

THESIS FOR THE DEGREE OF DOCTOR OF PHILOSOPHY

# **The Dirac Equation: Numerical and Asymptotic Analysis**

**Hasan Almasreh**

**CHALMERS** |  **UNIVERSITY OF GOTHENBURG**

Department of Mathematical Sciences  
Chalmers University of Technology and University of Gothenburg  
Gothenburg, Sweden 2012

The Dirac Equation: Numerical and Asymptotic Analysis

Hasan Almasreh

ISBN 978-91-628-8593-9

©Hasan Almasreh, 2012

Division of Mathematics–Physics Platform ( $\mathcal{MP}^2$ )

Department of Mathematical Sciences

Chalmers University of Technology and University of Gothenburg

SE - 412 96 Gothenburg

Sweden

Telephone: + 46 (0)31 772 1000

Printed at the Department of Mathematical Sciences

Gothenburg, Sweden 2012

# The Dirac Equation: Numerical and Asymptotic Analysis

**Hasan Almasreh**

*Department of Mathematical Sciences  
Chalmers University of Technology  
and University of Gothenburg*

## Abstract

The thesis consists of three parts, although each part belongs to a specific subject area in mathematics, they are considered as subfields of the perturbation theory. The main objective of the presented work is the study of the Dirac operator; the first part concerns the treatment of the spurious eigenvalues in the computation of the discrete spectrum. The second part considers G-convergence theory for positive definite parts of a family of Dirac operators and general positive definite self-adjoint operators. The third part discusses the convergence of wave operators for some families of Dirac operators and for general self-adjoint operators.

In the first and main part, a stable numerical scheme, using finite element and Galerkin-based  $hp$ -cloud methods, is developed to remove the spurious eigenvalues from the computational solution of the Dirac eigenvalue problem. The scheme is based on applying a Petrov-Galerkin formulation to introduce artificial diffusivity to stabilize the solution. The added diffusion terms are controlled by a stability parameter which is derived for the particular problem. The derivation of the stability parameter is the main part of the scheme, it is obtained for specific basis functions in the finite element method and then generalized for any set of admissible basis functions in the  $hp$ -cloud method.

In the second part, G-convergence theory is applied to positive definite parts of the Dirac operator perturbed by  $h$ -dependent abstract potentials, where  $h$  is a parameter allowed to grow to infinity. After shifting the perturbed Dirac operator so that the point spectrum is positive definite, the spectral measure is used to obtain projected positive definite parts of the operator, in particular the part that is restricted to the point spectrum. Using the general definition of G-convergence, G-limits, as  $h$  approaches infinity, are proved for these projected parts under suitable conditions on the perturbations. Moreover, G-convergence theory is also discussed for some positive definite self-adjoint  $h$ -dependent operators. The purpose of applying G-convergence is to study the asymptotic behavior of the corresponding eigenvalue problems. In this regard, the eigenvalue

problems for the considered operators are shown to converge, as  $h$  approaches infinity, to the eigenvalue problems of their associated G-limits.

In the third part, scattering theory is studied for the Dirac operator and general self-adjoint operators with classes of  $h$ -dependent perturbations. For the Dirac operator with different power-like decay  $h$ -dependent potentials, the wave operators exist and are complete. In our study, strong convergence, as  $h$  approaches infinity, of these wave operators is proved and their strong limits are characterized for specific potentials. For general self-adjoint operators, the stationary approach of scattering theory is employed to study the existence and convergence of the stationary and time-dependent  $h$ -dependent wave operators.

**Keywords:** Dirac operator, eigenvalue problem, finite element method, spurious eigenvalues, Petrov-Galerkin, cubic Hermite basis functions, stability parameter, meshfree method,  $hp$ -cloud, intrinsic enrichment, G-convergence,  $\Gamma$ -convergence, scattering theory, identification, wave operator, stationary approach.

## **Acknowledgements**

The first and greatest thank goes to my supervisor Professor Nils Svanstedt who passed away in the last year of my doctoral study. Thank you Nils for the time you spent in teaching me, guiding me, and encouraging me.

I want to thank Professor Mohammad Asadzadeh for his appreciated supervision, support, and guidance during the last period of my study. I would like to thank Professor Stig Larsson for his help and overall support. Also a thank goes to my assistant advisor Professor Sten Salomonson.

I would like to thank University of Gothenburg for the financial support. I also thank the staff at the mathematics department.

I thank the colleagues for creating a pleasant work environment. Thanks to the friends I met in Sweden who have made the life more enjoyable.

Many thanks to my mother, brothers, sisters, and parents-in-law for encouraging me and giving me the support I need during my study.

The last and deepest thank is forwarded to my wife May and to my son Adam for their love and the wonderful time they have created.

Hasan Almanasreh  
Gothenburg, November 2012



To  
May and Adam





## List of appended papers

The thesis is based on the following appended papers:

- I **Hasan Almasreh**, Sten Salomonson, and Nils Svanstedt, *Stabilized finite element method for the radial Dirac equation*. (Submitted)
- II **Hasan Almasreh**, *hp-Cloud approximation of the Dirac eigenvalue problem: The way of stability*. (To be submitted)
- III **Hasan Almasreh** and Nils Svanstedt, *G-convergence of Dirac operators*. (Published in Journal of Function Spaces and Applications)
- IV **Hasan Almasreh**, *On G-convergence of positive self-adjoint operators*. (Submitted)
- V **Hasan Almasreh**, *Strong convergence of wave operators for a family of Dirac operators*. (Submitted)
- VI **Hasan Almasreh**, *Existence and asymptotics of wave operators for self-adjoint operators*. (To be submitted)



# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>The Dirac equation</b>	<b>3</b>
<b>3</b>	<b>Computation of the eigenvalues of the Dirac operator</b>	<b>5</b>
3.1	Spurious eigenvalues in the computation . . . . .	5
3.1.1	Spuriousity in general eigenvalue problems . . . . .	6
3.1.2	Spuriousity in the Dirac eigenvalue problem . . . . .	8
3.1.3	More on spuriousity in the Dirac eigenvalue problem . . .	10
3.2	Stable computation of the eigenvalues . . . . .	11
3.2.1	The FEPG approximation . . . . .	14
3.2.2	The $hp$ -CPG approximation . . . . .	17
<b>4</b>	<b>G-convergence and eigenvalue problems</b>	<b>25</b>
4.1	Elliptic operators . . . . .	25
4.1.1	An overview . . . . .	25
4.1.2	A one dimensional example . . . . .	27
4.1.3	The definition . . . . .	28
4.1.4	Convergence of elliptic eigenvalue problems . . . . .	29
4.2	Positive definite self-adjoint operators . . . . .	30
4.2.1	The definition . . . . .	30
4.2.2	G-convergence of positive definite self-adjoint operators	31
4.3	Families of Dirac operators . . . . .	32
4.3.1	The Dirac operator with perturbation $(\tilde{\mathcal{H}}_h)$ . . . . .	33
4.3.2	G-convergence of projected parts of $\tilde{\mathcal{H}}_h$ . . . . .	33
<b>5</b>	<b>The wave operators for <math>h</math>-dependent self-adjoint operators</b>	<b>35</b>
5.1	A simple overview . . . . .	35
5.2	The time-dependent WO . . . . .	36
5.2.1	The strong time-dependent WO . . . . .	36
5.2.2	The weak time-dependent WO . . . . .	38
5.3	The stationary WO . . . . .	38
5.4	Pseudo-differential operators . . . . .	40
5.5	A family of Dirac operators . . . . .	40
5.5.1	An $h$ -dependent perturbation and the WO . . . . .	40
5.5.2	The asymptotics of the WOs and some particular cases .	42
5.6	Self-adjoint $h$ -dependent operators . . . . .	45
<b>6</b>	<b>Conclusion and future work</b>	<b>49</b>



# 1 Introduction

In quantum mechanics the Dirac equation is a wave equation that provides a description of the relativistic motion of the electrons as well the positrons, while the corresponding eigenvalue problem determines their energies (eigenvalues). The computation of the Dirac operator eigenvalues for single-electron systems is thoughtfully considered in the last decades in order to obtain stable solution that can be used as a basis in approximating the eigenvalues of the electron in some simple many-electron systems. The difficulty in computing the Dirac operator eigenvalues for a single-electron system is the presence of unphysical (spurious) eigenvalues among the genuine ones. Also, another challenging task is the study of the asymptotic behavior of the spectrum, in particular the eigenvalues, of families of perturbed Dirac operators.

The need for stable approximation for the Dirac operator eigenvalues with Coulomb interaction for single-electron systems makes the construction of a stable numerical scheme the main concern of this thesis. Here, we classify the spuriousity in two categories; the so-called instilled spurious eigenvalues and the spuriousity caused by the unphysical coincidence phenomenon. We provide a stable scheme to compute the Dirac operator eigenvalues implementing two different numerical methods; the finite element method (FEM) and the Galerkin-based  $hp$ -cloud method. The scheme relies on appropriate choices of the computational space that meets the properties of the Dirac wave function. On the other hand, it mainly relies on adding artificial stability terms controlled by a stability parameter. The stability parameter is derived for particular finite element basis functions in the FEM, and generalized to arbitrary basis functions in the  $hp$ -cloud approximation method. The stability scheme is computationally unexpensive, simple to apply, and guarantees complete removal of the spurious eigenvalues of both categories.

We also study the asymptotic behavior, as  $h \rightarrow \infty$ , of the eigenvalues of a family of perturbed Dirac operators by  $h$ -dependent potentials using the theory of G-convergence. We prove the G-limit operators of positive definite parts of this family under suitable assumptions on the perturbations. In particular we prove that the corresponding eigenvalues and the eigenvalue problem of the operator restricted to the point spectrum of the perturbed Dirac operators converge respectively to the eigenvalues and eigenvalue problem of the G-limit operator. Apart from this, we start employing  $\Gamma$ -convergence together with G-convergence to study the G-limits of some positive definite self-adjoint operators, and discuss the convergence of their corresponding eigenvalue problems.

Regarding the absolutely continuous part of the spectrum, we study scattering theory for a family of Dirac operators and general self-adjoint operators. For the Dirac operator with different power-like decay  $h$ -dependent potentials, the strong time-dependent wave operator (WO) exists and is complete. We prove the strong convergence, as  $h \rightarrow \infty$ , of this WO under suitable conditions on the assumed potentials. If the added potentials are of short-range type, the convergence study of the WOs is equivalent to the convergence study of the perturbed Dirac operator in the strong resolvent sense. For the Dirac operator with long-range potentials, we consider two simplified WOs for which the study of the asymptotic behavior is easier. Depending on the power of decay of the assumed potentials, the simplified WOs are obtained by considering two particular identifications. One of these identifications is an  $h$ -free operator, thus the study of the asymptotic behavior of the WOs is also reduced to the study of the convergence of the perturbed Dirac operator in the strong resolvent sense. The other identification still has the  $h$ -dependency, but the convergence of the WOs with this identification becomes easier to study. For general  $h$ -dependent self-adjoint operators, the existence and convergence, as  $h \rightarrow \infty$ , of the weak and strong time-dependent WOs and of the stationary WO are studied more extensively.

An outline of this work is as follows: In §2, we give preliminaries and introduce some elementary properties of the Dirac operator. In §3, we explain the occurrence of the spurious eigenvalues caused by using the projection method in the numerical approximation. Also we discuss the causes of the spuriousity of both categories in the computation of the Dirac eigenvalue problem. We continue with the discussion on the stability scheme and the stability parameters, where we also provide numerical examples using the FEM and  $hp$ -cloud method. In §4 we give basic preliminaries on G-convergence including a general overview, a one dimensional example, and some definitions. Likewise,  $\Gamma$ -convergence and its connection to G-convergence are stated. We also discuss G-convergence of elliptic and positive definite self-adjoint operators. Further, we apply G-convergence theory to positive definite parts of a family of Dirac operators. In §5 we provide a general overview of scattering theory and state the definitions of the time-dependent and stationary WOs and their properties. We also study the strong convergence of the WOs for a family of Dirac operators, and discuss the simplified WOs. Finally, we discuss and prove the existence and convergence of the time-dependent and stationary WOs for the general  $h$ -dependent self-adjoint operators. We conclude the work by giving a brief summary of the appended papers as well pointing out some future work.

## 2 The Dirac equation

The free Dirac equation describes the free motion of an electron (or a positron) with no external fields or presence of other particles. It is derived from the relativistic relation between energy and momentum

$$\lambda^2 = \mathbf{p}^2 c^2 + m^2 c^4, \quad (1)$$

where  $\lambda$  is the total electron energy,  $\mathbf{p}$  is the electron kinetic momentum,  $c$  is the speed of light, and  $m$  is the electron rest mass. The corresponding wave equation of quantum mechanics is obtained from the classical equation of motion (1) by replacing the energy  $\lambda$  and the momentum  $\mathbf{p}$  by their quanta

$$\lambda = i\hbar \frac{\partial}{\partial t} \quad \text{and} \quad \mathbf{p} = -i\hbar \nabla, \quad (2)$$

where  $t$  denotes the time,  $\hbar$  is the Planck constant divided by  $2\pi$ , and  $\nabla = (\frac{\partial}{\partial x_1}, \frac{\partial}{\partial x_2}, \frac{\partial}{\partial x_3})$ . Using (2), equation (1) can be written in the form

$$i\hbar \frac{\partial}{\partial t} u(x, t) = \sqrt{-c^2 \hbar^2 \Delta + m^2 c^4} u(x, t). \quad (3)$$

The problem with the existence of the Laplace operator under the square root was solved by Paul Dirac who derived the well-known Dirac equation that provides a description of the electron motion consistent with both the principles of quantum mechanics and the theory of special relativity. The free Dirac space-time equation (see [24] for more details) has the form

$$i\hbar \frac{\partial}{\partial t} u(x, t) = \mathbf{H}_0 u(x, t), \quad (4)$$

where  $\mathbf{H}_0 : H^1(\mathbb{R}^3; \mathbb{C}^4) \rightarrow L^2(\mathbb{R}^3; \mathbb{C}^4)$  is the free Dirac operator given as

$$\mathbf{H}_0 = -i\hbar c \boldsymbol{\alpha} \cdot \nabla + mc^2 \beta, \quad (5)$$

the symbols  $\boldsymbol{\alpha} = (\alpha_1, \alpha_2, \alpha_3)$  and  $\beta$  are the  $4 \times 4$  Dirac matrices given by

$$\alpha_j = \begin{pmatrix} 0 & \sigma_j \\ \sigma_j & 0 \end{pmatrix} \quad \text{and} \quad \beta = \begin{pmatrix} I & 0 \\ 0 & -I \end{pmatrix}.$$

Here  $I$  and  $0$  are the  $2 \times 2$  unity and zero matrices respectively, and  $\sigma_j$ 's are the  $2 \times 2$  Pauli matrices

$$\sigma_1 = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}, \quad \sigma_2 = \begin{pmatrix} 0 & -i \\ i & 0 \end{pmatrix}, \quad \text{and} \quad \sigma_3 = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}.$$

In the sequel we shall use the following notations;  $\mathbf{D}$ ,  $\mathbf{R}$ , and  $\mathbf{N}$  to denote respectively the domain, range, and null spaces of a given operator. The notations  $\sigma$ ,  $\sigma_p$ ,  $\sigma_{ac}$ , and  $\sigma_{ess}$  will denote respectively to the spectrum, point spectrum,

absolutely continuous spectrum, and essential spectrum of operators. For simplicity, we define  $X = H^1(\mathbb{R}^3, \mathbb{C}^4)$  and  $Y = L^2(\mathbb{R}^3, \mathbb{C}^4)$ . Separation of variables in (4) yields the free Dirac eigenvalue problem

$$\mathbf{H}_0 u(x) = \lambda u(x). \quad (6)$$

The free operator  $\mathbf{H}_0$  is essentially self-adjoint on  $C_0^\infty(\mathbb{R}^3; \mathbb{C}^4)$  and self-adjoint on  $X$ , moreover  $\sigma(\mathbf{H}_0) = (-\infty, -mc^2] \cup [mc^2, +\infty)$ . The free Dirac operator with an additional field  $V$  is given by

$$\mathbf{H} = \mathbf{H}_0 + V. \quad (7)$$

where  $V$  is a  $4 \times 4$  matrix-valued function acting as a multiplication operator in  $Y$ . The operator  $\mathbf{H}$  is essentially self-adjoint on  $C_0^\infty(\mathbb{R}^3; \mathbb{C}^4)$  and self-adjoint on  $X$  provided that the function  $V$  is Hermitian and for all  $x \in \mathbb{R}^3 \setminus \{0\}$  and  $i, j = 1, 2, 3, 4$ , satisfies  $|V_{ij}(x)| \leq a \frac{c}{2|x|} + b$ , where  $c$  is the speed of light,  $a < 1$ , and  $b > 0$ , see e.g. [51]. From now on, the function  $V$  will be considered as the Coulomb potential which has the form  $V(x) = \frac{-Z}{|x|} I$ , here  $I$  is the  $4 \times 4$  identity matrix ( $I$  will be dropped from the definition of the Coulomb potential for simplicity), and  $Z \in \{1, 2, \dots, 137\}$  is the electric charge number. The spectrum of the Dirac operator with Coulomb potential is  $(-\infty, -mc^2] \cup \{\lambda^k\}_{k \in \mathbb{N}} \cup [mc^2, +\infty)$ , where  $\{\lambda^k\}_{k \in \mathbb{N}}$  is a discrete sequence of eigenvalues.

For simple computations, to obtain the eigenvalues of the Dirac operator with Coulomb potential, the radial part of the operator is considered. Before proceeding, from now on, for simplicity, by the radial Dirac operator (eigenvalue problem) we shall mean the radial Dirac operator (eigenvalue problem) with Coulomb potential. The radial Dirac eigenvalue problem is obtained by separation of variables of the radial and angular parts, i.e., by assuming  $u(x) = \frac{1}{r} \begin{pmatrix} f(r) \mathcal{L}_{\kappa, m}(\varpi, \theta) \\ i g(r) \mathcal{L}_{-\kappa, m}(\varpi, \theta) \end{pmatrix}$ , where  $r$  represents the radial variable,  $f$  and  $g$  are the Dirac radial functions referred to as the large and small components respectively, and  $\mathcal{L}_{\cdot, m}$  is the angular part of the wave function  $u$ . The radial Dirac eigenvalue problem is then given by

$$H_\kappa \varphi(r) = \lambda \varphi(r), \quad \text{where} \quad (8)$$

$$H_\kappa = \begin{pmatrix} mc^2 + V(r) & c \left( -\frac{d}{dr} + \frac{\kappa}{r} \right) \\ c \left( \frac{d}{dr} + \frac{\kappa}{r} \right) & -mc^2 + V(r) \end{pmatrix} \quad \text{and} \quad \varphi(r) = \begin{pmatrix} f(r) \\ g(r) \end{pmatrix}. \quad (9)$$

As defined before,  $\lambda$  is the relativistic energy,  $V(r) = -Z/r$  is the radial Coulomb potential, and  $\kappa$  is the spin-orbit coupling parameter defined as  $\kappa = (-1)^{j+\ell+\frac{1}{2}} (j + \frac{1}{2})$ , where  $j$  and  $\ell$  are the total and orbital angular momentum numbers respectively.



### 3 Computation of the eigenvalues of the Dirac operator

Accurate and stable computation of the electron energies (eigenvalues) in single-electron systems (Hydrogen-like ions) is of vital interest in many applications. Approximation of the electron eigenvalues in many-electron systems, as in Helium-like ions, is based on studying quantum electrodynamic effects (QED-effects). QED-effects are known as a perturbation procedure which mainly concerns the interactions between the existing electrons in the system where these interactions are used to measure the electron correlation. An approach for calculating QED-effects, see [32, 41], is based on a basis set of eigenstates of Hydrogen-like ions (the radial Dirac operator). The main difficulty in computing the eigenvalues of the radial Dirac operator is the presence of unphysical values (eigenvalues that do not match the physical observations) among the genuine eigenvalues. These values are considered as a pollution to the spectrum and known as spurious eigenvalues. The spurious eigenvalues result in oscillations in the wave functions and the emergence of states that originally do not exist. In many cases, this will substantially reduce the computation reliability of the basis set (partially or may be completely) in the practical atomic calculations.

The spuriousity problem in the computation of the radial Dirac operator eigenvalues is a challenging issue which makes obtaining accurate and stable computation for these eigenvalues a field of study per se. Spurious eigenvalues are reported in most computational methods of eigenvalue problems, whether it is the finite element method (FEM), the finite difference method (FDM), the spectral domain approach (SDA), the boundary element method (BEM), the point matching method (PMM), or, further, the meshfree methods (MMs). Thus, spuriousity is an effect of the numerical methods and is found in the computational solution of many problems, rather than the Dirac eigenvalue problem [1, 39, 44], such as electromagnetic problems [35, 43] and general eigenvalue problems [61].

Below we present a classification of the spuriousity in the computation of the radial Dirac operator eigenvalues and its causes, we also explain the occurrence of spuriousity in the computation of general eigenvalue problems. We present two stable approaches for accurate computations with complete removal of spurious eigenvalues.

#### 3.1 Spurious eigenvalues in the computation

We classify the spuriousity in the computation of the eigenvalues of the radial Dirac operator in two categories

- (i) The instilled spuriousity.

(ii) The unphysical coincidence phenomenon.

The first category consists of those spurious eigenvalues that may occur within the genuine eigenvalues (they occur between the true energy levels). This type of spuriousity occurs for all values of the quantum number  $\kappa$ . The second type is the unphysical assigning of almost the same first eigenvalue (or almost the same entire set of eigenvalues) for  $2s_{1/2}(\kappa = -1)$  and  $2p_{1/2}(\kappa = 1)$ ,  $3p_{3/2}(\kappa = -2)$  and  $3d_{3/2}(\kappa = 2)$ ,  $4d_{5/2}(\kappa = -3)$  and  $4f_{5/2}(\kappa = 3)$ , and so on. To clarify, consider the computation of the electron eigenvalues in the Hydrogen atom using the FEM with linear basis functions (hat functions) given in Table 1, see Paper I in the appendix.

Table 1: The first computed eigenvalues, given in atomic unit, of the electron in the Hydrogen atom for point nucleus.

Level	$\kappa = 1$	$\kappa = -1$	Exact, $\kappa = -1$
1	-0.50000665661	-0.50000665659	-0.50000665659
2	-0.12500208841	-0.12500208839	-0.12500208018
3	-0.05555631532	-0.05555631532	-0.05555629517
$\Rightarrow$	-0.03141172061	-0.03141172060	Spurious Eigenvalue
4	-0.03118772526	-0.03118772524	-0.03125033803
5	-0.01974434510	-0.01974434508	-0.02000018105

The shaded value in the first level of Table 1 is what meant by the unphysical coincidence phenomenon, and the other two shaded values are the so-called instilled spuriousity. The right column contains the exact eigenvalues for  $\kappa = -1$  obtained using the relativistic formula.

### 3.1.1 Spuriousity in general eigenvalue problems

The numerical computation of the eigenvalue problems that is based on the projection method onto finite dimensional subspaces is often polluted by the presence of spurious eigenvalues [7]. The spurious eigenvalues appear particularly in the computation for those problems with eigenvalues in gaps of their essential spectrum. To understand why the projection method generates spurious eigenvalues, consider a self-adjoint operator  $T$  defined on a Hilbert space  $\mathcal{S}$ , and consider an orthogonal projection  $\Pi : \mathcal{S} \rightarrow \mathcal{L}$ , where  $\mathcal{L}$  is a finite dimensional subspace of  $\mathbf{D}(T)$ . Let  $z \in \mathbb{C}$  and define

$$\Theta(z) = \min_{\substack{f \in \mathcal{L} \\ f \neq 0}} \frac{\|\Pi(z - T)f\|_{\mathcal{S}}}{\|f\|_{\mathcal{S}}}. \quad (10)$$

If  $\Theta(\mu) = 0$ , then  $\mu = \mu(\mathcal{L})$  is a solution to the Rayleigh-Ritz problem

$$\mu = \min_{\substack{\dim(S)=k \\ S \subseteq \mathcal{L}}} \max_{g \in S} \mathcal{R}(g) = \min_{\substack{\dim(S)=k \\ S \subseteq \mathcal{L}}} \max_{g \in S} \frac{\langle Tg, g \rangle}{\|g\|_{\mathcal{T}}^2}, \quad (11)$$

where the opposite of the assertion is also true. Moreover, by assuming  $\Theta(\mu) = 0$ , we conclude that there exists  $f_0 \in \mathcal{L}$  such that

$$(\mu - T)f_0 \perp \mathcal{L}, \quad (12)$$

which particularly means that  $\mathcal{R}(f_0) = \mu$ . Thus  $\mu$  is close to the point spectrum  $\sigma_p(T)$ . But, generally, as  $\|(\mu - T)f\|_{\mathcal{T}}/\|f\|_{\mathcal{T}}$  is not necessarily small for  $f = f_0$ , any other  $f \in \mathcal{L}$ , or any  $f \in \mathcal{T}$ , then (12) does not guarantee that  $\mu$  is close to the spectrum  $\sigma(T)$  of  $T$ .

To verify the above theory, consider the following operator, see [7],

$$(Tf)(x) = \operatorname{sgn}(x)f(x), \quad (13)$$

defined on  $\mathcal{T} = L^2(-\pi, \pi)$ , where  $\operatorname{sgn}(x) = x/|x|$ . Since  $\|T\| = 1$ , then  $\sigma(T) \subseteq [-1, 1]$ , but for  $\mu \in (-1, 1)$ , the resolvent operator  $(T - \mu)^{-1}$  is well-defined and bounded, therefore  $\sigma(T) \subseteq \{-1, 1\}$ . However, it is easy to show that  $\pm 1$  are eigenvalues of  $T$ , these two eigenvalues are of infinite multiplicity, i.e.,  $\mathbf{N}(\mu - T)$  is infinite,  $\mu = \pm 1$ . Thus these two eigenvalues belong to  $\sigma_{ess}(T)$ . On the other hand, if  $\mathcal{L} \subset \mathcal{T}$  is spanned by the set of Fourier basis  $\{\varphi_{-n}, \varphi_{-n+1}, \dots, \varphi_{n-1}, \varphi_n\}$ , given by

$$\varphi_j(x) = \frac{1}{\sqrt{2\pi}} e^{-ijx}, \quad j = -n, -n+1, \dots, n-1, n, \quad (14)$$

then, the Galerkin approximation applied to  $T$  in the finite dimensional subspace  $\mathcal{L}$  implies that  $\mu_j(T, \mathcal{L})$  are the eigenvalues of the  $(2n+1) \times (2n+1)$  matrix  $A$  with entries  $(a_{jk})$  defined as

$$a_{jk} = \int_{-\pi}^{\pi} \operatorname{sgn}(x) \varphi_j(x) \varphi_k(-x) dx = \begin{cases} 0, & \text{for } k-j \text{ even,} \\ \frac{-2i}{\pi(k-j)}, & \text{for } k-j \text{ odd.} \end{cases} \quad (15)$$

The matrix  $A$  looks like

$$A = \begin{pmatrix} 0 & N & 0 & N & \dots & 0 \\ N & 0 & N & 0 & \dots & N \\ 0 & N & 0 & N & \dots & 0 \\ N & 0 & N & \ddots & \dots & N \\ \vdots & & & & \ddots & \vdots \\ 0 & N & 0 & N & \dots & 0 \end{pmatrix}, \quad (16)$$

here the letter  $N$  is used just to refer to a quantity different from zero (i.e., a number) and does not mean that these quantities are equal. It is clear that  $A$  consists of  $n + 1$  columns (the first set) whose odd entries are zero, and  $n$  columns (the second set) whose even entries are zero. If we disregard the zero entries (which are only  $n + 1$  entries) in each element of the first set, then we end with a set  $V = \{v_1, v_2, \dots, v_{n+1}\}$  where  $v_i \in \mathbb{R}^n, i = 1, 2, \dots, n+1$ . The set  $V$  is clearly linearly dependent, therefore the columns of the first set of the matrix  $A$  is linearly dependent, hence  $0 \in \sigma(A)$ . This, of course, violates the fact that  $0 \in \text{Res}(T)$ , where  $\text{Res}$  denotes the resolvent set. In this case, we conclude that  $0$  is a spurious eigenvalue that appears in the computed spectrum of the operator  $T$  caused by applying the projection method onto the finite dimensional subspace  $\mathcal{L}$ .

### 3.1.2 Spuriousity in the Dirac eigenvalue problem

The occurrence of the instilled spurious eigenvalues is a general phenomenon of the projection method in the numerical computations, thus the previous discussion can be considered as a good explanation for this type of spuriousity. Below we discuss the unphysical coincidence phenomenon as explained in [53].

Consider the radial Dirac eigenvalue problem (8), after applying the shift by  $-mc^2$  and assuming  $m = 1$ , it can be rewritten in the same form as (8) but with

$$H_\kappa = \begin{pmatrix} V(r) & c\left(-\frac{d}{dr} + \frac{\kappa}{r}\right) \\ c\left(\frac{d}{dr} + \frac{\kappa}{r}\right) & -2c^2 + V(r) \end{pmatrix}, \quad (17)$$

where the eigenvalues are also shifted but kept denoted as  $\lambda$ . Define the following transformation

$$\mathcal{U}_\kappa = \begin{pmatrix} 1 & \mathcal{U}_\kappa \\ \mathcal{U}_\kappa & 1 \end{pmatrix}, \quad (18)$$

where  $\mathcal{U}_\kappa = \frac{-Z\kappa}{c|\kappa|(|\kappa| + \varsigma)}$ , with  $\varsigma = \sqrt{\kappa^2 - Z^2/c^2}$ . We apply the above transformation to the radial function  $\varphi(r)$  given by (9) to get

$$\tilde{\varphi}_\kappa(r) = \mathcal{U}_\kappa \begin{pmatrix} f(r) \\ g(r) \end{pmatrix} = \begin{pmatrix} f(r) + \mathcal{U}_\kappa g(r) \\ g(r) + \mathcal{U}_\kappa f(r) \end{pmatrix} =: \begin{pmatrix} \tilde{f}_\kappa(r) \\ \tilde{g}_\kappa(r) \end{pmatrix}. \quad (19)$$

Using this transformation one can write

$$\mathcal{U}_\kappa^{-1} H_\kappa \mathcal{U}_\kappa^{-1} \tilde{\varphi}_\kappa(r) = \lambda_\kappa \mathcal{U}_\kappa^{-2} \tilde{\varphi}_\kappa(r). \quad (20)$$

By (20), and after adding the term  $c^2(1 - \frac{|\kappa|}{\zeta})\mathcal{U}_\kappa^{-2}\tilde{\varphi}_\kappa(r)$  to its both sides, the radial Dirac eigenvalue problem, (8), can be written in the form

$$H_{\kappa,\mu}\tilde{\varphi}_\kappa(r) = \mu_\kappa\mathcal{U}_\kappa^{-2}\tilde{\varphi}_\kappa(r), \quad (21)$$

where the operator  $H_{\kappa,\mu}$  is defined by

$$H_{\kappa,\mu} = \mathcal{U}_\kappa^{-1}H_\kappa\mathcal{U}_\kappa^{-1} + \Delta\mu_\kappa\mathcal{U}_\kappa^{-2} = \begin{pmatrix} 0 & cB_\kappa^+ \\ cB_\kappa & -2c^2 \end{pmatrix}, \quad (22)$$

and where  $\Delta\mu_\kappa = c^2(1 - \frac{|\kappa|}{\zeta})$ ,  $\mu_\kappa = \lambda_\kappa + \Delta\mu_\kappa$ ,  $B_\kappa = d/dr + \zeta\kappa/(|\kappa|r) - Z/\kappa$ , and  $B_\kappa^+ = -B_{-\kappa}$ . The same projection  $\mathcal{U}_\kappa$  can be applied to the Galerkin formulation of the radial Dirac eigenvalue problem in a finite dimensional subspace. In other words, if both radial functions  $f$  and  $g$  are expanded in a finite orthonormal basis set (orthonormal is assumed for simplicity, and it is not a requirement, since we can normalize, by a suitable linear transformation, any set of basis functions without changing the spectrum), then the above transformation applied to the discretization of the Galerkin formulation of the radial Dirac eigenvalue problem yields

$$(H_{\kappa,\mu})_{ij}(\tilde{\varphi}_\kappa)_{ij} = \mu_\kappa\mathcal{U}_\kappa^{-2}(\tilde{\varphi}_\kappa)_{ij}. \quad (23)$$

Here we have used the notation  $(\ )_{ij}$  to denote for the matrices (regardless their sizes) obtained from the Galerkin formulation. The vector  $(\tilde{\varphi}_\kappa)_{ij} = ((\tilde{f}_\kappa)_{ij}, (\tilde{g}_\kappa)_{ij})^t$  is the unknowns, and the matrix  $(H_{\kappa,\mu})_{ij}$  is given by

$$(H_{\kappa,\mu})_{ij} = \begin{pmatrix} 0 & c(B_\kappa^+)_{ij} \\ c(B_\kappa)_{ij} & -2c^2 \end{pmatrix}, \quad (24)$$

where  $(B_\kappa)_{ij}$  is the matrix of elements resulted from the discretization of the Galerkin formulation on the finite basis set.

We multiplying (23) from left by the matrix

$$\mathcal{A}_\kappa = \begin{pmatrix} A_\kappa & 0 \\ 0 & -A_\kappa^+ \end{pmatrix}, \quad (25)$$

where  $A_\kappa = (B_\kappa)_{ij} - \mu_\kappa\kappa Z/(|\kappa|c^2\zeta)$  and  $A_\kappa^+ = (B_\kappa^+)_{ij} - \mu_\kappa\kappa Z/(|\kappa|c^2\zeta)$ , to get

$$(H_{-\kappa,\mu})_{ij}\mathcal{A}_\kappa(\tilde{\varphi}_\kappa)_{ij} = \mu_\kappa\mathcal{U}_{-\kappa}^{-2}\mathcal{A}_\kappa(\tilde{\varphi}_\kappa)_{ij}, \quad (26)$$

where we have used the fact that

$$\mathcal{A}_\kappa((H_{\kappa,\mu})_{ij} - \mu_\kappa\mathcal{U}_\kappa^{-2}) = ((H_{-\kappa,\mu})_{ij} - \mu_\kappa\mathcal{U}_{-\kappa}^{-2})\mathcal{A}_\kappa. \quad (27)$$

Now we define the normalization factor  $\mathcal{N}_\kappa$  as

$$\mathcal{N}_\kappa = \langle A_\kappa(\tilde{f}_\kappa)_{ij}, A_\kappa(\tilde{f}_\kappa)_{ij} \rangle + \langle A_\kappa^+(\tilde{g}_\kappa)_{ij}, A_\kappa^+(\tilde{g}_\kappa)_{ij} \rangle, \quad (28)$$

here  $\langle \cdot, \cdot \rangle$  is the scalar product of vectors in the Euclidean space. Then for  $\mu_\kappa \neq 0$ , the eigenfunctions of  $(H_{\kappa,\mu})_{ij}$  and  $(H_{-\kappa,\mu})_{ij}$  are related by the following equation

$$(\tilde{\varphi}_{-\kappa})_{ij} = \mathcal{A}_\kappa(\tilde{\varphi}_\kappa)_{ij} / \sqrt{\mathcal{N}_\kappa}. \quad (29)$$

Substituting (29) in (26) yields

$$(H_{-\kappa,\mu})_{ij}(\tilde{\varphi}_{-\kappa})_{ij} = \mu_\kappa \mathcal{U}_{-\kappa}^{-2}(\tilde{\varphi}_{-\kappa})_{ij}. \quad (30)$$

Thus by (23) and (30), the nonzero eigenvalues of  $H_\kappa$  and  $H_{-\kappa}$  would coincide in the finite basis set. Since  $(H_{-\kappa,\mu})_{ij}$  and  $(H_{-\kappa,\mu})_{ij}$  are of the same size, then the number of their zero eigenvalues is the same. To conclude, the eigenvalues of  $H_\kappa$  and  $H_{-\kappa}$  would coincide in the projection method onto the finite dimensional subspaces in the numerical computations.

### 3.1.3 More on spuriosity in the Dirac eigenvalue problem

Most of computational methods of the eigenvalues of the radial Dirac operator consent that incorrect balancing and symmetric treatment of the large and small components of the wave function are the core of the problem [1, 39, 44]. We relate the occurrence of spuriosity of both categories to unsuitable computational spaces and to the symmetric treatment of the trial and test functions in the weak formulation of the equation. To get more understanding, we rewrite (8) to obtain explicit formulae for the radial functions  $f$  and  $g$ , see Paper I in the appendix,

$$f''(x) + \gamma_1(x, \lambda)f'(x) + \gamma_2(x, \lambda)f(x) = 0, \quad (31)$$

$$g''(x) + \theta_1(x, \lambda)g'(x) + \theta_2(x, \lambda)g(x) = 0, \quad (32)$$

where

$$\begin{aligned} \gamma_1(x, \lambda) &= -\frac{V'(x)}{w^-(x) - \lambda}, & \theta_1(x, \lambda) &= -\frac{V'(x)}{w^+(x) - \lambda}, \\ \gamma_2(x, \lambda) &= \frac{(w^+(x) - \lambda)(w^-(x) - \lambda)}{c^2} - \frac{\kappa^2 + \kappa}{x^2} - \frac{\kappa V'(x)}{x(w^-(x) - \lambda)}, \\ \theta_2(x, \lambda) &= \frac{(w^+(x) - \lambda)(w^-(x) - \lambda)}{c^2} - \frac{\kappa^2 - \kappa}{x^2} + \frac{\kappa V'(x)}{x(w^+(x) - \lambda)}, \end{aligned}$$

and  $w^\pm(x) = \pm mc^2 + V(x)$ . It is a well-known fact that the numerical methods are not stable when they are applied to convection dominated problems causing the solution to be disturbed by spurious oscillations. The following two criteria are frequently used to determine whether a given equation is convection dominated. Let

$$Pe_j = \frac{|u_j|h_j}{2K} \quad \text{and} \quad Da_j = \frac{s_j h_j}{|u_j|}, \quad (33)$$

where  $Pe_j$  and  $Da_j$  are known as the grid Peclet and Damköhler numbers respectively,  $h_j$  is the size of the element interval  $I_j$ ,  $u_j$  and  $s_j$  are respectively the coefficients of the convection and reaction terms corresponding to  $I_j$ , and  $K$  is the diffusivity size. In general, when the convection coefficient or the source term is larger than the diffusion coefficient, i.e., when  $Pe_j > 1$  or  $2Pe_j Da_j = (s_j h_j^2 / K) > 1$ , then the associated equation is a convection dominated one.

For both (31) and (32), the quantity  $2PeDa$  admits very large values if small number of nodal points in the discretization of the weak form is considered regardless the sizes of  $|\lambda|$ ,  $Z$ , and  $\kappa$ . Even with mesh refinement (increasing the number of nodal points),  $2PeDa$  still admits very large values. For (31),  $Pe$  is always less than one. As for (32), even with mesh refinement,  $Pe$  admits a value greater than one, see Paper II in the appendix for more details. Therefore, (31) and (32) are convection dominated equations. This means that the approximated solutions,  $f$  and  $g$ , will be disturbed by unphysical oscillations, these oscillations in the eigenfunctions are the cause of spurious eigenvalues.

### 3.2 Stable computation of the eigenvalues

In the coming discussion we present mesh-based and meshfree stable approaches for the approximation of the radial Dirac operator eigenvalues. As a mesh-based approach we use the finite element method (FEM), and as a meshfree approach we apply the  $hp$ -cloud method [18, 62]. For the purpose of obtaining a stability scheme based on the Petrov-Galerkin formulation with stability parameters for the particular problem, the  $hp$ -cloud method applied in this work is based on the Galerkin formulation. This means a background mesh must be employed in evaluating the integrals in the weak form, hence, the  $hp$ -cloud method used here is not really a truly meshfree method (MM). Therefore, the FEM and Galerkin-based  $hp$ -cloud method are similar in principle, while the latter approach can be regarded as a generalization of the FEM.

Based on (31) and (32), the radial functions  $f$  and  $g$  are continuous and have continuous first derivatives. Thus, the suitable choice of computational spaces

for the radial Dirac eigenvalue problem should possess these properties. Then, with homogeneous Dirichlet boundary condition for both radial functions, the proposed space is  $\mathbb{H}(\Omega) := C^1(\Omega) \cap H_0^1(\Omega)$ . Note that, except the states  $1s_{1/2}$  and  $2p_{1/2}$ , the radial functions are vanishing on the boundary in a damping way, consequently homogeneous Neumann boundary condition should be taken into account. Meanwhile, the upper boundary of the states  $1s_{1/2}$  and  $2p_{1/2}$  is treated as the others, but the first derivative of these states at the lower boundary is not zero, see e.g. [42]. For simplicity and to avoid further remarks, in the discussion below, general boundary conditions are assumed for all states, that is, homogeneous Dirichlet boundary condition. Thus the space  $\mathbb{H}(\Omega)$  is considered. However, for better rate of convergence of the approximation of the radial Dirac operator eigenvalues, the suitable Neumann boundary conditions, as discussed above, should be also implemented.

In our computation using the FEM, we use cubic Hermite basis functions, these functions treat also the first derivative values of the approximated function at the nodal values. Therefore, homogeneous Neumann boundary condition can be easily implemented by omitting the two basis functions that treat the function first derivative at the boundary nodal points, see the discussion below. Hence, in the approximation of the eigenvalues of the radial Dirac operator using the FEM, homogeneous Neumann boundary condition, as well homogeneous Dirichlet boundary condition, is implemented for all states. For the approximation using the  $hp$ -cloud method, homogeneous Dirichlet boundary condition is only considered.

Since the radial Dirac eigenvalue problem is a convection dominated problem, the FEM and  $hp$ -cloud method for this problem will be unstable, thus the occurrence of spurious eigenvalues. To stabilize the computation and to get rid of spurious eigenvalues completely, finite element Petrov-Galerkin (FEPG) (called also Streamline Upwind Petrov-Galerkin (SUPG)) [2, 14, 28] and  $hp$ -cloud Petrov-Galerkin ( $hp$ -CPG) (a technique of the general meshfree local Petrov-Galerkin (MLPG) methods [3, 19, 31]) methods are used. Apart from mesh consideration, the principle of the FEPG method is similar to that of the  $hp$ -CPG method, while the two methods mainly vary in the set of basis functions. The FEPG and  $hp$ -CPG methods are used to introduce artificial diffusion terms in the weak formulation of the equation to stabilize the approximated solution in a consistent way so that the solution of the original problem is also a solution to the weak form. The size of the added diffusivity is controlled by a stability parameter that is derived for the particular problem we consider.



To set the scheme, let  $\mathcal{V}_h$  be a finite dimensional subspace spanned by a suitable  $C^1$ -basis set on a partition  $k_h$  of the domain  $\Omega$ , where exponentially distributed nodal points are assumed to get sufficient information about the behavior of the radial functions near the origin where they oscillate heavily compared to regions away from it. We consider the weak form of the radial Dirac eigenvalue problem

$$\int_{\Omega} \mathbf{u}^t H_{\kappa} \varphi dr = \lambda \int_{\Omega} \mathbf{u}^t \varphi dr, \quad (34)$$

for all test functions  $\mathbf{u}$  in an appropriate function space, where, we recall that,  $H_{\kappa}$  is the radial Dirac operator given by

$$H_{\kappa} = \begin{pmatrix} mc^2 + V(r) & c\left(-\frac{d}{dr} + \frac{\kappa}{r}\right) \\ c\left(\frac{d}{dr} + \frac{\kappa}{r}\right) & -mc^2 + V(r) \end{pmatrix}, \quad (35)$$

and  $\varphi$  is the radial function given by

$$\varphi(r) = \begin{pmatrix} f(r) \\ g(r) \end{pmatrix}. \quad (36)$$

The usual Galerkin formulation is to consider the test function  $\mathbf{u}$  in the weak form above as  $(v, 0)^t$  and  $(0, v)^t$ , where  $v$  as well  $f$  and  $g$  is an element of  $\mathcal{V}_h$ . The FEPG and  $hp$ -CPG methods are formulated by considering  $\mathbf{u}$  in (34) as  $(v, \tau v')^t$  and  $(\tau v', v)^t$ , where  $v'$  means  $dv/dr$  and  $\tau$  is the stability parameter that controls the size of the artificial diffusivity. The stability parameter  $\tau$  is the main challenge in constructing the stability scheme and its derivation is the major task.

The derivation of  $\tau$  assumes the operator limit as the radial variable  $r \rightarrow \infty$ . This presumable assumption is inevitable and justifiable: The derivation leads to an approximation of the limit point eigenvalue depending on  $\tau$  which can be compared to the theoretical limit [23]. Thus, minimizing the error between these two limits provides  $\tau$ . By considering the limit operator at infinity, we consider the part that includes the convection terms of the operator which are mostly needed to be stabilized. Besides that, the stability parameter should be applicable at all positions, particularly for the large values of  $r$ . The derivation also considers the dominant terms with respect to the speed of light,  $c$ , as another minor simplification.

### 3.2.1 The FEPG approximation

In the FEPG method we let  $\mathcal{V}_h$  be spanned by the cubic Hermite basis functions

$$\phi_{j,1}(x) = \begin{cases} \frac{1}{h_j^2}(x - x_{j-1})^2 - \frac{2}{h_j^3}(x - x_{j-1})^2(x - x_j), & x \in I_j, \\ 1 - \frac{1}{h_{j+1}^2}(x - x_j)^2 + \frac{2}{h_{j+1}^3}(x - x_j)^2(x - x_{j+1}), & x \in I_{j+1}, \end{cases}$$

$$\phi_{j,2}(x) = \begin{cases} \frac{1}{h_j^2}(x - x_{j-1})^2(x - x_j), & x \in I_j, \\ (x - x_j) - \frac{1}{h_{j+1}}(x - x_j)^2 + \frac{1}{h_{j+1}^2}(x - x_j)^2(x - x_{j+1}), & x \in I_{j+1}. \end{cases}$$

These functions are continuous and admit continuous first derivatives, so they satisfy the continuity properties of the space  $\mathbb{H}(\Omega)$ . Moreover, they consist of two different bases, one treats the function values and the other treats the function first derivative values at the nodal points, see Figure 1. Thus any function  $w \in \mathcal{V}_h$  can be written as

$$w(r) = \sum_{j=1}^n w_j \phi_{j,1}(r) + \sum_{j=1}^n w'_j \phi_{j,2}(r), \quad (37)$$

where  $w_j$  and  $w'_j$  are respectively the function and the function first derivative values at the node  $r_j$ , and  $n$  is the number of type one basis functions  $\phi_{.,1}$  (which is the same as the number of type two basis functions  $\phi_{.,2}$ ) in the basis set.

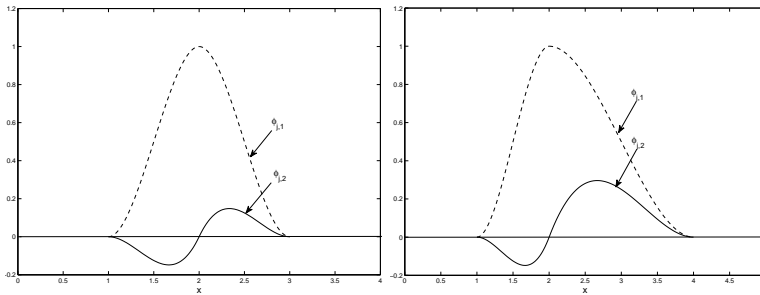


Figure 1: The CH basis functions with uniformly distributed nodal points (to the left), and non-uniformly distributed nodal points (to the right).

To treat the homogeneous Dirichlet boundary condition, the two basis functions of type  $\phi_{.,1}$  at the boundary nodal points are omitted. Also, for simplicity, we omit the two basis functions of type  $\phi_{.,2}$  at the boundary nodes, thus homogeneous Neumann boundary condition is also implemented. In the weak formulation (34), let  $v, f, g \in \mathcal{V}_h$ , this leads to the generalized eigenvalue problem

$$\mathbf{A}X = \lambda \mathbf{B}X. \quad (38)$$

The perturbed block matrices are given by  $\mathbf{A} = A + \tau \mathcal{A}$  and  $\mathbf{B} = B + \tau \mathcal{B}$ , where  $A$  and  $B$  are the matrices obtained from the FEM, and  $\mathcal{A}$  and  $\mathcal{B}$  are the

matrices obtained as a result of the correction part,  $\tau v'$ , in the test function. Note that  $\tau$  must be correlated with the size of the generated mesh, i.e., for a fine-structure mesh we expect  $\tau$  to be relatively small compared to a coarse mesh. On the other hand, to avoid the occurrence of complex eigenvalues,  $\tau$  should not be large compared to the mesh size. These properties are clear from the representation of  $\tau$  given by the following theorem, see Paper I in the appendix.

**Theorem 1** Considering the behavior of the eigenvalues as  $r$  tends to infinity, together with the dominant terms with respect to the speed of light, the mesh-dependent stability parameter,  $\tau_j$ , for an arbitrary  $j^{\text{th}}$  row of the matrices  $\mathcal{A}$  and  $\mathcal{B}$  in the generalized eigenvalue problem (38) has the form

$$\tau_j = \frac{9}{35} h_{j+1} \frac{(h_{j+1} - h_j)}{(h_{j+1} + h_j)}, \quad (39)$$

where  $h_j$  is the displacement between the nodes  $r_j$  and  $r_{j-1}$ .

Below, a numerical example of the computation of the eigenvalues of the electron in the Hydrogen-like Magnesium ion using the FEPG method is presented for  $\kappa = \pm 2$ . Note that, in all our computations, the eigenvalues are given in atomic unit. Table 2 shows the computation using the FEM with linear basis functions with 400 interior nodal points for point nucleus. Table 3 shows the same computation with the stability scheme.

Table 2: The first computed eigenvalues of the electron in the Hydrogen-like Magnesium ion using the FEM with linear basis functions for point nucleus.

Level	$\kappa = 2$	$\kappa = -2$	Exact, $\kappa = -2$
1	-18.0086349982	-18.0086349982	-18.0086349982
2	-8.00511829944	-8.00511829944	-8.00511739963
3	-4.50270135222	-4.50270135225	-4.50269856638
⇒	-2.88546212211	-2.88546212205	Spurious Eigenvalue
4	-2.88155295096	-2.88155295095	-2.88154739168
5	-2.00096852250	-2.00096852249	-2.00095939879
6	-1.47003410346	-1.47003410350	-1.47002066823
⇒	-1.13034880166	-1.13034880167	Spurious Eigenvalue
7	-1.12545691681	-1.12545691683	-1.12543844140
8	-0.889228944495	-0.889228944484	-0.889204706429
9	-0.720265553198	-0.720265553187	-0.720234829539
⇒	-0.600492562625	-0.600492562622	Spurious Eigenvalue
10	-0.595258516248	-0.595258516277	-0.595220579682
11	-0.500185771976	-0.500185772005	-0.500139887884
12	-0.426201311278	-0.426201311300	-0.426146735771

Table 3: The first computed eigenvalues of the electron in the Hydrogen-like Magnesium ion using the stability scheme for point nucleus.

Level	$\kappa = 2$	$\kappa = -2$	Exact, $\kappa = -2$
1		-18.0086349985	-18.0086349982
2	-8.00511739978	-8.00511740020	-8.00511739963
3	-4.50269856669	-4.50269856719	-4.50269856638
4	-2.88154739219	-2.88154739270	-2.88154739168
5	-2.00095939948	-2.00095939991	-2.00095939879
6	-1.47002066888	-1.47002066924	-1.47002066823
7	-1.12543844176	-1.12543844201	-1.12543844140
8	-.889204706068	-.889204706109	-.889204706429
9	-.720234827833	-.720234827687	-.720234829539
10	-.595220575840	-.595220575531	-.595220579682
11	-.500139880950	-.500139880357	-.500139887884
12	-.426146724530	-.426146723650	-.426146735771
13	-.367436809137	-.367436807839	-.367436826403
14	-.320073519367	-.320073498169	-.320073665658
15	-.281295132797	-.281293164731	-.281311119433

In Table 4, the computation is performed for extended nucleus using uniformly distributed charge with 397 interior nodal points, where 16 nodal points are considered in the domain  $[0, R]$  ( $R$  is the radius of the nucleus).

Table 4: The first computed eigenvalues of the electron in the Hydrogen-like Magnesium ion using the stability scheme for extended nucleus.

Level	$\kappa = 2$	$\kappa = -2$	Exact, $\kappa = -2$
1		-18.0086349986	-18.0086349982
2	-8.00511739975	-8.00511740015	-8.00511739963
3	-4.50269856673	-4.50269856733	-4.50269856638
4	-2.88154739230	-2.88154739279	-2.88154739168
5	-2.00095939956	-2.00095940014	-2.00095939879
6	-1.47002066903	-1.47002066934	-1.47002066823
7	-1.12543844179	-1.12543844207	-1.12543844140
8	-.889204706021	-.889204706003	-.889204706429
9	-.720234827640	-.720234827433	-.720234829539
10	-.595220575309	-.595220574883	-.595220579682
11	-.500139879906	-.500139879215	-.500139887884
12	-.426146722827	-.426146721812	-.426146735771
13	-.367436806543	-.367436805088	-.367436826403
14	-.320073514034	-.320073492344	-.320073665658
15	-.281294966822	-.281292979627	-.281311119433

Note that the exact eigenvalues in the tables above (as well in the computations below) are obtained, of course for point nucleus, using the relativistic formula.

### 3.2.2 The $hp$ -CPG approximation

The  $hp$ -cloud basis functions are obtained using moving least-squares (MLS) approximation method which allows polynomial enrichment and desired fundamental characters of the sought solution to be constructed in the approximation. The  $hp$ -cloud basis functions take the form

$$\psi_j(r) = P^t(r)M^{-1}(r)\varphi_j\left(\frac{r-r_j}{\rho_j}\right)P(r_j)\psi_j, \quad (40)$$

where  $M(r) = \sum_{j=1}^n \varphi_j\left(\frac{r-r_j}{\rho_j}\right)P(r_j)P^t(r_j)$  is the momentum matrix,  $P$  is a vector of intrinsic enrichments,  $\varphi_j$  is a weight function, and  $\rho_j$  is the dilation parameter that controls the support of the weight functions.

The weight function  $\varphi_j$  is the main feature in the definition of  $\psi_j$ , it is needed to be  $C^1$ -function in order to guarantee the continuity property of the space  $\mathbb{H}(\Omega)$ . For this purpose, we will consider quartic spline (which is a  $C^2$ -function) as a weight function defined as

$$\varphi(\check{r}) = \begin{cases} 1 - 6\check{r}^2 + 8\check{r}^3 - 3\check{r}^4, & \check{r} \leq 1, \\ 0, & \check{r} > 1, \end{cases} \quad (41)$$

where  $\check{r} = \frac{|r-r_j|}{\rho_j}$ . While the set of functions  $\{\psi_j\}_{j=1}^n$  builds a partition of unity (PU) ( $\sum_{j=1}^n \psi_j(r) = 1$ , for all  $r \in \Omega$ ), the set of their first derivatives  $\{\psi_{j,r}\}_{j=1}^n = \left\{\frac{d\psi_j(r)}{dr}\right\}_{j=1}^n$  builds a partition of nullity (PN) ( $\sum_{j=1}^n \psi_{j,r}(r) = 0$  for all  $r \in \Omega$ ), see Figure 2.

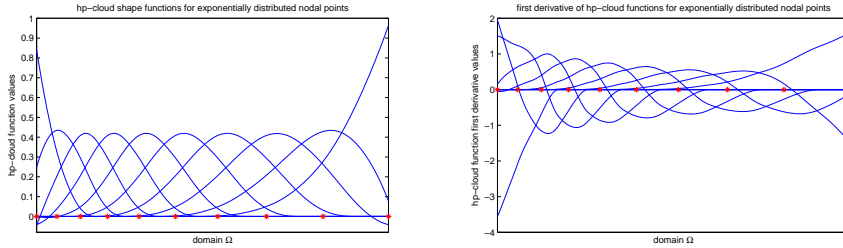


Figure 2: PU  $hp$ -clouds (to the left) and their PN first derivatives (to the right). Quartic spline is used as a weight function.

The invertibility of  $M$  depends on  $\rho_j$ , as  $\rho_j$  gets smaller as the matrix  $M$  has more tendency to be singular. So, in order to maintain the invertibility of  $M$ , it is necessarily to keep  $\rho_j$  sufficiently large. However,  $\rho_j$  can be chosen fixed

or arbitrary, in this work we consider (exponentially distributed nodal points are used)

$$\rho_j = \nu \cdot \max\{h_j, h_{j+1}\} = \nu h_{j+1}, \quad (42)$$

where  $\nu$  is the dimensionless size of the influence domain [30] which is chosen to be fixed in our computation. Note that the maximum in (42) is crucial to guarantee less possibility for singularity of  $M$ . The choices of  $\nu$  are constrained by two restrictions; the values of  $\nu$  should not be very small to ensure that any region is covered by at least two clouds, thus the invertibility of  $M$ . On the other hand, the values of  $\nu$  should not be very large to guarantee the local character of the approximation. Noting that as  $\nu \rightarrow 1$ , the  $hp$ -cloud,  $\psi_j$ , will act as a finite element basis function, and thus the features of the  $hp$ -cloud approximation are gradually lost. The optimal choices of  $\nu$  are left undetermined in general, but they can be individually specified for each problem by running numerical experiments [33, 60]. For the computation of the radial Dirac eigenvalue problem, for  $\nu \in [2.2, 2.7]$  good approximation is achieved, see Table 7 and Figure 5 below, with complete elimination of the spurious eigenvalues.

The intrinsic enrichment basis vector  $P$  is a very important ingredient in the construction of the  $hp$ -cloud functions. Using the vector  $P$ , all fundamental features of the sought solution as well as singularities and discontinuities can be inherited by the  $hp$ -cloud basis functions. This distinguishes the  $hp$ -cloud approximation by solving particular problems where much care is needed about the approximated solution such as solving equations with rough coefficients, problems with high oscillatory solutions, or eigenvalue problems that admit spurious eigenvalues. Note that yet another type of enrichment, called extrinsic enrichment, can be considered in the construction of the  $hp$ -cloud functions, but this type of enrichment is not adequate when applying the  $hp$ -CPG method [3]. Thus, in this work, extrinsic enrichment is not considered.

The number and type of the intrinsic enrichment functions in the basis set  $P$  can be chosen arbitrary for each cloud [21, 34], but for practical reasons (lowering both the condition number of  $M$  and the computational costs) we shall assume  $P(x) = [1, p_1(x)]$ , where the choices of  $p_1(x)$  follow two main properties; since  $\psi_j$  is needed to be a  $C^1$ -function, which is guaranteed only if both the weight function  $\varphi_j$  and the elements of  $P$  are also in  $C^1(\Omega)$ ,  $p_1(x)$  should be a  $C^1$ -function as well. Secondly,  $p_1(x)$  should possess the global behavior of the electron motion.

Slater type orbital functions (STOs) and Gaussian type orbital functions (GTOs) provide good description of the electron motion [10, 25]. But the

quadratic term in the exponent of the GTOs causes some numerical difficulty, in the sense that, the matrix  $M$  rapidly becomes poorly conditioned, this is also what is observed when applying quadratic basis enrichments, see [5]. Consequently, the STOs are considered as the intrinsic enrichment of the  $hp$ -cloud functions, thus  $p_1(x)$  can have, e.g., one of the following forms

$$\exp(-x), x \exp(-x/2), x(1-x/2) \exp(-x/2), \dots \text{ etc.}$$

Other possible intrinsic enrichments for the computation of the radial Dirac operator eigenvalues can be found in Paper II in the appendix. In the computations presented below, we consider  $p_1(x) = x(1-x/2) \exp(-x/2)$ .

The boundary conditions need special treatment: For the computation of the radial Dirac eigenvalue problem we assume homogeneous Dirichlet boundary condition, while it is well-known that imposing essential boundary conditions (EBCs) in MMs, in general, is a difficulty which needs to be treated with care. The reason is that the meshfree basis functions lack the Kronecker delta property ( $\psi_j(r_i) \neq \delta_{ji}$ ), thus EBCs are not directly imposed as for the FEM. To circumvent this difficulty, a coupling with finite element basis functions is considered, see Figure 3. By coupling with finite element basis functions at the lower and upper boundaries, the imposition of the homogeneous Dirichlet boundary condition is straightforward, e.g., by eliminating the two finite element basis functions at the boundary nodes.

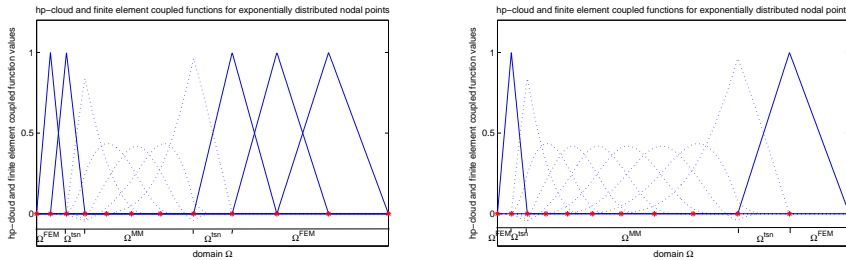


Figure 3: Coupled  $hp$ -cloud and finite element functions: general coupling (to the left), and coupling for the purpose of imposing EBCs (to the right) (two finite element shape functions are sufficient). Linear functions are used as finite element functions, and quartic spline as a weight function in the  $hp$ -clouds.

Two efficient approaches of coupling MMs with the FEM are coupling with Ramp functions [4] and coupling with reproducing conditions [26]. Using the former one, the derivative of the coupled approximation function on the boundary of the interface region,  $\Omega^{!sn}$  in Figure 3, is discontinuous, for this reason we

use the latter coupling approach, see e.g. [20]. The coupled  $hp$ -cloud and finite element function with the reproducing conditions is given as

$$\begin{aligned} \psi_j(r) = & \left( P^t(r) - \mathcal{G}_j(r) P^t(r_j) \chi_{\Omega^{\text{FEM}}}(r_j) \right) M^{-1}(r) \varphi_j\left(\frac{r-r_j}{\rho_j}\right) \times \\ & \times P(r_j) \chi_{\Omega^{\text{MM}}}(r_j) \psi_j + \mathcal{G}_j(r) \chi_{\Omega^{\text{FEM}}}(r_j) \psi_j, \end{aligned} \quad (43)$$

where  $\chi_{\Omega^{\text{FEM}}}$  and  $\chi_{\Omega^{\text{MM}}}$  are respectively the characteristic functions of the domains  $\Omega^{\text{FEM}}$  and  $\Omega^{\text{MM}}$ , see Figure 3, and  $\mathcal{G}_j$  is the finite element function.

To enhance the stability of the computation and to maintain the accuracy that may be affected or lost due to the round-off error, and also to get a lower condition number for the matrix  $M$ , the origin should be shifted to the evaluation point in the meshfree basis functions in general [20, 27, 30].

After constructing the  $hp$ -cloud basis functions, the  $hp$ -CPG method is formulated by assuming the weak form (34) where  $u$ , as before, takes the forms  $(v, \tau v')$  and  $(\tau v', v)$ , and  $v, f, g \in \mathcal{V}_h$ , where  $\mathcal{V}_h$  is now spanned by a set of functions of the form (43). This yields similar generalized eigenvalue problem as of (38). The stability parameter,  $\tau$ , is now different from the one given by Theorem 1, and can be considered as a generalization of it. The same principle as in Theorem 1 is used in deriving  $\tau$  by using the  $hp$ -cloud basis functions. The following theorem provides the representation of  $\tau$  which will be still denoted by the same notation.

**Theorem 2** Let  $M_{000}$  and  $M_{100}$  be the  $n \times n$  matrices ( $n$  is the number of  $hp$ -cloud basis functions) defined as

$$(M_{000})_{ij} = \int_{\Omega} \psi_j \psi_i dr, \text{ and } (M_{100})_{ij} = \int_{\Omega} \psi_j \psi_i' dr, \quad (44)$$

and let  $\sigma_{ji}$  and  $\eta_{ji}$  be the corresponding entries respectively. Define  $\vartheta$  as

$$\vartheta_{ji} = \begin{cases} -\sum_{k=i+1}^j h_k, & i < j, \\ 0, & i = j, \\ \sum_{k=j+1}^i h_k, & i > j, \end{cases}$$

where  $h_k$  is the displacement between the adjacent nodes  $r_k$  and  $r_{k-1}$ . Then the stability parameter,  $\tau_j$ , for an arbitrary  $j^{\text{th}}$  row of the matrices  $\mathcal{A}$  and  $\mathcal{B}$  in the generalized eigenvalue problem (38) is given by

$$\tau_j = \left| \frac{\sum_{i=1}^n \sigma_{ji} \vartheta_{ji}}{\sum_{i=1}^n \eta_{ji} \vartheta_{ji}} \right|. \quad (45)$$



The advantage of the  $hp$ -CPG stability parameter (45) is that it can be applied for general basis functions and not for particular ones as of the FEFG stability parameter (39).

**Remark 1** To capture the behavior of the radial functions near the origin where they oscillate heavily compared to regions away from it, the computation of the radial Dirac operator eigenvalues requires, as we mentioned before, exponentially distributed nodal points. In this regard, the following formula is used to discretize  $\Omega$

$$r_i = \exp \left( \ln(I_a + \varepsilon) + \left( \frac{\ln(I_b + \varepsilon) - \ln(I_a + \varepsilon)}{n} \right) i \right) - \varepsilon, \quad i = 0, 1, 2, \dots, n, \quad (46)$$

where  $n$  is the total number of nodal points,  $I_a$  and  $I_b$  are the lower and upper boundaries of  $\Omega$ , and  $\varepsilon \in [0, 1]$  is the nodes intensity parameter. The main role of  $\varepsilon$  is to control the intensity of the nodal points close to the origin ( $I_a$ ). As  $\varepsilon$  gets smaller as more nodes are dragged to the origin. In Paper II in the appendix, a study is carried out concerning the suitable choices of  $\varepsilon$ , it is shown that the most appropriate values for  $\varepsilon$  that provide good results are in the interval  $[10^{-6}, 10^{-4}]$ .

The results of the computation using the  $hp$ -CPG method with the stability parameter (45) are presented in Tables 5 and 6. In Table 5, the approximated eigenvalues of the electron in the Hydrogen-like Ununoctium ion are obtained using the usual and the stabilized  $hp$ -cloud methods. The computation is obtained at  $\rho_j = 2.2h_{j+1}$ ,  $\varepsilon = 10^{-5}$ , and  $n = 600$ . In the  $hp$ -cloud method, the instilled spurious eigenvalues appear for both positive and negative  $\kappa$  (the two shaded values in the fourteenth level). Also the the so-called unphysical coincidence phenomenon occurs for the positive  $\kappa$  (the shaded value in the first level). Note that these spurious eigenvalues are removed by the stability scheme.

Table 6 represents the stabilized  $hp$ -cloud approximation of the electron in the Hydrogen-like Ununoctium ion with different numbers of nodal points. The convergence rate of the first five eigenvalues is studied in Figure 4. In Figure 4,  $h$  is the maximum of the distances between the adjacent nodes, which equals to  $h_n = r_n - r_{n-1}$  for exponentially distributed nodal points. It can be verified from the figure that the convergence rates of the approximation of the first five eigenvalues,  $\lambda_1, \lambda_2, \dots, \lambda_5$ , are nearly 3.09, 2.66, 2.62, 2.59, and 2.56 respectively.

Table 5: The first computed eigenvalues of the electron in the Hydrogen-like Ununoctium ion using the  $hp$ -cloud and  $hp$ -CPG methods for point nucleus.

Level	$hp$ -cloud $\kappa = 2$	$hp$ -cloud $\kappa = -2$	Exact $\kappa = -2$	$hp$ -CPG $\kappa = -2$	$hp$ -CPG $\kappa = 2$
1	-1829.6307	-1829.6307	-1829.6307	-1829.6283	
2	-826.76981	-826.76981	-826.76835	-826.77147	-826.77388
3	-463.12149	-463.12149	-463.11832	-463.12471	-463.12611
4	-294.45523	-294.45523	-294.45098	-294.45915	-294.46006
5	-203.24689	-203.24689	-203.24195	-203.25115	-203.25179
6	-148.55882	-148.55882	-148.55344	-148.56324	-148.56372
7	-113.25360	-113.25360	-113.24791	-113.25808	-113.25845
8	-89.163854	-89.163854	-89.157945	-89.168323	-89.168622
9	-72.004533	-72.004533	-71.998465	-72.008947	-72.009194
10	-59.354813	-59.354813	-59.348624	-59.359134	-59.359342
11	-49.764290	-49.764290	-49.758009	-49.768490	-49.768669
12	-42.321471	-42.321471	-42.315117	-42.325523	-42.325679
13	-36.430396	-36.430396	-36.423983	-36.434277	-36.434414
14	-33.965028	-33.965028	-31.681730	-31.691878	-31.692001
15	-31.688189	-31.688189	-27.808134	-27.818109	-27.818219

Table 6: The first computed eigenvalues of the electron in the Hydrogen-like Ununoctium ion for  $\kappa = -2$  for point nucleus with different number of nodes, where  $\nu = 2.2$  and  $\varepsilon = 10^{-5}$  are used.

Level	$n = 200$	$n = 400$	$n = 600$	$n = 800$	$n = 1000$	Exact, $\kappa = -2$
1	-1829.5628	-1829.6224	-1829.6283	-1829.6297	-1829.6302	-1829.6307
2	-826.82670	-826.77726	-826.77147	-826.76987	-826.76923	-826.76835
3	-463.23292	-463.13630	-463.12471	-463.12146	-463.12016	-463.11832
4	-294.59147	-294.47367	-294.45915	-294.45503	-294.45336	-294.45098
5	-203.39386	-203.26721	-203.25115	-203.24654	-203.24466	-203.24195
6	-148.70878	-148.58009	-148.56324	-148.55835	-148.55635	-148.55344
7	-113.40170	-113.27527	-113.25808	-113.25304	-113.25096	-113.24791
8	-89.306709	-89.185557	-89.168323	-89.163201	-89.161076	-89.157945
9	-72.139617	-72.026008	-72.008947	-72.003802	-72.001653	-71.998465
10	-59.480154	-59.375861	-59.359134	-59.354006	-59.351849	-59.348624
11	-49.878353	-49.784751	-49.768490	-49.763410	-49.761256	-49.758009
12	-42.423104	-42.341207	-42.325523	-42.320517	-42.318374	-42.315117
13	-36.518814	-36.449288	-36.434277	-36.429365	-36.427242	-36.423983
14	-31.762955	-31.706134	-31.691878	-31.687081	-31.684984	-31.681730
15	-27.875610	-27.831538	-27.818109	-27.813442	-27.811376	-27.808134

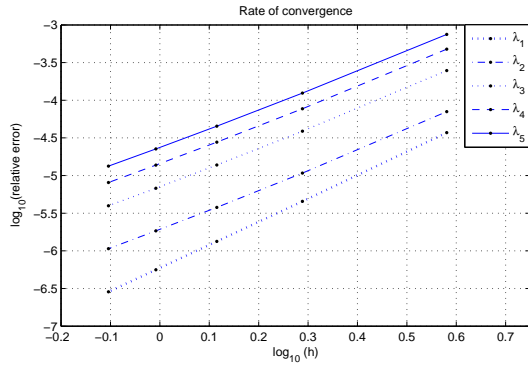


Figure 4: Studying the convergence rate of the first five eigenvalues in Table 6.

In Table 7, the computation of the eigenvalues of the electron in the Hydrogen-like Ununoctium ion is provided, the computation is obtained for  $\kappa = -2$  for point nucleus at  $n = 600$  and  $\varepsilon = 10^{-5}$ .

Table 7: The first computed eigenvalues of the electron in the Hydrogen-like Ununoctium ion for  $\kappa = -2$  for point nucleus with different values of  $\nu$ , where  $n = 600$  and  $\varepsilon = 10^{-5}$  are used.

Level	$\nu = 2.0$	$\nu = 2.2$	$\nu = 2.5$	$\nu = 2.7$	Exact solution
1	-1829.6287	-1829.6283	-1829.6276	-1829.6270	-1829.6307
2	-826.77119	-826.77147	-826.77197	-826.77233	-826.76835
3	-463.12417	-463.12471	-463.12567	-463.12638	-463.11832
4	-294.45850	-294.45915	-294.46033	-294.46120	-294.45098
5	-203.25046	-203.25115	-203.25244	-203.25340	-203.24195
6	-148.56255	-148.56324	-148.56460	-148.56562	-148.55344
7	-113.25741	-113.25808	-113.25949	-113.26054	-113.24791
8	-89.167688	-89.168323	-89.169756	-89.170831	-89.157945
9	-72.008358	-72.008947	-72.010396	-72.011489	-71.998465
10	-59.358602	-59.359134	-59.360592	-59.361700	-59.348624
11	-49.768025	-49.768490	-49.769950	-49.771070	-49.758009
12	-42.325133	-42.325523	-42.326981	-42.328113	-42.315117
13	-36.433970	-36.434277	-36.435728	-36.436870	-36.423983
14	-31.691663	-31.691878	-31.693318	-31.694472	-31.681730
15	-27.817992	-27.818109	-27.819533	-27.820699	-27.808134

In Figure 5, we study the rate of convergence of the approximation for the first five eigenvalues in Table 7.

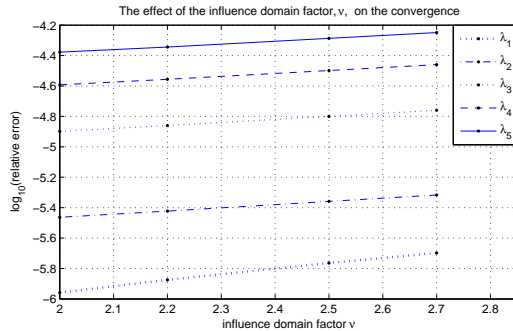


Figure 5: Studying the convergence rate with respect to  $\nu$ .

Figure 5 shows that smaller  $\nu$  gives better approximation. However, as we have mentioned, the appropriate values of  $\nu$  lie in  $[2.2, 2.7]$ , while other smaller values of  $\nu$  cause spurious eigenvalues, see Paper II in the appendix. This is seen evident since for small values of  $\nu$  the clouds are not stretched enough to capture the behavior of the sought solution. Also for small  $\nu$ , some regions of  $\Omega$  are covered only with one cloud function, which makes the momentum matrix  $M$  singular.

The stability of our computations is measured by the complete elimination of the spurious eigenvalues. This is tested using the relativistic formula which is defined as

$$\lambda_{n_r, \kappa} = \frac{mc^2}{\sqrt{1 + \frac{Z^2 \alpha^2}{(n_r - 1 + \sqrt{\kappa^2 - Z^2 \alpha^2})^2}}}, \quad (47)$$

where  $m$  and  $c$  are, as defined before, the electron rest mass and the speed of light respectively,  $\alpha$  is the fine structure constant which is equal to  $1/c$ , and  $n_r = 1, 2, \dots$  is the orbital level number. Note that, the relativistic formula only applies in the case of point nucleus. Since  $Z \in \{1, 2, \dots, 137\}$  and  $\kappa \in \mathbb{Z} \setminus \{0\}$ , this makes testing the computations with the derived stability scheme for all  $Z$  and  $\kappa$  a tedious work, but it is performed for all  $Z$  and  $\kappa = \pm 1, \pm 2, \dots, \pm 30$ . Our computations report no spurious eigenvalues, thus the numerical scheme is stable.

It is worth to mention that the FEPG method has a convergence rate higher than that of the  $hp$ -CPG method. Further, the  $hp$ -CPG method is more expensive due to the time consumption in evaluating the cloud functions that demand more integration points as  $\nu$  gets larger, which is the main disadvantage of MMs in general.

## 4 G-convergence and eigenvalue problems

In this part, we study the convergence of the eigenvalues and the corresponding eigenvalue problems for families of positive definite self-adjoint operators using the theory of G-convergence. First we discuss G-convergence of a class of elliptic and bounded positive definite self-adjoint operators. Then we consider G-convergence of a family of Dirac operators. Using the spectral measure, we consider projected positive definite parts of this family and then apply the theory of G-convergence.

The theory of G-convergence was introduced in the late 1960's [13, 45, 46, 47] for linear elliptic and parabolic operators with symmetric coefficient matrices. The concept was further extended to non-symmetric coefficient matrices [36, 48, 49, 50] and referred to as H-convergence. The theory was then generalized to positive definite self-adjoint operators [11] under the name G-convergence. The study of G-convergence of positive definite self-adjoint operators is often connected to the study of convergence of the associated quadratic forms in the calculus of variations via the notion of  $\Gamma$ -convergence which was introduced in the mid 1970's [12]. The monographs [8, 11] contain comprehensive material on the topic, where [11] deals with the connection to G-convergence. In this work, we will use the name G-convergence for the case of non-symmetric matrices as well.

### 4.1 Elliptic operators

#### 4.1.1 An overview

Let  $\Omega$  be an open bounded set in  $\mathbb{R}^N$ ,  $N \geq 1$ . To present the idea of G-convergence, a heat conduction example is considered. The  $h$ -dependent stationary heat equation with heat source  $f(x) \in H^{-1}(\Omega)$  and periodic heat conductivity matrix  $\mathbf{A}_h(x) = \mathbf{A}(hx)$ ,  $\mathbf{A}$  is  $\mathbf{Y}$ -periodic, is given by

$$\begin{cases} -\frac{\partial}{\partial x_i}((\mathbf{A}_h(x))_{ij} \frac{\partial u_h}{\partial x_j}) = f(x) & \text{in } \Omega, \\ u_h = 0 & \text{on } \partial\Omega. \end{cases} \quad (48)$$

The operator  $-\frac{\partial}{\partial x_i}((\mathbf{A}_h(x))_{ij} \frac{\partial}{\partial x_j})$  is defined on  $L^2(\Omega)$  with domain  $H_0^1(\Omega)$ ,  $h \in \mathbb{N}$  is a parameter that tends to infinity, and  $L^\infty(\Omega)^{N \times N} \ni (\mathbf{A}_h(x))_{ij}$  is positive definite and bounded.

The difficulty arises when  $h$  tends to infinity, where the highly oscillating coefficient matrix,  $\mathbf{A}_h$ , makes (48) hard to solve with direct numerical methods

with good accuracy. The idea we will advocate is to consider instead the limit equation as  $h \rightarrow \infty$  where the material is expected to behave as a homogeneous one. In other words, we are interested in finding the properties of a homogeneous equation that gives the same overall response as the heterogeneous one. This means that we look for the global macroscopic behavior of the solution. The limit problem of (48), as  $h \rightarrow \infty$ , can formally be written as

$$\begin{cases} -\frac{\partial}{\partial x_i}((\mathbf{B}(x))_{ij} \frac{\partial u}{\partial x_j}) = f(x) & \text{in } \Omega, \\ u = 0 & \text{on } \partial\Omega. \end{cases} \quad (49)$$

In a successful approach, the problem (49) contains no oscillations and hence is easier to be treated numerically. Thus, the task is to characterize the matrix  $\mathbf{B}$ .

The way of specifying the limit matrix  $\mathbf{B}$  is to let  $h \rightarrow \infty$  in the weak form of (48): Find  $u_h \in H_0^1(\Omega)$  such that

$$\langle \mathbf{A}_h(x) \nabla u_h(x), \nabla v(x) \rangle = \langle f(x), v(x) \rangle, \quad \forall v \in H_0^1(\Omega). \quad (50)$$

By the boundedness and coercivity of  $\mathbf{A}_h$ , the existence and uniqueness of a solution  $u_h$  to (48) are guaranteed by the Lax-Milgram theorem. Also these assumptions imply the boundedness of  $u_h$  and  $\nabla u_h$  in  $H_0^1(\Omega)$  and  $L^2(\Omega)$  respectively. Therefore, up to a subsequence of  $u_h$  still denoted by  $u_h$ ,

$$u_h(x) \rightharpoonup u(x) \quad \text{in } H_0^1(\Omega), \quad (51)$$

$$\nabla u_h(x) \rightharpoonup \nabla u(x) \quad \text{in } L^2(\Omega)^N, \quad (52)$$

where the notation  $\rightharpoonup$  refers to the weak convergence. Since  $\mathbf{A}_h$  is an element of  $L^\infty(\Omega)^{N \times N}$ , then up to a subsequence denoted by  $\mathbf{A}_h$ ,

$$\mathbf{A}_h \overset{*}{\rightharpoonup} M(\mathbf{A}), \quad \text{in } L^\infty(\Omega)^{N \times N}, \quad (53)$$

where  $M(\mathbf{A}) = \frac{1}{|\Upsilon|} \int_{\Upsilon} \mathbf{A}(y) dy$  is the average of  $\mathbf{A}_h$ , and  $\overset{*}{\rightharpoonup}$  refers to the weak\* convergence. We recall that (53) is also true in the  $L^2(\Omega)^{N \times N}$  sense, this is because  $L^2$  is continuously embedded in  $L^1$ , which implies  $L^\infty = (L^1)^* \subset \subset (L^2)^* = L^2$  ( $\star$  refers to the duality), hence the same topology on  $L^\infty(\Omega)$  can be also defined on  $L^2(\Omega)$ . Thus, we have two sequences,  $\nabla u_h$  and  $\mathbf{A}_h$ , which converge only weakly. This is the intricate task that we face to pass to limit as  $h \rightarrow \infty$  in (50) as nothing can be concluded about the limit of the product of two sequences that are only weakly convergent, and generally the following result is not true

$$\mathbf{A}_h(x) \nabla u_h(x) \rightharpoonup M(\mathbf{A}) \nabla u(x), \quad \text{in } L^2(\Omega)^N. \quad (54)$$

Here, another technique is employed to study the existence and characterization of the asymptotic limit of  $\mathbf{A}_h(x) \nabla u_h(x)$ , namely the theory of G-convergence.

### 4.1.2 A one dimensional example

Consider the following Dirichlet boundary value problem, see [15],

$$\begin{cases} -\frac{d}{dx}(\mathbf{A}_h(x)\frac{du_h}{dx}) = f & \text{in } \Omega = (x_0, x_1) \subset \mathbb{R}, \\ u_h \in H_0^1(\Omega), \end{cases} \quad (55)$$

where  $f \in L^2(\Omega)$ , and  $L^\infty(\Omega) \ni \mathbf{A}_h(x) = \mathbf{A}(hx)$ , with  $\mathbf{A} : \mathbb{R} \rightarrow \mathbb{R}$  is a  $Y$ -periodic function satisfying, for  $\alpha, \beta \in \mathbb{R}$ ,  $0 < \alpha \leq \mathbf{A}(x) \leq \beta < \infty$  a.e on  $\mathbb{R}$ . The bounds for  $\mathbf{A}_h$  give the existence and uniqueness of a solution  $u_h$  to (55). Moreover, the a priori estimate  $\|u_h\|_{H_0^1(\Omega)} \leq C$  implies that the sequence  $u_h$  is uniformly bounded in  $H_0^1(\Omega)$ . Hence by Rellich-Kondrachov compactness theorem, up to a subsequence still denoted by  $u_h$ , there exists  $u \in H_0^1(\Omega)$  such that

$$u_h \rightharpoonup u \text{ in } H_0^1(\Omega). \quad (56)$$

By the periodicity assumption on  $\mathbf{A}$  we have

$$\mathbf{A}_h \overset{*}{\rightharpoonup} M(\mathbf{A}) \text{ in } L^\infty(\Omega), \text{ (also weakly in } L^2(\Omega)). \quad (57)$$

One may hastily conclude that the asymptotic limit of (55) is

$$\begin{cases} -\frac{d}{dx}(M(\mathbf{A})\frac{du}{dx}) = f & \text{in } \Omega, \\ u \in H_0^1(\Omega). \end{cases} \quad (58)$$

But this is not the case in general, since the weak limit of the product of two sequences that are only weakly convergent is not the product of their individual weak limits. Here, the role of G-convergence theory comes in, it gives a strategy of identifying the correct limit of the problem.

In order to get the correct limit problem, we define

$$\xi_h = \mathbf{A}_h(x)\frac{du_h}{dx}. \quad (59)$$

By the boundedness of  $\mathbf{A}_h$  and the estimate  $\|u_h\|_{H_0^1(\Omega)} \leq C$ ,  $\xi_h$  is uniformly bounded in  $L^2(\Omega)$ . Since  $-\frac{d\xi_h}{dx} = f \in L^2(\Omega)$ , we conclude that  $\xi_h$  is uniformly bounded in  $H_0^1(\Omega)$ . By the compact embedding of  $H_0^1(\Omega)$  in  $L^2(\Omega)$ , up to a subsequence still denoted by  $\xi_h$ , we have

$$\xi_h \rightarrow \xi \text{ in } L^2(\Omega), \quad (60)$$

for some  $\xi$ , consequently

$$\frac{d\xi_h}{dx} \rightharpoonup \frac{d\xi}{dx} \text{ in } L^2(\Omega). \quad (61)$$

Note that  $0 < \frac{1}{\beta} < \frac{1}{\mathbf{A}_h(x)} < \frac{1}{\alpha} < \infty$ , hence

$$\frac{1}{\mathbf{A}_h} \overset{*}{\rightharpoonup} M\left(\frac{1}{\mathbf{A}}\right) \text{ in } L^\infty(\Omega), \text{ (also weakly in } L^2(\Omega)). \quad (62)$$

Since  $\frac{1}{\mathbf{A}_h(x)}\xi_h = \frac{du_h}{dx}$ , by (56) and (62) one gets

$$\frac{du}{dx} = M\left(\frac{1}{\mathbf{A}}\right)\xi. \quad (63)$$

By the result (61) together with (63) and a limit passage of  $-\frac{d\xi_h}{dx} = f$ , we conclude that  $u$  is the solution to the limit problem

$$\begin{cases} -\frac{d}{dx}\left(\frac{1}{M\left(\frac{1}{\mathbf{A}}\right)}\frac{du}{dx}\right) = f & \text{in } \Omega, \\ u \in H_0^1(\Omega). \end{cases} \quad (64)$$

By the uniqueness of the solution  $u$  to (64), and using Urysohn property, it follows that the whole sequence  $u_h$  converges weakly to  $u$ . The above conclusion can be summarized as  $\mathbf{A}_h$  G-converges to  $(1/M(1/\mathbf{A}))$  which is known as the harmonic mean of  $A_h$ .

It is important to point out that in the previous example the weak\* limit of  $\frac{1}{\mathbf{A}_h}$  characterizes the limit problem. This is only true for one dimensional problem, and it is not the case in higher dimensions ( $\mathbb{R}^N$ ,  $n \geq 2$ ). For more discussion on this issue we refer to [36].

### 4.1.3 The definition

Let  $\alpha$  and  $\beta$  be two real numbers such that  $0 < \alpha \leq \beta < \infty$ , and let  $S(\alpha, \beta, \Omega)$  be defined as  $S(\alpha, \beta, \Omega) = \{\mathbf{A} \in L^\infty(\Omega)^{N \times N}; (\mathbf{A}(x, \xi), \xi) \geq \alpha|\xi|^2 \text{ and } |\mathbf{A}(x, \xi)| \leq \beta|\xi|, \forall \xi \in \mathbb{R}^N \text{ and a.e. } x \in \Omega\}$ .

**Definition 1** The sequence  $\mathbf{A}_h \subset S(\alpha, \beta, \Omega)$  is said to be G-convergent to  $\mathbf{A} \in S(\alpha, \beta, \Omega)$ , denoted by  $\mathbf{A}_h \xrightarrow{G} \mathbf{A}$ , if for every  $f \in H^{-1}(\Omega)$ , the sequence  $u_h$  of solutions to the equation

$$\begin{cases} -\operatorname{div}(\mathbf{A}_h(x, Du_h)) = f & \text{in } \Omega, \\ u_h \in H_0^1(\Omega) \end{cases} \quad (65)$$

satisfies

$$\begin{aligned} u_h &\rightharpoonup u \text{ in } H_0^1(\Omega), \\ \mathbf{A}_h(\cdot, Du_h) &\rightharpoonup \mathbf{A}(\cdot, Du) \text{ in } L^2(\Omega)^N, \end{aligned}$$

where  $u$  is the unique solution of the problem

$$\begin{cases} -\operatorname{div}(\mathbf{A}(x, Du)) = f & \text{in } \Omega, \\ u \in H_0^1(\Omega). \end{cases} \quad (66)$$



G-convergence possesses the compactness property, i.e., if  $\mathbf{A}_h \in S(\alpha, \beta, \Omega)$ , then there exists a subsequence, denoted by  $\mathbf{A}_h$ , and  $\mathbf{A} \in S(\alpha, \beta, \Omega)$ , such that  $\mathbf{A}_h \xrightarrow{G} \mathbf{A}$ . The G-limit is unique and local, also if  $\mathbf{A}_h \xrightarrow{G} \mathbf{A}$ , then  $\mathbf{A}_h^t \xrightarrow{G} \mathbf{A}^t$ , here  $t$  denotes the transpose operator.

#### 4.1.4 Convergence of elliptic eigenvalue problems

For elliptic boundary value problems with source function  $f_h$  we have the following result, see e.g. [15] and Paper IV in the appendix.

**Theorem 3** Consider the Dirichlet boundary value problem

$$\begin{cases} -\operatorname{div}(\mathbf{A}_h(x)\nabla u_h) = f_h & \text{in } \Omega, \\ u_h \in H_0^1(\Omega). \end{cases} \quad (67)$$

If  $\mathbf{A}_h \in S(\alpha, \beta, \Omega)$  and  $f_h$  converges in  $H^{-1}(\Omega)$  to  $f$ , then the sequence  $u_h$  of solutions to (67) is weakly convergent in  $H_0^1(\Omega)$  to the solution of the problem

$$\begin{cases} -\operatorname{div}(\mathbf{A}(x)\nabla u) = f & \text{in } \Omega, \\ u \in H_0^1(\Omega), \end{cases} \quad (68)$$

where  $\mathbf{A}$  is the G-limit of  $\mathbf{A}_h$ .

The strength of G-convergence can be evidently seen by applying the concept to elliptic eigenvalue problems. Consider the linear elliptic eigenvalue problem

$$\begin{cases} -\operatorname{div}(\mathbf{A}_h(x)\nabla u_h^k) = \lambda_h^k u_h^k & \text{in } \Omega, \\ u_h^k \in H_0^1(\Omega), \end{cases} \quad (69)$$

where  $\mathbf{A}_h \in S(\alpha, \beta, \Omega)$  is symmetric and positive definite. Then, the set of eigenvalues  $\{\lambda_h^k\}$  is bounded and  $0 < \lambda_h^1 \leq \lambda_h^2 \leq \lambda_h^3 \leq \dots$ , also the multiplicity of each  $\lambda_h^k$  is finite. For the eigenvalue problem (69), the following result is formulated.

**Theorem 4** The sequences of eigenvalues  $\lambda_h^k$  and the corresponding eigenfunctions  $u_h^k$  of (69) converge to  $\lambda^k$  in  $\mathbb{R}$  and weakly to  $u^k$  in  $H_0^1(\Omega)$  respectively, where the eigencouple  $\{\lambda^k, u^k\}$  is the solution to the G-limit problem

$$\begin{cases} -\operatorname{div}(\mathbf{A}(x)\nabla u^k) = \lambda^k u^k & \text{in } \Omega, \\ u^k \in H_0^1(\Omega). \end{cases} \quad (70)$$

## 4.2 Positive definite self-adjoint operators

### 4.2.1 The definition

Let  $\mathcal{H}$  be a Hilbert space and let  $\lambda \geq 0$  be a real number, by  $\mathcal{P}_\lambda(\mathcal{H})$  we denote the class of self-adjoint operators  $A$  on a closed linear subspace  $\overline{\mathbf{D}(A)}$  of  $\mathcal{H}$  such that  $\langle Au, u \rangle \geq \lambda \|u\|_{\mathcal{H}}^2$ ,  $\forall u \in \mathbf{D}(A)$ .

**Definition 2** Let  $\lambda \geq 0$ , and let  $A_h \in \mathcal{P}_\lambda(\mathcal{H})$ . If  $\lambda > 0$ , we say that  $A_h \xrightarrow{G} A \in \mathcal{P}_\lambda(\mathcal{H})$  in  $\mathcal{H}$  if  $A_h^{-1}P_h u \rightarrow A^{-1}Pu$  in  $\mathcal{H}$ ,  $\forall u \in \mathcal{H}$ , where  $P_h$  and  $P$  are the orthogonal projections onto  $\overline{\mathbf{D}(A_h)}$  and  $\overline{\mathbf{D}(A)}$  respectively. If  $\lambda = 0$ , we say that  $A_h \in \mathcal{P}_0(\mathcal{H})$  converges to  $A \in \mathcal{P}_0(\mathcal{H})$  in the strong resolvent sense (SRS) if  $(\mu I + A_h) \xrightarrow{G} (\mu I + A)$  in  $\mathcal{H}$ ,  $\forall \mu > 0$ .

G-convergence of positive definite self-adjoint operators can be studied using  $\Gamma$ -convergence of the corresponding quadratic forms [11], where, generally, proving  $\Gamma$ -limits is simpler than proving G-limits. Below we define  $\Gamma$ -convergence and discuss its relation to G-convergence. First we need the following definitions.

**Definition 3** A function  $F : \mathcal{H} \rightarrow [0, \infty]$  is said to be lower semi-continuous (*lsc*) at  $u \in \mathcal{H}$ , if

$$F(u) \leq \sup_{U \in \mathcal{N}(u)} \inf_{v \in U} F(v),$$

where  $\mathcal{N}(u)$  is the set of all open neighborhoods of  $u$  in  $\mathcal{H}$ .

**Definition 4** A function  $F$  in  $\mathcal{H}$  is called a quadratic form if there exists a linear dense subspace  $\mathcal{X}$  of  $\mathcal{H}$  and a symmetric bilinear form  $B : \mathcal{X} \times \mathcal{X} \rightarrow [0, \infty)$  such that

$$F(u) = \begin{cases} B(u, u), & \forall u \in \mathcal{X}, \\ \infty, & \forall u \in \mathcal{H} \setminus \mathcal{X}. \end{cases}$$

Let  $F$  and  $B$  be as in the above definition, where  $\mathbf{D}(F) = \{u \in \mathcal{H}; F(u) < \infty\}$ . Then the operator associated to  $F$  is the linear operator  $A$  on  $\overline{\mathbf{D}(F)}$  with the domain being the set of all  $u \in \mathbf{D}(F)$  such that there exists  $v \in \overline{\mathbf{D}(F)}$  satisfying  $B(u, f) = \langle v, f \rangle$ ,  $\forall f \in \mathbf{D}(F)$  and  $Au = v$ ,  $\forall u \in \mathbf{D}(A)$ . If  $f = u$  then  $F(u) = \langle Au, u \rangle$ ,  $\forall u \in \mathbf{D}(A)$ .

Let  $\lambda \geq 0$ , by  $\tilde{\mathcal{Q}}_\lambda(\mathcal{H})$  we denote the class of quadratic forms  $F : \mathcal{H} \rightarrow [0, \infty]$  such that  $F(u) \geq \lambda \|u\|_{\mathcal{H}}^2$ . And by  $\mathcal{Q}_\lambda(\mathcal{H})$  we denote the subset of  $\tilde{\mathcal{Q}}_\lambda(\mathcal{H})$  whose elements are *lsc*.

**Definition 5** A sequence of functionals  $F_h : \mathcal{H} \rightarrow \overline{\mathbb{R}}$  is said to  $\Gamma$ -converge to  $F : \mathcal{H} \rightarrow \overline{\mathbb{R}}$ , written as  $F(u) = \Gamma - \lim_{h \rightarrow \infty} F_h(u)$  and denoted by  $F_h \xrightarrow{\Gamma} F$  if

$$F(u) = \Gamma - \liminf_{h \rightarrow \infty} F_h(u) = \Gamma - \limsup_{h \rightarrow \infty} F_h(u),$$

where  $\Gamma - \liminf_{h \rightarrow \infty} F_h(u) = \sup_{U \in \mathcal{N}(u)} \liminf_{h \rightarrow \infty} \inf_{v \in U} F_h(v)$  and  $\Gamma - \limsup_{h \rightarrow \infty} F_h(u) = \sup_{U \in \mathcal{N}(u)} \limsup_{h \rightarrow \infty} \inf_{v \in U} F_h(v)$ .

Note that if  $\mathcal{H}$  satisfies the first axiom of countability (the neighborhood system of every point in  $\mathcal{H}$  has a countable base), then  $F_h \xrightarrow{\Gamma} F$  in  $\mathcal{H}$  if and only if the following two conditions (called respectively the lim inf-inequality and lim-equality) are satisfied

$$(i) \quad \forall u \in \mathcal{H} \text{ and } \forall u_h \text{ converging to } u, F(u) \leq \liminf_{h \rightarrow \infty} F_h(u_h).$$

$$(ii) \quad \forall u \in \mathcal{H}, \exists u_h \text{ converging to } u \text{ such that } F(u) = \lim_{h \rightarrow \infty} F_h(u_h).$$

It is worth to mention that  $\Gamma$ -limit is always *lsc* and unique, also  $\Gamma$ -limit of non-negative quadratic form is a non-negative quadratic form.  $\Gamma$ -convergence possesses the compactness property, that is, if  $\mathcal{H}$  is a separable metric space, then every sequence  $F_h : \mathcal{H} \rightarrow \overline{\mathbb{R}}$  has a  $\Gamma$ -convergent subsequence.

The following theorem demonstrates the relation between G-convergence of operators of the class  $\mathcal{P}_\lambda(\mathcal{H})$  for  $\lambda \geq 0$  and  $\Gamma$ -convergence of the associated quadratic forms of the class  $\mathcal{Q}_\lambda(\mathcal{H})$ .

**Theorem 5** Let  $F_h$  and  $F$  be elements of  $\mathcal{Q}_0(\mathcal{H})$ , and let  $A_h, A \in \mathcal{P}_0(\mathcal{H})$  be the associated operators respectively. Then  $F_h \xrightarrow{\Gamma} F$  if and only if  $A_h \rightarrow A$  in the SRS. Also, for  $\mu > 0$ , if  $F_h, F \in \mathcal{Q}_\mu(\mathcal{H})$ , and  $A_h, A \in \mathcal{P}_\mu(\mathcal{H})$  are the associated operators respectively, then  $F_h \xrightarrow{\Gamma} F$  if and only if  $A_h \xrightarrow{G} A$ .

#### 4.2.2 G-convergence of positive definite self-adjoint operators

Let  $H_0$  be a positive definite bounded self-adjoint operator defined on  $L^2(\Omega)$  and let  $\mathbf{D}(H_0) = H_0^{-1}(\Omega)$ . Consider the perturbed operator  $H_h = H_0 + V_h$  where  $V_h(x)$  is a positive bounded real-valued multiplication operator in  $L^2(\Omega)$ . Using G-convergence together with  $\Gamma$ -convergence, we state the following results, see Paper IV in the appendix.

**Theorem 6** Let  $V_h$  be a sequence in  $L^\infty(\Omega)$  that converges weakly\* to  $V$ , then  $H_h$  G-converges to  $H = H_0 + V$ .

**Theorem 7** If  $V_h$  is a weakly convergent sequence in  $L^p(\Omega)$  for  $2 \leq p < \infty$  with a weak limit denoted by  $V$ , then  $H_h$  G-converges to  $H = H_0 + V$ .

Let  $\mathcal{X}$  and  $\mathcal{H}$  be two Hilbert spaces, and let  $\mathcal{B}(\mathcal{H})$  be the set of bounded linear operators on  $\mathcal{H}$ . As a generalization of Theorem 4, below we state the relation between the eigenvalue problems of an operator and its G-limit of the class  $\mathcal{P}_\lambda(\mathcal{H})$  for  $\lambda \geq 0$ , see Paper III in the appendix.

**Theorem 8** Let  $\lambda > 0$ , let  $A_h$  be a sequence in  $\mathcal{P}_\lambda(\mathcal{H})$  G-converging to  $A \in \mathcal{P}_\lambda(\mathcal{H})$ , and let  $\{\mu_h, u_h\}$  be the solution of the eigenvalue problem  $A_h u_h = \mu_h u_h$ . If  $\{\mu_h, u_h\} \rightarrow \{\mu, u\}$  in  $\mathbb{R} \times \mathcal{H}$ , then the limit couple  $\{\mu, u\}$  is the solution of the eigenvalue problem  $Au = \mu u$ .

It is clear that the assertion of Theorem 8 is also true if the sequence  $A_h \in \mathcal{P}_0(\mathcal{H})$  is convergent in the SRS to  $A \in \mathcal{P}_0(\mathcal{H})$ .

Note that if a sequence  $A_h$  is convergent in the SRS (or strongly convergent) to  $A$ , then every  $\lambda \in \sigma(A)$  is the limit of a sequence  $\lambda_h \in \sigma(A_h)$ , but not the limit of every sequence  $\lambda_h \in \sigma(A_h)$  lies in the spectrum of  $A$ , see [55]. Despite of this fact, the following theorem provides conditions by which G-convergence of an operator in  $\mathcal{P}_\lambda(\mathcal{H})$  (consequently the strong resolvent convergence in  $\mathcal{P}_0(Y)$ ) implies the convergence of the corresponding eigenvalues, see Paper III in the appendix.

**Theorem 9** Let  $\mathcal{X}$  be compactly and densely embedded in  $\mathcal{H}$ , and let  $A_h$  be a family of operators in  $\mathcal{P}_\lambda(\mathcal{H})$ ,  $\lambda > 0$ , with domain  $\mathcal{X}$ . If  $A_h$  G-converges to  $A \in \mathcal{P}_\lambda(\mathcal{H})$ , then  $A_h^{-1}$  converges in the norm of  $\mathcal{B}(\mathcal{H})$  to  $A^{-1}$ . Moreover, the  $k^{\text{th}}$  eigenvalue  $\mu_h^k$  of  $A_h$  converges to the  $k^{\text{th}}$  eigenvalue  $\mu^k$  of  $A$ ,  $\forall k \in \mathbb{N}$ .

Theorem 9 implies that, for those perturbations considered in Theorems 6 and 7, the eigenvalues of  $H_h$  converge to the eigenvalues of the G-limit operator  $H$ . Moreover, Theorem 8 guarantees that the eigenvalue problem  $H_h u_h = \mu_h u_h$  converges to the limit problem  $Hu = \mu u$ , where  $u$  is the limit of  $u_h$  in  $L^2(\Omega)$ .

**Remark 2** Let  $E^{H_h}$  and  $E^H$  be the spectral measures of  $H_h$  and  $H$  respectively, then G-convergence of  $H_h$  to  $H$  implies that  $E^{H_h}(\lambda) \rightarrow E^H(\lambda)$  strongly for all  $\lambda \in \mathbb{R}$  such that  $E^H(\lambda) = E^H(-\lambda)$ .

### 4.3 Families of Dirac operators

Here we consider an  $h$ -dependent perturbation added to the Dirac operator with Coulomb potential. The purpose is to apply G-convergence theory for positive definite parts of the perturbed operator and to investigate the asymptotic behavior of the corresponding eigenvalues in the gap.

### 4.3.1 The Dirac operator with perturbation ( $\tilde{\mathcal{H}}_h$ )

Let  $\mathcal{H}_h$  be defined as

$$\mathcal{H}_h = \mathbf{H} + V_h, \quad (71)$$

where  $V_h = V_h(x)$  is a  $4 \times 4$  matrix-valued function and, as defined before,  $\mathbf{H} = \mathbf{H}_0 + V$ , where again  $\mathbf{H}_0$  and  $V$  are respectively the free Dirac operator and the Coulomb potential. We recall here the spaces  $X = H^1(\mathbb{R}^3, \mathbb{C}^4)$  and  $Y = L^2(\mathbb{R}^3, \mathbb{C}^4)$ .

Recall that a function  $F$  is called homogeneous of degree  $p$  if for any nonzero scalar  $a$ ,  $F(ax) = a^p F(x)$ . The next theorem is of profound importance [54, 56].

**Theorem 10** Let, for  $h > 0$ ,  $V_h$  be a measurable  $(-1)$ -homogeneous Hermitian  $4 \times 4$  matrix-valued function with entries in  $L^p_{loc}(\mathbb{R}^3)$ ,  $p > 3$ . Then  $\mathcal{H}_h$  is essentially self-adjoint on  $C_0^\infty(\mathbb{R}^3; \mathbb{C}^4)$  and self-adjoint on  $X$ . Moreover,  $\sigma(\mathcal{H}_h) = (-\infty, -mc^2] \cup \{\lambda_h^k\}_{k \in \mathbb{N}} \cup [mc^2, +\infty)$ , where  $\{\lambda_h^k\}_{k \in \mathbb{N}}$  is a discrete sequence of  $h$ -dependent eigenvalues corresponding to the Dirac eigenvalue problem  $\mathcal{H}_h u_h = \lambda_h u_h$ .

We assume further that the  $4 \times 4$  matrix-valued function  $V_h$  is of the form  $V_h(x) = V_1(x)V_2(hx)$ , where  $V_1$  is  $(-1)$ -homogeneous and where the entries of  $V_2(y)$  are 1-periodic in  $y$ , i.e.,

$$V_2^{ij}(y+k) = V_2^{ij}(y), \quad k \in \mathbb{Z}^3.$$

We also assume that the entries of  $V_2$  belong to  $L^\infty(\mathbb{R}^3)$ . It is then well-known that

$$V_2^{ij}(hx) \ast \rightarrow M(V_2^{ij}) = \int_{\mathbb{T}^3} V_2^{ij}(y) dy, \quad \text{in } L^\infty(\mathbb{R}^3), \quad (72)$$

where  $\mathbb{T}^3$  is the unit torus in  $\mathbb{R}^3$ . It easily follows from this mean-value property that

$$V_h \rightarrow V_1 M(V_2), \quad \text{in } L^p(\mathbb{R}^3), \quad p > 3.$$

In the sequel, we consider a shifted family of Dirac operators denoted by  $\tilde{\mathcal{H}}_h$  and defined as  $\tilde{\mathcal{H}}_h = \tilde{\mathbf{H}} + V_h$ , where  $\tilde{\mathbf{H}} = \mathbf{H} + mc^2 \mathbf{I}$ . Also without loss of generality we set  $\tilde{h} = c = m = 1$ . By Theorem 10, for  $h > 0$ , we then get  $\sigma(\tilde{\mathcal{H}}_h) = (-\infty, 0] \cup \{\tilde{\lambda}_h^k\}_{k \in \mathbb{N}} \cup [2, \infty)$ .

### 4.3.2 G-convergence of projected parts of $\tilde{\mathcal{H}}_h$

Let  $X$  and  $Y$  be defined as before, and let  $E^{\tilde{\mathcal{H}}_h}$  and  $E^{\tilde{\mathbf{H}}}$  be the spectral measures of the families  $\tilde{\mathcal{H}}_h$  and  $\tilde{\mathbf{H}}$  respectively, by the spectral theorem

$$\tilde{\mathcal{H}}_h = \int_{\sigma(\tilde{\mathcal{H}}_h)} \lambda dE^{\tilde{\mathcal{H}}_h}(\lambda). \quad (73)$$

Define  $X_h^p = \bigoplus_{k \in \mathbb{N}} \mathcal{N}_h^k$  where  $\mathcal{N}_h^k = \{u_h \in X; \tilde{\mathcal{H}}_h u_h = \lambda_h^k u_h\}$ . Note that  $X_h^p$  is a closed subspace of  $Y$  invariant with respect to  $\tilde{\mathcal{H}}_h$ . Then we have the following theorem, see Paper III in the appendix.

**Theorem 11** Let  $E^{\tilde{\mathcal{H}}_h, p}$  be the point measure of  $\tilde{\mathcal{H}}_h$ , and consider the restriction  $\tilde{\mathcal{H}}_h^p$  of  $\tilde{\mathcal{H}}_h$  to  $X_h^p$  defined as

$$\tilde{\mathcal{H}}_h^p = \sum_{\lambda \in \sigma_p(\tilde{\mathcal{H}}_h)} \lambda E^{\tilde{\mathcal{H}}_h, p}(\lambda). \quad (74)$$

The operator  $\tilde{\mathcal{H}}_h^p$  is positive definite and self-adjoint on  $X$  with compact inverse  $(\tilde{\mathcal{H}}_h^p)^{-1}$ . Then there exists a positive definite self-adjoint operator  $\tilde{\mathcal{H}}^p$  such that, up to a subsequence,  $\tilde{\mathcal{H}}_h^p$  G-converges to  $\tilde{\mathcal{H}}^p$ . The operator  $\tilde{\mathcal{H}}^p$  is given by  $(\tilde{\mathbf{H}} + V_1 M(V_2))|_{X^p}$ , where  $\mathbf{D}(\tilde{\mathcal{H}}^p) = X^p = \bigoplus_{k \in \mathbb{N}} \mathcal{N}^k$  and  $\mathcal{N}^k = \{u \in X; \tilde{\mathcal{H}}^p u = \lambda^k u\}$ .

Now we can apply Theorem 9 to conclude that the sequence of  $k^{th}$  eigenvalues  $\lambda_h^k$  associated to  $\tilde{\mathcal{H}}_h$  converges to the  $k^{th}$  eigenvalue  $\lambda^k$  of  $\tilde{\mathcal{H}}^p$ .

For the absolutely continuous part of the operator  $\tilde{\mathcal{H}}_h$ , we let first  $X_h^{ac} = X_h^{ac,+} \oplus X_h^{ac,-}$ , where  $X_h^{ac,+}$  and  $X_h^{ac,-}$  are the closed subspaces, invariant with respect to  $\tilde{\mathcal{H}}_h$ , corresponding respectively to the absolutely continuous spectra  $\sigma_{ac}^+(\tilde{\mathcal{H}}_h) = [2, +\infty)$  and  $\sigma_{ac}^-(\tilde{\mathcal{H}}_h) = (-\infty, 0]$ . Let  $E^{\tilde{\mathcal{H}}_h, ac,+}(\lambda)$  be the absolutely continuous spectral measure corresponding to  $\tilde{\mathcal{H}}_h^{ac,+}$  and define

$$\tilde{\mathcal{H}}_h^{ac,+} = \int_{\lambda \in \sigma_{ac}^+(\tilde{\mathcal{H}}_h)} \lambda dE^{\tilde{\mathcal{H}}_h, ac,+}(\lambda). \quad (75)$$

By this construction, the operator  $\tilde{\mathcal{H}}_h^{ac,+}$  is the restriction of  $\tilde{\mathcal{H}}_h$  to  $X_h^{ac,+}$ , thus it is positive definite and self-adjoint on  $X$ . Therefore, there exists a subsequence of  $\tilde{\mathcal{H}}_h^{ac,+}$ , still denoted by  $\tilde{\mathcal{H}}_h^{ac,+}$ , which G-converges to a positive definite self-adjoint operator  $\tilde{\mathcal{H}}^{ac,+}$ . Moreover, convergence in the SRS can be drawn for  $-\tilde{\mathcal{H}}_h^{ac,-}$ ,

$$\tilde{\mathcal{H}}_h^{ac,-} = \int_{\lambda \in \sigma_{ac}^-(\tilde{\mathcal{H}}_h)} \lambda dE^{\tilde{\mathcal{H}}_h, ac,-}(\lambda), \quad (76)$$

where  $E^{\tilde{\mathcal{H}}_h, ac,-}(\lambda)$  is the absolutely continuous spectral measure corresponding to the operator  $\tilde{\mathcal{H}}_h^{ac,-}$ .

## 5 The wave operators for $\hbar$ -dependent self-adjoint operators

Scattering theory is a frame for comparing the dynamic behaviors of two quantum systems, and is well-known as perturbation theory of self-adjoint operators on the absolutely continuous spectrum. More specifically, scattering theory concerns studying the behavior, for large times, of the absolutely continuous solution of the convolution equation  $i\partial u/\partial t = Hu = (H_0 + \text{Interaction})u$  in terms of the absolutely continuous solution of the simple convolution equation  $i\partial u_0/\partial t = H_0u_0$ . Here  $H_0$  and  $H$  are self-adjoint operators acting on Hilbert spaces  $\mathcal{H}_0$  and  $\mathcal{H}$  respectively. That is, for a given initial solution  $f$  to the equation with interaction above, if  $f$  is an eigenvector corresponding to an eigenvalue  $\mu$ , then  $u(t) = \exp(-i\mu t)f$ , so that the time behavior is clear. But if  $f \in \mathcal{H}^{(ac)}$  (the absolutely continuous subspace of  $H$ ), it is not possible, in general, to calculate  $u(t)$  explicitly. Using scattering theory, one may study the asymptotic behavior of  $u(t) = \exp(-iHt)f$  as  $t \rightarrow \pm\infty$ ,  $f \in \mathcal{H}^{(ac)}$ , in terms of  $u_0(t) = \exp(-iH_0t)f_0$  for  $f_0 \in \mathcal{H}_0^{(ac)}$  (the absolutely continuous subspace of  $H_0$ ).

### 5.1 A simple overview

Consider a self-adjoint operator  $H_0$  in a Hilbert space  $\mathcal{H}_0$ , and assume that its absolutely continuous spectrum can be identified. Let  $H$  be another self-adjoint operator in a Hilbert space  $\mathcal{H}$  so that  $H$  is close to  $H_0$  in a certain sense. Scattering theory concerns the study of the absolutely continuous spectrum of the operator  $H$  and its connection to that of  $H_0$ . It is generally assumed that  $H = H_0 + V$ , where  $V$  is, in a particular measure, small compared to  $H_0$ , and thus the deduction of the spectral properties of the absolutely continuous spectrum of  $H$  depends on the presumed knowledge of the absolutely continuous spectrum of  $H_0$ .

Consider the free evolution problem

$$\begin{cases} i\frac{\partial}{\partial t}u_0(x, t) = H_0u_0(x, t), \\ u_0(x, 0) = u_0^0(x) \end{cases} \quad (77)$$

which has the solution  $u_0(t) = e^{-iH_0t}u_0^0$ . Let now

$$\begin{cases} i\frac{\partial}{\partial t}u(x, t) = Hu(x, t), \\ u(x, 0) = u^0(x) \end{cases} \quad (78)$$

be the evolution problem of the perturbed operator  $H = H_0 + V$ , which has the solution  $u(t) = e^{-iHt}u^0$ . The main task of scattering theory is to study the conditions under which, for all  $u^0 \in \mathcal{H}^{(ac)}$ , there exist  $u_0^{0,\pm} \in \mathcal{H}_0^{(ac)}$ , such that

$$\lim_{t \rightarrow \pm\infty} \|u(t) - \mathcal{J}u_0(t)\|_{\mathcal{H}} = 0, \quad (79)$$

for a bounded operator  $\mathcal{J}$ , where  $u_0(t) = e^{-iH_0t}u_0^{0,\pm}$ . Equivalently, scattering theory concerns the study of existence and completeness of the wave operator (WO)  $W_{\pm}(H, H_0; \mathcal{J})$ ,

$$W_{\pm}(H, H_0; \mathcal{J}) = s\text{-}\lim_{t \rightarrow \pm\infty} e^{iHt}\mathcal{J}e^{-iH_0t}u_0^{0,\pm}, \quad (80)$$

where the letter  $s$  refers to the strong sense convergence.

For comprehensive materials on scattering theory we refer to the monographs [40, 57]. Following the general notation in scattering theory, below we use  $s\text{-}\lim$  and  $w\text{-}\lim$  to denote the strong and weak limits respectively. Let  $H$  and  $H_0$  be self-adjoint operators in  $\mathcal{H}$  and  $\mathcal{H}_0$  with spectral families  $E$  and  $E_0$  respectively, below we define the time-dependent and stationary WOs.

## 5.2 The time-dependent WO

There are two time-dependent WOs, the strong and weak WOs. In what follows, the strong time-dependent WO will be referred to just as WO.

### 5.2.1 The strong time-dependent WO

The (modified or generalized) strong time-dependent WO  $W_{\pm}$  is defined as follows

**Definition 6** Let  $\mathcal{J} : \mathcal{H}_0 \rightarrow \mathcal{H}$  be a bounded operator (identification), the WO  $W_{\pm} = W_{\pm}(H, H_0; \mathcal{J})$  for  $H$  and  $H_0$  is the operator

$$W_{\pm}(H, H_0; \mathcal{J}) = s\text{-}\lim_{t \rightarrow \pm\infty} U(-t)\mathcal{J}U_0(t)P_0^{(ac)}, \quad (81)$$

provided that the corresponding strong limits exist ( $s$  refers to the strong sense convergence), where  $P_0^{(ac)}$  is the orthogonal projection onto the absolutely continuous subspace  $\mathcal{H}_0^{(ac)}$  of  $H_0$ ,  $U(t) = e^{-iHt}$ , and  $U_0(t) = e^{-iH_0t}$ . If  $\mathcal{H} = \mathcal{H}_0$  and  $\mathcal{J}$  is the identity operator, then the WO is denoted by  $W_{\pm}(H, H_0)$ .



The WO  $W_{\pm} = W_{\pm}(H, H_0; \mathcal{J})$  is bounded, and possesses the intertwining property, that is, for any bounded Borel function  $\phi$ ,

$$\phi(H)W_{\pm}(H, H_0; \mathcal{J}) = W_{\pm}(H, H_0; \mathcal{J})\phi(H_0), \quad (82)$$

also for any Borel set  $\Delta \subset \mathbb{R}$ ,

$$E(\Delta)W_{\pm}(H, H_0; \mathcal{J}) = W_{\pm}(H, H_0; \mathcal{J})E_0(\Delta). \quad (83)$$

The WO  $W_{\pm}$  admits the chain rule, i.e., if  $W_{\pm}(H, H_1; \mathcal{J}_1)$  and  $W_{\pm}(H_1, H_0; \mathcal{J}_0)$  exist, then the WO  $W_{\pm}(H, H_0; \mathcal{J}_{1,0}) = W_{\pm}(H, H_1; \mathcal{J}_1)W_{\pm}(H_1, H_0; \mathcal{J}_0)$  also exists, where  $\mathcal{J}_{1,0} = \mathcal{J}_1\mathcal{J}_0$ .

Note that the operator  $W_{\pm}(H, H_0)$  is isometric. To prove that  $W_{\pm}(H, H_0; \mathcal{J})$  is isometric, it is equivalent to prove that for any  $u \in \mathcal{H}_0^{(ac)}$ ,

$$\lim_{t \rightarrow \pm\infty} \|\mathcal{J}U_0(t)u\|_{\mathcal{H}} = \|u\|_{\mathcal{H}_0}.$$

The following remark states the equivalence between WOs with different identifications.

**Remark 3** Assume that, with an identification  $\mathcal{J}_1$ , the WO  $W_{\pm}(H, H_0; \mathcal{J}_1)$  exists, and suppose that  $\mathcal{J}_2$  is another identification such that  $\mathcal{J}_1 - \mathcal{J}_2$  is compact, then the WO  $W_{\pm}(H, H_0; \mathcal{J}_2)$  exists and  $W_{\pm}(H, H_0; \mathcal{J}_1) = W_{\pm}(H, H_0; \mathcal{J}_2)$ . Moreover, the condition that  $\mathcal{J}_1 - \mathcal{J}_2$  is compact can be replaced by  $s\text{-}\lim_{t \rightarrow \pm\infty} (\mathcal{J}_1 - \mathcal{J}_2)U_0(t)P_0^{(ac)} = 0$ .

Assume the existence of the WO  $W_{\pm}$ , another task that is not less important is to show the completeness of  $W_{\pm}$ .

**Definition 7** The WO  $W_{\pm}$  is said to be complete if  $\mathbf{R}(W_{\pm}) = \mathcal{H}^{(ac)}$ .

If the WO  $W_{\pm}$  is complete, then the absolutely continuous operators  $H^{(ac)}$  and  $H_0^{(ac)}$  are unitary equivalent. Since, by the chain rule,  $P^{(ac)} = W_{\pm}(H, H) = W_{\pm}(H, H_0)W_{\pm}^*(H, H_0)$  where  $P^{(ac)}$  is the orthogonal projection onto the absolutely continuous subspace  $\mathcal{H}^{(ac)}$  of  $H$ , to prove the completeness of the WO  $W_{\pm}(H, H_0)$  is equivalent to prove the existence of the WO  $W_{\pm}^*(H, H_0) = W_{\pm}(H_0, H)$ . On the other hand, the completeness of the WO  $W_{\pm}(H, H_0; \mathcal{J})$  is equivalent to the existence of  $W_{\pm}(H_0, H; \mathcal{J}^*)$  and that the identification  $\mathcal{J}$  is boundedly invertible.

### 5.2.2 The weak time-dependent WO

The weak time-dependent WO  $\widetilde{W}_\pm$  is defined as follows

**Definition 8** Let  $\mathcal{J} : \mathcal{H}_0 \rightarrow \mathcal{H}$  be a bounded identification, the weak WO  $\widetilde{W}_\pm(H, H_0; \mathcal{J})$  for  $H$  and  $H_0$  is the operator

$$\widetilde{W}_\pm(H, H_0; \mathcal{J}) = w\text{-}\lim_{t \rightarrow \pm\infty} P^{(ac)} U(-t) \mathcal{J} U_0(t) P_0^{(ac)}, \quad (84)$$

provided that the corresponding weak limits exist ( $w$  refers to the weak sense convergence).

Note that the boundedness and intertwining properties of the WO  $W_\pm$  are preserved for the WO  $\widetilde{W}_\pm$ , whereas the chain rule property is not valid for the weak WO. This is evident since the weak limit of the product of two sequences that are only weakly convergent is not necessarily the product of their weak limits, even more this weak limit does not necessarily exist. On contrast to  $W_\pm$ , if the weak WO  $\widetilde{W}_\pm(H, H_0; \mathcal{J})$  exists, then it is necessarily that  $\widetilde{W}_\pm(H_0, H; \mathcal{J}^*)$  also exists.

### 5.3 The stationary WO

Let  $R(z)$  and  $R_0(z)$  be the resolvent operators of  $H$  and  $H_0$  respectively, and let  $M_0$  and  $M$  be dense sets in  $\mathcal{H}_0$  and  $\mathcal{H}$  respectively.

Let  $\epsilon > 0$ , and let  $\theta(\lambda, \epsilon)$  be defined as

$$\theta(\lambda, \epsilon) = (2\pi i)^{-1} (R(\lambda + i\epsilon) - R(\lambda - i\epsilon)) = \pi^{-1} \epsilon R(\lambda + i\epsilon) R(\lambda - i\epsilon). \quad (85)$$

Further, let  $\mathfrak{H}$  be an auxiliary Hilbert space, the concept  $H$ -smoothness in the strong and weak senses is defined as follows

**Definition 9** An  $H$ -bounded operator,  $A : \mathcal{H} \rightarrow \mathfrak{H}$ , is called  $H$ -smooth (in the strong sense) if one of the following bounds is satisfied

$$\sup_{\|v\|_{\mathcal{H}}=1, v \in \mathbf{D}(H)} \int_{-\infty}^{\infty} \|Ae^{-iHt}v\|_{\mathfrak{H}}^2 dt < \infty.$$

$$\sup_{\epsilon > 0, \mu \in \mathbb{R}} \|AR(\mu \pm i\epsilon)\|^2 < \infty.$$

**Definition 10** An  $H$ -bounded operator,  $A : \mathcal{H} \rightarrow \mathfrak{H}$ , is called  $H$ -smooth in the weak sense if

$$w\text{-}\lim_{\epsilon \rightarrow \infty} A\theta(\lambda, \epsilon)A^* \quad (86)$$

exists for a.e.  $\lambda \in \mathbb{R}$ .

Equivalent conditions for the weak  $H$ -smoothness are stated by the following remark (other conditions can be found in [57]).

**Remark 4** An operator  $A : \mathcal{H} \rightarrow \mathfrak{H}$  is weakly  $H$ -smooth if and only if any of the following two conditions is satisfied

$$\|A\theta(\lambda, \epsilon)A^*\| \leq C(\lambda), \quad \text{a.e. } \lambda \in \mathbb{R}. \quad (87)$$

$$\epsilon^{1/2}\|AR(\lambda \pm i\epsilon)\| \leq C(\lambda), \quad \text{a.e. } \lambda \in \mathbb{R}. \quad (88)$$

To define the stationary WO, we first define the following

$$\mathcal{G}_{\pm}(H, H_0; \mathcal{J}) = \lim_{\epsilon \rightarrow 0} \pi^{-1} \epsilon \langle \mathcal{J}R_0(\lambda \pm i\epsilon)u_0, R(\lambda \pm i\epsilon)u \rangle. \quad (89)$$

Let, for all  $u_0 \in M_0$  and  $u \in M$ , the limit (89) exist for a.e.  $\lambda \in \mathbb{R}$ , then the stationary WO  $\mathcal{W}_{\pm} = \mathcal{W}_{\pm}(H, H_0; \mathcal{J})$  for the operators  $H$  and  $H_0$  with the identification  $\mathcal{J}$  is the operator on  $M_0 \times M$  defined by the following sesquilinear form

$$\langle \mathcal{W}_{\pm}u_0, u \rangle = \int_{-\infty}^{\infty} \mathcal{G}_{\pm}(H, H_0; \mathcal{J})d\lambda. \quad (90)$$

The WO  $\mathcal{W}_{\pm}$  is bounded, satisfying the intertwining property, and  $\mathbf{R}(\mathcal{W}_{\pm}) \subseteq \mathcal{H}^{(ac)}$ . Moreover, by the existence of  $\mathcal{W}_{\pm}(H, H_0; \mathcal{J})$ , then the adjoint WO  $\mathcal{W}_{\pm}^*(H, H_0; \mathcal{J})$  also exists and given by

$$\mathcal{W}_{\pm}^*(H, H_0; \mathcal{J}) = \mathcal{W}_{\pm}(H_0, H; \mathcal{J}^*). \quad (91)$$

Further more, on the relation to the weak time-dependent WOs, if both of  $\widetilde{\mathcal{W}}_{\pm}(H, H_0; \mathcal{J})$  and  $\mathcal{W}_{\pm}(H, H_0; \mathcal{J})$  exist, then they coincide with each other. This statement is also true if  $(H, H_0; \mathcal{J})$  is replaced by any of the collections  $(H_0, H; \mathcal{J}^*)$ ,  $(H, H; \mathcal{J}\mathcal{J}^*)$ , or  $(H_0, H_0; \mathcal{J}^*\mathcal{J})$ .

The importance of the stationary approach in scattering theory can be summarized as:

Let the WOs  $\widetilde{\mathcal{W}}_{\pm}(H, H_0; \mathcal{J})$  and  $\widetilde{\mathcal{W}}_{\pm}(H_0, H_0; \mathcal{J}^*\mathcal{J})$  exist, and let

$$\mathcal{W}_{\pm}^*(H, H_0; \mathcal{J})\mathcal{W}_{\pm}(H, H_0; \mathcal{J}) = \mathcal{W}_{\pm}(H_0, H_0; \mathcal{J}^*\mathcal{J}) \quad (92)$$

be satisfied, then the WO  $\mathcal{W}_{\pm}(H, H_0; \mathcal{J})$  exists.

Below we define a particular class of pseudo-differential operators (PSDOs) that are necessary to state the results on the convergence of the WOs for a family of Dirac operators.

## 5.4 Pseudo-differential operators

The class  $\mathcal{S}_{\rho,\delta}^r(\mathbb{R}^3, \mathbb{R}^3)$  of symbols is defined as follows

**Definition 11** The class  $\mathcal{S}_{\rho,\delta}^r(\mathbb{R}^3, \mathbb{R}^3)$  is the vector space of all smooth functions  $\mathcal{P}(x, \zeta) : \mathbb{R}^3 \times \mathbb{R}^3 \rightarrow \mathbb{C}$  such that for all multi-indices  $\alpha$  and  $\gamma$

$$|\partial_x^\alpha \partial_\zeta^\gamma \mathcal{P}(x, \zeta)| \leq c_{\alpha,\gamma} \langle x \rangle^{r-\rho|\alpha|+\delta|\gamma|}, \quad (93)$$

where  $r \in \mathbb{R}$ ,  $\rho > 0$ ,  $\delta < 1$ , and  $\langle x \rangle = (1 + |x|^2)^{1/2}$ . The function  $\mathcal{P}$  is called the symbol of the PSDO and  $r$  is called the order of  $\mathcal{P}$ .

Let  $\mathcal{P}(x, \zeta) \in \mathcal{S}_{\rho,\delta}^r(\mathbb{R}^3, \mathbb{R}^3)$ , the associated PSDO,  $\mathcal{P}$ , to  $\mathcal{P}$  is defined by the following integral

$$(\mathcal{P}f)(x) = (2\pi)^{-3/2} \int_{\mathbb{R}^3} e^{ix \cdot \zeta} \mathcal{P}(x, \zeta) \hat{f}(\zeta) d\zeta, \quad (94)$$

where  $f \in \mathcal{H}$  and  $\hat{f}(\zeta) = (2\pi)^{-3/2} \int_{\mathbb{R}^3} e^{-ix \cdot \zeta} f(x) dx$  is the Fourier transform of  $f$ .

## 5.5 A family of Dirac operators

Consider the free Dirac operator  $\mathbf{H}_0$ , and let  $\mathbf{V}$  be a short-range potential (decaying faster than the Coulomb potential), then the WO  $W_\pm(\mathbf{H}_0 + \mathbf{V}, \mathbf{H}_0)$  exists and is complete. The proofs of existence and completeness of  $W_\pm(\mathbf{H}_0 + \mathbf{V}, \mathbf{H}_0)$  are similar to that of the Schrödinger operator. For  $\mathbf{V}$  being the Coulomb potential, the WO  $W_\pm = W_\pm(\mathbf{H}_0 + \mathbf{V}, \mathbf{H}_0; \mathcal{J})$ , with a bounded identification  $\mathcal{J}$ , has been studied in [16, 17]. If  $\mathbf{V}$  is of long-range type (decaying as the Coulomb potential or slower), the existence and completeness of the WO  $W_\pm$  have been studied in [22, 37, 38, 52]. The asymptotic behavior of the WO  $W_\pm$  with respect to the speed of light,  $c$ , as  $c \rightarrow \infty$ , has been discussed for short-range potentials in [58] and for long-range potentials in [59].

### 5.5.1 An $h$ -dependent perturbation and the WO

Consider the free Dirac operator  $\mathbf{H}_0$ , and let  $\mathbf{V}_h$  be an  $h$ -dependent potential. We define the following family of Dirac operators

$$\mathbf{H}_h = \mathbf{H}_0 + \mathbf{V}_h. \quad (95)$$

We assume that the potential  $\mathbf{V}_h$  is real and bounded, thus the operators  $\mathbf{H}_h$  and  $\mathbf{H}_0$  have the same domain  $X$  and that  $\mathbf{H}_h$  is self-adjoint on  $X$ , for  $h > 0$ . Also,

for simplicity, we let  $\hbar = c = 1$ .

We assume further that  $\mathbf{V}_h$  is of long-range type for all  $h > 0$ , that is, for all multi-index  $\alpha$ ,  $\mathbf{V}_h$  fulfills the following condition

$$|\partial^\alpha \mathbf{V}_h(x)| \leq C \langle x \rangle^{-\rho-|\alpha|}, \quad \text{for all } h > 0, \text{ and } \rho \in (0, 1], \quad (96)$$

where we recall that  $\langle x \rangle = (1 + |x|^2)^{1/2}$ , and  $C$  is a constant independent of  $x$  and  $h$ .

Let  $P_h^{(ac)}$  be the orthogonal projection onto the absolutely continuous subspace of  $\mathbf{H}_h$ , and define  $U_h(t) = e^{-i\mathbf{H}_h t}$  and  $U_0(t) = e^{-i\mathbf{H}_0 t}$ . Now, by (96), and according to [22], the WOs  $W_{\pm, h}$  and  $W_{\pm, h}^*$ , defined as

$$W_{\pm, h} = W_{\pm}(\mathbf{H}_h, \mathbf{H}_0; \mathcal{J}_{\pm, h}) = s\text{-}\lim_{t \rightarrow \pm\infty} U_h(-t) \mathcal{J}_{\pm, h} U_0(t) \quad (97)$$

and

$$W_{\pm, h}^* = W_{\pm}(\mathbf{H}_0, \mathbf{H}_h; \mathcal{J}_{\pm, h}^*) = s\text{-}\lim_{t \rightarrow \pm\infty} U_0(-t) \mathcal{J}_{\pm, h}^* U_h(t) P_h^{(ac)}, \quad (98)$$

exist, moreover the WO  $W_{\pm, h}$  is complete. The identification  $\mathcal{J}_{\pm, h}$  is defined by the following PSDO

$$(\mathcal{J}_{\pm, h} g)(x) = (2\pi)^{-3/2} \int_{\mathbb{R}^3} e^{ix \cdot \zeta + i\Phi_{\pm, h}(x, \zeta)} \mathcal{P}_{\pm, h}(x, \zeta) \mathcal{C}_{\pm}(x, \zeta) \psi(|\zeta|^2) \hat{g}(\zeta) d\zeta, \quad (99)$$

where  $\psi \in C_0^\infty(\mathbb{R}_+)$  is introduced to localize  $\mathcal{J}_{\pm, h}$  in compact intervals of the positive part of the absolutely continuous spectrum,  $(m, \infty)$ , and where  $\mathcal{C}_{\pm}(x, \zeta)$  is a cut-off function defined as

$$\mathcal{C}_{\pm}(x, \zeta) = \theta(x) \omega_{\pm}(\langle \tilde{x}, \tilde{\zeta} \rangle), \quad \text{for all } y \in \mathbb{R}^3 \setminus \{0\}, \tilde{y} = y/|y|. \quad (100)$$

The function  $\theta$  is smooth and is introduced to avoid the singularity of  $\tilde{x}$  at  $x = 0$ , and  $\omega_{\pm}(\tau) = 1$  near  $\pm 1$  and  $\omega_{\pm}(\tau) = 0$  near  $\mp 1$ . Thus the cut-off function  $\mathcal{C}_{\pm}$  is supported in the cone

$$\Xi_{\pm}(\varrho) = \{(x, \zeta) \in \mathbb{R}^6 : \pm \langle x, \zeta \rangle \geq \varrho |x| |\zeta|\}, \quad \varrho \in (-1, 1). \quad (101)$$

Below, in a chain of definitions, we give the construction of the phase function  $\Phi_{\pm, h}(x, \zeta)$  and the amplitude function  $\mathcal{P}_{\pm, h}(x, \zeta)$ . The function  $\Phi_{\pm, h}(x, \zeta)$  is defined as follows

$$\Phi_{\pm, h}(x, \zeta) = \sum_{n=1}^N \Phi_{\pm, h}^{(n)}(x, \zeta), \quad x \in \Xi_{\pm}(\varrho), \quad (102)$$

where  $N$  satisfies  $(N + 1)\rho > 1$ , and for  $n \geq 0$ ,  $\Phi_{\pm,h}^{(n+1)}(x, \zeta) = Q_{\pm}(\zeta)F_{\pm,h}^{(n)}$  which is defined as

$$(Q_{\pm}(\zeta)F)(x) = \pm \int_0^{\infty} (F(x \pm t\zeta, \zeta) - F(\pm t\zeta, \zeta)) dt. \quad (103)$$

The functions  $F_{\pm,h}^{(n)}$  are defined as

$$F_{\pm,h}^{(0)}(x, \zeta) = \eta(\zeta)\mathbf{V}_h(x) - \frac{1}{2}\mathbf{V}_h^2(x), \quad F_{\pm,h}^{(1)}(x, \zeta) = \frac{1}{2}|\nabla\Phi_{\pm,h}^{(1)}(x, \zeta)|^2, \quad (104)$$

and for  $n \geq 2$

$$F_{\pm,h}^{(n)}(x, \zeta) = \sum_{k=1}^{n-1} \langle \nabla\Phi_{\pm,h}^{(k)}(x, \zeta), \nabla\Phi_{\pm,h}^{(n)}(x, \zeta) \rangle + \frac{1}{2}|\nabla\Phi_{\pm,h}^{(n)}(x, \zeta)|^2. \quad (105)$$

The amplitude function  $\mathcal{P}_{\pm,h}(x, \zeta)$  is defined as

$$\mathcal{P}_{\pm,h}(x, \zeta) = (I - S_{\pm,h}(x, \zeta))^{-1}p_0(\zeta), \quad x \in \Xi_{\pm}(\varrho), \quad (106)$$

where  $p_0(\zeta) = p_{+,0}(\zeta)$ , and

$$p_{\pm,0}(\zeta) = \frac{1}{2}(I \pm \eta^{-1}(\zeta)(\zeta_{\alpha} + mc^2\beta)), \quad (107)$$

with  $\eta(\zeta) = \sqrt{|\zeta|^2 + m^2c^4}$  and  $\zeta_{\alpha} = \alpha \cdot \zeta = \sum_{k=1}^3 \alpha_k \zeta_k$ . Finally,  $S_{\pm,h}(x, \zeta)$  is given by

$$S_{\pm,h}(x, \zeta) = (2\eta(\zeta))^{-1} \left( \mathbf{V}_h(x) + \sum_{k=1}^3 \partial_k \Phi_{\pm,h}(x, \zeta) \alpha_k \right), \quad x \in \Xi_{\pm}(\varrho). \quad (108)$$

Note that the WOs defined above are for the positive part of the absolutely continuous spectrum,  $(m, \infty)$ . For the negative part of the absolutely continuous spectrum,  $(-\infty, -m)$ , the WOs can be defined in a similar way with minor modifications, see [22]. The asymptotic study carried out below can be applied for the WOs on the negative part as well.

### 5.5.2 The asymptotics of the WOs and some particular cases

Define the WOs  $W_{\pm}^{\dagger} := s\text{-}\lim_{h \rightarrow \infty} W_{\pm,h}$  and  $W_{\pm}^{\dagger,*} := s\text{-}\lim_{h \rightarrow \infty} W_{\pm,h}^*$ , where  $W_{\pm,h}$  and  $W_{\pm,h}^*$  are given respectively by (97) and (98). Let the perturbed Dirac operator  $\mathbf{H}_h$  converge in the SRS to  $\mathbf{H}_{\infty}$ , and assume that the identification  $\mathcal{J}_{\pm,h}$

converges strongly to  $\mathcal{J}_{\pm, \infty}$ , then the WOs  $W_{\pm}^{\dagger}$  and  $W_{\pm}^{\dagger, *}$  exist. The task now is to characterize the limits, as  $h \rightarrow \infty$ , of the WOs  $W_{\pm, h}$  and  $W_{\pm, h}^*$ , which is equivalent to the problem of interchanging  $s\text{-}\lim_{h \rightarrow \infty}$  and  $s\text{-}\lim_{t \rightarrow \pm\infty}$ . To this end, we state the following two lemmas.

**Lemma 1** Define the function  $\mathcal{K}_{u_0, h}^{(1)}$  as

$$\mathcal{K}_{u_0, h}^{(1)}(t) = \left\| \left( \mathbf{H}_h \phi(\mathbf{H}_h) \mathcal{J}_{\pm, h} \phi(\mathbf{H}_0) - \phi(\mathbf{H}_h) \mathcal{J}_{\pm, h} \mathbf{H}_0 \phi(\mathbf{H}_0) \right) U_0(t) u_0 \right\|_Y,$$

where  $u_0 \in X$ . Then for some continuous function  $\phi : \mathbb{R} \rightarrow \mathbb{R}$  such that  $x\phi(x)$  is bounded on  $\mathbb{R}$  and for any  $\varepsilon > 0$  there exist  $D_1, D_2 \in \mathbb{R}$  such that

$$\int_{D_1}^{\infty} \mathcal{K}_{u_0, h}^{(1)}(t) dt \leq \varepsilon \text{ and } \int_{-\infty}^{D_2} \mathcal{K}_{u_0, h}^{(1)}(t) dt \leq \varepsilon \text{ for all } h > 0.$$

**Lemma 2** Define the function  $\mathcal{K}_{u, h}^{(2)}$  as

$$\mathcal{K}_{u, h}^{(2)}(t) = \left\| \left( \mathbf{H}_0 \phi(\mathbf{H}_0) \mathcal{J}_{\pm, h}^* \phi(\mathbf{H}_h) - \phi(\mathbf{H}_0) \mathcal{J}_{\pm, h}^* \mathbf{H}_h \phi(\mathbf{H}_h) \right) U_h(t) u \right\|_Y,$$

where  $u_0 \in X$ . Then for some continuous function  $\phi : \mathbb{R} \rightarrow \mathbb{R}$  such that  $x\phi(x)$  is bounded on  $\mathbb{R}$  and for any  $\varepsilon > 0$  there exist  $D_3, D_4 \in \mathbb{R}$  such that

$$\int_{D_3}^{\infty} \mathcal{K}_{u, h}^{(2)}(t) dt \leq \varepsilon \text{ and } \int_{-\infty}^{D_4} \mathcal{K}_{u, h}^{(2)}(t) dt \leq \varepsilon \text{ for all } h > 0.$$

By Lemmas 1 and 2, and according to [9], the limits  $s\text{-}\lim_{h \rightarrow \infty}$  and  $s\text{-}\lim_{t \rightarrow \pm\infty}$  in the definition of the WOs  $W_{\pm}^{\dagger}$  and  $W_{\pm}^{\dagger, *}$  are interchangeable. Thus we have the following result.

**Theorem 12** Let the WOs  $W_{\pm, h}$  and  $W_{\pm, h}^*$  be defined by (97) and (98) respectively. Suppose that, as  $h \rightarrow \infty$ , the Dirac operator  $\mathbf{H}_h$  converges to  $\mathbf{H}_{\infty}$  in the SRS, and the identification  $\mathcal{J}_{\pm, h}$  converges strongly to  $\mathcal{J}_{\pm, \infty}$ . Then the WOs  $W_{\pm}^{\dagger}$  and  $W_{\pm}^{\dagger, *}$  exist,

$$W_{\pm}^{\dagger} = W_{\pm}(\mathbf{H}_{\infty}, \mathbf{H}_0; \mathcal{J}_{\pm, \infty}),$$

and

$$W_{\pm}^{\dagger, *} = W_{\pm}(\mathbf{H}_0, \mathbf{H}_{\infty}; \mathcal{J}_{\pm, \infty}^*).$$

**Remark 5** In Theorem 12 we assume that  $\mathcal{J}_{\pm, h}$  converges strongly to  $\mathcal{J}_{\pm, \infty}$ , this also implies that  $\mathcal{J}_{\pm, h}^*$  converges strongly to  $\mathcal{J}_{\pm, \infty}^*$ . However, in general, the strong convergence of an operator does not imply the strong convergence of its adjoint operator to the adjoint of its strong limit. Hence, in order to study

the convergence of the adjoint WO in the strong sense for other self-adjoint operators, we should assume if necessary, the strong convergence of the identifications adjoint operators as well.

In what follows we assume the hypotheses of Theorem 12 and study different cases of the identification  $\mathcal{J}_{\pm,h}$ . Note that, in the first case we also consider short-range potentials, so the identification is just the identity operator. On the other hand, the other two cases are stated briefly, where a dwell-upon discussion is available in Paper V in the appendix.

**The case  $\rho > 1$ .**

In this case we can set  $\mathcal{J}_{\pm,h} = I$ , this is due to the fact that for short-range potentials, the WOs  $W_{\pm}(\mathbf{H}_h, \mathbf{H}_0)$  and  $W_{\pm}(\mathbf{H}_0, \mathbf{H}_h)$  exist and are complete. Therefore, the limits  $s\text{-}\lim_{h \rightarrow \infty}$  and  $s\text{-}\lim_{t \rightarrow \pm\infty}$  are interchangeable in the definitions of the WOs  $W_{\pm}^{\dagger}$  and  $W_{\pm}^{\dagger,*}$ . Thus, if the perturbed Dirac operator  $\mathbf{H}_h$  is convergent to  $\mathbf{H}_{\infty}$  in the SRS, then

$$W_{\pm}^{\dagger} = W_{\pm}(\mathbf{H}_{\infty}, \mathbf{H}_0) \quad (109)$$

and

$$W_{\pm}^{\dagger,*} = W_{\pm}(\mathbf{H}_0, \mathbf{H}_{\infty}). \quad (110)$$

**The case  $\rho = 1$ .**

Let  $\Phi_{\pm}(x, \zeta)$  be an  $h$ -free function satisfying

$$|\partial_x^{\alpha} \partial_{\zeta}^{\gamma} \Phi_{\pm}(x, \zeta)| \leq c_{\alpha,\gamma} \langle x \rangle^{1-\rho-|\alpha|}, \quad x \in \Xi_{\pm}(\varrho), \quad (111)$$

and let  $\mathcal{J}_{\pm}^{(1)}$  (with adjoint operator denoted by  $\mathcal{J}_{\pm}^{(1),*}$ ) be defined as

$$(\mathcal{J}_{\pm}^{(1)} g)(x) = (2\pi)^{-3/2} \int_{\mathbb{R}^3} e^{ix \cdot \zeta + i\Phi_{\pm}(x, \zeta)} p_0(\zeta) \mathcal{C}_{\pm}(x, \zeta) \psi(|\zeta|^2) \hat{g}(\zeta) d\zeta. \quad (112)$$

Then

$$W_{\pm}^{\dagger} = W_{\pm}(\mathbf{H}_{\infty}, \mathbf{H}_0; \mathcal{J}_{\pm}^{(1)}) \quad (113)$$

and

$$W_{\pm}^{\dagger,*} = W_{\pm}(\mathbf{H}_0, \mathbf{H}_{\infty}; \mathcal{J}_{\pm}^{(1),*}). \quad (114)$$

**The case  $\rho \in (1/2, 1)$ .**

Let  $\mathcal{J}_{\pm,h}^{(2)}$  (with adjoint operator denoted by  $\mathcal{J}_{\pm}^{(2),*}$ ) be given by (99), but with  $\Phi_{\pm,h}(x, \zeta)$  defined as

$$\Phi_{\pm,h}(x, \zeta) = \pm\eta(\zeta) \int_0^{\infty} (\mathbf{V}_h(x \pm t\zeta) - \mathbf{V}_h(\pm t\zeta)) dt \quad (115)$$



and with  $p_0(\zeta)$  instead of  $\mathcal{P}_{\pm,h}(x, \zeta)$ . Assume further that  $\mathbf{V}_h$  is given so that  $\mathbf{H}_h = \mathbf{H}_0 + \mathbf{V}_h$  and  $\Phi_{\pm,h}(x, \zeta)$  converge in the SRS respectively to  $\mathbf{H}_\infty = \mathbf{H}_0 + \mathbf{V}_\infty$  and

$$\Phi_{\pm,\infty}(x, \zeta) = \pm\eta(\zeta) \int_0^\infty (\mathbf{V}_\infty(x \pm t\zeta) - \mathbf{V}_\infty(\pm t\zeta)) dt. \quad (116)$$

Then

$$W_\pm^\dagger = W_\pm(\mathbf{H}_\infty, \mathbf{H}_0; \mathcal{J}_{\pm,\infty}^{(2)}) \quad (117)$$

and

$$W_\pm^{\dagger,*} = W_\pm(\mathbf{H}_0, \mathbf{H}_\infty; \mathcal{J}_{\pm,\infty}^{(2),*}), \quad (118)$$

where

$$(\mathcal{J}_{\pm,\infty}^{(2)}g)(x) = (2\pi)^{-3/2} \int_{\mathbb{R}^3} e^{ix \cdot \zeta + i\Phi_{\pm,\infty}(x, \zeta)} p_0(\zeta) \mathcal{C}_\pm(x, \zeta) \psi(|\zeta|^2) \hat{g}(\zeta) d\zeta. \quad (119)$$

## 5.6 Self-adjoint $h$ -dependent operators

Let  $\mathcal{H}_0$  and  $\mathcal{H}$  be two Hilbert spaces, and let  $M_0$  and  $M$  be dense sets in  $\mathcal{H}_0$  and  $\mathcal{H}$  respectively. Let  $H_0$  and  $H_h$  be two self-adjoint operators in  $\mathcal{H}_0$  and  $\mathcal{H}$  respectively, with  $\mathbf{D}(H_0) = \mathcal{X}_0$  and  $\mathbf{D}(H_h) = \mathcal{X}$ , and with corresponding resolvent operators  $R_0$  and  $R_h$  respectively. Let also  $P_0^{(ac)}$  and  $P_h^{(ac)}$  be respectively the orthogonal projections onto the absolutely continuous subspaces of  $H_0$  and  $H_h$ . Assume that  $H_h = H_0 + V_h$ , where  $V_h$  admits the following factorization

$$V_h = H_h \mathcal{J}_h - \mathcal{J}_h H_0 = A_h^* A_0, \quad (120)$$

where  $\mathcal{J}_h : \mathcal{H}_0 \rightarrow \mathcal{H}$  is a bounded identification, and  $A_h : \mathcal{H} \rightarrow \mathfrak{H}$  and  $A_0 : \mathcal{H}_0 \rightarrow \mathfrak{H}$  are respectively  $H_h$ -bounded, for all  $h > 0$ , and  $H_0$ -bounded operators, where  $\mathfrak{H}$  is an auxiliary Hilbert space. Note that (120) is understood as the equalities of the corresponding sesquilinear forms.

Define the time-dependent WO  $W_\pm^\dagger(H, H_0; \mathcal{J})$  as

$$\begin{aligned} W_\pm^\dagger(H, H_0; \mathcal{J}) &= s\text{-}\lim_{h \rightarrow \infty} W_\pm(H_h, H_0; \mathcal{J}_h) \\ &= s\text{-}\lim_{h \rightarrow \infty} s\text{-}\lim_{t \rightarrow \pm\infty} U_h(-t) \mathcal{J}_h U_0(t) P_0^{(ac)}, \end{aligned} \quad (121)$$

where  $U_h(t) = e^{-iH_h t}$ ,  $U_0(t) = e^{-iH_0 t}$ , and  $H$  and  $\mathcal{J}$  are some limit operators in appropriate sense of  $H_h$  and  $\mathcal{J}_h$  respectively.

Let  $\mathcal{G}_\pm^\dagger(H, H_0; \mathcal{J})$  be defined as

$$\mathcal{G}_\pm^\dagger(H, H_0; \mathcal{J}) = \lim_{h \rightarrow \infty} \lim_{\epsilon \rightarrow 0} \pi^{-1} \epsilon \langle \mathcal{J}_h R_0(\lambda \pm i\epsilon) u_0, R_h(\lambda \pm i\epsilon) u \rangle, \quad (122)$$

where  $u_0 \in M_0$  and  $u \in M$ . We define the stationary WO  $\mathcal{W}_\pm^\dagger = \mathcal{W}_\pm^\dagger(H, H_0; \mathcal{J})$  on  $M_0 \times M$  by the sesquilinear form

$$\langle \mathcal{W}_\pm^\dagger u_0, u \rangle = \int_{-\infty}^{\infty} \mathcal{G}_\pm^\dagger(H, H_0; \mathcal{J}) d\lambda. \quad (123)$$

We also define the weak WO  $\widetilde{\mathcal{W}}_\pm^\dagger(H, H_0; \mathcal{J})$  as

$$\begin{aligned} \widetilde{\mathcal{W}}_\pm^\dagger(H, H_0; \mathcal{J}) &= w\text{-}\lim_{h \rightarrow \infty} \widetilde{\mathcal{W}}_\pm(H_h, H_0; \mathcal{J}_h) \\ &= w\text{-}\lim_{h \rightarrow \infty} w\text{-}\lim_{t \rightarrow \pm\infty} P_h^{(ac)} U_h(-t) \mathcal{J}_h U_0(t) P_0^{(ac)}. \end{aligned} \quad (124)$$

In the coming discussion we state some results regarding the existence of the WOs  $\mathcal{W}_\pm^\dagger$ ,  $\widetilde{\mathcal{W}}_\pm^\dagger$ ,  $W_\pm^\dagger$ , and their adjoint operators that are denoted respectively by  $\mathcal{W}_\pm^{\dagger,*}$ ,  $\widetilde{\mathcal{W}}_\pm^{\dagger,*}$ , and  $W_\pm^{\dagger,*}$ . These results are briefly stated, where the details are available in Paper VI in the appendix.

**Theorem 13** Assume the following

- (i)  $A_0$  is weakly  $H_0$ -smooth.
- (ii) For all  $h > 0$ ,  $A_h R_h(\lambda \pm i\epsilon)$  is strongly convergent as  $\epsilon \rightarrow 0$  for a.e.  $\lambda \in \mathbb{R}$ .
- (iii) If  $T_h$  is the strong limit of  $A_h R_h(\lambda \pm i\epsilon)$  as  $\epsilon \rightarrow 0$  obtained in (ii),  $T_h$  converges weakly to some  $T_\infty$  for a.e.  $\lambda \in \mathbb{R}$ .
- (iv)  $\mathcal{J}_h$  converges weakly to  $\mathcal{J}_\infty$ .

Then the WO  $\mathcal{W}_\pm^\dagger(H, H_0; \mathcal{J})$  exists, also  $\mathcal{W}_\pm^\dagger(H_0, H; \mathcal{J}^*)$  exists and

$$\mathcal{W}_\pm^{\dagger,*}(H, H_0; \mathcal{J}) = \mathcal{W}_\pm^\dagger(H_0, H; \mathcal{J}^*). \quad (125)$$

Note that the assertions of Theorem 13 remain unchanged if its first three hypotheses are replaced as: For a.e.  $\lambda \in \mathbb{R}$ , as  $\epsilon \rightarrow 0$ , the operator  $A_0 \theta_0(\lambda, \epsilon)$  is strongly convergent and  $T_{h,\epsilon} := A_h R_h(\lambda \pm i\epsilon)$  is weakly convergent to some  $T_{h,0}$  for all  $h > 0$ , and the resulting limit  $T_{h,0}$  converges weakly to some  $T_{\infty,0}$  as  $h \rightarrow \infty$ .

Similar assertions as of Theorem 13 can be formulated as in the following theorem.

**Theorem 14** Assume the following

- (i) For all  $h > 0$ ,  $A_h$  is weakly  $H_h$ -smooth.
- (ii) The operator  $A_0 R_0(\lambda \pm i\epsilon)$  is strongly convergent as  $\epsilon \rightarrow 0$  for a.e.  $\lambda \in \mathbb{R}$ .
- (iii) If  $T_h$  is the weak limit of  $A_h \theta_h(\lambda, \epsilon)$  as  $\epsilon \rightarrow 0$  obtained in (i),  $T_h$  converges weakly to some  $T_\infty$  for a.e.  $\lambda \in \mathbb{R}$ .
- (iv) If  $E_h$  is the spectral family of  $H_h$ , then  $E_h(\lambda)$  and  $\mathcal{J}_h$  converge weakly to  $E_\infty(\lambda)$  and  $\mathcal{J}_\infty$  respectively for a.e.  $\lambda \in \mathbb{R}$ .

Then the WO  $\mathcal{W}_\pm^\dagger(H, H_0; \mathcal{J})$  exists, also  $\mathcal{W}_\pm^\dagger(H_0, H; \mathcal{J}^*)$  exists and

$$\mathcal{W}_\pm^{\dagger,*}(H, H_0; \mathcal{J}) = \mathcal{W}_\pm^\dagger(H_0, H; \mathcal{J}^*). \quad (126)$$

The assertions of Theorem 14 are also true if its first three hypotheses are replaced by the following: For a.e.  $\lambda \in \mathbb{R}$ , as  $\epsilon \rightarrow 0$ , the operator  $A_0 R_0(\lambda \pm i\epsilon)$  is weakly convergent and  $S_{h,\epsilon} := A_h \theta_h(\lambda, \epsilon)$  is strongly convergent to some  $S_{h,0}$  for all  $h > 0$ , and the resulting limit  $S_{h,0}$  converges weakly to some  $S_{\infty,0}$  as  $h \rightarrow \infty$ .

The existence of  $\mathcal{W}_\pm^\dagger(H, H; \mathcal{J}\mathcal{J}^*)$  and  $\mathcal{W}_\pm^\dagger(H_0, H_0; \mathcal{J}^*\mathcal{J})$  is proved in Theorems 15 and 16 respectively.

**Theorem 15** Let the hypotheses of Theorem 13 be satisfied, and let further  $\mathcal{J}_h^*$  and  $R_h$  be strongly convergent. Then the WO  $\mathcal{W}_\pm^\dagger(H, H; \mathcal{J}\mathcal{J}^*)$  exists, moreover we have

$$\mathcal{W}_\pm^\dagger(H, H_0; \mathcal{J}) \mathcal{W}_\pm^{\dagger,*}(H, H_0; \mathcal{J}) = \mathcal{W}_\pm^\dagger(H, H; \mathcal{J}\mathcal{J}^*). \quad (127)$$

**Theorem 16** Let the hypotheses of Theorem 14 be satisfied, and let  $\mathcal{J}_h$  be strongly convergent. Then the WO  $\mathcal{W}_\pm^\dagger(H_0, H_0; \mathcal{J}^*\mathcal{J})$  exists and

$$\mathcal{W}_\pm^{\dagger,*}(H, H_0; \mathcal{J}) \mathcal{W}_\pm^\dagger(H, H_0; \mathcal{J}) = \mathcal{W}_\pm^\dagger(H_0, H_0; \mathcal{J}^*\mathcal{J}). \quad (128)$$

Similarly to the coincidence between the usual stationary and weak time-dependent WOs, we have the coincidence between the stationary WO  $\mathcal{W}_\pm^\dagger$  and the weak time-dependent WO  $\widetilde{\mathcal{W}}_\pm^\dagger$ , that is, if both of  $\mathcal{W}_\pm^\dagger(H, H_0; \mathcal{J})$  and  $\widetilde{\mathcal{W}}_\pm^\dagger(H, H_0; \mathcal{J})$  exist, then they coincide with each other. The same assertion can be concluded for any of the collections  $(H_0, H; \mathcal{J}^*)$ ,  $(H_0, H_0; \mathcal{J}^*\mathcal{J})$ , and  $(H, H; \mathcal{J}\mathcal{J}^*)$  as well.

Also, by the hypotheses of Theorem 13 (equivalently the hypotheses of Theorem 14), the WO  $\widetilde{W}_{\pm}^{\dagger}(H, H_0; \mathcal{J})$  exists, consequently  $\widetilde{W}_{\pm}^{\dagger}(H_0, H; \mathcal{J}^*)$  exists and

$$\widetilde{W}_{\pm}^{\dagger,*}(H, H_0; \mathcal{J}) = \widetilde{W}_{\pm}^{\dagger}(H_0, H; \mathcal{J}^*). \quad (129)$$

For the WOs  $\widetilde{W}_{\pm}^{\dagger}(H, H; \mathcal{J}\mathcal{J}^*)$  and  $\widetilde{W}_{\pm}^{\dagger}(H_0, H_0; \mathcal{J}^*\mathcal{J})$ , we have the following two theorems.

**Theorem 17** Suppose the hypotheses of Theorem 15 are satisfied, then the WO  $\widetilde{W}_{\pm}^{\dagger}(H, H; \mathcal{J}\mathcal{J}^*)$  exists.

**Theorem 18** Suppose the hypotheses of Theorem 16 are satisfied, then the WO  $\widetilde{W}_{\pm}^{\dagger}(H_0, H_0; \mathcal{J}^*\mathcal{J})$  exists.

The existence of the time-dependent WOs  $W_{\pm}^{\dagger}(H, H_0; \mathcal{J})$  and  $W_{\pm}^{\dagger}(H_0, H; \mathcal{J}^*)$  is summarized in the following theorems.

**Theorem 19** If the hypotheses of Theorem 15 are satisfied, then  $W_{\pm}^{\dagger}(H_0, H; \mathcal{J}^*)$  exists.

**Theorem 20** If the hypotheses of Theorem 16 are satisfied, then  $W_{\pm}^{\dagger}(H, H_0; \mathcal{J})$  exists.

After proving the existence of the WOs  $W_{\pm}^{\dagger}(H, H_0; \mathcal{J})$  and  $W_{\pm}^{\dagger}(H_0, H; \mathcal{J}^*)$ , we would like to study the asymptotic behavior, as  $h \rightarrow \infty$ , of the WOs  $W_{\pm}(H_h, H_0; \mathcal{J}_h)$  and  $W_{\pm}(H_0, H_h; \mathcal{J}_h^*)$ . The problem of finding these asymptotic limits is reduced, as we mentioned before, to the problem of interchanging  $s\text{-}\lim_{h \rightarrow \infty}$  and  $s\text{-}\lim_{t \rightarrow \pm\infty}$ . By the existence of  $W_{\pm}^{\dagger}(H, H_0; \mathcal{J})$  and  $W_{\pm}^{\dagger}(H_0, H; \mathcal{J}^*)$ , Lemmas 1 and 2 are satisfied for the collections  $(H_h, H_0, \mathcal{J}_h, \mathcal{H})$  and  $(H_0, H_h, \mathcal{J}_h^*, \mathcal{H}_0)$  respectively. This implies that, according to [9], in the definitions of  $\mathcal{W}_{\pm}^{\dagger}(H, H_0; \mathcal{J})$  and  $\mathcal{W}_{\pm}^{\dagger}(H_0, H; \mathcal{J}^*)$ , the limits  $s\text{-}\lim_{h \rightarrow \infty}$  and  $s\text{-}\lim_{t \rightarrow \pm\infty}$  are interchangeable. Therefore, if  $H_h$  converges to  $H_{\infty}$  in the SRS, and  $\mathcal{J}_h$  and  $\mathcal{J}_h^*$  converge strongly to  $\mathcal{J}_{\infty}$  and  $\mathcal{J}_{\infty}^*$  respectively, then

$$s\text{-}\lim_{h \rightarrow \infty} W_{\pm}(H_h, H_0; \mathcal{J}_h) = W_{\pm}(H_{\infty}, H_0; \mathcal{J}_{\infty}) \quad (130)$$

and

$$s\text{-}\lim_{h \rightarrow \infty} W_{\pm}(H_0, H_h; \mathcal{J}_h^*) = W_{\pm}(H_0, H_{\infty}; \mathcal{J}_{\infty}^*). \quad (131)$$

## 6 Conclusion and future work

The numerical scheme we provide for approximating the eigenvalues of the Dirac operator is of vital importance in the sense that it results to the complete removal of the spurious eigenvalues. The stability approach for the FEM in Paper I yields a better rate of convergence compared to the Galerkin-based  $hp$ -cloud approach, Paper II. Though the computation using the  $hp$ -cloud method is more time consuming compared to the FEM, but as the method is applied for the first time to the Dirac operator makes it a considerable novel effort toward getting faster computation. An ongoing work concerns improving the rate of convergence of the approximations using these two approaches, where the focus is on the computational aspects of the  $hp$ -cloud approach. Also, as a future work, we will consider developing of approximation methods for the eigenvalues of the electron in the Helium-like ion systems.

In Paper III, we apply G-convergence theory for positive definite parts of the Dirac operator where our purpose is to study the behavior of the eigenvalues of a family of perturbed Dirac operators by abstract  $h$ -dependent potentials. That the Dirac operator is not bounded affects construction of operator convergence methods. So, the results in Paper III can be viewed as a modest progress for certain abstract potentials that may not appear in applications. Our concluding paper in G-convergence, Paper IV, gives a simplified and brief knowledge on the theory of G-convergence. It includes some review material about G-convergence of elliptic as well as some positive definite self-adjoint operators. To extend the asymptotic analysis of the perturbed Dirac operator, we intend to use suitable variational convergence techniques. We believe that this can lead to useful results, but they are not completely formalized and still under scrutiny.

In Paper V, we apply scattering theory for the free Dirac operator with long-range  $h$ -dependent potentials and study the strong convergence of the wave operator (WO). Though the application of this study is not seemingly evident for the Dirac operator, but it is manifest for other differential operators, in particular the Schrödinger operator. In Paper VI, we provide a general asymptotic study for the WOs for general self-adjoint  $h$ -dependent operators. This study has not come with new results regarding the strong convergence of the time-dependent WO, but it conforms to the existence result in Paper V and the references therein. In addition, in Paper VI, we extend the study and prove the convergence of the weak time-dependent and stationary WOs in a general setting. As a future work, we intend to study the asymptotics of the WOs for other self-adjoint operators. Also we like to study other classes of perturbations with different factorizations.

## References

- [1] E. Ackad and M. Horbatsch, *Numerical solution of the Dirac equation by a mapped fourier grid method*, J. Phys. A: Math. Gen., 38(2005), pp. 3157-3171.
- [2] R. C. Almeida and R. S. Silva, *A stable Petrov-Galerkin method for convection-dominated problems*, Comput. Methods Appl. Mech. Engng., 140(1997).
- [3] S. N. Atluri and S. Shen, *The meshless local Petrov-Galerkin (MLPG) method: A simple and less-costly alternative to the finite element and boundary element methods*, CMES, 3(2002), pp. 11-51.
- [4] T. Belytschko, D. Organ, and Y. Krongauz, *A coupled finite element-element-free Galerkin method*, Comput. Mech., 17(1995), pp. 186-195.
- [5] T. Belytschko, Y. Krongauz, D. Organ, M. Fleming, and P. Krysl, *Meshless methods: An overview and recent developments*, Comput. Methods Appl. Mech. Engng., 139(1996), pp. 3-47.
- [6] M. S. Birman and M. Z. Solomjak, *Spectral Theory of Self-adjoint Operators in Hilbert Space*, D. Reidel Publishing Company, Dordrecht, Holland, 1987.
- [7] L. Boulton and M. Levitin, *Trends and Tricks in Spectral Theory*, Ediciones IVIC, Caracas, Venezuela, 2007.
- [8] A. Braides,  *$\Gamma$ -convergence for Beginners*, Oxford University Press Inc., New York, 2002.
- [9] E. Brüning and F. Gesztesy, *Continuity of wave and scattering operators with respect to interactions*, J. Math. Phys., 24(1983), pp. 1516-1528.
- [10] J. S. Chen, W. Hu, and M. A. Puso, *Orbital hp-clouds for solving Schrödinger equation in quantum mechanics*, Comput. Methods Appl. Mech. Engng., 196(2007), pp. 3693-3705.
- [11] G. Dal Maso, *An Introduction to  $\Gamma$ -convergence*, Birkhäuser, Boston, 1993.

- [12] E. De Giorgi and T. Franzoni, *Su un tipo convergenza variazionale*, Atti Accad. Naz. Lincei Rend. Cl. Sci. Mat. 58(1975), pp. 842-850.
- [13] E. De Giorgi and S. Spagnolo, *Sulla convergenze degli integrali dell'energia per operatori ellittici del secondo ordine*, Boll. Un. Mat. Ital., 8(1973), pp. 391-411.
- [14] P. A. B. De Sampaio, *A Petrov-Galerkin/modified operator formulation for convection-diffusion problems*, Int. J. Numer. methods Engng., 30(1990).
- [15] A. Defranceschi, *An introduction to homogenization and G-convergence*, School on homogenization, ICTP, Trieste, September 6-17, 1993.
- [16] J. Dollard, *Asymptotic convergence and Coulomb interaction*, J. Math. Phys., 5(1964), pp. 729-738.
- [17] J. Dollard and G. Velo, *Asymptotic behaviour of a Dirac particle in a Coulomb field*, Nuovo Cimento, 45(1966), pp. 801-812.
- [18] C. A. Duarte and J. T. Oden, *H-p clouds—An h-p meshless method*, Numer. Meth. Part. D. E., 12(1996), pp. 673-705.
- [19] T. Fries and H. Matthies, *A review of Petrov-Galerkin stabilization approaches and an extension to meshfree methods*, Institute of scientific computing, Technical University Braunschweig, Brunswick, Germany, 2004.
- [20] T. Fries and H. Matthies, *Classification and overview of Meshfree Methods*, Institute of scientific computing, Technical University Braunschweig, Brunswick, Germany, 2004.
- [21] O. Garcia, E. A. Fancello, C. S. de Barcellos, and C. A. Duarte, *hp-Clouds in Mindlin's thick plate model*, Int. J. Numer. methods Engng., 47(2000), pp. 1381-1400.
- [22] Y. Gätel and D. R. Yafaev, *Scattering theory for the Dirac operator with a long-range electromagnetic potential*, J. Fun. Anal., 184(2001), pp. 136-176.
- [23] M. Griesemer and J. Lutgen, *Accumulation of Discrete Eigenvalues of the Radial Dirac Operator*, J. Funct. Anal., 162(1999).

- [24] H. Haken and H. C. Wolf, *The Physics of Atoms and Quanta, Introduction to Experiments and Theory*, Springer-Verlag, Berlin, 1996.
- [25] W. J. Hehre, L. Radom, P. Schleyer, and J. Pople, *Ab Initio Molecular Orbital Theory*, John Willey & Sons, New York, 1986.
- [26] A. Huerta and S. Fernández-Méndez, *Coupling element free Galerkin and finite element methods*, ECCOMAS, Barcelona, 2000.
- [27] A. Huerta and S. Fernández-Méndez, *Enrichment and coupling of the finite element and meshless methods*, Int. J. Numer. methods Engng., 48(2000), pp. 1615-1636.
- [28] S. Idelsohn, N. Nigro, M. Storti, and G. Buscaglia, *A Petrov-Galerkin formulation for advection-reaction-diffusion problems*, Comput. Methods Appl. Mech. Engng., 136(1996).
- [29] T. Kato, *Perturbation Theory for Linear Operators*, Springer-Verlag, Berlin Heidelberg, 1976.
- [30] C. V. Le, *Novel numerical procedures for limit analysis of structures: Mesh-free methods and mathematical programming*, Ph.D thesis, University of Sheffield, England, 2009.
- [31] H. Lin and S. N. Atluri, *Meshless local Petrov-Galerkin (MLPG) method for convection-diffusion problems*, CMES, 1(2000), pp. 45-60.
- [32] I. Lindgren, S. Salomonson, and B. Åsén, *The covariant-evolution-operator method in bound-state QED*, Physics Reports, 389(2004), pp. 161-261.
- [33] G. R. Liu, *Mesh Free Methods: Moving Beyond the Finite Element Method*, CRC press, 2003.
- [34] P. de T. R. Mendonça, C. S. de Barcellos, A. Duarte, *Investigations on the hp-cloud method by solving Timoshenko beam problems*, Comput. Mech., 25(2000), pp. 286-295.
- [35] G. Mur, *On the causes of spurious solutions in electromagnetics*, Electromagnetic, 22(2002), pp. 357-367.
- [36] F. Murat, *H-convergence*, Séminaire d'Analyse Fonctionnelle et Numérique de l'Université d'Alger, 1977.



- [37] Pl. Muthuramalingam, *Scattering theory by Enss' method for operator valued matrices: Dirac operator in an electric field*, J. Math. Soc. Japan, 37(1985), pp. 415-432.
- [38] Pl. Muthuramalingam and K. B. Sinha, *Existence and completeness of wave operators for the Dirac operator in an electromagnetic field with long range potentials*, J. Indian Math. Soc. (N.S.), 50(1986), pp. 103-130.
- [39] G. Pestka, *Spurious roots in the algebraic Dirac equation*, Chem. Phys. Lett. 376(2003), pp. 659-661.
- [40] M. Reed and B. Simon, *Methods of Modern Mathematical Physics: Scattering Theory*, Vol. 3, Academic Press, San Diego, CA, 1979.
- [41] L. Rosenberg, *Virtual-pair effects in atomic structure theory*, Phys. Rev. A, 39(1989), pp. 4377-4386.
- [42] S. Salomonson and P. Öster, *Relativistic all-order pair functions from a discretized single-particle Dirac Hamiltonian*, Phys. Rev. A, 40(1989), pp. 5548-5558.
- [43] W. Schroeder and I. Wolf, *The origin of spurious modes in numerical solutions of electromagnetic field eigenvalue problems*, IEEE Tran. on Micr. Theory and Tech., 42(1994), pp. 644-653.
- [44] V. M. Shabaev, I. I. Tupitsyn, V. A. Yerokhin, G. Plunien, and G. Soff, *Dual kinetic balance approach to basis-set expansions for the Dirac equation*, Phys. Rev. Lett., 93(2004).
- [45] S. Spagnolo, *Sul limite delle soluzioni di problemi di Cauchy relativi all'equazione del calore*, Ann. Sc. Norm. Sup. Pisa cl. Sci., 21(1967), pp. 657-699.
- [46] S. Spagnolo, *Sulla convergenza delle soluzioni di equazioni paraboliche ed ellittiche*, Ann. Sc. Norm. Sup. Pisa cl. Sci., 22(1968), pp. 571-597.
- [47] S. Spagnolo, *Convergence in energy for elliptic operators*, Proc. Third Symp. Numer. Solut. Partial Diff. Equat. (College Park, 1975), pp. 469-498, Academic Press, San Diego, 1976.
- [48] L. Tartar, *Quelques remarques sur l'homogénéisation*, Proc. of the Japan-France seminar 1976, "Functional analysis and numer-

- ical analysis”, pp. 469-482, Japan society for the promotion of science, 1978.
- [49] L. Tartar, *Convergence d’opérateurs différentiels*, Analisi convessa, Roma, 1977.
- [50] L. Tartar, *Cours peccot au collège de France*, Paris, 1977.
- [51] B. Thaller, *The Dirac Equation*, Springer-verlag, Berlin, 1993.
- [52] B. Thaller and V. Enss, *Asymptotic observables and Coulomb scattering for the Dirac equation*, Ann. Inst. H. Poincaré Phys. Theor., 45(1986), pp. 147-171.
- [53] I. I. Tupitsyn and V. M. Shabaev, *Spurious states of the Dirac equation in a finite basis set*, Optika i Spektroskopiya, 105(2008), pp. 203-209.
- [54] J. Weidman, *Linear Operators in Hilbert Spaces*, Springer-verlag, New York, 1980.
- [55] J. Weidman, *Strong operator convergence and spectral theory of ordinary differential operators*, Univ. Lagel. Acta Math. No., 34(1997), pp. 153-163.
- [56] J. Weidman, *Lineare Operatoren in Hilberträumen*, Teubner verlag, Wiesbaden, 2003.
- [57] D. R. Yafaev, *Mathematical Scattering Theory: General Theory*, Amer. Math. Soc., Providence, Rhode Island, 1992.
- [58] K. Yajima, *Nonrelativistic limit of the Dirac theory, scattering theory*, J. Fac. Sci. Univ. Tokyo Sect. 1A, Math., 23(1976), pp. 517-523.
- [59] O. Yamada, *The nonrelativistic limit of modified wave operators for Dirac operators*, Proc. Japan Acad., Ser. A, 59(1983), pp. 71-74.
- [60] Y. You, J. Chen, and H. Lu *Filters, reproducing kernel, and adaptive meshfree method*, Comput. Mech., 31(2003), pp. 316-326.
- [61] S. Zhao, *On the spurious solutions in the high-order finite difference methods for eigenvalue problems*, Comp. Meth. Appl. Mech. Engng., 196(2007), pp. 5031-5046.
- [62] C. Zuppa, *Modified Taylor reproducing formulas and h-p clouds*, Math. Comput., 77(2008), pp. 243-264.