mathbf{u}\|_{\mathbf{H}^1} \|\mathbf{v}\|_{\mathbf{H}^1}} \geq \kappa.$$

From that we can infer the solvability of (5.5) for all  $\mathbf{f} \in \mathbf{V}^*$ .

To define a Galerkin discretization of (5.5), let  $(S_{h,0})_{h \in (0,1]}$  be a family of finite-dimensional subspaces of  $\mathbf{H}_0^1(\Omega)$ , which satisfy the approximation property

$$\inf_{u_h \in S_{h,0}} \|\nabla(u - u_h)\|_{L^2} \leq C_{\text{apx}} h |u|_{\mathbf{H}^2} \quad \text{for all } u \in \mathbf{H}_0^1(\Omega) \cap \mathbf{H}^2(\Omega), \quad (5.6)$$

for every  $h$ , where  $C_{\text{apx}}$  is independent of  $h$ . We then define the approximation space

$$\mathbf{V}_h = \{\mathbf{v}_h \in S_{h,0}^m : (\mathbf{g}, \mathbf{v}_h) = 0 \text{ for } i = 1, \dots, m\}.$$

The Galerkin discretization of (5.5) is given by

$$a(\mathbf{u}_h, \mathbf{v}) = (\mathbf{f}, \mathbf{v}) \quad \forall \mathbf{v} \in \mathbf{V}_h. \quad (5.7)$$

Note that it may occur that  $\mathbf{V}_h = S_{h,0}^m$ . However, it follows immediately from (5.6) that, for sufficiently small  $h$ , the co-dimension of  $\mathbf{V}_h$  in  $S_{h,0}^m$  is also one.

Our main result of this appendix ensures solvability of the Galerkin discretization (5.7), provided that  $h$  is sufficiently small. The proof parallels Theorem 5.7.6 in Brenner & Scott.<sup>3</sup>

**Theorem 5.1.** *There exists  $h_0 > 0$  such that, for every  $h \in (0, h_0]$  and for every  $\mathbf{f} \in L^2(\Omega; \mathbb{R}^m)$ , there is a unique solution  $\mathbf{u}_h$  of the Galerkin discretization (5.7), satisfying*

$$\|\nabla(\mathbf{u} - \mathbf{u}_h)\|_{L^2} \leq C_1 h |\mathbf{u}|_{\mathbf{H}^2} \quad \text{and} \quad \|\mathbf{u} - \mathbf{u}_h\|_{L^2} \leq C_2 h^2 |\mathbf{u}|_{\mathbf{H}^2},$$

where  $\mathbf{u}$  is the exact solution of (5.5). Furthermore, there exists  $\kappa_d > 0$  such that,

$$\inf_{\mathbf{v}_h \in \mathbf{V}_h} \sup_{\mathbf{w}_h \in \mathbf{V}_h} \frac{a(\mathbf{v}_h, \mathbf{w}_h)}{\|\nabla \mathbf{v}_h\|_{L^2} \|\nabla \mathbf{w}_h\|_{L^2}} \geq \kappa_d \quad \forall h \in (0, h_0]. \quad (5.8)$$

**Remark 5.1.** We remark that Theorem 5.1 as well as the following auxiliary Lemmas hold for any finite number of linear constraints. For example, if  $\mathbf{V}$  is given by

$$\mathbf{V} = \{\mathbf{v} \in \mathbf{H}_0^1(\Omega; \mathbb{R}^m) : (\mathbf{g}_i, \mathbf{v}) = 0, i = 1, \dots, n\},$$

where  $\mathbf{g}_i \in L^2(\Omega; \mathbb{R}^m)$ , and if the functions  $\mathbf{g}_i, i = 1, \dots, n$ , are linearly independent, then either minor modifications of the proofs, or simply a successive application of



24 *Bernhard Langwallner, Christoph Ortner, and Endre Süli*

the results for a single constraint can establish the results for a finite number of constraints.  $\square$

The proof of Theorem 5.1 requires two auxiliary results, which we provide in the following two lemmas. The first result shows that the constrained variational problem (5.5) inherits the  $H^2$ -regularity of the Laplace operator.

**Lemma 5.1.** *Let  $\mathbf{f} \in L^2(\Omega; \mathbb{R}^m)$  and let  $\mathbf{u} \in \mathbf{V}$  be the solution of (5.5), then  $\mathbf{u} \in H^2(\Omega; \mathbb{R}^m)$  and*

$$|\mathbf{u}|_{H^2} \leq C'_{\text{reg}} \|\mathbf{f}\|_{L^2}.$$

**Proof.** The result follows by explicitly computing a representation of  $\mathbf{u}$  in terms of the solution  $\tilde{\mathbf{u}}$  of a Poisson problem in the entire space  $H_0^1(\Omega; \mathbb{R}^m)$ :

$$(\nabla \tilde{\mathbf{u}}, \nabla \mathbf{v}) = (\mathbf{f} - \mathbf{M}\mathbf{u}, \mathbf{v}) \quad \forall \mathbf{v} \in H_0^1(\Omega; \mathbb{R}^m).$$

Let  $\mathbf{e}$  be the Riesz representation of  $\mathbf{g}$  in  $H_0^1(\Omega)$ ,

$$(\nabla \mathbf{e}, \nabla \mathbf{v}) = (\mathbf{g}, \mathbf{v}) \quad \forall \mathbf{v} \in H_0^1(\Omega; \mathbb{R}^m);$$

then it follows that the function

$$\tilde{\mathbf{u}} - t\mathbf{e}, \quad \text{where } t = \frac{(\nabla \tilde{\mathbf{u}}, \nabla \mathbf{e})}{\|\nabla \mathbf{e}\|^2},$$

belongs to  $\mathbf{V}$ , and that

$$(\nabla(\tilde{\mathbf{u}} - t\mathbf{e}), \nabla \mathbf{v}) = (\nabla \tilde{\mathbf{u}}, \nabla \mathbf{v}) - t(\mathbf{g}, \mathbf{v}) = (\mathbf{f} - \mathbf{M}\mathbf{u}, \mathbf{v}) \quad \forall \mathbf{v} \in \mathbf{V}.$$

Hence, we deduce that

$$\mathbf{u} = \tilde{\mathbf{u}} - t\mathbf{e}.$$

Assumption (5.1) ensures that  $\tilde{\mathbf{u}}$  and  $\mathbf{e}$  belong to  $H^2(\Omega; \mathbb{R}^m)$ , and therefore, we can deduce that  $\mathbf{u} \in H^2(\Omega)$ . The existence of  $C'_{\text{reg}}$  follows from the open mapping theorem and the fact that  $|\mathbf{u}|_{H^2}$  is a norm on  $\mathbf{V}$ .  $\square$

Our second auxiliary result shows that the constrained subspace  $\mathbf{V}_h$  inherits the approximation property (5.6) of  $S_{h,0}$ .

**Lemma 5.2.** *There exists  $h_0 > 0$  and a constant  $C'_{\text{apx}}$  such that*

$$\inf_{\mathbf{u}_h \in \mathbf{V}_h} \|\nabla(\mathbf{u} - \mathbf{u}_h)\|_{L^2} \leq C'_{\text{apx}} h |\mathbf{u}|_{H^2} \quad \text{for all } \mathbf{u} \in \mathbf{V} \cap H^2(\Omega; \mathbb{R}^m) \quad \text{and for } h \in (0, h_0].$$

**Proof.** Let  $\mathbf{e}_h \in S_{h,0}^m$  be the solution of

$$(\nabla \mathbf{e}_h, \nabla \mathbf{v}) = (\mathbf{g}, \mathbf{v}) \quad \forall \mathbf{v} \in S_{h,0}^m.$$

It is not difficult to verify, for  $h$  sufficiently small, say  $h \in (0, h_0]$ , that there exists  $\mathbf{v} \in S_{h,0}^m$  such that  $(\mathbf{g}, \mathbf{v}) \neq 0$  and hence  $\mathbf{e} \neq 0$  for  $h \in (0, h_0]$ .

Let  $\tilde{\mathbf{u}}_h \in S_{h,0}^m$  be the Ritz projection of  $\mathbf{u}$ , i.e.,

$$(\nabla \tilde{\mathbf{u}}_h, \nabla \mathbf{v}) = (\nabla \mathbf{u}, \nabla \mathbf{v}) \quad \forall \mathbf{v} \in S_{h,0}^m.$$

We construct the final approximant  $\mathbf{u}_h$ , as in the proof of Lemma 5.1,

$$\mathbf{u}_h = \tilde{\mathbf{u}}_h - \frac{(\nabla \tilde{\mathbf{u}}_h, \nabla \mathbf{e}_h)}{\|\nabla \mathbf{e}_h\|^2} \mathbf{e}_h \in \mathbf{V}_h.$$

Since  $(\mathbf{g}, \mathbf{u}) = 0$ , the error  $\|\nabla(\mathbf{u}_h - \tilde{\mathbf{u}}_h)\|_{L^2}$  can be estimated as follows:

$$\|\nabla(\mathbf{u}_h - \tilde{\mathbf{u}}_h)\|_{L^2} = \frac{|(\nabla \mathbf{e}_h, \nabla \tilde{\mathbf{u}}_h)|}{\|\nabla \mathbf{e}_h\|_{L^2}} = \frac{|(\mathbf{g}, \tilde{\mathbf{u}}_h - \mathbf{u})|}{\|\nabla \mathbf{e}_h\|_{L^2}} \leq \frac{\|\mathbf{g}\|_{H^{-1}}}{\|\nabla \mathbf{e}_h\|_{L^2}} \|\nabla(\mathbf{u} - \tilde{\mathbf{u}}_h)\|_{L^2}.$$

Since  $\mathbf{e}_h$  converges in  $H_0^1(\Omega)$  to the Riesz representation of  $\mathbf{g}$ , the term  $\|\mathbf{g}\|_{H^{-1}} / \|\nabla \mathbf{e}_h\|_{L^2}$  converges to one as  $h \rightarrow 0$ , and in particular is uniformly bounded on  $(0, h_0]$ , provided  $h_0$  is chosen sufficiently small. Invoking the approximation property (5.6) we arrive at the desired approximation result.  $\square$

### Proof of Theorem 5.1.

*Step 1: The Schatz argument.* The main part of the proof follows an argument originally given by Schatz.<sup>32</sup> Our presentation is largely analogous to Theorem 5.7.6 in Brenner & Scott.<sup>3</sup> First, let us simply assume the existence of a discrete solution  $\mathbf{u}_h$ . From Galerkin orthogonality we get  $a(\mathbf{u} - \mathbf{u}_h, \mathbf{u} - \mathbf{u}_h) = a(\mathbf{u} - \mathbf{u}_h, \mathbf{u} - \mathbf{v})$  for every  $\mathbf{v} \in \mathbf{V}_h$ . Then, using (5.4) we deduce that

$$\alpha \|\nabla(\mathbf{u} - \mathbf{u}_h)\|_{L^2}^2 \leq C_a \|\nabla(\mathbf{u} - \mathbf{u}_h)\|_{L^2} \|\nabla(\mathbf{u} - \mathbf{v})\|_{L^2} + K \|\mathbf{u} - \mathbf{u}_h\|_{L^2}^2, \quad (5.9)$$

for every  $\mathbf{v} \in \mathbf{V}_h$ , where we have used the continuity of  $a$ . Considering the adjoint problem

$$a(\mathbf{z}, \mathbf{w}) = (\mathbf{u} - \mathbf{u}_h, \mathbf{z}) \quad \forall \mathbf{z} \in \mathbf{V}$$

we can show that

$$\|\mathbf{u} - \mathbf{u}_h\|_{L^2}^2 \leq C'_{\text{reg}} h \|\nabla(\mathbf{u} - \mathbf{u}_h)\|_{L^2} \|\mathbf{w}\|_{H^2} \leq C'_{\text{reg}} h \|\nabla(\mathbf{u} - \mathbf{u}_h)\|_{L^2} \|\mathbf{u} - \mathbf{u}_h\|_{L^2},$$

where Lemma 5.1 was used to obtain  $H^2$ -regularity for  $\mathbf{w}$ . This results in

$$\|\mathbf{u} - \mathbf{u}_h\|_{L^2} \leq C'_{\text{reg}} h \|\nabla(\mathbf{u} - \mathbf{u}_h)\|_{L^2}. \quad (5.10)$$

Applying this bound to (5.9) and choosing  $h$  sufficiently small, we obtain the bound

$$\|\nabla(\mathbf{u} - \mathbf{u}_h)\|_{L^2} \leq C \|\nabla(\mathbf{u} - \mathbf{v})\|_{L^2} \quad \forall \mathbf{v} \in \mathbf{V}_h. \quad (5.11)$$

$\mathbf{V}_h$  is a finite-dimensional space, so existence of a solution  $\mathbf{u}_h$  for arbitrary right-hand sides, and its uniqueness are equivalent. Suppose, for  $\mathbf{f} = 0$ , the discrete problem had a nontrivial solution  $\mathbf{u}_h \neq 0$ . Then, equation (5.11) would produce a contradiction for  $h$  sufficiently small because  $\mathbf{u} = 0$ . Hence, there is a unique solution  $\mathbf{u}_h$ .

26 *Bernhard Langwallner, Christoph Ortner, and Endre Süli*

Taking the infimum over  $\mathbf{v} \in \mathbf{V}_h$  in (5.11) and using Lemma 5.2 yields the first error bound stated in the theorem. Combining this bound with (5.10) provides the second error bound.

*Step 2. Uniform Inf-Sup Constant.* Unique solvability of (5.5) for  $\mathbf{f} = \mathbf{0}$  implies that  $a$  satisfies the inf-sup condition

$$\inf_{\mathbf{u} \in \mathbf{V}} \sup_{\mathbf{v} \in \mathbf{V}} \frac{a(\mathbf{u}, \mathbf{v})}{\|\nabla \mathbf{u}\|_{L^2} \|\nabla \mathbf{v}\|_{L^2}} = \kappa > 0. \quad (5.12)$$

Our aim now is to prove the validity of a corresponding condition for  $\mathbf{V}_h \times \mathbf{V}_h$ , which is uniform in  $h$ .

To obtain  $\mathcal{O}(h)$ -convergence of  $\mathbf{u}_h$  to  $\mathbf{u}$  it had to be assumed that  $\mathbf{f} \in L^2(\Omega)^m$ . However, both the continuous and discrete solution  $\mathbf{u} \in \mathbf{V}$ , respectively  $\mathbf{u}_h \in \mathbf{V}_h$ , exist and are unique if  $\mathbf{f} \in H^{-1}(\Omega)^m$ ,

$$a(\mathbf{u}, \mathbf{v}) = \langle \mathbf{f}, \mathbf{v} \rangle \quad \forall \mathbf{v} \in \mathbf{V} \quad \text{and} \quad a(\mathbf{u}_h, \mathbf{v}) = \langle \mathbf{f}, \mathbf{v} \rangle \quad \forall \mathbf{v} \in \mathbf{V}_h.$$

We note that, to prove (5.11), we only used that  $\mathbf{u} - \mathbf{u}_h \in L^2(\Omega; \mathbb{R}^m)$  but not  $\mathbf{f} \in L^2(\Omega; \mathbb{R}^m)$ . Hence, choosing  $\mathbf{v} = \mathbf{0}$  in (5.11) and using a triangle inequality and the inf-sup condition (5.12) we deduce that,

$$\|\nabla \mathbf{u}_h\|_{L^2} \leq (1 + C) \|\nabla \mathbf{u}\|_{L^2} \leq (1 + C) \kappa^{-1} \|\mathbf{f}\|_{\mathbf{V}^*} \quad \forall h \in (0, h_0], \quad (5.13)$$

where  $h_0 > 0$  is chosen sufficiently small, and where  $C$  is the constant from (5.11). If  $\mathbf{f} \in \mathbf{V}_h^*$  then, by the Hahn–Banach theorem,  $\mathbf{f}$  can be extended to an element of  $\mathbf{V}^*$  while preserving its norm, and hence, we obtain

$$\|\nabla L_h^{-1} \mathbf{f}\|_{L^2} \leq (1 + C) \kappa^{-1} \|\mathbf{f}\|_{\mathbf{V}_h^*} \quad \forall \mathbf{f} \in \mathbf{V}_h^* \quad \forall h \in (0, h_0],$$

where  $L_h^{-1}$  denotes the solution operator for (5.7). This statement is equivalent to the uniform inf-sup condition (5.8) with  $\kappa_d = \kappa/(1 + C)$ .  $\square$

### Acknowledgment

The authors were supported by the EPSRC project ‘New Frontiers in the Mathematics of Solids’.

### References

1. R. A. Adams and J. F. Fournier. *Sobolev Spaces*. Pure and Applied Mathematics Series. Academic Press, Amsterdam, second edition, 2003. International Series in Pure and Applied Mathematics.
2. X. Blanc, C. Le Bris, and P.-L. Lions. From molecular models to continuum mechanics. *Arch. Ration. Mech. Anal.*, 164(4):341–381, 2002.
3. S. C. Brenner and L. R. Scott. *The mathematical theory of finite element methods*, volume 15 of *Texts in Applied Mathematics*. Springer-Verlag, New York, second edition, 2002.
4. F. Brezzi and M. Fortin. *Mixed and hybrid finite element methods*, volume 15 of *Springer Series in Computational Mathematics*. Springer-Verlag, New York, 1991.

5. F. Brezzi, J. Rappaz, and P.-A. Raviart. Finite-dimensional approximation of nonlinear problems. I. Branches of nonsingular solutions. *Numer. Math.*, 36(1):1–25, 1980/81.
6. E. Cancès, R. Chakir, and Y. Maday. Numerical analysis of nonlinear eigenvalue problems. *arXiv:0905.1645v2*, 2009.
7. E. Cancès, R. Chakir, and Y. Maday. Numerical analysis of the planewave discretization of orbital-free and Kohn–Sham models part i: The Thomas–Fermi–von Weizsäcker model. *arXiv:0909.1464v1*, 2009.
8. I. Catto, C. Le Bris, and P.-L. Lions. *The mathematical theory of thermodynamic limits: Thomas–Fermi type models*. Oxford Mathematical Monographs. The Clarendon Press Oxford University Press, New York, 1998.
9. D. M. Ceperley and B. J. Alder. Ground state of the electron gas by a stochastic method. *Phys. Rev. Lett.*, 45(7):566–569, Aug 1980.
10. B. Dacorogna. *Direct methods in the calculus of variations*, volume 78 of *Applied Mathematical Sciences*. Springer, New York, second edition, 2008.
11. M. Dobrowolski and R. Rannacher. Finite element methods for nonlinear elliptic systems of second order. *Math. Nachr.*, 94:155–172, 1980.
12. L. C. Evans. *Partial differential equations*, volume 19 of *Graduate Studies in Mathematics*. American Mathematical Society, Providence, RI, 1998.
13. M. Finnis. *Interatomic Forces in Condensed Matter*. Oxford University Press, USA, 2003.
14. C. J. García-Cervera. An efficient real space method for orbital-free density-functional theory. *Commun. Comput. Phys.*, 2(2):334–357, 2007.
15. C. J. García-Cervera. A remark on “an efficient real space method for orbital-free density-functional theory”. *Commun. Comput. Phys.*, 3(4):968–972, 2008.
16. C. J. García-Cervera, J. Lu, and W. E. A sub-linear scaling algorithm for computing the electronic structure of materials. *Commun. Math. Sci.*, 5(4):999–1026, 2007.
17. V. Gavini, K. Bhattacharya, and M. Ortiz. Quasi-continuum orbital-free density-functional theory: a route to multi-million atom non-periodic DFT calculation. *J. Mech. Phys. Solids*, 55(4):697–718, 2007.
18. V. Gavini, J. Knap, K. Bhattacharya, and M. Ortiz. Non-periodic finite-element formulation of orbital-free density functional theory. *J. Mech. Phys. Solids*, 55(4):669–696, 2007.
19. S. K. Ghosh and A. K. Dhara. Density-functional theory of two-dimensional electron gas in a magnetic field. *Physical Review A*, 40(10):6103–6106, 1989.
20. P. Grisvard. *Elliptic problems in nonsmooth domains*, volume 24 of *Monographs and Studies in Mathematics*. Pitman (Advanced Publishing Program), Boston, MA, 1985.
21. P. Hohenberg and W. Kohn. Inhomogeneous electron gas. *Phys. Rev.*, 136(3B):B864–B871, 1964.
22. W. Kohn. Nobel Lecture: Electronic structure of matterwave functions and density functionals. *Reviews of Modern Physics*, 71(5):1253–1266, 1999.
23. W. Kohn and L. J. Sham. Self-consistent equations including exchange and correlation effects. *Phys. Rev.*, 140(4A):A1133–A1138, 1965.
24. B. Langwallner, C. Ortner, and E. Süli. Optimal convergence rates for the finite element discretization of an electronic density functional. *In preparation*.
25. E. H. Lieb. Thomas–Fermi and related theories of atoms and molecules. *Rev. Modern Phys.*, 53(4):603–641, 1981.
26. L. Lignère and A. C. Carter. An introduction to orbital-free density functional theory. In S. Yip, editor, *Handbook of Materials Modeling*, volume 15 of *Texts in Applied Mathematics*, chapter 1.8. Springer-Verlag, second edition, 2005.

28 *Bernhard Langwallner, Christoph Ortner, and Endre Süli*

27. H. Maurer and J. Zowe. First and second order necessary and sufficient optimality conditions for infinite-dimensional programming problems. *Math. Programming*, 16(1):98–110, 1979.
28. D. Negrut, M. Anitescu, A. El-Azab, and P. Zapol. Quasicontinuum-like reduction of density functional theory calculations of nanostructures. *Journal of Nanoscience and Nanotechnology*, 8(7):3729–3740, 2008.
29. D. Negrut, M. Anitescu, T. Munson, and P. Zapol. Simulating nanoscale processes in solids using DFT and the quasicontinuum method (IMECE2005-81755). *Proceedings of IMECE*, 2005.
30. R.G. Parr and W. Yang. *Density-Functional Theory of Atoms and Molecules*. Oxford University Press, USA, 1989.
31. J. P. Perdew and A. Zunger. Self-interaction correction to density-functional approximations for many-electron systems. *Phys. Rev. B*, 23(10):5048–5079, May 1981.
32. A. H. Schatz. An observation concerning Ritz–Galerkin methods with indefinite bilinear forms. *Math. Comp.*, 28:959–962, 1974.
33. Y. A. Wang, N. Govind, and E. A. Carter. Orbital-free kinetic-energy functionals for the nearly free electron gas. *Physical Review B*, 58(20):13465–13471, 1998.
34. Y. A. Wang, N. Govind, and E. A. Carter. Orbital-free kinetic-energy density functionals with a density-dependent kernel. *Physical Review B*, 60(24):16350–16358, 1999.
35. S. C. Watson and E. A. Carter. Linear-scaling parallel algorithms for the first principles treatment of metals. *Computer Physics Communications*, 128(1-2):67–92, 2000.
36. E. Zeidler. *Nonlinear functional analysis and its applications. I*. Springer-Verlag, New York, 1986. Fixed-point theorems, Translated from the German by Peter R. Wadsack.