

This item was submitted to [Loughborough's Research Repository](#) by the author.  
Items in Figshare are protected by copyright, with all rights reserved, unless otherwise indicated.

## Analysis of Schrodinger operators with inverse square potentials II: FEM and approximation of eigenfunctions in the periodic case

PLEASE CITE THE PUBLISHED VERSION

<http://dx.doi.org/10.1002/num.21861>

PUBLISHER

© Wiley Periodicals, Inc.

VERSION

AM (Accepted Manuscript)

PUBLISHER STATEMENT

This work is made available according to the conditions of the Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International (CC BY-NC-ND 4.0) licence. Full details of this licence are available at:  
<https://creativecommons.org/licenses/by-nc-nd/4.0/>

LICENCE

CC BY-NC-ND 4.0

REPOSITORY RECORD

Hunsicker, Eugenie, Hengguang Li, Victor Nistor, and Ville Uski. 2019. "Analysis of Schrodinger Operators with Inverse Square Potentials II: FEM and Approximation of Eigenfunctions in the Periodic Case". figshare. <https://hdl.handle.net/2134/17169>.

# ANALYSIS OF SCHRÖDINGER OPERATORS WITH INVERSE SQUARE POTENTIALS II: FEM AND APPROXIMATION OF EIGENFUNCTIONS IN THE PERIODIC CASE

EUGENIE HUNSICKER, HENGGUANG LI, VICTOR NISTOR, AND VILLE USKI

**ABSTRACT.** In this paper, we consider the problem of optimal approximation of eigenfunctions of Schrödinger operators with isolated inverse square potentials and of solutions to equations involving such operators. It is known in this situation that the finite element method performs poorly with standard meshes. We construct an alternative class of graded meshes, and prove and numerically test optimal approximation results for the finite element method using these meshes. Our numerical tests are in good agreement with our theoretical results.

## CONTENTS

1. Introduction and statement of main results	1
1.1. Notation and Results	3
Acknowledgements	6
2. Background results	7
3. Approximation and mesh refinement	9
3.1. Construction of the meshes	9
3.2. Proof of Theorem 3.1	11
4. Applications to Finite Element Methods	14
4.1. The condition number of the stiffness matrix	15
5. Numerical tests of the finite element method	18
References	21

## 1. INTRODUCTION AND STATEMENT OF MAIN RESULTS

Schrödinger type operators of the form  $H = -\Delta + V$  with inverse square potentials  $V$  arise in a variety of interesting contexts motivated by continuum mechanics, by quantum physics, by theoretical numerical analysis considerations, and by questions in other areas. The purpose of this paper is to develop numerical approximation tools for studying the spectra of such operators. Let us denote by

---

*Date:* December 17, 2013.

V.N. was partially supported by the NSF Grants OCI-0749202 and DMS-1016556. Manuscripts available from <http://www.math.psu.edu/nistor/>. H.L. was partially supported by the NSF Grant DMS-1158839. V.U. was supported by and EH was supported in part by Leverhulme Trust grant J11695.

$\rho$  the smoothed distance to the set  $\mathcal{S}$  of singular points of  $V$ . The standard example of a Schrödinger operator with  $c/\rho$  potential is a special case of the inverse square potentials we consider, where the function  $\rho^2 V$  vanishes to order 1 at the singularity, and the results in this paper do apply to such operators. However, in addition, Hamiltonians with true inverse square potentials arise in relativistic quantum mechanics from the square of the Dirac operator coupled with an interaction potential. They arise also in the interaction of a polar molecule with an electron. See [39, 43] for further applications of inverse square potentials in physics. See also [3, 14, 15, 23, 29, 35, 38] for related results on operators with singular coefficients.

Consider a Hamiltonian operator  $H := -\Delta + V$  that is periodic on  $\mathbb{R}^3$  with periodicity lattice  $\Lambda$ . The standard approach to studying the spectrum of  $H$  on complex function spaces is by considering the associated operators on Bloch waves for every vector in the first Brillouin zone for the lattice. Mathematically speaking, this associates to  $H$  a family of Bloch operators,  $H_{\mathbf{k}} := -\sum_{j=1}^3 (\partial_j + i\mathbf{k}_j)^2 + V$ , parametrized over vectors  $\mathbf{k}$  in the fundamental domain of the dual lattice to  $\Lambda$ , which act on  $\Lambda$ -periodic functions in  $\mathbb{R}^3$ , and thus have discrete spectra. (The definition of the operators  $H_{\mathbf{k}}$  will be discussed in more detail below.) Equivalently, we can consider the Bloch operators as acting on functions on the 3-torus obtained by identifying opposite sides of the fundamental domain of the lattice.

When the potential is smooth, the eigenfunctions of the operators  $H_{\mathbf{k}}$  are smooth, and therefore the associated eigenvalue problems on the torus can be approximately solved using a standard mesh with the finite element method. Further, in this case, again because the eigenfunctions are all smooth, the convergence rate for the finite element approximations in terms of standard Sobolev spaces can be made as high as desired by choosing to work with elements of sufficiently high polynomial degree on the tetrahedra of the mesh. However, when the potential,  $V$ , has singularities, the eigenfunctions of  $H_{\mathbf{k}}$  have singularities at the singularities of the potential (which can even blow up like  $\rho^\alpha$  for  $\alpha \in (-1/4, 0)$ ), and in particular have limited Sobolev regularity. Thus the convergence rates for these methods are less than optimal (suboptimal) if quasi-uniform meshes or other classical approximation methods are used.

For problems with singular solutions, the phenomenon of the classical finite element methods exhibiting suboptimal convergence rates has been observed by many authors. For instance, in the setting of reentrant corners, the problem was studied by many authors, including by Apel, Nicaise, and Schöberl in [2], by Babuška, Kellogg, and Pitkäranta in [5], by Bacuta, Bramble, and Xu in [20], by Demkowicz, Monk, Schwab, and Vardapetyan in [24], by Mazzucato, Li, and Nistor in [37], and by Wahlbin in [42]. A consequence of their research is that the framework of quasi-uniform meshes leads to suboptimal rates of convergence in the Finite Element Method for problems on non-convex domains. We can make the ideas of “optimal” and “suboptimal” rates of convergence precise as follows. Let  $N$  be the number of degrees of freedom in the finite element space. By “optimal rates of convergence” we mean the rates of convergence of the order  $N^{-m/3}$  obtained using quasi-uniform meshes and continuous piecewise polynomials of degree  $m$  when the solution is in

$H^{m+1}$  (see [17, 22] and [41]). Any rate of convergence that is less than optimal will be called *suboptimal*.

In this paper, we present a modification of the finite element method for approximating eigenvalues and eigenfunctions of the Bloch operators,  $H_{\mathbf{k}}$ , in the case that the potential  $V$  has inverse square singularities. The modification uses graded meshes near the singularities. The workhorse theorem in this paper is Theorem 3.1, which is an approximation theorem for functions in a family of weighted Sobolev spaces using the modified finite elements. By Theorem 2.1 from [30], all eigenfunctions of the operators  $H_{\mathbf{k}}$  can be decomposed into the sum of a well-understood singular part and a part that lies in all weighted Sobolev spaces in this family. As a corollary of these two theorems, we get our main theorems, Theorems 1.1 and 1.2, which give convergence results in terms of the weighted Sobolev spaces for the eigenvalue problem for  $H_{\mathbf{k}}$  and for solutions to inhomogeneous equations of the form  $(L + H_{\mathbf{k}})u = f$  using the finite element method with these meshes. The results show that the convergence rates can again be made as high as desired by choosing sufficiently high degree polynomials on the tetrahedra of the mesh. In particular, in the case of linear elements, we recover the convergence rate obtained for linear elements in the smooth setting. In the last section, we carry out numerical tests using linear elements that show good agreement with these theoretical results.

**1.1. Notation and Results.** Before we can state our approximation results, we must fix some notation and state the assumptions we make about our Hamiltonian operators. As above, consider a Hamiltonian operator  $H := -\Delta + V$  that is periodic on  $\mathbb{R}^3$  with periodicity lattice  $\Lambda$ . Its fundamental domain is a parallelepiped whose faces can be identified under the symmetries of  $H$  to form the torus  $\mathbb{T} = \mathbb{R}^3/\Lambda$ , which is how we will denote this fundamental domain in the remainder of this paper. Let  $\rho(x)$  be a continuous function on  $\mathbb{T}$  that is given by  $\rho(x) = |x - p|$  for  $x$  close to  $p$ , is smooth except at the points of  $\mathcal{S}$ , and may be assumed to be equal to one outside a neighborhood of  $\mathcal{S}$ .

We need two assumptions about the potentials  $V$  that we will consider in this paper. First, we assume that  $V$  is smooth except at a *finite* set of points  $\mathcal{S} \subset \mathbb{T} = \mathbb{R}^3/\Lambda$ , near which it has singularities of the form  $Z/\rho^2$ , where  $Z$  is continuous on  $\mathbb{T}$  and smooth in *generalized* spherical coordinates  $(r, x')$ ,  $r \in [0, \infty)$ ,  $x' \in S^2$  around  $p \in \mathcal{S}$ . We want to make this more precise. Let  $O$  be the origin in  $\mathbb{R}^3$ . By *blowing up  $O$  in  $\mathbb{R}^3$*  we shall mean replacing  $\mathbb{R}^3$  with  $S^2 \times [0, \infty)$  such that  $(r, x')$  in the blown up space corresponds to  $rx' \in \mathbb{R}^3$ . In this way, the point  $O$  was replaced by the copy  $S^2 \times \{0\}$  of the two dimensional sphere  $S^2$ . A function  $f : \mathbb{R}^3 \setminus \{0\} \rightarrow \mathbb{R}$  will be called *smooth in generalized polar coordinates* if it extends to a smooth function on  $S^2 \times [0, \infty)$ . Let us denote then by  $\mathbb{M}$  the smooth manifold with boundary obtained by replacing each point  $p \in \mathcal{S}$  with the two dimensional sphere  $p \times S^2$ , that is, by blowing up each singular point of  $V$ , according to the procedure just explained. Smoothness in polar coordinates around each singular point then means smoothness on  $\mathbb{M}$ :

$$(1) \quad \textbf{Assumption 1 :} \quad Z := \rho^2 V \in \mathcal{C}(\mathbb{T}) \cap \mathcal{C}^\infty(\mathbb{M}).$$

(While the local structure of  $\mathbb{M}$  is the one just explained, it is difficult to give a global description of  $\mathbb{M}$ , as a set. However,  $\mathbb{M} := \mathbb{T} \setminus \mathcal{S} \cup \mathcal{S} \times S^2$ , where the union is disjoint, with the smooth structure defined above.) Assumption 1, more precisely the continuity of  $Z$  at  $\mathcal{S}$ , allows us to formulate our second assumption. Namely,

$$(2) \quad \textbf{Assumption 2} : \eta := \min_{p \in \mathcal{S}} \sqrt{1/4 + Z(p)} > 0.$$

We will see that the value of  $\eta$  determines the strength of the singularity in the eigenfunctions associated to the Hamiltonians  $H_{\mathbf{k}}$ . In particular, we assume that for all  $p \in \mathcal{S}$ ,  $Z(p) > -1/4$ . These assumptions are sharp in the sense that the analysis yields fundamentally different results if either one fails. In particular, the value  $Z(p) = -1/4$  corresponds to the critical coupling for an isolated inverse square potential in  $\mathbb{R}^3$  where the system undergoes a transition between the conformal and non-conformal regimes [39]. If the first assumption fails, then the available analytic techniques are much weaker, see for instance [26, 27]. In either case, the approximation theorems in this paper fail if either assumption is violated. More details of this are included in [30], and a study of the analysis when these assumptions are relaxed will be examined in a forthcoming paper.

We are interested in understanding the spectrum and eigenfunctions of the operators  $H_{\mathbf{k}}$ . As mentioned above, we do this by studying Bloch waves. Recall that if  $\mathbf{k}$  is an element of the first Brillouin zone of  $\Lambda$ , that is, is an element of the fundamental domain of the dual lattice of  $\Lambda$ , then a Bloch wave with wave vector  $\mathbf{k}$  is a function in  $L^2_{loc}(\mathbb{R}^3)$  that satisfies the semi-periodicity condition

$$(3) \quad \psi_{\mathbf{k}}(x + X) = e^{i\mathbf{k} \cdot X} \psi_{\mathbf{k}}(x) \quad \forall X \in \Lambda.$$

It is well known that such a Bloch wave can be written as

$$(4) \quad \psi_{\mathbf{k}}(x) = e^{i\mathbf{k} \cdot x} u_{\mathbf{k}}(x)$$

for a function  $u_{\mathbf{k}}$  that is truly periodic with respect to  $\Lambda$  and thus can be considered as living on the three-torus  $\mathbb{T}$  [25]. We define the  $\mathbf{k}$ -Hamiltonian  $H_{\mathbf{k}}$  on  $L^2(\mathbb{T})$  by

$$(5) \quad H_{\mathbf{k}} := - \sum_{j=1}^3 (\partial_j + i\mathbf{k}_j)^2 + V.$$

Then we have further that if a Bloch wave  $\psi_{\mathbf{k}}$  satisfies  $H\psi_{\mathbf{k}} = \lambda\psi_{\mathbf{k}}$ , then the function  $u_{\mathbf{k}} := e^{-i\mathbf{k} \cdot x} \psi_{\mathbf{k}}(x)$  is a standard  $L^2$ -eigenfunction of  $H_{\mathbf{k}}$  with eigenvalue  $\lambda$ .

Let  $\lambda_j$ ,  $j \geq 1$ , be the eigenvalues of  $H_{\mathbf{k}}$ , arranged in increasing order,  $\dots \leq \lambda_j \leq \lambda_{j+1} \leq \dots$ , and repeated according to their multiplicities. We know that  $H_{\mathbf{k}}$  is self-adjoint and has an orthonormal basis of eigenvectors by the results of [30]. In particular, in this case, by the multiplicity of an eigenvalue  $\lambda$  we shall mean the dimension of the corresponding eigenspace  $E(\lambda)$ . We fix an eigenbasis  $(u_j)$  of  $H_{\mathbf{k}}$ . One of our goals is to approximate the eigenvalues  $\lambda_j$  and the corresponding eigenfunctions  $u_j$ .

As usual, for our finite element approximation results, we consider a sequence  $S_n$  of finite dimensional subspaces of the domain of  $H_{\mathbf{k}}$  and project onto  $S_n$  to obtain a discrete formulation. To define the appropriate projection, we use Theorem 2.1, which says that for sufficiently large  $L \geq 0$ ,  $L + H_{\mathbf{k}}$  is an isomorphism between

appropriate weighted Sobolev spaces. This ensures the coercivity of the natural sesquilinear form  $a$  defined as follows:

$$(6) \quad a(v, w) := ((L + H_{\mathbf{k}})v, w) = ((\nabla + i\mathbf{k})v, (\nabla + i\mathbf{k})w) + ((L + V)u, v),$$

where  $(\nabla + i\mathbf{k})u$  is the vector with components  $(\partial_j + ik_j)u$  and  $(v, w) := \int_{\mathbb{T}} v \bar{w} dx$  is the sesquilinear inner product on the complex Hilbert space  $L^2(\mathbb{T})$ . Now let  $R_n$  denote the projection onto  $S_n$  taken with respect to the form  $a(y, w)$ . The operator  $R_n$  will be called the associated Riesz projection, as usual. Let also  $H_{\mathbf{k},n} := R_n H_{\mathbf{k}} R_n$  be the associated finite element approximation of  $H_{\mathbf{k}}$  acting on  $S_n$ . Denote by  $\lambda_{j,n}$  the eigenvalues of the approximation  $H_{\mathbf{k},n}$ , again arranged in increasing order,  $\dots \leq \lambda_{j,n} \leq \lambda_{j+1,n} \leq \dots$ , and repeated according to their multiplicities. The spaces  $S_n$  we use for our theorems are defined in terms of a sequence of graded tetrahedral meshes  $\mathcal{T}_n := k^n(\mathcal{T}_0)$  on  $\mathbb{T}$  (sometimes called triangulations), given by sequential refinements, associated to a scaling parameter  $k$ , of an original tetrahedral mesh  $\mathcal{T}_0$ . We describe the meshing refinement procedure in detail in Section 3. We will take  $S_n = S(\mathcal{T}_n, m)$ , the finite element spaces associated to these meshes (*i.e.*, using continuous, piecewise polynomials of degree  $m$ ).

We now state our two main theorems, which will be proved in the main body of the paper. Both of these theorems are given in terms of weighted Sobolev spaces whose definition we now recall

$$(7) \quad \mathcal{K}_a^m(\mathbb{T}; \mathcal{S}) := \{v : \mathbb{T} \setminus \mathcal{S} \rightarrow \mathbb{C}, \rho^{|\alpha|-a} \partial^\alpha v \in L^2(\mathbb{T}), \forall |\alpha| \leq m\},$$

with semi-norms and norms

$$(8) \quad |v|_{\mathcal{K}_a^m(\mathbb{T}; \mathcal{S})}^2 := \sum_{|\alpha|=m} \|\rho^{|\alpha|-a} \partial^\alpha v\|_{L^2(\mathbb{T})}^2, \quad \|v\|_{\mathcal{K}_a^m(\mathbb{T}; \mathcal{S})}^2 = \sum_{|\alpha| \leq m} |v|_{\mathcal{K}_a^m(\mathbb{T}; \mathcal{S})}^2.$$

These spaces have been considered in many other papers, most notably in Kondratiev's groundbreaking paper [34]. Our first main theorem is a theoretical result for the finite element method approximation of eigenvalues and eigenfunctions of  $H_{\mathbf{k}}$  using tetrahedralisations with graded meshes:

**Theorem 1.1.** *Let  $\lambda_j$  be an eigenvalue of  $H_{\mathbf{k}}$  and fix  $a \leq m$ ,  $0 < a < \eta$ , with  $\eta = \min_{p \in \mathcal{S}} \sqrt{1/4 + Z(p)}$ , as in Assumption 2. Consider the spaces  $S_n$  associated to the nested sequence  $\mathcal{T}_n$  of meshes on  $\mathbb{T}$  defined by the scaling parameter  $k = 2^{-m/a}$  and piecewise polynomials of degree  $m$ . Then there exists a constant  $c(\lambda_j, a)$  with the following property. Let  $R_n$  be the associated Riesz projections. Denote by  $\lambda_{j,n}$  the eigenvalues of the approximation  $H_{\mathbf{k},n} := R_n H_{\mathbf{k}} R_n$  acting on  $S_n$ , again arranged in increasing order,  $\dots \leq \lambda_{j,n} \leq \lambda_{j+1,n} \leq \dots$ , and repeated according to their multiplicities. Then*

$$|\lambda_j - \lambda_{j,n}| \leq c(\lambda_j, a) \dim(S_n)^{-2m/3}.$$

Moreover, let  $E'_n(\lambda)$  be the sum of the eigenspaces  $E(\lambda_{j,n})$  for  $\lambda_j = \lambda$ . Then for each  $j$ , there exist suitable  $u_{j,n} \in E'_n(\lambda_j)$  such that

$$\|u_j - u_{j,n}\|_{H^1(\mathbb{T})} \leq \|u_j - u_{j,n}\|_{\mathcal{K}_1^1(\mathbb{T}; \mathcal{S})} \leq c(\lambda_j, a) \dim(S_n)^{-m/3}.$$

Recall that the analogous result in the setting of smooth potentials is the same as this result, except the mesh is not graded, which corresponds to the parameter value  $k = 0.5$ , and the weighted Sobolev space in the estimate is replaced with the standard Sobolev space  $H^1$ . Thus our result shows that if  $\eta \geq m$ , ungraded meshes give the optimal convergence rate for elements of degree  $m$ . Moreover, if  $\eta < m$ , our numerical tests seem to indicate that graded meshes are necessary to obtain optimal convergence.

For our second main theorem, we consider the finite element approximations of the equation

$$(9) \quad (L + H_{\mathbf{k}})v = f, \quad \text{for } L > C_0,$$

where  $C_0$  is the constant from Theorem 2.1 below. Let  $v$  be the solution of Equation (9) above. We then define the usual Galerkin finite element approximation  $v_n$  of  $v$  as the unique  $v_n \in S_n := S(\mathcal{T}_n, m)$  such that

$$(10) \quad a(v_n, w_n) := (f, w_n), \quad \text{for all } w_n \in S_n.$$

Here is our second main theorem.

**Theorem 1.2.** *The sequence  $\mathcal{T}_n := k^n(\mathcal{T}_0)$  of meshes on  $\mathbb{T}$  defined using the  $k$ -refinement, for  $k = 2^{-m/a}$ ,  $0 < a < \eta$ ,  $a \leq m$ , has the following property: The sequence  $v_n \in S_n := S(\mathcal{T}_n, m)$  of finite element (Galerkin) approximations of  $v$  from Equation (10) satisfies*

$$(11) \quad \|v - v_n\|_{\mathcal{K}_1^1(\mathbb{T}; \mathcal{S})} \leq C \dim(S_n)^{-m/3} \|f\|_{\mathcal{K}_{a-1}^{m-1}(\mathbb{T}; \mathcal{S})},$$

where  $C$  is independent of  $n$  and  $f$ .

The proofs of these theorems use the regularity results from [30], the approximation result of Theorem 3.1, and some general results (C  a's Lemma and results of Babu  ska and Osborn on the approximation of eigenvectors). We will recall the statements of the relevant regularity results from that paper and some additional background material in Section 2. Thus although as indicated by the title, [30] is the first part of an extended project of which this paper forms the second part, this paper may be read independently of [30].

The remainder of the paper is organized as follows. In Section 3, we first describe the  $k$ -refinement algorithm for the three dimensional tetrahedral meshes, which results in a sequence of meshes  $\mathcal{T}_n$ . We then prove a general interpolation approximation result for the sequence of finite element spaces associated to this sequence of meshes. In Section 4 we use this general approximation result to prove our main approximation results. This section includes in particular the proofs of Theorem 1.1 and Theorem 1.2, as well as an additional result about the condition number of the stiffness matrix associated to the finite element spaces  $S_n$ . In the last section, Section 5, we discuss results of numerical tests of the method for solving equations of the form  $(L + H_{\mathbf{k}})v = f$  and compare them to the theoretical results.

**Acknowledgements.** We would like to thank Bernd Ammann, Douglas Arnold, and Catarina Carvalho for useful discussions. We also thank the Leverhulme Trust whose funding supported the fourth author during this project. This project was

started while Hunsicker and Nistor were visiting the Max Planck Institute for Mathematics in Bonn, Germany, and we are grateful for its support. We would like to thank Serge Nicaise for pointing out an imprecision in our previous statement of our theorem on eigenvalue approximation. We would like to also thank two anonymous referees for their careful reading of the paper and their useful comments.

## 2. BACKGROUND RESULTS

In this section we recall some definitions and results from [30], as well as the classical approximation result for Lagrange interpolants (see [4, 17, 22, 41]), that will be used in the proofs of the approximation theorems above.

The first result that we recall guarantees the existence of solutions of equations of the form  $(L + H_{\mathbf{k}})v = f$  for  $L$  greater than some constant  $C_0$ , and identifies the natural domain of  $H_{\mathbf{k}}$ . Let us fix smooth functions  $\chi_p$  supported near points of  $\mathcal{S}$  such that the functions  $\chi_p$  have disjoint supports and  $\chi_p = 1$  in a small neighborhood of  $p \in \mathcal{S}$ . Then Theorem 1.1, Lemma 3.4, and Proposition 3.6 from [30] combine to give right away the following result.

**Theorem 2.1.** *Let  $V$  be a potential satisfying both Assumptions 1 and 2. Then there exists a constant  $C_0$  such that  $L + H_{\mathbf{k}} : \mathcal{K}_{a+1}^{m+1}(\mathbb{T}; \mathcal{S}) \rightarrow \mathcal{K}_{a-1}^{m-1}(\mathbb{T}; \mathcal{S})$  is an isomorphism for all  $m \in \mathbb{Z}_{\geq 0}$ , all  $|a| < \eta$ , and all  $L > C_0$ . In addition, recall the form  $a(\cdot, \cdot)$  from (6). Then, we have that  $a(\cdot, \cdot)$  is continuous on  $\mathcal{K}_1^1(\mathbb{T}; \mathcal{S}) \times \mathcal{K}_1^1(\mathbb{T}; \mathcal{S})$  and coercive. Namely, there is  $C > 0$ , such that, for any  $u \in \mathcal{K}_1^1(\mathbb{T}; \mathcal{S})$ ,*

$$(12) \quad a(u, u) \geq C \|u\|_{\mathcal{K}_1^1(\mathbb{T}; \mathcal{S})}^2.$$

Moreover, for any  $u \in \mathcal{K}_{a+1}^{m+1}(\mathbb{T}; \mathcal{S})$  satisfying  $(L + H_{\mathbf{k}})u = f \in H^{m-1}(\mathbb{T})$ , we can find constants  $a_p \in \mathbb{R}$  such that

$$u_{\text{reg}} := u - \sum_{p \in \mathcal{S}} a_p \chi_p \rho^{\sqrt{1/4 + Z(p)} - 1/2} \in \mathcal{K}_2^{m+1}(\mathbb{T}; \mathcal{S}),$$

with  $Z$  as in Assumption 1.

We obtain, in particular, that  $H_{\mathbf{k}}$  has a natural self-adjoint extension, the Friedrichs extension. Therefore, from now on, we shall extend  $H_{\mathbf{k}}$  to the domain of the Friedrichs extension of  $L + H_{\mathbf{k}}$ , as in the above Theorem. Let us denote by  $\mathcal{D}(H_{\mathbf{k}})$  its domain. Then Theorem 2.1 gives that  $\mathcal{D}(H_{\mathbf{k}}) = \mathcal{K}_2^2(\mathbb{T}; \mathcal{S})$  for  $\min_p Z(p) > 3/4$ , and, in general,

$$(13) \quad \mathcal{D}(H_{\mathbf{k}}) \subset \mathcal{K}_{a+1}^2(\mathbb{T}; \mathcal{S}), \quad \text{for } a < \eta := \min_{p \in \mathcal{S}} \sqrt{1/4 + Z(p)} \text{ and } a \leq 1$$

so that  $\mathcal{D}(H_{\mathbf{k}}) \subset \mathcal{K}_1^1(\mathbb{T}; \mathcal{S}) \subset H^1(\mathbb{T})$ , since we assumed that  $\min_p Z(p) > -1/4$ .

We can now state a regularity theorem for the eigenfunctions of  $H_{\mathbf{k}}$  near a point  $p \in \mathcal{S}$ , or equivalently, for Bloch waves associated to the wavevector  $\mathbf{k}$ .

**Theorem 2.2.** *Assume that  $V$  satisfies Assumptions 1 and 2 and let  $u \in \mathcal{D}(H_{\mathbf{k}})$  satisfy  $H_{\mathbf{k}}u = \lambda u$ , for some  $\lambda \in \mathbb{R}$ . Then we can find constants  $a_p \in \mathbb{R}$  such that*

$$u - \sum_{p \in \mathcal{S}} a_p \chi_p \rho^{\sqrt{1/4 + Z(p)} - 1/2} \in \mathcal{K}_{b+1}^{m+1}(\mathbb{T}; \mathcal{S}), \quad \forall b < \min_{p \in \mathcal{S}} \sqrt{9/4 + Z(p)}.$$



In particular,  $u \in \mathcal{K}_{a+1}^{m+1}(\mathbb{T}; \mathcal{S})$ , where  $a < \eta := \min_{p \in \mathcal{S}} \sqrt{1/4 + Z(p)}$  and  $m \in \mathbb{Z}_+$  are arbitrary.

See also [32, 33] for some related classical results in this area. Theorems 2.1 and 2.2 lead to an estimate for the distance from an element in the domain of  $H_{\mathbf{k}}$  to the approximation spaces that we construct using graded meshes. The key issue in the numerical approximation for the problem associated with  $H_{\mathbf{k}}$  is the effectiveness of the algorithm resolving singularities of the form  $\rho^\alpha$ , where  $\alpha > -1/2$  can be negative.

Next, recall the definition of Lagrange interpolants associated to a mesh. Let us choose  $\mathbb{P}$  to be a parallelepiped that is a fundamental domain of the lattice  $\Lambda$ . That is,  $\mathbb{R}^3 = \cup_{y \in \Lambda} (y + \mathbb{P})$  and all  $y + \mathbb{P}$  disjoint. Let  $\mathcal{T} = \{T_i\}$  be a *mesh* on  $\mathbb{P}$ , that is a mesh of  $\mathbb{P}$  with tetrahedra  $T_i$ . We can identify this  $\mathcal{T}$  with a mesh  $\mathcal{T}'$  of the fundamental region of the lattice  $\Lambda$  (that is, to the Brillouin zone of  $\Lambda$ ). Fix an integer  $m \in \mathbb{N}$  that will play the role of the order of approximation. We denote by  $S(\mathcal{T}, m)$  the finite element space associated to the degree  $m$  Lagrange tetrahedron. That is,  $S(\mathcal{T}, m)$  consists of all continuous functions  $\chi : \mathbb{P} \rightarrow \mathbb{R}$  such that  $\chi$  coincides with a polynomial of degree  $\leq m$  on each tetrahedron  $T \in \mathcal{T}$  and  $\chi$  is *periodic*. This means the values of  $\chi$  on opposite faces coincide, so  $\chi$  will have a continuous, periodic extension to the whole space, or alternatively, can be thought of as a continuous function on  $\mathcal{T}$ . We shall denote by  $w_I = w_{I, \mathcal{T}} \in S(\mathcal{T}, m)$  the Lagrange interpolant of  $w \in H^2(\mathbb{P})$ . Let us recall the definition of  $w_{I, \mathcal{T}}$ . First, given a tetrahedron  $T$ , let  $[t_0, t_1, t_2, t_3]$  be the barycentric coordinates on  $T$ . The nodes of the degree  $m$  Lagrange tetrahedron  $T$  are the points of  $T$  whose barycentric coordinates  $[t_0, t_1, t_2, t_3]$  satisfy  $mt_j \in \mathbb{Z}$ . The *degree  $m$  Lagrange interpolant*  $w_{I, \mathcal{T}}$  of  $w$  is the unique function  $w_{I, \mathcal{T}} \in S(\mathcal{T}, m)$  such that  $w = w_{I, \mathcal{T}}$  at the nodes of each tetrahedron  $T \in \mathcal{T}$ . The shorter notation  $w_I$  will be used when only one mesh is understood in the discussion.

The classical approximation result for Lagrange interpolants ([4, 17, 22, 41]) can now be stated.

**Theorem 2.3.** *Let  $\mathcal{T}$  be a mesh of a polyhedral domain  $D \subset \mathbb{R}^3$  with the property that all tetrahedra comprising  $\mathcal{T}$  have (plane and dihedral) angles  $\geq \alpha$  and edges  $\leq h$ . Then there exists a constant  $C(\alpha, m) > 0$  such that, for any  $u \in H^{m+1}(D)$ ,*

$$\|u - u_I\|_{H^1(D)} \leq C(\alpha, m) h^m \|u\|_{H^{m+1}(D)}.$$

Finally, we recall two properties of functions in the weighted Sobolev spaces  $\mathcal{K}_a^m(\mathbb{T}; \mathcal{S})$  that are useful for the analysis of the approximation scheme we use with graded meshes. The proofs of these lemmas are contained in [31] and are based on the definitions and straightforward calculations.

**Lemma 2.4.** *Let  $D$  be a small neighborhood of a point  $p \in \mathcal{S}$  such that on  $D$ ,  $\rho$  is given by distance to  $p$ . Let  $0 < \gamma < 1$  and denote by  $\gamma D$  the region obtained by radially shrinking around  $p$  by a factor of  $\gamma$ . Then*

$$\|w\|_{\mathcal{K}_a^m(D)} = \gamma^{a-3/2} \|w\|_{\mathcal{K}_a^m(\gamma D)}.$$

**Lemma 2.5.** *If  $m \geq m'$ ,  $a \geq a'$  and  $0 < \rho < \delta$  on  $D$ , then*

$$\|w\|_{\mathcal{K}_{a'}^{m'}(D)} \leq \delta^{a-a'} \|w\|_{\mathcal{K}_a^m(D)}.$$

We can now continue to the definition of the mesh refinement technique and the proof of the general approximation theorem underlying our two main theorems.

### 3. APPROXIMATION AND MESH REFINEMENT

Let  $\mathcal{T}$  be a mesh of  $\mathbb{T}$ , such that any point  $p \in \mathcal{S}$  is a node of  $\mathcal{T}$ . Note that the singular expansion of Theorem 2.2 shows that the value of an eigenfunction  $u$  of  $H_{\mathbf{k}}$  at a singular point in  $\mathcal{S}$  may not be defined. Therefore, we define the *modified* degree  $m$  Lagrange “interpolant”  $u_I = u_{I,\mathcal{T}}$  associated to the mesh  $\mathcal{T}$ , such that

$$(14) \quad \begin{cases} u_I(x) = u(x) \text{ for any node } x \notin \mathcal{S} \\ u_I(x) = 0 \text{ if } x \in \mathcal{S}. \end{cases}$$

Alternatively, we can take the modified Lagrange interpolant to be zero on the whole tetrahedron that contains a singular point.

Our two main theorems follow from standard results, such as Céa’s Lemma (for the proof of Theorem 1.2) and the results used in [6, 8, 7, 16, 40] (for the proof of Theorem 1.1), together with the following underlying approximation theorem:

**Theorem 3.1.** *There exists a sequence  $\mathcal{T}_n$  of meshes of  $\mathbb{T}$  that depends only on the choice of a parameter  $k = 2^{-m/a}$ ,  $0 < a < \eta$  and  $a \leq m$ , with the following property. If  $u \in \mathcal{K}_{a+1}^{m+1}(\mathbb{T}; \mathcal{S})$ , then the modified Lagrange interpolant  $u_{I,\mathcal{T}_n} \in S(\mathcal{T}_n, m)$  of  $u$  satisfies*

$$\|u - u_{I,\mathcal{T}_n}\|_{\mathcal{K}_1^1(\mathbb{T}; \mathcal{S})} \leq C \dim(S_n)^{-m/3} \|u\|_{\mathcal{K}_{a+1}^{m+1}(\mathbb{T}; \mathcal{S})},$$

where  $C$  depends only on  $m$  and  $a$  (so it is independent of  $n$  and  $u$ ).

In this section we will define the mesh refinement process and prove Theorem 3.1. The first step is to describe the refinement procedure that results in our sequence of meshes (or triangulations). This is based on the construction in [10] and in [13]. Thus we refer the reader to those papers for details, and here give only an outline and state the critical properties. The second step is to prove a sequence of simple lemmas used in the estimates. The third step is to prove the estimate separately on smaller regions. This uses the scaling properties of the meshes in Lemmas 2.4 and 2.5 together with Theorem 2.3.

**3.1. Construction of the meshes.** We continue to keep the approximation degree  $m$  fixed throughout this section. Fix a parameter  $a$  and let  $k = 2^{-m/a}$ . In our estimates, we will chose  $a$  such that  $a < \eta := \min_p \sqrt{1/4 + Z(p)}$  and  $a \leq m$ . Let  $l$  denote the smallest distance between the points in  $\mathcal{S}$ . Choose an initial mesh  $\mathcal{T}_0$  of  $\mathbb{P}$  with tetrahedra such that all singular points of  $V$  (*i.e.*, all points of  $\mathcal{S}$ ) are among the vertices of  $\mathcal{T}_0$  and no tetrahedron has more than one vertex in  $\mathcal{S}$ . We assume that this mesh is such that if  $F_1$  and  $F_2$  are two opposite faces of  $\mathbb{P}$ , which hence correspond to each other through periodicity, then the resulting triangulations of  $F_1$  and  $F_2$  will also correspond to each other, that is, they are congruent in an obvious sense.

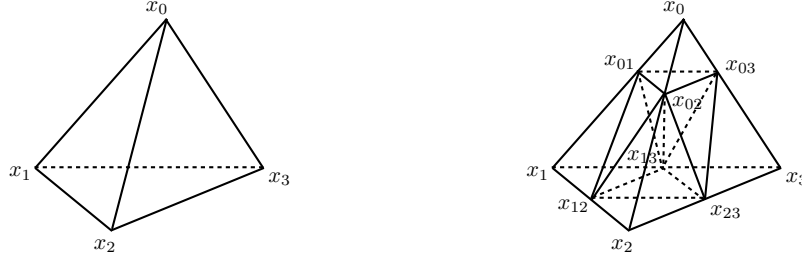


FIGURE 1. The initial tetrahedron  $\{x_0, x_1, x_2, x_3\}$  (left); eight sub-tetrahedra after one  $k$ -refinement (right),  $k = \frac{|x_0x_{01}|}{|x_0x_1|} = \frac{|x_0x_{02}|}{|x_0x_2|} = \frac{|x_0x_{03}|}{|x_0x_3|}$ .

We start with a special refinement of an arbitrary tetrahedron  $T$  that has one of the vertices in the set  $\mathcal{S}$ . Our assumptions then guarantee that all the other vertices of  $T$  will not be in  $\mathcal{S}$ . We use the refinement in [11, 31], which generalizes the 2D refinement introduced in [9] (this refinement was also used in [18, 19, 37]) and we thus introduce the  $k$ -refinement algorithm for a single tetrahedron that divides  $T$  into eight sub-tetrahedra as follows:

**Algorithm 3.2.  $k$ -refinement for a single tetrahedron:** Let  $\{x_0, x_1, x_2, x_3\}$  be the vertices of  $T$ . Suppose that  $x_0 \in \mathcal{S}$ . Therefore,  $x_0$  is the vertex around which a grading ratio  $k \in (0, 1/2]$  will be applied in the next refinement. We first generate new nodes  $x_{ij}$ ,  $0 \leq i < j \leq 3$ , on each edge of  $T$ , such that  $x_{ij} = (x_i + x_j)/2$  for  $1 \leq i < j \leq 3$  and  $x_{0j} = (1 - k)x_0 + kx_j$  for  $1 \leq j \leq 3$ . Note that the node  $x_{ij}$  is on the edge connecting  $x_i$  and  $x_j$ . Connecting these nodes  $x_{ij}$  on all the faces, we obtain 4 sub-tetrahedra and one octahedron. The octahedron then is cut into four tetrahedra using  $x_{13}$  as the common vertex. Therefore, after one refinement, we obtain eight sub-tetrahedra (Figure 1), namely, we obtain the tetrahedra with the following sets of vertices:

$$\begin{aligned} &\{x_0, x_{01}, x_{02}, x_{03}\}, \{x_1, x_{01}, x_{12}, x_{13}\}, \{x_2, x_{02}, x_{12}, x_{23}\}, \{x_3, x_{03}, x_{13}, x_{23}\} \\ &\{x_{01}, x_{02}, x_{03}, x_{13}\}, \{x_{01}, x_{02}, x_{12}, x_{13}\}, \{x_{02}, x_{03}, x_{13}, x_{23}\}, \{x_{02}, x_{12}, x_{13}, x_{23}\}. \end{aligned}$$

**Algorithm 3.3.  $k$ -refinement for a mesh:** Let  $\mathcal{T}$  be a triangulation of the domain  $\mathbb{P}$  such that all points in  $\mathcal{S}$  are among the vertices of  $\mathcal{T}$  and no tetrahedron contains more than one point in  $\mathcal{S}$  among its vertices. Then we divide each tetrahedron  $T$  of  $\mathcal{T}$  that has a vertex in  $\mathcal{S}$  using the  $k$ -refinement and we divide each tetrahedron  $T$  that has no vertices in  $\mathcal{S}$  using the  $1/2$ -refinement. The resulting mesh will be denoted by  $k(\mathcal{T})$ . We then define  $\mathcal{T}_n = k^n(\mathcal{T}_0)$ , where  $\mathcal{T}_0$  is the initial mesh of  $\mathbb{P}$ .

**Remark 3.4.** According to [13], when  $k = 1/2$ , which is the case when the tetrahedron under consideration is away from  $\mathcal{S}$ , the recursive application of Algorithm 3.2 on the tetrahedron generates tetrahedra within at most three similarity classes. On the other hand, if  $k < 1/2$ , the eight sub-tetrahedra of  $T$  are not necessarily similar. Thus, with one  $k$ -refinement, the sub-tetrahedra of  $T$  may belong to at most eight similarity classes. Note that the first sub-tetrahedron in Algorithm 3.2

is similar to the original tetrahedron  $T$  with the vertex  $x_0 \in \mathcal{S}$  and therefore, a further  $k$ -refinement on this sub-tetrahedron will generate eight children tetrahedra within the same eight similarity classes as sub-tetrahedra of  $T$ . Hence, successive  $k$ -refinements of a tetrahedron  $T$  in the initial triangulation  $\mathcal{T}_0$  will generate tetrahedra within at most three similarity classes if  $T$  has no vertex in  $\mathcal{S}$ . On the other hand, successive  $k$ -refinements of a tetrahedron  $T$  in the initial triangulation will generate tetrahedra within at most  $1 + 7 \times 3 = 22$  similarity classes if  $T$  has a point in  $\mathcal{S}$  as a vertex. Thus, our  $k$ -refinement is conforming and yields only non-degenerate tetrahedra, all of which will belong to only finitely many similarity classes.

*Remark 3.5.* Recall that our initial mesh  $\mathcal{T}_0$  has matching restrictions to corresponding faces. Since the singular points in  $\mathcal{S}$  are not on the boundary of  $\mathbb{P}$ , the refinement on opposite boundary faces of  $\mathbb{P}$  is obtained by the usual mid-point decomposition. Therefore, the same matching property will be inherited by  $\mathcal{T}_n$ . In particular, we can extend  $\mathcal{T}_n$  to a mesh in the whole space by periodicity. We will, however, not make use of this periodic mesh on the whole space.

For each point  $p \in \mathcal{S}$  and each  $j$ , we denote by  $\mathcal{V}_{pj}$  the union of all tetrahedra of  $\mathcal{T}_j$  that have  $p$  as a vertex. Thus  $\mathcal{V}_{pj}$  is obtained by scaling the tetrahedra in  $\mathcal{V}_{p0}$  by a factor of  $k^j$  with center  $p$ . In particular, the level  $n \geq j$  refinements of  $\mathcal{T}_0$  give rise to a mesh on  $\mathcal{R}_{pj} := \mathcal{V}_{p(j-1)} \setminus \mathcal{V}_{pj}$ . Define

$$\Omega := \mathbb{P} \setminus \cup_{p \in \mathcal{S}} \mathcal{V}_{p0}.$$

According to Algorithm 3.3,  $\Omega$  and  $\cup_{p \in \mathcal{S}} \mathcal{V}_{p0}$  are triangulated differently. For  $\mathcal{V}_{p0}$ , only the tetrahedra touching  $p$  are refined by the  $k$ -refinement ( $k < 0.5$ ) (Algorithm 3.2) for each refinement, while other tetrahedra are refined by the  $1/2$ -refinement. For  $\Omega$ , we use the  $1/2$ -refinement for each refinement, which is, of course, a uniform refinement. Then, we can decompose  $\mathbb{P}$  as the union

$$(15) \quad \mathbb{P} = \Omega \cup_{p \in \mathcal{S}} \left( \cup_{j=1}^n \mathcal{R}_{pj} \cup \mathcal{V}_{pn} \right),$$

where each set in the union is a union of tetrahedra in  $\mathcal{T}_n$ .

*Remark 3.6.* Note that the size of each simplex of  $\mathcal{T}_n$  contained in  $\Omega$  is  $\mathcal{O}(2^{-n})$ , the size of each simplex of  $\mathcal{T}_n$  contained in  $\mathcal{R}_{pj}$  is  $\mathcal{O}(k^j 2^{-(n-j)})$ , and the size of  $\mathcal{V}_{pn}$  is  $\mathcal{O}(k^n)$ . In addition, the number of tetrahedra in  $\mathcal{T}_n$  is  $\mathcal{O}(2^{3n})$  (see Algorithm 3.3).

Recall equation (10), where the finite element approximation  $v_n \in S(\mathcal{T}_n, m)$  to the equation  $(L + H_{\mathbf{k}})v = f$  is defined. In this case,  $\mathcal{T}_n$  is obtained by applying  $n$  times the  $k$ -refinements to  $\mathcal{T}_0$ , where  $k = 2^{-m/a}$ ,  $0 < a < \eta$ ,  $a \leq m$ , and  $L > C_0$  satisfies Theorem 2.1. Note that (12) and the continuity of  $a(\cdot, \cdot)$  give that the finite element solution  $v_n \in \mathcal{K}_1^1$  is well defined.

**3.2. Proof of Theorem 3.1.** By construction, the restriction of  $\mathcal{T}_n$  to  $\mathcal{R}_{pj}$  scales to the restriction of  $\mathcal{T}_{n-j+1}$  to  $\mathcal{R}_{p1}$ . We denote by  $u_{I,n} = u_{I,\mathcal{T}_n}$  the *modified* interpolation in (14) on  $\mathcal{T}_n$ . The following lemma is based on the definition of the  $k$ -refinement and the discussion in Remark 3.6.

**Lemma 3.7.** *For all  $x \in \mathcal{R}_{pj}$  and a function  $u(x)$ , define by scaling the new function  $\hat{u}(\psi_{-(j-1)}(x)) := u(x)$ , where  $\psi_{-(j-1)}(x) := p + (x-p)/k^{(j-1)}$  is the dilation with ratio  $k^{-(j-1)}$  and center  $p$ . Then,  $u_{I,n}(x) = \widehat{u_{I,n}}(\psi_{-(j-1)}(x)) = \hat{u}_{I,n-j+1}(\psi_{-(j-1)}(x))$ .*

Recall that  $\rho^2 V \in \mathcal{C}^\infty(\mathbb{M}) \cap \mathcal{C}(\mathbb{T})$  and  $\min_p Z(p) > -1/4$ . That is,  $V$  satisfies Assumptions 1 and 2.

We can now give the proof of Theorem 3.1.

*Proof.* Recall that  $\mathcal{V}_{p0}$  consists of the tetrahedra of the initial mesh  $\mathcal{T}_0$  that have  $p$  as a vertex and that all the regions  $\mathcal{V}_p$  are away from each other (they are closed and disjoint). We used this to define  $\Omega := \mathbb{P} \setminus \cup_p \mathcal{V}_{p0}$ . The region  $\mathcal{V}_{pj}$  is obtained by dilating  $\mathcal{V}_p$  with the ratio  $k^j < 1$  and center  $p$ . Finally, recall that  $\mathcal{R}_{pj} = \mathcal{V}_{p(j-1)} \setminus \mathcal{V}_{pj}$ . Let  $R$  be any of the regions  $\Omega$ ,  $\mathcal{R}_{pj}$ , or  $\mathcal{V}_{pn}$ . Since the union of these regions is  $\mathbb{P}$ , it is enough to prove that

$$\|u - u_{I,\mathcal{T}_n}\|_{\mathcal{K}_1^1(R \setminus \mathcal{S})} \leq C \dim(S_n)^{-m/3} \|u\|_{\mathcal{K}_{a+1}^{m+1}(R \setminus \mathcal{S})},$$

for a constant  $C$  independent of  $R$  and  $n$ . The result will follow by squaring all these inequalities and adding them up. In fact, since  $\dim(S_n)^{-m/3} = \mathcal{O}(2^{-nm})$ , it is enough to prove

$$(16) \quad \|u - u_{I,\mathcal{T}_n}\|_{\mathcal{K}_1^1(R \setminus \mathcal{S})} \leq C 2^{-nm} \|u\|_{\mathcal{K}_{a+1}^{m+1}(R \setminus \mathcal{S})},$$

again for a constant  $C$  independent of  $R$  and  $n$ .

If  $R = \Omega := \mathbb{P} \setminus \cup_p \mathcal{V}_{p0}$ , the estimate in (16) follows right away from Theorem 2.3. For the other estimates, recall that  $0 < k = 2^{-m/a}$ , where  $0 < a < \eta$  and  $a \leq m$ . We next establish the desired interpolation estimate on the region  $R = \mathcal{R}_{pj}$ , for any fixed  $p \in \mathcal{S}$  and  $j = 1, 2, \dots, n$ . Let  $\hat{u}(x) = u(\psi_{j-1}(x))$ , where  $\psi_{-(j-1)}(x) := p + (x-p)/k^{(j-1)}$  is the dilation with ratio  $k^{-(j-1)}$  and center  $p$ . From Lemmas 2.4 and 3.7, we have

$$\|u - u_{I,n}\|_{\mathcal{K}_1^1(\mathcal{R}_{pj})} = (k^{j-1})^{1/2} \|\hat{u} - \widehat{u_{I,n}}\|_{\mathcal{K}_1^1(\mathcal{R}_{p1})} = (k^{j-1})^{1/2} \|\hat{u} - \hat{u}_{I,n-j+1}\|_{\mathcal{K}_1^1(\mathcal{R}_{p1})}.$$

Since  $\mathcal{K}_a^m(\mathcal{R}_{p1})$  is equivalent to  $H^m(\mathcal{R}_{p1})$ , we can apply Theorem 2.3 with  $h = \mathcal{O}(2^{-(n-j+1)})$  to get

$$(17) \quad \|u - u_{I,n}\|_{\mathcal{K}_1^1(\mathcal{R}_{pj})} \leq C (k^{j-1})^{1/2} 2^{-m(n-j+1)} \|\hat{u}\|_{\mathcal{K}_{a+1}^{m+1}(\mathcal{R}_{p1})}.$$

Now applying Lemma 2.4 to scale back again and using also  $k = 2^{-m/a}$ , we get that the right hand side in (17)

$$\begin{aligned} C (k^{j-1})^{1/2} 2^{-m(n-j+1)} \|\hat{u}\|_{\mathcal{K}_{a+1}^{m+1}(\mathcal{R}_{p1})} &= C (k^{j-1})^a 2^{-m(n-j+1)} \|u\|_{\mathcal{K}_{a+1}^{m+1}(\mathcal{R}_{pj})} \\ &\leq C 2^{-mn} \|u\|_{\mathcal{K}_{a+1}^{m+1}(\mathcal{R}_{pj})}. \end{aligned}$$

This proves the estimate in (16) for  $R = \mathcal{R}_{pj}$ .

It remains to prove this estimate for  $R = \mathcal{V}_{pn}$ . For any function  $w$  on  $\mathcal{V}_{pn}$ , we let  $\hat{w}(x) = w(\psi_n(x))$  be a function on  $\mathcal{V}_{p0}$ . Therefore, by Lemma 2.4

$$(18) \quad \|u - u_{I,n}\|_{\mathcal{K}_1^1(\mathcal{V}_{pn})} = (k^n)^{1/2} \|\widehat{u - u_{I,n}}\|_{\mathcal{K}_1^1(\mathcal{V}_{p0})},$$

and by Lemma 3.7 (which follows from the definition of the meshes  $\mathcal{T}_k$  and from the fact that interpolation commutes with changes of variables),

$$(19) \quad (k^n)^{1/2} \|u - \widehat{u_{I,n}}\|_{\mathcal{K}_1^1(\mathcal{V}_{p0})} = (k^n)^{1/2} \|\hat{u} - \hat{u}_{I,0}\|_{\mathcal{K}_1^1(\mathcal{V}_{p0})}.$$

Now let  $\chi$  be a smooth cutoff function on  $\mathcal{V}_{p0}$  such that  $\chi = 0$  in a neighborhood of  $p$  and  $= 1$  at every other node of  $\mathcal{V}_{p0}$ .

Define  $\hat{v} := \hat{u} - \chi\hat{u}$ . Then, by (14),

$$(20) \quad \begin{aligned} (k^n)^{1/2} \|\hat{u} - \hat{u}_{I,0}\|_{\mathcal{K}_1^1(\mathcal{V}_{p0})} &= (k^n)^{1/2} \|\hat{v} + \chi\hat{u} - \hat{u}_{I,0}\|_{\mathcal{K}_1^1(\mathcal{V}_{p0})} \\ &\leq (k^n)^{1/2} (\|\hat{v}\|_{\mathcal{K}_1^1(\mathcal{V}_{p0})} + \|\chi\hat{u} - \hat{u}_{I,0}\|_{\mathcal{K}_1^1(\mathcal{V}_{p0})}) \\ &= (k^n)^{1/2} (\|\hat{v}\|_{\mathcal{K}_1^1(\mathcal{V}_{p0})} + \|\chi\hat{u} - (\chi\hat{u})_{I,0}\|_{\mathcal{K}_1^1(\mathcal{V}_{p0})}). \end{aligned}$$

Since  $\chi$  vanishes in the neighborhood of  $p$  we can consider multiplication by  $\chi$  as  $\mathcal{C}^\infty$  times a degree 0 b-operator. Thus it is a bounded operator on any weighted Sobolev space. Thus

$$(21) \quad \begin{aligned} \|\hat{v}\|_{\mathcal{K}_1^1(\mathcal{V}_{p0})} &\leq \|\hat{v}\|_{\mathcal{K}_1^{m+1}(\mathcal{V}_{p0})} \\ &\leq \|\hat{u}\|_{\mathcal{K}_1^{m+1}(\mathcal{V}_{p0})} + \|\chi\hat{u}\|_{\mathcal{K}_1^{m+1}(\mathcal{V}_{p0})} \leq C\|\hat{u}\|_{\mathcal{K}_1^{m+1}(\mathcal{V}_{p0})}, \end{aligned}$$

where  $C$  depends on  $m$  and, through  $\chi$ , the nodes in the triangulation.

Using (18), (19), (20), (21), Lemma 2.5, and Theorem 2.3, we have

$$\begin{aligned} \|u - u_{I,n}\|_{\mathcal{K}_1^1(\mathcal{V}_{pn})} &\leq C(k^n)^{1/2} (\|\hat{u}\|_{\mathcal{K}_1^{m+1}(\mathcal{V}_{p0})} + \|\chi\hat{u} - (\chi\hat{u})_{I,0}\|_{\mathcal{K}_1^1(\mathcal{V}_{p0})}) \\ &\leq C(k^n)^{1/2} (\|\hat{u}\|_{\mathcal{K}_1^{m+1}(\mathcal{V}_{p0})} + \|\chi\hat{u}\|_{H^{m+1}(\mathcal{V}_{p0})}) \\ &\leq C(k^n)^{1/2} (\|\hat{u}\|_{\mathcal{K}_1^{m+1}(\mathcal{V}_{p0})} + \|\hat{u}\|_{\mathcal{K}_1^{m+1}(\mathcal{V}_{p0})}) \\ &\leq C\|u\|_{\mathcal{K}_1^{m+1}(\mathcal{V}_{pn})} \leq Ck^{na}\|u\|_{\mathcal{K}_{a+1}^{m+1}(\mathcal{V}_{pn})} \leq C2^{-mn}\|u\|_{\mathcal{K}_{a+1}^{m+1}(\mathcal{V}_{pn})}. \end{aligned}$$

This proves the estimate of Equation (16) for  $R = V_{pn}$  and completes the proof of Theorem 3.1.  $\square$

*Remark 3.8.* Theorem 3.1 is obtained for the grading parameter  $k = 2^{-m/a}$  satisfying  $0 < a < \eta$  and  $a \leq m$ . We can also always decrease  $k$  because  $\eta$  can be decreased. Going in the opposite direction, that is, increasing  $k$ , will lead to weaker error estimates and convergence rates. For instance, we may find the upper bound of the interpolation error as follows. The estimates on  $\|u - u_{I,n}\|_{\mathcal{K}_1^1(\mathcal{V}_{pn})}$  in the proof above give

$$\|u - u_{I,n}\|_{\mathcal{K}_1^1(\mathcal{V}_{pn})} \leq Ck^{na}\|u\|_{\mathcal{K}_{a+1}^{m+1}(\mathcal{V}_{pn})} \leq C(\dim(S_n))^{-a \log_2(1/k)/3} \|u\|_{\mathcal{K}_{a+1}^{m+1}(\mathcal{V}_{pn})}.$$

Then, examining the estimates on  $\mathcal{R}_{pj}$  and on  $\Omega$ , for  $2^{-m/a} < k \leq 1/2$ , we have the following global upper bound for the error estimate in the case of “insufficient grading”

$$(22) \quad \|u - u_{I,n}\|_{\mathcal{K}_1^1(\mathbb{T};S)} \leq C\dim(S_n)^{-a \log_2(1/k)/3} \|u\|_{\mathcal{K}_{a+1}^{m+1}(\mathbb{T};S)}.$$

Note that insufficient grading still leads to reduction in the error, though with a slower rate than that for meshes with optimal grading parameter. This is numerically notable, especially when  $k$  is close to the upper bound of the optimal range

$2^{-m/a}$ . See Section 5 for a comparison of numerical tests for different values of the grading parameter  $k$  for good and for insufficient grading. The results of those numerical tests seem to be in agreement with Equation (22), but more tests would be needed for a firm confirmation. We will certainly investigate this aspect in the future when more computing power is available to us.

#### 4. APPLICATIONS TO FINITE ELEMENT METHODS

We can now turn to the proofs of the theorems stated in the introduction. First, Theorem 1.1 follows from our general approximation result, Theorem 3.1, and the standard results on approximations of eigenvalues and eigenvectors (eigenfunctions in our case) discussed, for instance, in [6, 8, 7, 16, 40]. Let  $\lambda_j$  be the eigenvalues of  $H_{\mathbf{k}}$  arranged in increasing order and repeated according to their multiplicities. Using the notation introduced in the introduction, we have the following. Let us denote by  $E(\lambda)$  the eigenspace of  $H_{\mathbf{k}}$  corresponding to the eigenvalue  $\lambda$  and by  $E_1(\lambda) \subset E(\lambda)$ , the subset consisting of functions of  $\mathcal{K}_1^1(\mathbb{T}; \mathcal{S})$ -norm one. Then the following result is well known (see for instance Equations (1.1) and (1.2) in [8]). We state it only for our operator  $H_{\mathbf{k}}$ , although it is valid for more general self-adjoint operators with compact resolvent.

Let  $S_n \subset \mathcal{K}_1^1(\mathbb{T}; \mathcal{S})$  be a finite dimensional subspace. Let us denote by  $R_n : \mathcal{K}_1^1(\mathbb{T}; \mathcal{S}) \rightarrow S_n$  the projection in the inner product defined by the bilinear form  $a$  of Equation (6) (the Riesz projection) and by  $\lambda_{j,n}$  the eigenvalues of  $R_n H_{\mathbf{k}} R_n$  arranged in increasing order and repeated according to their multiplicities.

**Theorem 4.1.** *For each  $j$ , there exists a constant  $C_j > 0$  with the following property. Let us denote  $\epsilon_j := \sup_{u \in E_1(\lambda_j)} \inf_{\chi \in S_n} \|u - \chi\|_{\mathcal{K}_1^1(\mathbb{T}; \mathcal{S})}$ . Then*

$$|\lambda_j - \lambda_{j,n}| \leq C_j \epsilon_j^2.$$

*Furthermore, let  $E'_n(\lambda)$  be the sum of eigenspaces  $E_n(\lambda_{j,n})$  of  $R_n H_{\mathbf{k}} R_n$  corresponding to  $\lambda_{j,n}$  with  $\lambda_j = \lambda$ . Then there exists  $w_{j,n} \in E'_n(\lambda_j)$  such that*

$$\|u_j - w_{j,n}\|_{\mathcal{K}_1^1(\mathbb{T}; \mathcal{S})} \leq C_j \epsilon_j.$$

The proof of Theorem 1.1 will then be obtained from Theorem 4.1 as follows.

*Proof.* (of Theorem 1.1). We need to estimate  $\sup_{u \in E_1(\lambda)} \inf_{\chi \in S_n} \|u - \chi\|_{\mathcal{K}_1^1(\mathbb{T}; \mathcal{S})}$ . To this end, let us notice that any  $u \in E(\lambda) \subset \mathcal{K}_1^1(\mathbb{T}; \mathcal{S})$  satisfies  $(L + H_{\mathbf{k}})u = (L + \lambda)u$ . Theorem 2.1 then gives  $\|u\|_{\mathcal{K}_{a+1}^{m+1}(\mathbb{T}; \mathcal{S})} \leq C_{m,\lambda} \|u\|_{\mathcal{K}_{a-1}^{m-1}(\mathbb{T}; \mathcal{S})}$  for a suitably large  $C_{m,\lambda}$  that depends on  $\lambda$  and  $a < \eta$ . A bootstrap argument then gives for any  $u \in E(\lambda)$  that  $\|u\|_{\mathcal{K}_{a+1}^{m+1}(\mathbb{T}; \mathcal{S})} \leq C'_{m,\lambda} \|u\|_{\mathcal{K}_1^1(\mathbb{T}; \mathcal{S})}$ . Theorem 3.1 then gives for  $u \in E_1(\lambda_j)$  (thus  $\|u\|_{\mathcal{K}_1^1(\mathbb{T}; \mathcal{S})} = 1$ ), the following:

$$\begin{aligned} \sup_{u \in E_1(\lambda)} \inf_{\chi \in S_n} \|u - \chi\|_{\mathcal{K}_1^1(\mathbb{T}; \mathcal{S})} &\leq \sup_{u \in E_1(\lambda)} \|u - u_{I, \mathcal{T}_n}\|_{\mathcal{K}_1^1(\mathbb{T}; \mathcal{S})} \\ &\leq C \sup_{u \in E_1(\lambda)} \dim(S_n)^{-m/3} \|u\|_{\mathcal{K}_{a+1}^{m+1}(\mathbb{T}; \mathcal{S})} \leq c(m, \lambda_j) \dim(S_n)^{-m/3}. \end{aligned}$$

The proof of Theorem 1.1 is now complete.  $\square$

Next, the proof of Theorem 1.2 follows from Theorem 3.1, Theorem 2.1, the Lax-Milgram lemma and Céa's lemma. We note some consequences of this theorem.

*Remark 4.2.* First, in the case  $f \in H^{m-1}(\mathbb{T})$ , by the estimate in Equation (11), we have

$$\|v - v_n\|_{\mathcal{K}_1^1(\mathbb{T};\mathcal{S})} \leq C \dim(S_n)^{-m/3} \|f\|_{\mathcal{K}_{a-1}^{m-1}(\mathbb{T};\mathcal{S})} \leq C \dim(S_n)^{-m/3} \|f\|_{H^{m-1}(\mathbb{T})},$$

as long as the index in Theorem 1.2 is chosen such that  $0 < a \leq 1$ .

As in the classical finite element method, a duality argument yields the following  $L^2$ -convergence result.

**Theorem 4.3.** *In addition to the assumptions and notation in Theorem 1.2, assume that  $0 < a \leq 1$ . Then the following  $L^2$  estimate holds*

$$\|v - v_n\|_{L^2(\mathbb{T})} \leq C \dim(S_n)^{(-m-1)/3} \|f\|_{H^{m-1}(\mathbb{T})}.$$

*Proof.* We sketch the proof by using the duality argument in weighted Sobolev spaces. Consider the equation

$$(23) \quad (L + H_{\mathbf{k}})w = v - v_n \quad \text{in } \mathbb{T}.$$

(So we use periodic boundary conditions on  $\mathbb{P}$ .) The definition of the Galerkin projection  $v_n$  of  $v$ , Equation (10), then gives

$$(v - v_n, v - v_n) = ((L + H_{\mathbf{k}})w, v - v_n) = ((L + H_{\mathbf{k}})(w - w_n), v - v_n),$$

where  $w_n$  is the finite element solution of Equation (23) on  $\mathcal{T}_n$ . We also have  $\|w\|_{\mathcal{K}_{a+1}^2(\mathbb{T};\mathcal{S})} \leq C \|v - v_n\|_{L^2(\mathbb{T})}$  by Theorem 2.1, since  $v - v_n \in L^2(\mathbb{T}) \subset \mathcal{K}_{a-1}^0(\mathbb{T};\mathcal{S})$ . Therefore, applying Theorem 1.2 to  $v - v_n \in L^2(\mathbb{T})$  and  $m = 1$ , we have

$$\begin{aligned} \|v - v_n\|_{L^2(\mathbb{T})} &\leq C \|w - w_n\|_{\mathcal{K}_1^1(\mathbb{T};\mathcal{S})} \|v - v_n\|_{\mathcal{K}_1^1(\mathbb{T};\mathcal{S})} / \|v - v_n\|_{L^2(\mathbb{T})} \\ &\leq C \dim(S_n)^{-1/3} \|v - v_n\|_{\mathcal{K}_1^1(\mathbb{T};\mathcal{S})} \leq C \dim(S_n)^{(-m-1)/3} \|f\|_{H^{m-1}(\mathbb{T})}. \end{aligned}$$

This completes the proof.  $\square$

**4.1. The condition number of the stiffness matrix.** It is important that the discrete system that we use is well-conditioned for us to be able to realize the theoretical approximation bounds in practice. Thus we need additionally to obtain upper and lower bounds on the eigenvalues of the stiffness matrix that arises in calculation.

Recall the standard nodal basis function  $\phi_j$  of the space  $S_n := S(\mathcal{T}_n, m)$ . It consists of functions that are equal to 1 at one node and equal to zero at all the other nodes. For convenience, we now instead consider the rescaled bases  $\varphi_j := h_j^{-1/2} \phi_j$ , where  $h_j$  is the diameter of the support patch for  $\phi_j$ . Then, we consider the scaled stiffness matrix

$$(24) \quad A_n := (a(\varphi_i, \varphi_j))$$

from our graded finite element discretization (10). In practice,  $A_n$  can be obtained from the usual stiffness matrix  $(a(\phi_i, \phi_j))$  by a diagonal preconditioning process.



We point out that similar scaled matrices were considered in [12, 36] to study the condition numbers of other Galerkin-based methods.

For a symmetric matrix  $A$ , we shall denote by  $\lambda_{\max}(A)$  the largest eigenvalue of  $A$  and by  $\lambda_{\min}(A)$  the smallest eigenvalue of  $A$ . Thus the spectrum of  $A$  is contained in  $[\lambda_{\min}(A), \lambda_{\max}(A)]$ , but is not contained in any smaller interval. We first prove the following estimates needed below.

**Lemma 4.4.** *Let  $T_i$  be a tetrahedron in the mesh  $\mathcal{T}_n$  and let  $\text{diam}(T_i)$  denote the diameter of  $T_i$ . Then, for any  $\psi_n \in S_n$  and  $\psi \in H^1(\mathbb{T})$ , there exists a constant  $C > 0$  independent of  $n$ ,  $\psi_n$  and  $\psi$ , such that*

$$(25) \quad \|\psi_n\|_{\mathcal{K}_1^1(T_i)} \leq C \text{diam}(T_i)^{1/2} \|\psi_n\|_{L^\infty(T_i)} \leq C \|\psi_n\|_{L^6(T_i)},$$

$$(26) \quad \|\psi\|_{L^6(\mathbb{T})} \leq C \|\psi\|_{H^1(\mathbb{T})}.$$

Furthermore, writing  $\psi_n = \sum c_j \varphi_j$ , where  $\varphi_j := h_j^{-1/2} \phi_j$  are the rescaled basis functions, we get

$$(27) \quad C^{-1} \sum_{j \in \text{node}(T_i)} c_j^2 \leq \text{diam}(T_i) \|\psi_n\|_{L^\infty(T_i)}^2 \leq C \sum_{j \in \text{node}(T_i)} c_j^2.$$

*Proof.* We shall show (25) and (27) since (26) is a particular case of the well known Sobolev embedding theorem, see [28] for example.

To prove the second estimate in (25), let us first recall that all the tetrahedra  $T_i$  belong to a finite class of shapes (or similarity classes) in our graded triangulation. The second estimate in (25) is then a direct consequence of the scaling argument in [17, 21]. Recall that this scaling argument, to be used also below, is to map an arbitrary tetrahedron  $T_i$  to a standard tetrahedron  $T_{ref}$  and then to use the equivalence of norms on finite dimensional spaces. The resulting constant  $C$  will then of course depend on the shape regularity of the mesh. We now turn to the proof of the first estimate in (25).

By the definition of the weighted space (8) and the usual scaling argument, we first have

$$(28) \quad |\psi_n|_{\mathcal{K}_1^1(T_i)} \leq \|\psi_n\|_{H^1(T_i)} \leq C \text{diam}(T_i)^{1/2} \|\psi_n\|_{L^\infty(T_i)}.$$

We shall consider the two possibilities when one of the vertices of  $T_i$ , call it  $Q$ , is in  $\mathcal{S}$  and when none of the vertices of  $T_i$  is in  $\mathcal{S}$ .

First, if the vertex  $Q$  of  $T_i$  is in  $\mathcal{S}$ , we use a new coordinate system, which is the translation of the old coordinate system, such that  $Q$  is the origin. In the new coordinate system, denote by  $T_\gamma := \{\gamma x, \forall x \in T_i\}$  the dilated tetrahedron with the constant  $\gamma = \text{diam}(T_i)^{-1}$ . Therefore  $\text{diam}(T_\gamma) \simeq 1$ . For a function  $v$  on  $T_i$ , we define for  $x \in T_\gamma$ ,  $v^\gamma(x) := v(\gamma^{-1}x)$ . Recall  $\rho$  in (7) is the distance to  $Q$  on  $T_i$  and therefore  $\rho(\gamma x) = \gamma \rho(x)$  for  $x \in T_i$ . Using the scaling argument and norm equivalence on finite dimensional spaces, we have

$$(29) \quad \begin{aligned} \|\psi_n\|_{\mathcal{K}_1^1(T_i)}^2 &= \int_{T_i} \rho^{-2}(x) \psi_n^2(x) dx = \int_{T_\gamma} \gamma^2 \rho^{-2}(\gamma x) [\psi_n^\gamma(\gamma x)]^2 \gamma^{-3} d(\gamma x) \\ &\leq C \gamma^{-1} \|\psi_n^\gamma\|_{L^\infty(T_\gamma)}^2 \leq \text{diam}(T_i) \|\psi_n\|_{L^\infty(T_i)}^2. \end{aligned}$$

On the other hand, if no vertex of  $T_i$  belongs to  $\mathcal{S}$ , the construction of our graded meshes shows that

$$\|\rho^{-1}(x)\|_{L^\infty(T_i)} \leq C \text{diam}(T_i)^{-1}.$$

Combining this with the above standard scaling argument, we have

$$(30) \quad \|\psi_n\|_{\mathcal{K}_1^0(T_i)} \leq C \text{diam}(T_i)^{-1} \|\psi_n\|_{L^2(T_i)} \leq C \text{diam}(T_i)^{1/2} \|\psi_n\|_{L^\infty(T_i)}.$$

Combining Equations (28), (29), and (30) completes the proof for the first estimate in (25) since the diameters  $\text{diam}(T_i)$  are bounded.

For the estimate in (27), let  $\hat{T}$  be the usual reference tetrahedron and  $F_i$  be the affine mapping such that  $F_i(T_i) = \hat{T}$ . For any function  $v$  on  $T_i$ , we denote by  $\hat{v} := v \circ F_i^{-1}$  the resulting function on  $\hat{T}$ . Let us also denote by  $\psi_n = \sum c_i \varphi_i = \sum \bar{c}_i \phi_i$ . Based on the definition of the basis function  $\varphi_i$ ,

$$(31) \quad C^{-1} \text{diam}(T_i)^{1/2} \bar{c}_i \leq c_i \leq C \text{diam}(T_i)^{1/2} \bar{c}_i.$$

Then, both  $\|\hat{\psi}_n\|_{L^\infty}$  and  $(\sum_{j \in \text{node}(\hat{T})} \bar{c}_j^2)^{1/2}$  are norms for the finite element function  $\hat{\psi}_n|_{\hat{T}}$ , where the summation on  $\bar{c}_j$  is for all the nodes in  $\hat{T}$ . Based on equivalence of all norms for a finite dimensional space, we have

$$C^{-1} \left( \sum_{j \in \text{node}(\hat{T})} \bar{c}_j^2 \right)^{1/2} \leq \|\hat{\psi}_n\|_{L^\infty(\hat{T})} \leq C \left( \sum_{j \in \text{node}(\hat{T})} \bar{c}_j^2 \right)^{1/2}.$$

This, together with (31), implies

$$C^{-1} \sum_{j \in \text{node}(T_i)} c_j^2 \leq \text{diam}(T_i) \|\psi_n\|_{L^\infty(T_i)}^2 \leq C \sum_{j \in \text{node}(T_i)} c_j^2,$$

which completes the proof.  $\square$

Therefore, we have the following estimates on the eigenvalues of the stiffness matrix.

**Lemma 4.5.** *Let  $A_n$  be the stiffness matrix from the finite element discretization corresponding to the rescaled nodal basis  $\varphi_j$  of the space  $S_n := S(\mathcal{T}_n, m)$  in Equation (24). Then,*

$$\lambda_{\max}(A_n) \leq M,$$

where the constant  $M$  is independent of the mesh level  $n$ .

*Proof.* Let us fix the mesh level  $n$ . All the constants below will be independent of  $n$ . Let  $\{T_i\}$  be the tetrahedra forming our mesh  $\mathcal{T}_n$ . Let  $\psi_n \in S_n$  be arbitrary and write  $\psi_n = \sum_j c_j \varphi_j$  and  $\mathbf{V} := (c_j)$ . By (25) and (27), we have

$$\begin{aligned} \mathbf{V}^T A_n \mathbf{V} &= a(\psi_n, \psi_n) \leq C \|\psi_n\|_{\mathcal{K}_1^1(\mathbb{T})}^2 = C \sum_i \|\psi_n\|_{\mathcal{K}_1^1(T_i)}^2 \\ &\leq C \sum_i \text{diam}(T_i) \|\psi_n\|_{L^\infty(T_i)}^2 \leq C \sum_j c_j^2 \leq C \mathbf{V}^T \mathbf{V}. \end{aligned}$$

This completes the proof.  $\square$

**Lemma 4.6.** *We use the same notation as in Lemma 4.5. Then smallest eigenvalue  $\lambda_{\min}(A_n)$  of the stiffness matrix  $A_n$  satisfies*

$$\lambda_{\min}(A_n) \geq C \dim(S_n)^{-2/3}.$$

*Proof.* For any  $\psi_n \in S_n$ , we use the notation  $\psi_n = \sum_j c_j \varphi_j$ ,  $\mathbf{V} := (c_j)$ , and  $\text{diam}(T_i)$  denotes the diameter of  $T_i$ , as in the proof of the previous lemma. In view of (27), the inverse estimate (25), Hölder's inequality, and the Sobolev embedding estimate (26), we then have

$$\begin{aligned} \mathbf{V}^T \mathbf{V} &= \sum_j c_j^2 \leq C \sum_i \text{diam}(T_i) \|\psi_n\|_{L^\infty(T_i)}^2 \leq C \sum_i \|\psi_n\|_{L^6(T_i)}^2 \\ &\leq C \left( \sum_i 1 \right)^{\frac{2}{3}} \left( \sum_i \|\psi_n\|_{L^6(T_i)}^6 \right)^{\frac{1}{3}} \leq C \dim(S_n)^{\frac{2}{3}} \|\psi_n\|_{L^6(\mathbb{T})}^2 \\ &\leq C \dim(S_n)^{\frac{2}{3}} \|\psi_n\|_{H^1(\mathbb{T})}^2 \leq C \dim(S_n)^{\frac{2}{3}} \mathbf{V}^T A_n \mathbf{V}. \end{aligned}$$

□

Then, we have the estimate on the condition number.

**Theorem 4.7.** *Let  $A_n = (a(\varphi_i, \varphi_j))$  be the stiffness matrix. Then the condition number  $\kappa(A_n)$  of  $A_n$  satisfies*

$$\kappa(A_n) \leq C \dim(S_n)^{2/3}.$$

*The constant  $C$  depends on the finite element space, but not on  $n$ .*

*Proof.* Using  $\kappa(A_n) = \lambda_{\max}(A_n)/\lambda_{\min}(A_n)$ , we obtain the estimate by Lemmas 4.5 and 4.6. □

*Remark 4.8.* Similar estimates on condition numbers have been derived by Bank and Scott [12]. We also mention that Apel and Heinrich [1] studied the condition number from 3D graded meshes for edge singularities. They also recommended the scaling of basis functions to precondition the matrix when the solution possesses severe edge singularities.

## 5. NUMERICAL TESTS OF THE FINITE ELEMENT METHOD

We now present the numerical tests for the finite element solution defined in (10) approximating possibly singular solutions to Equation (9).

To be more precise, suppose that our periodicity lattice is  $2\mathbb{Z}^3$  and we choose our fundamental domain  $\mathbb{P} = [-1, 1]^3$  to be a cube of side length 2. We impose periodic boundary condition on the following model problem

$$(32) \quad (L + H_0)v := (-\Delta + \delta\psi r^{-2} + L)v = 1 \quad \text{in } \Omega,$$

where  $r = |x|$ ,  $\delta > -1/4$ ,  $L \geq 0$ , and the cut-off function  $\psi := e^{r_c^2/(r^4 - r_c^2) + 1}$  for  $r^2 \leq r_c$  and  $\psi = 0$  for  $r^2 > r_c$ ; in the tests, we chose  $r_c = 0.25$ . Note that if  $\delta > 0$ , it is clear that the operator  $L + H_0$  is positive on  $\mathcal{K}_1^1$  (see Theorem 2.1). We use the continuous piecewise linear finite element method on triangulations graded

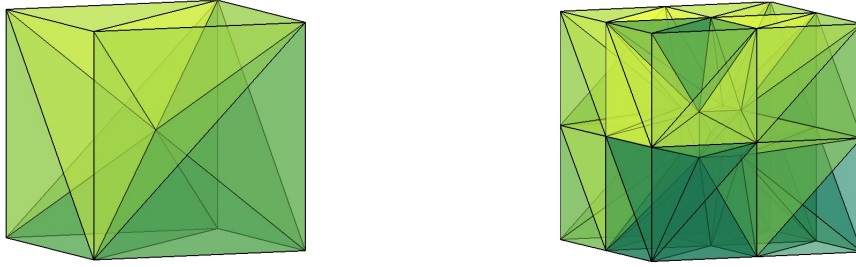


FIGURE 2. The initial mesh on the cube  $\mathbb{P} = [-1, 1]^3$  (left); the mesh after one  $k$  refinement for the origin,  $k = 0.2$  (right).

toward the origin with grading ratio  $k > 0$  (Recall that  $k = 0.5$  corresponds to the quasi-uniform refinement.)

To enforce the periodic boundary condition for the finite element functions, we use meshes where all the boundary nodal points are symmetric about the mid-plane between opposite faces of the cube. Any set of the symmetric nodes will be associated to the same shape function in the discretization. For example, nodes on edges of the cube generally have three mirror images over two mid-planes (two direct mirror images and the third is symmetric over the line of intersection of these two mid-planes), and these four points are associated to the same shape function. Consequently, the eight vertices of the cube are associated to the same shape function through symmetry. See Figure 2 for example. In particular, a mesh on  $\mathbb{T}$  identifies with a mesh with suitable properties on  $\mathbb{P}$ .

Our first tests are for Equation (32) with  $\delta = 4.0$  and  $L = 0$ . According to Theorem 1.2 in the case  $m = 1$ , the optimal rate of convergence for the finite element solution,  $\dim(S_n)^{-1/3}$ , should be obtained on triangulations with any  $k \leq 0.5$ , since  $\eta = \sqrt{1/4 + 4} > 1$ . The convergence rates  $e$  associated to triangulations with different values of  $k$  are listed in Table 1. Starting from an initial triangulation, we compute the rates based on the comparison of the numerical errors on triangulations with consecutive  $k$ -refinements,

$$(33) \quad e := \log_2 \frac{|v_{j-1} - v_j|_{\mathcal{K}_1^1}}{|v_j - v_{j+1}|_{\mathcal{K}_1^1}},$$

where  $v_j$  is the finite element solution on the mesh after  $j$   $k$ -refinements. Recall the dimension of the finite element space grows by a factor of 8 with one  $k$ -refinement. Thus, by Theorem 1.2, for a sequence of optimal meshes, the error  $|v - v_j|_{\mathcal{K}_1^1}$  is reduced by a factor of 2 for linear finite element approximations with each  $k$ -refinement. Thus,  $e \rightarrow 1$  implies that the optimal rate of convergence in Theorem 1.2 is achieved.

Table 1 shows that the convergence rates  $e$  approach 1 for all values of the grading parameter  $k$ . This is in agreement with our theory that the optimal rates of convergence are obtained for any triangulations with  $k \leq 0.5$ . Note that  $k = 0.5$  corresponds to a standard ungraded mesh, and the result for this value of  $k$  is the standard convergence result for such meshes. This result can be recovered in this

$j \backslash e$	$k = 0.1$	$k = 0.2$	$k = 0.3$	$k = 0.4$	$k = 0.5$
2	0.42	0.44	0.56	0.33	-0.20
3	0.48	0.68	0.75	0.79	0.70
4	0.78	0.81	0.86	0.88	0.85
5	0.91	0.92	0.94	0.95	0.93
6	0.97	0.97	0.98	0.99	0.98

TABLE 1. Convergence rates  $e$  of finite element solutions solving equation (32) with  $\delta = 4.0$  and  $L = 0$  on different graded tetrahedra.

$j \backslash e$	$k = 0.1$	$k = 0.2$	$k = 0.3$	$k = 0.4$	$k = 0.5$
2	0.20	0.30	0.33	0.11	-0.03
3	0.54	0.66	0.69	0.61	0.39
4	0.74	0.81	0.83	0.77	0.60
5	0.88	0.91	0.92	0.87	0.72
6	0.95	0.97	0.98	0.92	0.79

TABLE 2. Convergence rates  $e$  of finite element solutions solving equation (32) with  $\delta = 0.6$  and  $L = 0$  on different graded tetrahedra.

example since the singularity in the solution is not strong enough to be detectable for linear finite elements. That is, the solution has sufficient regularity in terms of regular Sobolev spaces for the standard mesh result to hold for linear elements. However, a graded mesh would be necessary to obtain the optimal convergence rate for elements of higher degree.

In the second test, we implemented our method solving equation (32) with  $\delta = 0.6$ ,  $L = 0$  and summarize the results in Table 2. Based on the upper bound  $\eta = \sqrt{1/4 + 0.6}$  given in Theorem 1.2, we expect the optimal rate of convergence for the numerical solution as long as the grading parameter  $k < 2^{-1/\eta} \approx 0.47$ . The convergence rates in Table 2 tend to 1 when  $k \leq 0.4$ , which implies the optimality of our finite element approximation on these meshes. However, when  $k = 0.5$ , the convergence rate is far less than 1 and there is a large gap between the rates corresponding to  $k = 0.4$  and  $k = 0.5$ . This further confirms our theory that the upper bound of the suitable range of  $k$  for an optimal finite element approximation lies in  $(0.4, 0.5)$ .

The third tests are for negative potentials in equation (32), where we set  $\delta = -0.1$  and  $L = 20$  to satisfy the positivity requirement in Theorem 2.1. Our theoretical results indicate that the singularity in the solution due to the singular potential is stronger in this case and the optimal rate can be achieved only if the grading parameter  $k < 2^{-1/\sqrt{1/4-0.1}} \approx 0.167$ . Because of the limitation of the computation power, we only display the convergence results up to the 7th refinement for various graded parameters  $k$  in Table 3. However, we still see the trend that appropriate gradings improve the convergence rate as predicted in Theorem 1.2. When  $k$  is close to the optimal value 0.167 (i.e.,  $k = 0.1$  and  $0.2$ ), we have remarkable improvements.

$j \backslash e$	$k = 0.1$	$k = 0.2$	$k = 0.3$	$k = 0.4$	$k = 0.5$
2	-0.10	-0.05	-0.09	-0.16	-0.03
3	0.32	0.37	0.30	0.19	0.07
4	0.51	0.52	0.44	0.32	0.18
5	0.67	0.64	0.53	0.40	0.26
6	0.80	0.72	0.59	0.45	0.32

TABLE 3. Convergence rates  $e$  of finite element solutions solving equation (32) with  $\delta = -0.1$  and  $L = 20$  on different graded tetrahedra.

In particular, for  $k = 0.1$ , based on Table 3, we expect that the optimal rate occurs with further refinements.

We have also implemented the method on graded meshes for the eigenvalue problem associated with equation (32), especially on the computation of the first eigenvalues. Namely,

$$H_0 u := (-\Delta + \delta \psi r^{-2})u = \lambda_1 u$$

on the cube  $\mathbb{P} = [-1, 1]^3$ , where  $\lambda_1$  is the first eigenvalue of the operator. Depending on the choice of  $\delta$ , the convergence rates for the numerical eigenvalues on graded meshes are roughly twice the rates for the numerical solutions of equation (32) (see Tables 1, 2, and 3), and present similar trends for different gradings.

All our numerical tests (Tables 1, 2, 3, and corresponding eigenvalue computations) verify Theorem 1.1 by comparing the rates of convergence for different singular potentials on different graded triangulations for the model operator in (32). The theoretical upper bounds  $2^{-1/\eta}$  of the optimal range for the grading parameter  $k$  are also demonstrated in these numerical results. In these tests, the initial triangulation of the cube consists of 12 tetrahedra and we consecutively refine the mesh using the  $k$ -refinements up to level 7 that includes  $12 \times 8^7 \approx 2.5 \times 10^7$  tetrahedra and roughly 4.2 million unknowns. Numerical experiments show that the condition numbers of our discrete systems grow by a factor of 4 for consecutive refinements, regardless of the value of  $k$ , which resembles the estimates given in [12] for the Laplace operator. However, the values of  $k$  affect the magnitude of the condition numbers. In general, smaller  $k$  leads to bad shapes for the tetrahedra and therefore results in larger condition numbers. The preconditioned conjugate gradient (PCG) method (using the inverse of the diagonal entries as the preconditioner) was used as the numerical solver for the discrete systems.

## REFERENCES

- [1] Th. Apel and B. Heinrich. Mesh refinement and windowing near edges for some elliptic problem. *SIAM J. Numer. Anal.*, 31(3):695–708, 1994.
- [2] Thomas Apel, Serge Nicaise, and Joachim Schöberl. Crouzeix-Raviart type finite elements on anisotropic meshes. *Numer. Math.*, 89(2):193–223, 2001.
- [3] D. Arroyo, A. Bespalov, and N. Heuer. On the finite element method for elliptic problems with degenerate and singular coefficients. *Math. Comp.*, 76(258):509–537 (electronic), 2007.

- [4] I. Babuška and A. K. Aziz. Survey lectures on the mathematical foundations of the finite element method. In *The mathematical foundations of the finite element method with applications to partial differential equations (Proc. Sympos., Univ. Maryland, Baltimore, Md., 1972)*, pages 1–359. Academic Press, New York, 1972. With the collaboration of G. Fix and R. B. Kellogg.
- [5] I. Babuška, R. B. Kellogg, and J. Pitkäranta. Direct and inverse error estimates for finite elements with mesh refinements. *Numer. Math.*, 33(4):447–471, 1979.
- [6] I. Babuška and J. Osborn. Eigenvalue problems. In *Handbook of numerical analysis, Vol. II*, Handb. Numer. Anal., II, pages 641–787. North-Holland, Amsterdam, 1991.
- [7] I. Babuška and J. E. Osborn. Estimates for the errors in eigenvalue and eigenvector approximation by Galerkin methods, with particular attention to the case of multiple eigenvalues. *SIAM J. Numer. Anal.*, 24(6):1249–1276, 1987.
- [8] I. Babuška and J. E. Osborn. Finite element-Galerkin approximation of the eigenvalues and eigenvectors of selfadjoint problems. *Math. Comp.*, 52(186):275–297, 1989.
- [9] C. Bacuta, V. Nistor, and L. Zikatanov. Improving the rate of convergence of ‘high order finite elements’ on polygons and domains with cusps. *Numerische Mathematik*, 100:165–184, 2005.
- [10] C. Bacuta, V. Nistor, and L. Zikatanov. Improving the rate of convergence of high-order finite elements on polyhedra. I. A priori estimates. *Numer. Funct. Anal. Optim.*, 26(6):613–639, 2005.
- [11] C. Bacuta, V. Nistor, and L. Zikatanov. Improving the rate of convergence of high-order finite elements on polyhedra. II. Mesh refinements and interpolation. *Numer. Funct. Anal. Optim.*, 28(7-8):775–824, 2007.
- [12] R.E. Bank and L.R. Scott. On the conditioning of finite element equations with highly refined meshes. *SIAM J. Numer. Anal.*, 26(6):1383–1394, 1989.
- [13] J. Bey. Tetrahedral grid refinement. *Computing*, 55(4):355–378, 1995.
- [14] S. Bidwell, M. E. Hassell, and C. R. Westphal. A weighted least squares finite element method for elliptic problems with degenerate and singular coefficients. *Math. Comp.*, 82(282):673–688, 2013.
- [15] V. Bonnaillie-Noël and M. Dauge. Asymptotics for the low-lying eigenstates of the Schrödinger operator with magnetic field near corners. *Ann. Henri Poincaré*, 7(5):899–931, 2006.
- [16] J. H. Bramble and J. E. Osborn. Rate of convergence estimates for nonselfadjoint eigenvalue approximations. *Math. Comp.*, 27:525–549, 1973.
- [17] S. Brenner and R. Scott. *The mathematical theory of finite element methods*, volume 15 of *Texts in Applied Mathematics*. Springer-Verlag, New York, second edition, 2002.
- [18] S.C. Brenner, J. Cui, T. Gudi, and L.-Y. Sung. Multigrid algorithms for symmetric discontinuous Galerkin methods on graded meshes. *Numer. Math.*, 119(1):21–47, 2011.
- [19] S.C. Brenner, J. Cui, and L.-Y. Sung. Multigrid methods for the symmetric interior penalty method on graded meshes. *Numer. Linear Algebra Appl.*, 16(6):481–501, 2009.
- [20] C. Băcuță, J.H. Bramble, and J. Xu. Regularity estimates for elliptic boundary value problems in Besov spaces. *Math. Comp.*, 72:1577–1595, 2003.
- [21] P. Ciarlet. *The Finite Element Method for Elliptic Problems*, volume 4 of *Studies in Mathematics and Its Applications*. North-Holland, Amsterdam, 1978.
- [22] P. Ciarlet. Basic error estimates for elliptic problems. In *Handbook of numerical analysis, Vol. II*, Handb. Numer. Anal., II, pages 17–352. North-Holland, Amsterdam, 1991.
- [23] M. Costabel, M. Dauge, and S. Nicaise. Corner singularities of Maxwell interface and eddy current problems. In *Operator theoretical methods and applications to mathematical physics*, volume 147 of *Oper. Theory Adv. Appl.*, pages 241–256. Birkhäuser, Basel, 2004.
- [24] L. Demkowicz, P. Monk, Ch. Schwab, and L. Vardapetyan. Maxwell eigenvalues and discrete compactness in two dimensions. *Comput. Math. Appl.*, 40(4-5):589–605, 2000.

- [25] M.S.P. Eastham. *The Spectral Theory of Periodic Differential Equations*. Scottish Academic Press, Edinburgh, 1973.
- [26] V. Felli, A. Ferrero, and S. Terracini. Asymptotic behavior of solutions to Schrödinger equations near an isolated singularity of the electromagnetic potential. *J. Eur. Math. Soc. (JEMS)*, 13(1):119–174, 2011.
- [27] V. Felli, E. Marchini, and S. Terracini. On the behavior of solutions to Schrödinger equations with dipole type potentials near the singularity. *Discrete Contin. Dyn. Syst.*, 21(1):91–119, 2008.
- [28] D. Gilbarg and N.S. Trudinger. *Elliptic partial differential equations of second order*. Springer-Verlag, Berlin, 1977. Grundlehren der Mathematischen Wissenschaften, Vol. 224.
- [29] X. Gong, L. Shen, D. Zhang, and A. Zhou. Finite element approximations for Schrödinger equations with applications to electronic structure computations. *J. Comput. Math.*, 26(3):310–323, 2008.
- [30] E. Hunsicker, H. Li, V. Nistor, and V. Uski. Analysis of Schrödinger operators with inverse square potentials I: regularity results in 3D. *Bull. Math. Soc. Sci. Math. Roumanie (N.S.)*, 55(103)(2):157–178, 2012.
- [31] E. Hunsicker, V. Nistor, and J. Sofo. Analysis of periodic Schrödinger operators: regularity and approximation of eigenfunctions. *J. Math. Phys.*, 49(8):083501, 21, 2008.
- [32] T. Kato. Fundamental properties of Hamiltonian operators of Schrödinger type. *Trans. Amer. Math. Soc.*, 70:195–211, 1951.
- [33] T. Kato. On the eigenfunctions of many-particle systems in quantum mechanics. *Comm. Pure Appl. Math.*, 10:151–177, 1957.
- [34] V. A. Kondrat'ev. Boundary value problems for elliptic equations in domains with conical or angular points. *Transl. Moscow Math. Soc.*, 16:227–313, 1967.
- [35] H. Li. A-priori analysis and the finite element method for a class of degenerate elliptic equations. *Math. Comp.*, 78:713–737, 2009.
- [36] H. Li. A note on the conditioning of a class of generalized finite element methods. *Appl. Numer. Math.*, 62:754–766, 2012.
- [37] H. Li, A. Mazzucato, and V. Nistor. Analysis of the finite element method for transmission/mixed boundary value problems on general polygonal domains. *Electron. Trans. Numer. Anal.*, 37:41–69, 2010.
- [38] H. Li and V. Nistor. Analysis of a modified Schrödinger operator in 2D: regularity, index, and FEM. *J. Comput. Appl. Math.*, 224(1):320–338, 2009.
- [39] S. Moroz and R. Schmidt. Nonrelativistic inverse square potential, scale anomaly, and complex extension. *Ann. Phys.*, 325(2):491–513, 2010.
- [40] J. Osborn. Spectral approximation for compact operators. *Math. Comput.*, 29:712–725, 1975.
- [41] C. Schwab. *P- And Hp- Finite Element Methods: Theory and Applications in Solid and Fluid Mechanics*. 1999.
- [42] L. Wahlbin. On the sharpness of certain local estimates for  $\mathring{H}^1$  projections into finite element spaces: influence of a re-entrant corner. *Math. Comp.*, 42(165):1–8, 1984.
- [43] H. Wu and D.W.L Sprung. Inverse-square potential and the quantum vortex. *Physical Review A*, 49:4305–4311, 1994.



EUGENIE HUNSICKER, DEPARTMENT OF MATHEMATICAL SCIENCES, LOUGHBOROUGH UNIVERSITY, LOUGHBOROUGH, LEICESTERSHIRE, LE11 3TU, UK

*E-mail address:* `E.Hunsicker@lboro.ac.uk`

HENGGUANG LI, DEPARTMENT OF MATHEMATICS, WAYNE STATE UNIVERSITY, DETROIT, MI 48202, USA

*E-mail address:* `hli@math.wayne.edu`

V. NISTOR, PENNSYLVANIA STATE UNIVERSITY, MATH. DEPT., UNIVERSITY PARK, PA 16802, USA, AND INST. MATH. ROMANIAN ACAD. PO BOX 1-764, 014700 BUCHAREST ROMANIA

*E-mail address:* `nistor@math.psu.edu`

VILLE USKI, DEPARTMENT OF MATHEMATICAL SCIENCES, LOUGHBOROUGH UNIVERSITY, LOUGHBOROUGH, LEICESTERSHIRE, LE11 3TU, UK

*E-mail address:* `V.Uski@lboro.ac.uk`