# COUPLED-CLUSTER THEORY REVISITED

## PART II: ANALYSIS OF THE SINGLE-REFERENCE COUPLED-CLUSTER EQUATIONS

Mihály A. Csirik[1,*] and Andre Laestadius[1,2]

**Abstract.** In a series of two articles, we propose a comprehensive mathematical framework for Coupled-Cluster-type methods. In this second part, we analyze the nonlinear equations of the single-reference Coupled-Cluster method using topological degree theory. We establish existence results and qualitative information about the solutions of these equations that also sheds light of the numerically observed behavior. In particular, we compute the topological index of the zeros of the single-reference Coupled-Cluster mapping. For the truncated Coupled-Cluster method, we derive an energy error bound for approximate eigenstates of the Schrödinger equation.

## 1. Introduction

The present article is concerned with the analysis of the single-reference Coupled-Cluster (SRCC) equations and its truncated variants. It is known that the *truncated* CC equations have a large number of solutions, some of which exhibit unphysical behavior. Recall that the truncated CC equations are *not* equivalent to the truncated CI equations (see Sect. 2.3 of Part I), hence, there is no obvious connection with the Schrödinger equation in general – truncated CC solutions can be very far from the desired Full CI solutions. Also, the SRCC method is typically unreliable for treating degenerate states.

Since the early days of the CC method, researchers in quantum chemistry have been interested in understanding the complicated behavior of the (truncated) CC equations [17–22, 30, 32, 44, 45]. This line of investigation mainly consisted in somehow "connecting" the truncated CC solutions to the untruncated ones, thereby identifying which truncated solutions can be considered "physical" and which should be discarded as "unphysical". The comprehensive investigation by K. Kowalski and P. Piecuch [32] classifies solutions based on a certain homotopy. While the universality of this approach is unclear (see Rem. 4.30 below), Kowalski and Piecuch [32] demonstrates without doubt that the solutions of the truncated CC equations exhibit rather complicated

[1] Hylleraas Centre for Quantum Molecular Sciences, Department of Chemistry, University of Oslo, P.O. Box 1033, Blindern, N-0315 Oslo, Norway.

[2] Department of Computer Science, Oslo Metropolitan University, P.O. Box 4, St. Olavs plass, NO-0130 Oslo, Norway.

*Corresponding author: m.a.csirik@kjemi.uio.no

behavior, which signals the need for a deeper analytical investigation. Since the CC equations are a system of quartic polynomial equations, the complete set of solutions can be computed by numerical means (Kowalski and Piecuch used the HOMPACK software [40]). However, in our analysis we view the CC equations as an abstract nonlinear equation and disregard their polynomial structure (this route was also taken by the seminal [37] and its follow-up works, discussed below). The main reason for this is that system is *very large*, and this seems prohibitive for the application of algebraic approaches.

In this article, we analyze the CC equations and the homotopy of Kowalski and Piecuch using a standard tool of nonlinear analysis, topological degree theory. In order to do this, we need to restrict ourselves to the finite-dimensional case (Rem. 3.1). The present article differs in flavor from the previous mathematical investigations, which where more standard "numerical analysis" approaches. However, a closer look at the qualitative properties of the (truncated) CC equations seems necessary, that, in addition to the ground-state also describes excited states.

## 1.1. Previous work

Similarly to Part I, our approach is based on the analysis of the single-reference CC method by R. Schneider [37]. In that seminal work, a thorough description of the basic building blocks of the method, namely excitation- and cluster operators, and their algebraic and functional-analytic properties are given. Under certain assumptions, the CC equations (a nonlinear system of equations consisting of quartic polynomials) are formulated in terms of a locally strongly monotone and locally Lipschitz operator defined on an appropriate space. Under further hypotheses, this establishes local existence and uniqueness of a (Galerkin projected) ground-state solution of the equation and moreover *quasi-optimality* of the projected CC solution (Thm. 5.8 ibid.). Perhaps the most important contribution of [37] is a quadratic energy error estimate (Thm. 6.3 ibid.). It is worth emphasizing that Schneider's analysis is a local one: *"[...] experience indicates that, in general, it cannot be expected that strong monotonicity always holds, or the constants might be extremely bad. Therefore we expect to get local existence results at best."* (ibid. p. 30)

Schneider's original analysis was carried out in the finite-dimensional case only. This was remedied in two subsequent articles by T. Rohwedder [35] and then by both of them [36]. The article [35] establishes important technical tools and rigorously proves that the untruncated CC problem is equivalent to the Full CI problem in the infinite-dimensional case, *i.e.* to essentially the Schrödinger equation. Using the said tools, the subsequent paper [36] also establishes local uniqueness and existence of a solution to the truncated CC equations in a neighborhood of the untruncated CC solution, under certain assumptions. Further, Rohwedder and Schneider [36] also extends the energy error estimates of [37] to the infinite-dimensional case.

This line of investigation was continued by S. Kvaal and A. Laestadius [24] for the Extended CC (ECC) method based on the "bivariational principle" [1]. In this case, local strong monotonicity for the ECC mapping can be established so that quasi-optimality follows along similar lines as previously done by Schneider and Rohwedder. In the ECC theory, the traditional CC theory is recovered as a special case.

Furthermore, the local strong monotonicity-based analysis was applied to a variant of the traditional CC method by F.M. Faulstich *et al.* [12], namely to the Tailored Coupled-Cluster (TCC) method. The TCC approach splits the computational task of evaluating the ground state energy into two parts: solving for the statically correlated wave function on a complete active space, and then on top of that accounting for the dynamical correlation using the CC method. Numerical investigations based on [12] were conducted in [13].

Finally, we mention the survey article [23] for more details on the use of local strong monotonicity-based methods in the analysis of CC methods.

## 1.2. Outline

While the knowledge of the notations and results of Part I is not strictly necessary for this Part II, some of the results in Section 4 do employ *excitation graphs* to a certain extent, introduced in Section 3.2 of Part I.

To make the present work as self-contained as possible, we define all the necessary concepts without the use of Part I; in other words, in the "traditional" second-quantized way.

In Section 2, we describe the setting of the quantum-mechanical problems the CC theory is aimed at. In Section 3 we state a few results that we use from finite-dimensional topological degree theory.

The analysis of the single-reference CC (SRCC) method begins in Section 4. Basic properties of the SRCC mapping are discussed in Section 4.1. After this, the local properties of the SRCC mapping, such as strong monotonicity and topological index in both the non-degenerate-, and in the degenerate case are considered in Section 4.2. We also look at the complex SRCC mapping in Section 4.3.

In Section 4.4, an important class of homotopies is defined that can be used for proving the existence of a solution to the truncated SRCC mappings. In Section 4.5 a homotopy is considered that was invented specifically to connect CC methods of different truncation levels. We prove an existence result and calculate the topological index of the homotopy. Finally, we derive an energy error estimate in Section 4.6 using the results of Appendix B.

## 2. Background

The usual notation $B(a, r)$ is used for the open ball of radius $r$ and center $a$, also $B^*(a, r) = B(a, r) \setminus \{a\}$ denotes the punctured ball.

The spectrum of a linear operator $A$ is written $\sigma(A)$, the elements of its discrete spectrum as $\mathcal{E}_n(A)$, where $n = 0, 1, 2, \ldots$, if $A$ is bounded from below. We use the usual notation $[A, B] = AB - BA$ for the commutator. The (conjugate) transpose of $A$ is denoted as $A^\dagger$. For normed spaces $V$ and $W$, the symbol $\mathcal{L}(V, W)$ denotes normed space of *bounded* linear mappings $V \to W$ endowed with the operator norm $\|\cdot\|_{\mathcal{L}(V,W)}$. Furthermore, $V^*$ denotes the (continuous) dual space. As usual, $[a, b]$ denotes the (closed) line segment between $a, b \in V$.

### 2.1. Second quantization

In this section, we first briefly review the framework of second quantization. The notations mainly follow [25, 39]. Next, we formulate the Schrödinger Hamiltonian in second quantization. Finally, we introduce the excitation-, and cluster operators and cluster amplitude spaces through the use of creation-, and annihilation operators.

#### 2.1.1. Fermionic Fock space

For simplicity, and because the Hamiltonian we will be considering is spin-independent, we neglect spin and consider $\mathfrak{h} = L^2(\mathbb{R}^3)$ as the *one-particle Hilbert space*. The tensor powers of $\mathfrak{h}$ are denoted by $\mathfrak{h}^N = \bigotimes^N \mathfrak{h}$ and antisymmetric powers of $\mathfrak{h}$ are denoted by $\mathfrak{h}_a^N = \bigwedge^N \mathfrak{h}$ for any $N \geq 0$, where we introduced the notation $\mathfrak{h}^0 = \mathbb{C}$ and $\mathfrak{h}_a^0 = \mathbb{C}$. The *fermionic Fock space* is then the Hilbert space given by the infinite direct sum $\mathfrak{F}_a := \bigoplus_{N \geq 0} \mathfrak{h}_a^N = \mathfrak{h}_a^0 \oplus \mathfrak{h}_a^1 \oplus \mathfrak{h}_a^2 \oplus \ldots$, in other words

$$\mathfrak{F}_a = \left\{ \Psi = (\psi^0, \psi^1, \ldots) : \psi^N \in \mathfrak{h}_a^N \ (N \geq 0), \ \|\Psi\|_{\mathfrak{F}}^2 := \sum_{N \geq 0} \|\psi^N\|_{\mathfrak{h}^N}^2 < \infty \right\},$$

endowed with the inner product $\langle \Psi_1, \Psi_2 \rangle_{\mathfrak{F}} = \sum_{N \geq 0} \langle \psi_1^N, \psi_2^N \rangle_{\mathfrak{h}^N}$, for any $\Psi_1 = (\psi_1^0, \psi_1^1, \ldots) \in \mathfrak{F}_a$ and $\Psi_2 = (\psi_2^0, \psi_2^1, \ldots) \in \mathfrak{F}_a$. The *vacuum state* is the distinguished element $\Omega = (1, 0, \ldots) \in \mathfrak{F}_a$. The space $\mathfrak{h}_a^N$ can be identified as a subspace of $\mathfrak{F}_a$, in this context it is called the *$N$-particle sector of $\mathfrak{F}$*. For $\Psi_1 \in \mathfrak{h}_a^{N_1}$ and $\Psi_2 \in \mathfrak{h}_a^{N_2}$ define $\Psi_1 \wedge \Psi_2 \in \mathfrak{h}_a^{N_1+N_2}$ *via*

$$\Psi_1 \wedge \Psi_2(\mathbf{x}_1, \ldots, \mathbf{x}_{N_1+N_2}) = C_{N_1,N_2} \sum_{\sigma \in \mathfrak{S}_{N_1+N_2}} (\text{sgn } \sigma) \Psi_1 \big( \mathbf{x}_{\sigma(1)}, \ldots, \mathbf{x}_{\sigma(N_1)} \big) \Psi_2 \big( \mathbf{x}_{\sigma(N_1+1)}, \ldots, \mathbf{x}_{\sigma(N_1+N_2)} \big),$$

with the normalization constant $C_{N_1,N_2} = (N_1! N_2! (N_1 + N_2)!)^{-1/2}$. Here, $\mathfrak{S}_n$ denotes the permutation group of $\{1, \ldots, n\}$ and $\text{sgn}\sigma$ the sign of the permutation $\sigma$.

Let $\{\varphi_p\}_{p=1}^{\infty} \subset \mathfrak{h}$ be an orthonormal basis and define the *Slater determinant* $\Phi_\alpha := \varphi_{\alpha_1} \wedge \ldots \wedge \varphi_{\alpha_N} \in \mathfrak{h}_a^N$ for every multiindex $\alpha \in \mathbb{N}_0^N$ with $\alpha_1 < \ldots < \alpha_N$. Then $\{\Phi_\alpha\}_\alpha$ is an orthonormal basis of $\mathfrak{h}_a^N$.

### 2.1.2. Creation-, and annihilation operators

For fixed $\varphi \in \mathfrak{h}$, the (fermionic) *creation operator* $a^\dagger(\varphi) : \mathfrak{F}_a \to \mathfrak{F}_a$ is defined as $a^\dagger(\varphi)\Psi = \varphi \wedge \Psi$ for any $\Psi \in \mathfrak{h}_a^N$ and extended boundedly and linearly to the whole space. Also, for fixed $\varphi \in \mathfrak{h}$, define the (fermionic) *annihilation operator* $a(\varphi) : \mathfrak{F}_a \to \mathfrak{F}_a$ as the adjoint of $a^\dagger(\varphi)$, *i.e.* $\langle \Psi_1, a^\dagger(\varphi)\Psi_2 \rangle_{\mathfrak{F}} = \langle a(\varphi)\Psi_1, \Psi_2 \rangle_{\mathfrak{F}}$ for all $\Psi_1, \Psi_2 \in \mathfrak{F}_a$. It is an easy calculation to show that $(a(\varphi)\Psi)(\mathbf{x}_1, \ldots, \mathbf{x}_{N-1}) = N^{-1/2} \int \varphi(\mathbf{x})\Psi(\mathbf{x}, \mathbf{x}_1, \ldots, \mathbf{x}_{N-1}) \, d\mathbf{x}$ for all $\Psi \in \mathfrak{h}_a^N$ ($N \geq 1$) and $a(\varphi)\Omega = 0$. The *canonical anticommutation relations (CAR)*

$$a^\dagger(\varphi)a(\varphi') + a(\varphi')a^\dagger(\varphi) = \langle \varphi, \varphi' \rangle_{\mathfrak{h}} 1_{\mathfrak{F}_a},$$
$$a^\dagger(\varphi)a^\dagger(\varphi') + a^\dagger(\varphi')a^\dagger(\varphi) = 0,$$
$$a(\varphi)a(\varphi') + a(\varphi')a(\varphi) = 0,$$

hold true for any $\varphi, \varphi' \in \mathfrak{h}$, where $1_{\mathfrak{F}_a} : \mathfrak{F}_a \to \mathfrak{F}_a$ denotes the identity map on $\mathfrak{F}_a$. Since both $a^\dagger(\varphi)$ and $a(\varphi)$ are linear and bounded in $\varphi$, it is enough to specify them on an orthonormal basis $\{\varphi_p\}_{p\geq 1} \subset \mathfrak{h}$, *i.e.* $a_p^\dagger := a^\dagger(\varphi_p)$ and $a_p := a(\varphi_p)$ ($p \geq 1$) completely determines the families $\{a^\dagger(\varphi)\}_{\varphi \in \mathfrak{h}}$ and $\{a(\varphi)\}_{\varphi \in \mathfrak{h}}$. It is then clear that any Slater determinant $\Phi_\alpha \in \mathfrak{h}_a^N$ can be "created" from the vacuum state as $\Phi_\alpha = a_{\alpha_1}^\dagger \cdots a_{\alpha_N}^\dagger \Omega$.

### 2.1.3. First-, and second quantization of one-, and two-body operators

Let $h : \mathfrak{h} \to \mathfrak{h}$ be a linear operator. Define its *first quantization* as $\sum_{1 \leq j \leq N} h_j : \mathfrak{h}_a^N \to \mathfrak{h}_a^N$ where

$$h_j = 1 \otimes \ldots \otimes 1 \otimes \underbrace{h}_{j\text{th}} \otimes 1 \otimes \ldots \otimes 1.$$

Here, $1$ denotes the identity on $\mathfrak{h}$. The *second quantization* of $h$ is defined as $0 \oplus \bigoplus_{N \geq 1} \sum_{1 \leq j \leq N} h_j : \mathfrak{F}_a \to \mathfrak{F}_a$.

Analogously, one can define the first quantization of a two-body operator $W : \mathfrak{h}_a^2 \to \mathfrak{h}_a^2$ as $\sum_{1 \leq i,j \leq N} W_{ij} : \mathfrak{h}_a^N \to \mathfrak{h}_a^N$, where

$$W_{ij} = 1 \otimes \ldots \otimes 1 \otimes \underbrace{W}_{i\text{th}} \otimes 1 \ldots 1 \otimes \underbrace{W}_{j\text{th}} \otimes 1 \otimes \ldots \otimes 1.$$

The second quantization of $W$ is $0 \oplus 0 \oplus \bigoplus_{N \geq 2} \sum_{1 \leq i < j \leq N} W_{ij} : \mathfrak{F}_a \to \mathfrak{F}_a$. Both the one-, and the two-body operators can be expressed using creation-, and annihilation operators as follows.

**Theorem 2.1.** *Let $\{\varphi_p\}_{p\geq 1} \subset \mathfrak{h}$ be an orthonormal basis. Then*

$$0 \oplus \bigoplus_{N \geq 1} \sum_{1 \leq j \leq N} h_j = \sum_{p,q \geq 1} \langle h\varphi_p, \varphi_q \rangle a_p^\dagger a_q.$$

*and*

$$0 \oplus 0 \oplus \bigoplus_{N \geq 2} \sum_{1 \leq i < j \leq N} W_{ij} = \sum_{\substack{1 \leq p \leq q \\ 1 \leq r \leq s}} \langle W\varphi_p \wedge \varphi_q, \varphi_r \wedge \varphi_s \rangle a_p^\dagger a_q^\dagger a_r a_s.$$

## 2.2. Schrödinger Hamiltonian in second-quantized form

For convenience and in accordance to Part I, we define the $N$-particle fermionic spaces

$$\mathfrak{L}^2 := \mathfrak{h}_a^N = \bigwedge^N L^2(\mathbb{R}^3), \quad \text{and} \quad \mathfrak{H}^1 := \mathfrak{L}^2 \cap H^1(\mathbb{R}^{3N}),$$

endowed with usual inner products $\langle \cdot, \cdot \rangle$ and $\langle \cdot, \cdot \rangle_{\mathfrak{H}^1}$, and norms $\|\cdot\|$ and $\|\cdot\|_{\mathfrak{H}^1}$, see Section 2.1 in Part I. We also set $\mathfrak{H}^2 := \mathfrak{L}^2 \cap H^2(\mathbb{R}^{3N})$, and $\mathfrak{H}^{-1} := (\mathfrak{H}^1)^*$ as usual. Henceforth, we employ the same convention as explained after Remark 2.1 in Part I.

Recall (see Sect. 2.2 of Part I) the Schrödinger Hamiltonian $\mathcal{H} : \mathfrak{L}^2 \to \mathfrak{L}^2$, which is the self-adjoint operator

$$\mathcal{H} = \sum_{1 \leq i \leq N} \left[ -\frac{1}{2} \triangle_{\mathbf{x}_i} + V(\mathbf{x}_i) \right] + \sum_{1 \leq i < j \leq N} w(\mathbf{x}_i - \mathbf{x}_j),$$

with domain $D(\mathcal{H}) = \mathfrak{H}^2$, where $V, w : \mathbb{R}^3 \to \mathbb{R}$ are Kato class potentials: $V, w \in L^{3/2}(\mathbb{R}^3) + L^{\infty}_{\varepsilon}(\mathbb{R}^3)$ and $w$ is even. The quadratic form corresponding to $\mathcal{H}$ is denoted as $\mathcal{E}(\Psi) = \langle \mathcal{H}\Psi, \Psi \rangle$, and has form domain $Q(\mathcal{E}) = \mathfrak{H}^1$. Recall that there is a constant $M > 0$, such that

$$\langle \mathcal{H}\Psi, \Phi \rangle \leq M \|\Psi\|_{\mathfrak{H}^1} \|\Phi\|_{\mathfrak{H}^1} \tag{2.1}$$

for all $\Psi, \Phi \in \mathfrak{H}^1$. Thus, $\mathcal{H}$ can be extended to a bounded mapping $\mathfrak{H}^1 \to \mathfrak{H}^{-1}$, which we denote with the same symbol. We say that $\Psi \in \mathfrak{H}^1$ and $\mathcal{E} \in \mathbb{R}$ satisfy the *weak Schrödinger equation* if $\langle \mathcal{H}\Psi, \Phi \rangle = \mathcal{E}\langle \Psi, \Phi \rangle$ for all $\Phi \in \mathfrak{H}^1$.

According to Theorem 2.1, $\mathcal{H}$ can also be given in the second quantized form

$$\mathcal{H} = \sum_{p,q \geq 1} h_{pq} a_p^{\dagger} a_q + \sum_{\substack{1 \leq p \leq q \\ 1 \leq r \leq s}} W_{pq,rs} a_p^{\dagger} a_q^{\dagger} a_r a_s,$$

where $h_{pq} = \frac{1}{2}\langle \boldsymbol{\nabla}\varphi_p, \boldsymbol{\nabla}\varphi_q \rangle + \langle V\varphi_p, \varphi_q \rangle$ and $W_{pq,rs} = \langle W\varphi_p \wedge \varphi_q, \varphi_r \wedge \varphi_s \rangle$.

### 2.2.1. Truncation of the orbital basis

Our analysis of the CC method is restricted to the finite-dimensional case due to reasons described later (Rem. 3.1). Let $K \geq N$ and assume that an $L^2$-orthonormal *(spin-)orbital set* $\{\varphi_p\}_{1 \leq p \leq K} \subset H^1(\mathbb{R}^3)$ is given. We define

$$\mathfrak{B}_K = \{\Phi_\alpha \in \mathfrak{H}^1 : 1 \leq \alpha_1 < \ldots < \alpha_N \leq K\},$$

as in Section 2.1 of Part I. Then $\mathfrak{B}_K$ is $\mathfrak{L}^2$-orthonormal. Set

$$\mathfrak{H}^1_K = \operatorname{Span} \mathfrak{B}_K \subset \mathfrak{H}^1.$$

We define $\mathcal{H}_K : \mathfrak{H}^1_K \to (\mathfrak{H}^1_K)^*$ to be the projection of $\mathcal{H} : \mathfrak{H}^1 \to \mathfrak{H}^{-1}$ onto the finite-dimensional subspace $\mathfrak{H}^1_K \subset \mathfrak{H}^1$, *i.e.* $\langle \mathcal{H}_K \Phi, \Psi \rangle = \langle \mathcal{H}\Phi, \Psi \rangle$ for all $\Phi, \Psi \in \mathfrak{H}^1_K$. Then $\Psi \in \mathfrak{H}^1_K$ and $\mathcal{E} \in \mathbb{R}$ are said to satisfy the *projected Schrödinger equation* if

$$\langle \mathcal{H}_K \Psi, \Phi \rangle = \mathcal{E}\langle \Psi, \Phi \rangle \quad \text{for all} \quad \Phi \in \mathfrak{H}^1_K.$$

For a convergence theory for this eigenvalue problem, see *e.g.* Theorem 5.11 of [41]. Vaguely speaking, if $\bigcup_{K > N} \{\varphi_p\}_{1 \leq p \leq K}$ is dense in $H^1(\mathbb{R}^3)$, then convergence of the projected eigenvalues to the discrete spectrum of $\mathcal{H}$ is guaranteed as $K \to \infty$ (see [37]).

From the computational standpoint, orbital basis truncation is always necessary when solving the Schrödinger equation. Of course, the spectral structure changes (the essential spectrum disappears in finite dimensions), and this means that we replaced the original problem with one having different qualitative properties. However, the CI and CC methods are mainly used to approximate the ground-state-, or the first few excited energies (and other quantities related to the lowest eigenstates). With a sufficiently good choice the orbitals $\{\varphi_p\}_{1 \leq p \leq K}$, these energies can be approximated rather well in a lot of situations. Therefore, it is safe to regard the rest of the projected spectrum as "junk", and simply ignore it. We can conclude that the restriction to finite dimensions allows for a meaningful theory that is still relevant to CC methods.

### 2.3. The Hartree–Fock method

In practice, the CC method usually takes the HF orbitals as an input and builds up correction based on them, as described in Section 2.3 of Part I. Here, we collect the basic facts about the HF method that will be used later on.

The HF method[1] is based on the minimiziation of the energy over Slater determinants. It is fairly easy to see by direct calculation that the energy $\mathcal{E}(\Phi) = \langle \mathcal{H}\Phi, \Phi \rangle$ of a Slater determinant $\Phi = \varphi_1 \wedge \ldots \wedge \varphi_N \in \mathfrak{H}^1$, with $\{\varphi_j\}_{j=1}^N \subset H^1(\mathbb{R}^3)$ being $L^2$-orthonormal is given by

$$\mathcal{E}(\varphi_1, \ldots, \varphi_N) := \mathcal{E}(\Phi) = \frac{1}{2}\sum_{i=1}^N \int_{\mathbb{R}^3} |\boldsymbol{\nabla}\varphi_i|^2 + \sum_{i=1}^N \int_{\mathbb{R}^3} V|\varphi_i|^2$$
$$+ \frac{1}{2}\iint_{\mathbb{R}^3 \times \mathbb{R}^3} \left[ \sum_{i,j=1}^N |\varphi_i(\mathbf{x})|^2|\varphi_j(\mathbf{y})|^2 - \left| \sum_{i=1}^N \varphi_i(\mathbf{x})\overline{\varphi_i(\mathbf{y})} \right|^2 \right] w(\mathbf{x} - \mathbf{y})\,\mathrm{d}\mathbf{x}\mathrm{d}\mathbf{y},$$

see *e.g.* [7]. Hence the task is to determine the *HF energy*

$$\mathcal{E}_{\mathrm{HF}} = \mathcal{E}(\Phi_{\mathrm{HF}}) = \min_{\substack{\varphi_1, \ldots, \varphi_N \in H^1(\mathbb{R}^3) \\ \langle \varphi_i, \varphi_j \rangle = \delta_{ij}}} \mathcal{E}(\varphi_1, \ldots, \varphi_N), \tag{2.2}$$

along with a *HF minimizer* (or *HF determinant*) $\Phi_{\mathrm{HF}}$. The existence of a HF minimizer $\Phi_{\mathrm{HF}}$ is guaranteed for the case of positive ions and neutral atoms and molecules (corresponding to the electronic molecular Hamiltonian).

**Theorem 2.2** ([27]). *If $N < Z + 1$, then there exists a minimizer $\Phi_{\mathrm{HF}}$ to* (2.2).

Following the seminal work [27], much more has been discovered about the mathematical structure of the HF energy functional [14, 26, 28]. As usual, a minimizer satisfies the corresponding Euler–Lagrange equations (which are called *Hartree–Fock equations* in this context). In practice, it is this system of nonlinear integro-differential equations which is discretized and solved. To describe the HF equations, we make a definition that will be convenient for later purposes. Fix $\Phi = \varphi_1 \wedge \ldots \wedge \varphi_N$, where $\{\varphi_p\}_{p=1}^N \subset H^1(\mathbb{R}^3)$ is $L^2$-orthonormal. Define the lower semibounded, self-adjoint operator $F_\Phi : L^2(\mathbb{R}^3) \to L^2(\mathbb{R}^3)$ with domain $D(F_\Phi) = H^2(\mathbb{R}^3)$ *via* the instruction

$$(F_\Phi \psi)(\mathbf{x}) = -\frac{1}{2}\triangle\psi(\mathbf{x}) + V(\mathbf{x})\psi(\mathbf{x}) + \left( \sum_{i=1}^N \int_{\mathbb{R}^3} w(\mathbf{x} - \mathbf{y})|\varphi_i(\mathbf{y})|^2 \right)\psi(\mathbf{x})$$
$$- \sum_{i=1}^N \left( \int_{\mathbb{R}^3} w(\mathbf{x} - \mathbf{y})\overline{\varphi_i(\mathbf{y})}\psi(\mathbf{y})\,\mathrm{d}\mathbf{y} \right)\varphi_i(\mathbf{x})$$

for all $\psi \in D(F_\Phi)$ and all $\mathbf{x} \in \mathbb{R}^3$. The operator $F_\Phi$ is called the *mean-field operator*[2]. The form domain of $F_\Phi$ is $H^1(\mathbb{R}^3)$. The essential spectrum of $F_\Phi$ is $[0, +\infty)$. We summarize the basic properties of the mean-field operator in the next theorem. Let

$$\mu_n(F_\Phi) = \min_{\substack{\psi_1, \ldots, \psi_n \in H^1(\mathbb{R}^3) \\ \langle \varphi_i, \varphi_j \rangle = \delta_{ij}}} \max_{\substack{\psi \in \mathrm{Span}\{\psi_1, \ldots, \psi_n\} \\ \|\psi\|=1}} \langle F_\Phi \psi, \psi \rangle$$

denote the min-max values of $F_\Phi$.

---

[1]a.k.a. Self-Consistent Field (SCF) method

[2]It is sometimes called the Fock operator, but we reserve that name for its $N$-particle version, *i.e.* its first quantization.

**Theorem 2.3.** *Assume that there exists a HF minimizer $\Phi_{\mathrm{HF}} = \varphi_1 \wedge \ldots \wedge \varphi_N$.*

(i) *(Hartree–Fock equations) There exists a unitary matrix $\mathbf{U} \in \mathbb{C}^{N \times N}$ so that $\widetilde{\varphi}_i$ are eigenfunctions of $F_\Phi$ corresponding to its $N$ lowest eigenvalues $\lambda_1 \leq \ldots \leq \lambda_N$,*

$$F_\Phi \widetilde{\varphi}_j = \lambda_j \widetilde{\varphi}_j, \quad \text{for all} \quad j = 1, \ldots, N, \tag{2.3}$$

*where $(\widetilde{\varphi}_1, \ldots, \widetilde{\varphi}_N) = \mathbf{U}(\varphi_1, \ldots, \varphi_N)$.*

(ii) *(Aufbau principle) If $\mu_{N+1}(F_\Phi)$ is an eigenvalue of $F_\Phi$, then $\lambda_N = \mu_N(F_\Phi) < \mu_{N+1}(F_\Phi) \leq 0$.*

*Proof.* See, for example [4]. □

The eigenvalue $\lambda_N$ corresponds to the *highest occupied molecular orbital* (HOMO) and $\lambda_{N+1}$ to the *lowest unoccupied molecular orbital* (LUMO). Their difference, $\varepsilon_{\min} := \lambda_{N+1} - \lambda_N$ is called the *HOMO-LUMO gap*, which is an important quantity in quantum chemistry [2].

The first quantization of the mean-field operator $F_\Phi$ is called the *Fock operator* $\mathcal{F}_\Phi : \mathfrak{L}^2 \to \mathfrak{L}^2$ with domain $D(\mathcal{F}_\Phi) = \bigwedge^N D(F_\Phi) = \mathfrak{H}^2$,

$$\mathcal{F}_\Phi = F_\Phi \otimes I \otimes \ldots \otimes I + \ldots + I \otimes \ldots \otimes I \otimes F_\Phi.$$

Henceforth we omit $\Phi$ from the notation, and let $\mathcal{F} := \mathcal{F}_\Phi$ whenever $\Phi$ is clear from the context. It is immediate that the HF determinant $\Phi_0 := \Phi_{\mathrm{HF}}$ is an eigenfunction of $\mathcal{F}$,

$$\mathcal{F}\Phi_0 = \Lambda_0 \Phi_0, \quad \text{with} \quad \Lambda_0 = \sum_{i=1}^{N} \lambda_i.$$

The Fock operator gives rise to a splitting of the molecular Hamiltonian

$$\mathcal{H} = \mathcal{F} + \mathcal{W}, \quad \text{where} \quad \mathcal{W} := \mathcal{H} - \mathcal{F} \tag{2.4}$$

is called the *fluctuation operator*.

For the rest of the section, we consider the finite-dimensional case. In practice, the Galerkin projection of the Hartree–Fock equations equations (2.3) are solved to obtain the orbitals $\{\varphi_i\}_{i=1}^{N}$. Since the mean-field operator $F_\Phi$ is self-adjoint, its eigenfunctions can be used to extend these orbitals to an orthonormal basis $\{\varphi_i\}_{i=1}^{K} \subset H_K^1(\mathbb{R}^3)$. Similarly to $\mathcal{H}_K$ we introduce the projected versions of $\mathcal{F}$ and $\mathcal{W}$, denoted as $\mathcal{F}_K$ and $\mathcal{W}_K$. In the orbital basis $\{\varphi_i\}_{i=1}^{K}$, the Fock operator takes the diagonal form $\mathcal{F}_K = \sum_{i=1}^{K} \lambda_i a_i^\dagger a_i$.

Furthermore, if $\Phi_\alpha = \varphi_{\alpha_1} \wedge \ldots \wedge \varphi_{\alpha_N}$ with $\alpha_1 < \ldots < \alpha_N$ that is obtained from $\Phi_0 = \varphi_1 \wedge \ldots \wedge \varphi_N$ by swapping $r$ orbitals $\varphi_{I_j}$ with $\varphi_{A_j}$ where $I_j \in \{1, \ldots, N\}$, $A_j \notin \{1, \ldots, N\}$ $(j = 1, \ldots, r \leq N)$, then

$$\mathcal{F}_K \Phi_\alpha = \Lambda_\alpha \Phi_\alpha, \quad \text{with} \quad \Lambda_\alpha = \sum_{i=1}^{N} \lambda_{\alpha_i} = \Lambda_0 + \varepsilon_\alpha, \quad \varepsilon_\alpha = \sum_{j=1}^{r} \lambda_{A_j} - \lambda_{I_j}. \tag{2.5}$$

We would like to emphasize that in the finite-dimensional case, it is possible to have $\lambda_N \leq 0 < \lambda_{N+1}$[3].

Unfortunately, in the infinite-dimensional case, it might not be possible to construct a complete eigenbasis $\{\varphi_i\}_{i=1}^{\infty} \subset H^1(\mathbb{R}^3)$ for the mean-field operator $F_\Phi$.

---

[3] The addition of diffuse functions of the form $e^{-\beta|\mathbf{x}|^2}$ for $\beta \ll 1$ can significantly reduce the HOMO-LUMO gap $\varepsilon_{\min}$.

## 2.4. Excitation-, and cluster operators

We now briefly recall the basic properties excitation-, and cluster operators using the second-quantized language. Alternatively, the reader may consult Section 3.2 of Part I. Here, we introduce the concepts in arbitrary dimensions for completeness, although we only need the finite-dimensional case in the sequel.

Suppose that $\Phi_0 := \Phi_{\mathrm{HF}} = \varphi_1 \wedge \ldots \wedge \varphi_N$ is a HF minimizer for some orthonormal $\{\varphi_p\}_{p=1}^N$, as discussed in the preceding section. The wavefunction $\Phi_0$ is called the *reference determinant* or simply the *reference*. The functions $\{\varphi_p\}_{p=1}^N$ are called *occupied orbitals* and are extended to a complete orthonormal basis of $H_K^1(\mathbb{R}^3)$ with the *virtual orbitals* $\{\varphi_p\}_{p=N+1}^K$ (here, $K = \infty$ is allowed). According to Section 2.1.2, any Slater determinant $\Phi_\alpha$ may be obtained from the Fock vacuum *via* a string of creation operators, in particular $\Phi_0 = a_1^\dagger \cdots a_N^\dagger \Omega$. Here, and henceforth, the symbol 0 denotes $(1, \ldots, N)$ as multiindex.

As in the preceding section, we may partition the indices $\alpha_1 < \ldots < \alpha_N$ into occupied-, and virtual ones by letting $I_{r+1} < \ldots < I_N$ be the $\alpha_j$'s which lie in $\{1, \ldots, N\}$ and $A_1 < \ldots < A_r$ be the $\alpha_j$'s which lie in $\{N+1, N+2, \ldots\}$, where $r$ is the number of virtual indices in $\alpha$, also called the *rank of* $\alpha$. Now let $I_1 < \ldots < I_r$ be given by $\{I_1, \ldots, I_r\} = \{1, \ldots, N\} \smallsetminus \{I_{r+1}, \ldots, I_N\}$, then we can write

$$\Phi_\alpha = a_{A_1}^\dagger a_{I_1} \cdots a_{A_r}^\dagger a_{I_r} \Phi_0,$$

since the string of creation and annihilation operators acting on $\Phi_0$ clearly swaps the occupied orbital $\varphi_{I_j}$ with $\varphi_{A_j}$ in $\Phi_0$ for all $j = 1, \ldots, r$, yielding precisely $\Phi_\alpha$. Then $X_\alpha = a_{A_1}^\dagger a_{I_1} \cdots a_{A_r}^\dagger a_{I_r}$ is called an *excitation operator*. We say that the Slater determinant $\Phi_\alpha$ (or the multiindex $\alpha$), is *singly-* (S), *doubly-* (D), *triply-* (T) etc. *excited* if the rank of $\alpha$ is 1, 2, 3, etc.

This way we constructed a family of bounded operators $X_\alpha$, where $\alpha$ runs over all possible *nonzero* multiindices. The adjoint of an excitation operator $X_\alpha$, denoted $X_\alpha^\dagger$, is called a *de-excitation operator*. Notice that both excitation-, and de-excitation operators conserve the particle number.

**Theorem 2.4.** *The following properties hold true.*

(i) *The product of (de-)excitation operators is an (de-)excitation operator, or, zero.*
(ii) *(commutativity)* $X_\alpha X_\beta = X_\beta X_\alpha$ *and* $X_\alpha^\dagger X_\beta^\dagger = X_\beta^\dagger X_\alpha^\dagger$.
(iii) *(nilpotency)* $X_\alpha^2 = 0$ *and* $(X_\alpha^\dagger)^2 = 0$.

Linear combinations of excitation operators are called *cluster operators*, *i.e.* operators of the form $C = \sum_{\alpha \neq 0} c_\alpha X_\alpha$. Clearly, there is a one-to-one correspondence between $C$ and the so-called *cluster amplitude* $(c_\alpha)_{\alpha \neq 0}$. We use the standard convention that capital letters $C, T, S, \ldots$ denote cluster operators and small letters $c, t, s, \ldots$ their corresponding (unique) cluster amplitudes. Clearly, there is a one-to-one correspondence between functions $\Psi \in \mathfrak{H}^1$ with $\langle \Psi, \Phi_0 \rangle = 0$ and the cluster operators $C_\Psi$ defined as

$$C_\Psi = \sum_{\alpha \neq 0} c_\alpha X_\alpha, \quad \text{where} \quad c_\alpha = \langle \Psi, \Phi_\alpha \rangle. \tag{2.6}$$

**Theorem 2.5** ([35]). *Fix* $\Psi \in \mathfrak{H}^1$, *such that* $\langle \Psi, \Phi_0 \rangle = 0$. *Then the following properties hold true.*

(i) *The vector space of cluster operators extended with the identity is a closed, commutative, nilpotent subalgebra of* $\mathcal{L}(\mathfrak{H}^1)$. *In particular,* $T^{N+1} = 0$ *for any* $T$.
(ii) *The cluster operator* $C_\Psi$ *satisfies* $C_\Psi \in \mathcal{L}(\mathfrak{H}^1, \mathfrak{H}^1)$. *Furthermore, there is a constant* $b > 0$ *independent of* $\Psi$ *such that* $\|\Psi\|_{\mathfrak{H}^1} \leq \|C_\Psi\|_{\mathcal{L}(\mathfrak{H}^1, \mathfrak{H}^1)} \leq b \|\Psi\|_{\mathfrak{H}^1}$.
(iii) $C_\Psi^\dagger \in \mathcal{L}(\mathfrak{H}^1, \mathfrak{H}^1)$, *and there is a constant* $b' > 0$ *independent of* $\Psi$ *such that* $\|C_\Psi^\dagger\|_{\mathcal{L}(\mathfrak{H}^1, \mathfrak{H}^1)} \leq b' \|\Psi\|_{\mathfrak{H}^1}$, *and there cannot be a uniform lower bound in terms of* $\|\Psi\|_{\mathfrak{H}^1}$.
(iv) $C_\Psi$ *can be extended to* $\mathcal{L}(\mathfrak{H}^{-1}, \mathfrak{H}^{-1})$.

According to Lemma 5.2 of [35], for any cluster operator $C \in \mathcal{L}(\mathfrak{H}^1, \mathfrak{H}^1)$, there exists a cluster operator $T \in \mathcal{L}(\mathfrak{H}^1, \mathfrak{H}^1)$, such that $I + C = e^T$. This $T$ is uniquely given by $\log(I + C)$. We would like to point out, that the exponential representation is only true in the *untruncated case*, *i.e.* when *all* cluster amplitudes $(t_\alpha)$ and all excitation operators $X_\alpha$ are considered. If some truncation is in effect (see below), the set $\{e^T \Phi_0\}$ is unknown.

## 2.5. Cluster amplitude spaces

As mentioned in the previous section, the linear combination coefficients $t = (t_\alpha)$ of a cluster operator expansion $T = \sum_{\alpha \neq 0} t_\alpha X_\alpha$ are called cluster amplitudes. If some *truncation scheme* is present, some of the $t_\alpha$'s and $X_\alpha$'s will be missing (again we refer to Sect. 3.5 of Part I for the precise formulation). This truncation procedure is crucial for a feasible implementation of the CC method, and most commonly only the singly- (S), or singly and doubly- (SD), or the singly, doubly and triply (SDT) excited amplitudes are kept.

The vector space of (possibly truncated) cluster amplitudes form a subspace of $\ell^2$. We define the *(cluster) amplitude space*

$$\mathbb{V} = \left\{ t \in \ell^2 : \|T\Phi\|_{\mathfrak{H}^1} < \infty \right\},$$

endowed with the inner product $\langle t, s \rangle_{\mathbb{V}} = \langle T\Phi_0, S\Phi_0 \rangle_{\mathfrak{H}^1}$. Analogously to $\mathfrak{H}^1 \subset \mathcal{L}^2 \subset \mathfrak{H}^{-1}$, the spaces $\mathbb{V} \subset \ell^2 \subset \mathbb{V}^*$ form a Gelfand triple, so Remark 2.1 of Part I applies. To $\mathbb{V}$, there corresponds the *functional amplitude space*

$$\mathfrak{V} = \{T\Phi_0 \in \mathfrak{H}^1 : t \in \mathbb{V}\}.$$

In accordance to the concepts developed in Part I, we denote the *full*, *i.e.* untruncated amplitude space and the corresponding functional amplitude space by $\mathbb{V}(G^{\text{full}})$ and $\mathfrak{V}(G^{\text{full}})$, respectively.

**Remark 2.6.** In the case $K < \infty$, and assuming that HOMO-LUMO gap satisfies $\varepsilon_{\min} > 0$ (no matter how small), following [37], we introduce the norm

$$\vert\!\vert\!\vert t \vert\!\vert\!\vert^2 := \sum_{\alpha \neq 0} \varepsilon_\alpha |t_\alpha|^2 \quad \text{for all} \quad t \in \mathbb{V},$$

where $\varepsilon_\alpha$ denotes the eigenvalues of the Fock operator $\mathcal{F}_K$, see (2.5). It was shown in [37] that under certain assumptions, the constants in the norm equivalence $\vert\!\vert\!\vert \cdot \vert\!\vert\!\vert \sim \| \cdot \|_{\mathbb{V}}$ are independent of $K$. Finally, we also define the norm $\vert\!\vert\!\vert \cdot \vert\!\vert\!\vert$ on $\mathfrak{V}$ *via* $\vert\!\vert\!\vert T\Phi_0 \vert\!\vert\!\vert := \vert\!\vert\!\vert t \vert\!\vert\!\vert$ for any $t \in \mathbb{V}$.

## 3. TOPOLOGICAL DEGREE THEORY

In this section, we give a short introduction to the basic notions of *finite-dimensional* topological degree theory, and state the results that we will need in the forthcoming analysis of the CC method. While topological degree theory is well-known in the field of nonlinear analysis and dynamical systems, it is not so ubiquitous in mathematical physics (a nice application is furnished by Ginzburg–Landau theory, see *e.g.* [9], Sect. 2.4). Therefore, an inclusion of a brief summary of this important and powerful tool in the present article seems justified.

Here, we mainly follow the set of notes by Dinca and Mawhin [9] (which was released as a book recently [10]), where the proofs may be found. The concept of the topological (a.k.a. mapping-, or Brouwer-) degree of a mapping goes back to Kronecker, Poincaré and Brouwer; Leray and Schauder extended the concept to infinite dimensions and demonstrated its usefulness in the theory of partial differential equations. (See Chap. 23 of [9] for a fascinating historical perspective.) The standard textbooks on the topic are [8, 29, 43].

**Remark 3.1.** Unfortunately, there cannot be a "topological degree theory" for general continuous mappings in infinite dimensions. The reason for this, roughly speaking, is that such a theory would imply the Brouwer fixed point theorem, but that does not hold in Hilbert space as pointed out by Kakutani ([11], Example 5.1.7). There are many generalizations of topological degree theory to infinite dimensions, all requiring some stringent assumptions on the class of mappings considered [11, 29, 31, 38].

## 3.1. Construction and basic properties

The first result gives an axiomatic characterization of the topological degree as the *unique* additive homotopy invariant that can be attached to continuous mappings (up to normalization).

**Theorem 3.2** (Existence and uniqueness of the degree). *There exists a unique integer-valued function*

$$(\mathcal{A}, D, z) \mapsto \deg(\mathcal{A}, D, z),$$

*where $D \subset \mathbb{R}^n$ is an open, bounded and nonempty subset, $\mathcal{A} \in C(\overline{D}, \mathbb{R}^n)$ and $z \notin \mathcal{A}(\partial D)$, such that the following properties hold true:*

(i) *(Normalization) If $z \in D$, then we have $\deg(\mathrm{id}, D, z) = 1$.*
(ii) *(Additivity) Let $D_1, D_2 \subset D$ be disjoint open subsets such that $z \notin \mathcal{A}(\overline{D} \smallsetminus (D_1 \cup D_2))$. Then $\deg(\mathcal{A}, D, z) = \deg(\mathcal{A}, D_1, z) + \deg(\mathcal{A}, D_2, z)$.*
(iii) *(Homotopy invariance) If $\mathcal{K} \in C(\overline{D} \times [0,1], \mathbb{R}^n)$ and $z \notin \mathcal{K}(\partial D \times [0,1])$, then $\lambda \mapsto \deg(\mathcal{K}(\cdot, \lambda), D, z)$ is constant.*

Among the many useful properties of the degree we highlight the following few.

**Corollary 3.3.**

(i) *(Excision property) Let $D \subset \mathbb{R}^n$ be an open, bounded and $D' \subset D$ open. If $z \notin \mathcal{A}(\overline{D} \smallsetminus D')$, then $\deg(\mathcal{A}, D, z) = \deg(\mathcal{A}, D', z)$.*
(ii) *(Existence property) If $z \notin \mathcal{A}(\overline{D})$, then $\deg(\mathcal{A}, D, z) = 0$. Said differently, if $\deg(\mathcal{A}, D, z) \neq 0$ for some open, bounded $D \subset \mathbb{R}^n$ such that $z \notin \mathcal{A}(\partial D)$, then there is a solution $u \in D$ to $\mathcal{A}(u) = z$.*
(iii) *(Additivity property) Suppose that $\{D_j\} \subset D$ is a sequence of open and disjoint sets. If $z \notin \mathcal{A}(\overline{D} \smallsetminus \bigcup D_j)$, then $\deg(\mathcal{A}, D_j, z) = 0$ for all but finitely many $j$, and $\deg(\mathcal{A}, D, z) = \sum_j \deg(\mathcal{A}, D_j, z)$.*

The following formula is essential for the practical computation of the degree.

**Proposition 3.4.** *Let $\mathcal{A} \in C(\overline{D}, \mathbb{R}^n) \cap C^1(D, \mathbb{R}^n)$ such that $z \notin \mathcal{A}(\partial D)$ and $\det \mathcal{A}'(u) \neq 0$ for all $u \in \mathcal{A}^{-1}(z)$. Then $\deg(\mathcal{A}, D, z) = \sum_{u \in \mathcal{A}^{-1}(z)} \mathrm{sgn} \det \mathcal{A}'(u)$.*

Based on this, we call a point $u$ *non-degenerate* if $\det \mathcal{A}'(u) \neq 0$ and *degenerate* otherwise. When talking about an *isolated solution* $u$ of $\mathcal{A}$, *i.e.* a point $u \in \mathcal{A}^{-1}(z)$ such that there is a ball $B(u,r)$ with $\mathcal{A}^{-1}(z) \cap B(u,r) = \{u\}$, the concept of the *index* is useful.

**Definition 3.5.** Let $\mathcal{A} \in C(\overline{D}, \mathbb{R}^n)$ and let $u$ be an isolated solution. Then the *(topological) index* of $\mathcal{A}$ at $u$ is defined as $i(\mathcal{A}, u) = \deg(\mathcal{A}, B(u,r), z)$, where $\mathcal{A}^{-1}(z) \cap B(u,r) = \{u\}$.

It follows from the excision property that $i(\mathcal{A}, u)$ is independent of $r$ (for sufficiently small values of $r > 0$), so the definition makes sense.

**Proposition 3.6.**

(i) *Let $\mathcal{A} \in C(\overline{D}, \mathbb{R}^n)$ such that $z \notin \mathcal{A}(\partial D)$ and $\mathcal{A}^{-1}(z)$ is finite. Then $\deg(\mathcal{A}, D, z) = \sum_{u \in \mathcal{A}^{-1}(z)} i(\mathcal{A}, u)$.*
(ii) *Let $u \in \mathbb{R}^n$ be a zero of $\mathcal{A} \in C(\overline{D}, \mathbb{R}^n)$, where $D$ an open neighborhood of $u$. If $\mathcal{A}$ is differentiable at $u$ with $\det \mathcal{A}'(u) \neq 0$, then $i(\mathcal{A}, u) = \mathrm{sgn} \det \mathcal{A}'(u)$.*

The topological degree (resp. index) can be naturally extended to continuous mappings of type $\mathcal{A} : X \to Y$, where $X$ and $Y$ are $n$-dimensional oriented topological vector spaces and the degree (resp. index) is independent of the choice of the basis. We refer the reader to Chapter 6 from [9] for details.

**Theorem 3.7** ([10], Thm. 1.3.1 and [9], Thm. 6.3.1). *Let $\mathcal{A} : D \to Y$ be continuous, where $X$ and $Y$ are $n$-dimensional topological vector spaces and $D \subset X$ is an open neighborhood of $u \in X$. Let $h : X \to \mathbb{R}^n$ and $g : Y \to \mathbb{R}^n$ be linear homeomorphisms. Suppose that $\mathcal{A}$ is differentiable at $u$ and that $\mathcal{A}'(u) : X \to Y$ is invertible. Then $i(\mathcal{A}, u) = \mathrm{sgn} \det g\mathcal{A}'(u)h^{-1}$.*

The following result says the degree is stable under (almost all) small perturbations of the right-hand side of $\mathcal{A}(u) = z$, and that the degree provides a lower bound for the number of solutions of the perturbed equation.

**Theorem 3.8** ([8], Cor. 7.4). *Let $D \subset \mathbb{R}^n$ be a bounded open set and let $\mathcal{A} : D \to \mathbb{R}^n$ be a $C^1$ mapping. If $z \notin \mathcal{A}(\partial D)$, then there is a $\delta > 0$ such that if $z' \in B(z, \delta) \setminus E_{z,\delta}$, then*

$$\deg(\mathcal{A}, D, z) = \deg(\mathcal{A}, D, z'),$$

*where $E_{z,\delta} \subset \{w \in B(z, \delta) : \det \mathcal{A}'(u) = 0, \ \mathcal{A}(u) = w\}$. Furthermore, $\mathcal{A}^{-1}(z')$ consists of a finite number $m$ of points, where $|\deg(\mathcal{A}, D, z)| \leq m$ and $m \equiv \deg(\mathcal{A}, D, z) \pmod 2$.*

The proof is based on Sard's theorem, which says that $E_{z,\delta}$ has $n$-dimensional Lebesgue measure zero. Next, we recall a tool that is very useful for the computation of the degree in the degenerate case.

**Theorem 3.9** ([9], Thm. 6.4.2). *Let $X$ and $Z$ be $n$-dimensional topological vector spaces. Let $L : X \to Z$ be a linear mapping with $\ker L \neq \{0\}$ and $V \subset Z$ be a vector space such that $Z = V \oplus \operatorname{ran} L$. Let $D \subset X$ open and bounded, $r : \overline{D} \to V$ continuous with $0 \notin \mathcal{A}(\partial D)$, where $\mathcal{A} = L + r$. Then, for each invertible linear map $J : \ker L \to V$, and each projector $P$ with $\operatorname{ran} P = \ker L$, the relation*

$$\deg(\mathcal{A}, D, 0) = i(L + JP, 0) \deg\big(J^{-1}r|_{\ker L}, D \cap \ker L, 0\big)$$

*called* Leray's second reduction formula *holds true.*

**Corollary 3.10.** *Let $L : X \to Z$ be a linear mapping with $\ker L \neq \{0\}$ and let $Q$ be a projector with $\ker Q = \operatorname{ran} L$. Also, let $\mathcal{N} : D \times [0, 1] \to Z$ be continuous, with $D \subset X$ open and bounded. Furthermore, assume the following hold true.*

(i) $Lu + \lambda \mathcal{N}(u, \lambda) \neq 0$ *for all $u \in \partial D$ and $\lambda \in (0, 1]$.*
(ii) $Q\mathcal{N}(u, 0) \neq 0$ *for each $u \in \partial D$.*

*Then*
$$\deg(L + \mathcal{N}(\cdot, 1), D, 0) = i(L + Q) \deg(Q\mathcal{N}(\cdot, 0)|_{\ker L}, D \cap \ker L, 0).$$

*Proof.* We adapt the the proof of Theorem 2.2.3 from [10]. Define the homotopy

$$\mathcal{A}(u, \lambda) = Lu + (1 - \lambda)Q\mathcal{N}(u, \lambda) + \lambda\mathcal{N}(u, \lambda).$$

Fix $\lambda \in (0, 1]$. Projecting the equation $\mathcal{A}(u, \lambda) = 0$ with $Q$ and $I - Q$ (*i.e.* onto the complementary spaces $\operatorname{ran} Q$ and $\operatorname{ran}(I - Q) = \ker Q = \operatorname{ran} L$), we get that $\mathcal{A}(u, \lambda) = 0$ is equivalent to the system $Q\mathcal{N}(u, \lambda) = 0$, $Lu + \lambda(I - Q)\mathcal{N}(u, \lambda) = 0$. But this is equivalent to $Lu + \lambda\mathcal{N}(u, \lambda) = 0$. By assumption (i), $\mathcal{A}(u, \lambda) \neq 0$ for all $u \in \partial D$ and $\lambda \in (0, 1]$. Further, $\mathcal{A}(u, 0) = 0$ is equivalent to the system $u \in \ker L$, $Q\mathcal{N}(u, 0) = 0$, hence by assumption (ii), $\mathcal{A}(u, 0) \neq 0$ for all $u \in \partial D$. Using the homotopy invariance and Leray's second reduction formula with $V = \ker L$, $J = I$ and $r = \mathcal{N}(\cdot, 0)$, we get

$$\deg(\mathcal{A}(\cdot, 1), D, 0) = \deg(\mathcal{A}(\cdot, 0), D, 0) = i(L + Q, 0) \deg(Q\mathcal{N}(\cdot, 0)|_{\ker L}, D \cap \ker L, 0).$$

$\square$

## 3.2. Orientation-preserving mappings

An important class of mappings for which the degree behaves rather nicely is the following.

**Definition 3.11.** Let $U \subset \mathbb{R}^n$ be open.

(i) A linear map $L : \mathbb{R}^n \to \mathbb{R}^n$ is said to be *orientation-preserving* if $\det L \geq 0$.

(ii) A mapping $\mathcal{A} \in C^2(U, \mathbb{R}^n)$ is said to be *strictly orientation-preserving in $U$* if
    (a) $\mathcal{A}'(u)$ is orientation-preserving for all $u \in U$, and
    (b) the set $\{u \in U : \det \mathcal{A}'(u) = 0\}$ is nowhere dense in $U$.

(iii) A mapping $\mathcal{A} \in C(U, \mathbb{R}^n)$ is said to be *orientation-preserving* if for every $V \subset U$ open with $\overline{V} \subset U$ is compact, there exists a sequence $\{\mathcal{A}_j\}$ such that $\mathcal{A}_j$ is strictly orientation preserving for all $j$ and $\mathcal{A}_j \to \mathcal{A}$ uniformly on $V$.

**Theorem 3.12** ([10], Thm. 7.2.1 and [9], Thm. 19.3.1). *Let $U \subset \mathbb{R}^n$ be open, $D \subset U$ an open bounded set with $\overline{D} \subset U$ and let $\mathcal{A} : U \to \mathbb{R}^n$ be orientation-preserving. If $z \notin \mathcal{A}(\partial D)$, then the following hold true.*

  (i) $\deg(\mathcal{A}, D, z) \geq 0$.
 (ii) $\deg(\mathcal{A}, D, z) > 0$ *if and only if* $z \in \mathcal{A}(D)$.
(iii) *If* $\deg(\mathcal{A}, D, z) = 1$, *then* $\mathcal{A}^{-1}(z) \cap D$ *is connected.*
(iv) *If* $\mathcal{A}(D)$ *contains a point of some component $C$ of $\mathbb{R}^n \setminus \mathcal{A}(\partial D)$, then $C \subset \mathcal{A}(D)$.*

The preceding notions are related to the well-known monotone-type mappings.

**Theorem 3.13** ([10], Prop. 7.2.1 and [9], Prop. 19.4.1). *Let $U \subset \mathbb{R}^n$ be open, and suppose that $\mathcal{A} : U \to \mathbb{R}^n$ is monotone, i.e.*

$$\langle \mathcal{A}(u) - \mathcal{A}(v), u - v \rangle_{\mathbb{R}^n} \geq 0$$

*for all $u, v \in U$. Then $\mathcal{A}$ is orientation-preserving.*

### 3.3. Holomorphic mappings

Next, we consider the complex case.

**Definition 3.14.** Let $U \subset \mathbb{C}^n$ be open. A complex mapping $\mathcal{A} : U \to \mathbb{C}^n$ is said to be *holomorphic at $a \in \mathbb{C}^n$* if there is a $\mathbb{C}$-linear mapping $L_a : \mathbb{C}^n \to \mathbb{C}^n$ and a mapping $r_a : U \setminus \{a\} \to \mathbb{C}^n$ with $r_a = o(1)$, such that

$$\mathcal{A}(a + h) = \mathcal{A}(a) + L_a h + r_a(h)$$

for all $h \in \mathbb{C}$ such that $a + h \in U$. In this case, $L_a = \mathcal{A}'(a)$, where $\mathcal{A}'(a)h = h_1 \partial_{z_1} \mathcal{A}(a) + \ldots + h_n \partial_{z_n} \mathcal{A}(a)$ and $\partial_{z_k} \mathcal{A}(a)$ denotes the $k^{\text{th}}$ complex partial derivative of $\mathcal{A}$ at $a$.

We now remind the reader of some elementary linear algebra [34]. Every complex vector space $V$ is also vector space over $\mathbb{R}$. This space will be denoted as $V_{\mathbb{R}}$ and called the *realification* of $V$. The realification of a linear operator $A : V \to W$ over $\mathbb{C}$ is the linear map $A_{\mathbb{R}} : V_{\mathbb{R}} \to W_{\mathbb{R}}$. If $\{e_1, \ldots, e_n\}$ is a basis in $V$, then $\{e_1, \ldots, e_n, ie_1, \ldots, ie_n\}$ is a basis in $V_{\mathbb{R}}$. Further, if $\{e'_1, \ldots, e'_m\}$ is a basis in $W$, and $A = B + iC$ with some real matrices $B$ and $C$, then the matrix of $A_{\mathbb{R}}$ with respect to the bases $\{e_1, \ldots, e_n, ie_1, \ldots, ie_n\}$ and $\{e'_1, \ldots, e'_n, ie'_1, \ldots, ie'_n\}$ is

$$\begin{pmatrix} B & -C \\ C & B \end{pmatrix}.$$

The determinant of the realification $A_{\mathbb{R}}$ obeys the following important rule:

$$\det A_{\mathbb{R}} = |\det A|^2. \tag{3.1}$$

This motivates that we define the realification of a mapping $\mathcal{A} : U \to \mathbb{C}^n$ with $\mathcal{A} = (\mathcal{A}_1, \ldots, \mathcal{A}_n)$ as $\mathcal{A}_{\mathbb{R}} : U \to \mathbb{R}^{2n}$ *via*

$$\mathcal{A}_{\mathbb{R}}(x, y) = (\operatorname{Re} \mathcal{A}_1(z), \operatorname{Im} \mathcal{A}_1(z), \ldots, \operatorname{Re} \mathcal{A}_n(z), \operatorname{Im} \mathcal{A}_n(z)),$$

where $x = (x_1, \ldots, x_n)$, $y = (y_1, \ldots, y_n)$ and $z = x + iy$. Notice that we do not explicitly denote the realification of the set $U$.

**Definition 3.15.** Let $U \subset \mathbb{C}^n$ be open, $\mathcal{A} : U \to \mathbb{C}^n$ be a holomorphic mapping and $D \subset U$ a domain. The *degree of $\mathcal{A}$ in $D$* is defined as the degree of the realification, $\deg(\mathcal{A}, D, z) := \deg(\mathcal{A}_{\mathbb{R}}, D, z)$.

It can be proved using (3.1), that if $\mathcal{A} : U \to \mathbb{C}^n$ is holomorphic, then $\det \mathcal{A}'_{\mathbb{R}}(x, y) = |\det \mathcal{A}'(z)|^2$ for all $z \in U$. This innocent-looking identity has striking consequences.

**Theorem 3.16** ([9], Sect. 19.5)**.** *Let $U \subset \mathbb{C}^n$ be open, $\mathcal{A} : U \to \mathbb{C}^n$ be a holomorphic mapping and $D \subset U$ bounded and open. Then the following statements hold true.*

(i) *$\mathcal{A}_{\mathbb{R}} : U \to \mathbb{R}^{2n}$ is orientation-preserving. In particular, $\deg(\mathcal{A}, D, z) \geq 0$ and $\deg(\mathcal{A}, D, z) > 0$ if and only if $z \in \mathcal{A}(D)$.*

(ii) *Let $\overline{D} \subset U$ and $z \notin \mathcal{A}(\partial D)$. If $\deg(\mathcal{A}, D, z) = k$, then $\mathcal{A}(u) = z$ has at most $k$ solutions in $D$. If $\deg(\mathcal{A}, D, z) = 1$, then $\mathcal{A}(u) = z$ has a unique solution in $D$.*

(iii) *Let $\zeta \in D$ be a solution of $\mathcal{A}(\zeta) = z$ and denote by $C$ the connected component of $\mathcal{A}^{-1}(z)$ containing $\zeta$. Then either $\zeta$ is an isolated solution of $\mathcal{A}(\zeta) = z$ (and hence $C = \{\zeta\}$) or $C \cap G \neq \emptyset$ for any neighborhood $G$ of $\partial D$.*

(iv) *Let $\overline{D} \subset U$ and $z \notin \mathcal{A}(\partial D)$. Then $\deg(\mathcal{A}, D, z) = 1$ if and only if there is a unique non-degenerate solution $\zeta \in D$ of $\mathcal{A}(\zeta) = z$.*

(v) *([42], Thm. 42(b)) The number of zeros in $D$ is finite.*

From (iv) we can conclude that if $\zeta$ is a degenerate isolated solution of $\mathcal{A}(\zeta) = z$, then necessarily $i(\mathcal{A}, \zeta) \geq 2$. Further in the holomorphic case, Theorem 3.8 can sharpened as follows.

**Theorem 3.17.** *Let $\mathcal{A} : D \to \mathbb{C}^n$ be a holomorphic mapping and $D \subset \mathbb{C}^n$ bounded and open. If $z \notin \mathcal{A}(\partial D)$, then there is a $\delta > 0$ such that if $z' \in B(z, \delta) \setminus E$, then*

$$\deg(\mathcal{A}, D, z) = \deg(\mathcal{A}, D, z'),$$

*where $E \subset \{w \in B(z, \delta) : \det \mathcal{A}'(u) = 0, \ \mathcal{A}(u) = w\}$ has measure 0. Furthermore, $\mathcal{A}^{-1}(z')$ consists of a finite number $m$ of points, where $m = \deg(\mathcal{A}, D, z)$.*

*Proof.* From Theorem 3.8, we have $\deg(\mathcal{A}, D, z) \leq m$ and the converse inequality $m \leq \deg(\mathcal{A}, D, z') = \deg(\mathcal{A}, D, z)$ follows from Theorem 3.16(ii). $\qquad\square$

The following result is concerned with holomorphic extensions.

**Theorem 3.18** ([42], Thm. 45)**.** *Let $\mathcal{A} : \mathbb{R}^n \to \mathbb{R}^n$ be a real analytic mapping that can be extended to a holomorphic mapping $\widetilde{\mathcal{A}} : \mathbb{C}^n \to \mathbb{C}^n$. Let $D \subset \mathbb{R}^n$ be an open bounded set such that $0 \notin \mathcal{A}(\partial D)$. Let*

$$D_\varepsilon = \{x \in \mathbb{C}^n : \mathrm{Re}\, x \in D, \quad |\,\mathrm{Im}\, x| < \varepsilon\}.$$

*Then, for sufficiently small $\varepsilon > 0$,*

$$|\deg(\mathcal{A}, D, 0)| \leq \deg\left(\widetilde{\mathcal{A}}, D_\varepsilon, 0\right), \quad \deg(\mathcal{A}, D, 0) = \deg\left(\widetilde{\mathcal{A}}, D_\varepsilon, 0\right) \mod 2.$$

## 4. Analysis of the SRCC method

As stressed earlier, we consider the finite-dimensional case only, by which we mean that the cardinality of the orbital basis $\{\varphi_p\}_{p=1}^K$, $K$ is finite. Recall the definition of $\mathfrak{H}_K^1$ from Section 2.2 and the interpretation of the projected Hamiltonian from Section 2.2.1.

## 4.1. Definitions and basic properties

Let $\mathbb{V} := \mathbb{V}(G)$ be the real amplitude space corresponding to some consistent excitation graph $G$ (see Sect. 3.5 of Part I). In particular, one can take $\mathbb{V}$ to be the S, D, SD, SDT, etc. truncated amplitude space, or the full amplitude space, see Section 2.5. Let $\mathfrak{V}$ be the corresponding functional amplitude space, *i.e.* the set of wavefunctions of the form $T\Phi_0$ for any $t \in \mathbb{V}$, where $T$ is the cluster operator corresponding to $t$ – this convention will be used throughout.

Define the mapping $\mathcal{A} : \mathbb{V} \to \mathbb{V}^*$ *via* the instruction

$$\langle \mathcal{A}(t), s \rangle := \langle e^{-T}\mathcal{H}_K e^T \Phi_0, S\Phi_0 \rangle, \tag{4.1}$$

for any $t, s \in \mathbb{V}$. Then, $\mathcal{A}$ is well-defined, because (4.1) can be rewritten as $\langle \mathcal{A}(t), s \rangle = \left\langle \mathcal{H}_K e^T \Phi_0, e^{-T^\dagger} S\Phi_0 \right\rangle = \left\langle \mathcal{H} e^T \Phi_0, e^{-T^\dagger} S\Phi_0 \right\rangle$, and here $e^T\Phi_0 \in \mathfrak{H}^1_K$ and $e^{-T^\dagger} S\Phi_0 \in \mathfrak{H}^1_K$.

Recall the definition (see (2.6) or Thm. 4.4 of Part I) of the CC energy,

$$\mathcal{E}_{\mathrm{CC}}(t) := \langle e^{-T}\mathcal{H}_K e^T \Phi_0, \Phi_0 \rangle = \langle \mathcal{H}_K e^T \Phi_0, \Phi_0 \rangle,$$

where the second equality follows from $(e^{-T})^\dagger \Phi_0 = \Phi_0$. Note that $\mathcal{E}_{\mathrm{CC}}(t) \in \mathbb{R}$ since the amplitude space $\mathbb{V}$ is assumed to be real. The similarity-transformed Hamiltonian occurs often in the forthcoming discussion, so we introduce the notation

$$\mathcal{H}_K(t) := e^{-T}\mathcal{H}_K e^T : \mathfrak{H}^1_K \to (\mathfrak{H}^1_K)^*, \tag{4.2}$$

which is a bounded map for any fixed cluster amplitude $t \in \mathbb{V}$. Furthermore, for a given bounded map $\mathcal{T} : \mathfrak{H}^1_K \to (\mathfrak{H}^1_K)^*$, we define the operator $\mathcal{T}_{\mathfrak{V}} : \mathfrak{V} \to \mathfrak{V}$ *via* $\langle \mathcal{T}_{\mathfrak{V}} \Psi, \Psi' \rangle = \langle \mathcal{T}\Psi, \Psi' \rangle$ for all $\Psi, \Psi' \in \mathfrak{V}$. For instance, we will often encounter the projected similarity-transformed Hamiltonian $\mathcal{H}_K(t)_{\mathfrak{V}}$.

We will also use a notation analogous to (4.2) for the similarity-transformed fluctuation operator $\mathcal{W}_K$ (see (2.4)), *i.e.* $\mathcal{W}_K(t) = e^{-T}\mathcal{W}_K e^T$. In our finite-dimensional setting, the similarity-transformed Fock operator can be given explicitly as

$$e^{-T}\mathcal{F}_K e^T = \mathcal{F}_K + [\mathcal{F}_K, T], \quad \text{and} \quad [\mathcal{F}_K, X_\alpha] = \varepsilon_\alpha X_\alpha, \tag{4.3}$$

for any $t \in \mathbb{V}$, see *e.g.* Lemma 15 from [12]. In particular,

$$[[\mathcal{H}_K(t), U], V] = [[\mathcal{W}_K(t), U], V], \tag{4.4}$$

for any $t, u, v \in \mathbb{V}$.

The following simple observation shows the equivalence of the (strong) Schrödinger equation with the Full CC method.

**Lemma 4.1.** *Assume that the Slater determinant basis satisfies $\Phi_\alpha \in \mathfrak{H}^2_K$. Suppose that $\mathbb{V} = \mathbb{V}(G^{\mathrm{full}})$, and that $\mathcal{A}(t_*) = 0$. Then the function $\Psi = (c_0 I + C)\Phi_0 \in \mathfrak{H}^2_K$ satisfies the Schrödinger equation $\mathcal{H}_K \Psi = \mathcal{E}\Psi$ if and only if $e^{-T_*}(c_0 I + C) = r_0 I + R$, where ($c_0 = r_0$ and)*

$$\left. \begin{array}{r} \mathcal{E}_{\mathrm{CC}}(t_*) r_0 + \langle \mathcal{H}_K(t_*) R\Phi_0, \Phi_0 \rangle = \mathcal{E} r_0 \\ \mathcal{H}_K(t_*)_{\mathfrak{V}} R\Phi_0 = \mathcal{E} R\Phi_0 \end{array} \right\}. \tag{4.5}$$

*Furthermore,*

$$\sigma(\mathcal{H}_K) = \{\mathcal{E}_{\mathrm{CC}}(t_*)\} \cup \sigma(\mathcal{H}_K(t_*)_{\mathfrak{V}}). \tag{4.6}$$

*Proof.* Using the splitting $\mathrm{Span}\{\Phi_0\} \oplus \mathfrak{V}$, the similarity-transformed Hamiltonian is block upper triangular in the Slater basis,

$$\mathcal{H}_K(t_*) = \begin{pmatrix} \mathcal{E}_{\mathrm{CC}}(t_*) & \langle \mathcal{H}_K(t_*) \cdot, \Phi_0 \rangle \\ 0 & \mathcal{H}_K(t_*)_{\mathfrak{V}} \end{pmatrix},$$

due to $\langle \mathcal{H}_K(t_*)\Phi_0, \Phi_\alpha \rangle = 0$. The proof now follows by noting that the eigenvalues of $\mathcal{H}_K(t_*)$ and $\mathcal{H}_K$ are the same, and the eigenvectors of $\mathcal{H}_K(t_*)$ are of the form $e^{-T_*}\Phi$, where $\mathcal{H}_K\Phi = \mathcal{E}'\Phi$. Formula (4.6) follows by the fact that the spectrum of a block triangular matrix is the union of the spectra of the blocks in the diagonal. $\qquad\square$

We distinguish two cases. Obviously, $r_0 \neq 0$ and $R = 0$ is a solution to the system (4.5) if and only if $\mathcal{E}_{\mathrm{CC}}(t_*) = \mathcal{E}$. In this case, $\Psi = e^{T_*}\Phi_0$ is a solution.

If, however, $\mathcal{E}_{\mathrm{CC}}(t_*) \neq \mathcal{E}$, then $R$ cannot be 0 (because that would imply $\Psi = 0$). In this case, $\Psi = (r_0 I + R)e^{T_*}\Phi_0$, where $r_0 = \langle \mathcal{H}_K(t_*)R\Phi_0, \Phi_0 \rangle / (\mathcal{E} - \mathcal{E}_{\mathrm{CC}}(t_*))$. Note that it is possible to have $r_0 = 0$, in which case $\langle \Psi, \Phi_0 \rangle = 0$. We return to this latter case in Remark 4.16 and discuss the case $\mathcal{E}_{\mathrm{CC}}(t_*) = \mathcal{E}$ below.

**Remark 4.2.** If $\mathcal{E}_{\mathrm{CC}}(t_*) = \mathcal{E}$, then (4.5) reduces to

$$\left.\begin{array}{r}\langle \mathcal{H}_K(t_*)R\Phi_0, \Phi_0 \rangle = 0 \\ \mathcal{H}_K(t_*)_{\mathfrak{V}}R\Phi_0 = \mathcal{E}R\Phi_0\end{array}\right\}.$$

Suppose that this system has $\mu$ linearly independent solutions $R_1, \ldots, R_\mu$. Then it is easy to see, using $\Psi = (r_0 I + R_\mu)e^{T_*}\Phi_0$, that the wavefunctions

$$\left\{ e^{T_*}\Phi_0, R_1 e^{T_*}\Phi_0, \ldots, R_\mu e^{T_*}\Phi_0 \right\}$$

span the eigenspace $\ker(\mathcal{H}_K - \mathcal{E})$. In particular, we have $\dim \ker(\mathcal{H}_K - \mathcal{E}) = \mu + 1$. Note also that in this case $\sigma(\mathcal{H}_K) = \sigma(\mathcal{H}_K(t_*)_{\mathfrak{V}})$.

Let us now recall that the CC equation $\mathcal{A}(t_*) = 0$ can be cast in a form that closely resembles the CI eigenvalue problem (see (2.3) of Part I) (although it is *not* equivalent to it in general).

**Lemma 4.3** ([37], Thm. 5.6). *Let $G$ be excitation complete, and let $\mathbb{V} = \mathbb{V}(G)$ be the corresponding amplitude space. Then the "linked" CC equation $\mathcal{A}(t_*) = 0$ is equivalent to the "unlinked" (a.k.a. "energy-dependent") CC equation*

$$\langle \mathcal{H}_K e^{T_*}\Phi_0, S\Phi_0 \rangle = \mathcal{E}_{\mathrm{CC}}(t_*)\langle e^{T_*}\Phi_0, S\Phi_0 \rangle \quad \text{for all} \quad s \in \mathbb{V}. \tag{4.7}$$

*Proof.* We have

$$\langle (\mathcal{H}_K - \mathcal{E}_{\mathrm{CC}}(t_*))e^{T_*}\Phi_0, S\Phi_0 \rangle = \langle e^{-T_*}(\mathcal{H}_K - \mathcal{E}_{\mathrm{CC}}(t_*))e^{T_*}\Phi_0, (e^{T_*})^\dagger S\Phi_0 \rangle$$

$$= \langle e^{-T_*}(\mathcal{H}_K - \mathcal{E}_{\mathrm{CC}}(t_*))e^{T_*}\Phi_0, \Pi_{\mathfrak{V}}(e^{T_*})^\dagger S\Phi_0 \rangle + \left\langle e^{-T_*}(\mathcal{H}_K - \mathcal{E}_{\mathrm{CC}}(t_*))e^{T_*}\Phi_0, \underbrace{\Pi_{\Phi_0}(e^{T_*})^\dagger S\Phi_0}_{\text{const}\cdot\Phi_0} \right\rangle$$

$$= \langle e^{-T_*}\mathcal{H}_K e^{T_*}\Phi_0, \Pi_{\mathfrak{V}}(e^{T_*})^\dagger S\Phi_0 \rangle,$$

where second term on the right-hand side of the penultimate equality vanishes by the definition of $\mathcal{E}_{\mathrm{CC}}(t_*)$. The proof is completed by recalling that $\Pi_{\mathfrak{V}}(e^{T_*})^\dagger : \mathfrak{V} \to \mathfrak{V}$ is surjective due to Proposition 3.30 of Part I. $\qquad\square$

The "unlinked" form is less useful in practice, because the expansion of $\mathcal{H}_K e^T$ does not terminate like the Baker–Campbell–Hausdorff series

$$\mathcal{H}_K(t) = \sum_{j=0}^{4} \frac{1}{j!}[\mathcal{H}_K, T]_{(j)}. \tag{4.8}$$

More generally, the *doubly*[4] similarity-transformed Hamiltonian $\mathcal{H}_K(t+s) = e^{-S}\mathcal{H}_K(t)e^S$ can also be expanded using the Baker–Campbell–Hausdorff series but in this case

$$\mathcal{H}_K(t+s) = \sum_{j=0}^{2N} \frac{1}{j!}[\mathcal{H}_K(t), S]_{(j)}, \tag{4.9}$$

---

[4]The doubly similarity-transformation above differs from the one considered in Arponen's Extended CC (ECC) theory [1].

*i.e.* the series terminates at $2N$. To see this, simply note that $[\mathcal{H}_K(t), S]_{(2N+1)}$ consists of terms of the form $S^i \mathcal{H}_K(t) S^k$, where $i + k = 2N + 1$, so $i, k \geq N + 1$, which, using the nilpotency of the cluster operators implies that all terms for $j \geq 2N + 1$ vanish.

## 4.2. Local properties – real case

Next, we look at the local behavior of the CC mapping $\mathcal{A} : \mathbb{V} \to \mathbb{V}^*$ for general (real) amplitude spaces $\mathbb{V}$. For fixed $t \in \mathbb{V}$, define the *modified similarity-transformed Hamiltonian*,

$$\widehat{\mathcal{H}_K}(t) := \mathcal{H}_K(t) - \sum_{\alpha \in \Xi(G)^c} \langle \mathcal{H}_K(t) \Phi_0, \Phi_\alpha \rangle X_\alpha, \tag{4.10}$$

where $\Xi(G)^c = \Xi(G^{\text{full}}) \smallsetminus \Xi(G)$ and the *set of excitations* $\Xi(G)$ was defined in (3.3) of Part I. For example, if SD truncation is in effect, then $\Xi(G)$ consists of all singly-, and doubly excited multiindices (see Sect. 2.4). Also, if $\mathbb{V}$ is the full amplitude space, then $\Xi(G)$ consists of all multiindices $\alpha \neq 0$.

**Definition 4.4.** The amplitude space $\mathbb{V}(G)$ is said to be *rank-regular*, if

$$\langle X_\alpha \Phi_\beta, \Phi_\gamma \rangle = 0 \quad \text{for all } \beta, \gamma \in \Xi(G) \text{ and } \alpha \in \Xi(G)^c.$$

We immediately get that $\widehat{\mathcal{H}_K}(t)_{\mathfrak{V}} = \mathcal{H}_K(t)_{\mathfrak{V}}$ if $\mathbb{V}(G)$ is rank-regular. The next proposition shows that the truncated subgraphs typically used in practice – such as ones coming from S, D, SD, SDT, etc. truncations – are rank-regular. We refer the reader to Section 3.2 of Part I for details on truncated subgraphs of the excitation graph.

**Proposition 4.5.** *Suppose that the excitation graph $G$ is a rank-truncated subgraph of the form $G = G(1, 2, \ldots, \rho)$, for some $\rho = 1, \ldots, N$ or $G = G(D)$. Then $\mathbb{V}(G)$ is rank-regular.*

*Proof.* The set $\Xi(G)^c$ consists of elements of rank $\rho + 1, \ldots, N$ (or empty), so that $\langle X_\alpha \Phi_\beta, \Phi_\gamma \rangle = 0$ for all rk $\alpha \notin \{1, 2, \ldots, \rho\}$ and all rk $\beta$, rk $\gamma \in \{1, 2, \ldots, \rho\}$, due to the fact that $X_\alpha \Phi_\beta$ is of rank $\text{rk}(\alpha) + \text{rk}(\beta) \notin \{1, 2, \ldots, \rho\}$. The proof of the case $G = G(D)$ is similar. $\qquad \square$

Next, we compute the derivative of the mapping $\mathcal{A}$, which explains the definition of $\widehat{\mathcal{H}_K}(t)$ and of rank-regularity.

**Lemma 4.6.** *Let $t_*$ be a zero of $\mathcal{A} : \mathbb{V} \to \mathbb{V}^*$. Then the derivative $\mathcal{A}'(t_*) \in \mathcal{L}(\mathbb{V}, \mathbb{V}^*)$ is given by*

$$\langle \mathcal{A}'(t_*) u, v \rangle = \left\langle \left( \widehat{\mathcal{H}_K}(t_*) - \mathcal{E}_{\text{CC}}(t_*) \right) U \Phi_0, V \Phi_0 \right\rangle \tag{4.11}$$

*for all $u, v \in \mathbb{V}$.*

*Proof.* The derivative $\mathcal{A}' : \mathbb{V} \to \mathcal{L}(\mathbb{V}, \mathbb{V}^*)$ is readily computed as

$$\begin{aligned}
\frac{\mathrm{d}}{\mathrm{d}h} \langle \mathcal{A}(t + hu), v \rangle \Big|_{h=0} &= \frac{\mathrm{d}}{\mathrm{d}h} \langle e^{-T-hU} \mathcal{H}_K e^{T+hU} \Phi_0, V \Phi_0 \rangle \Big|_{h=0} \\
&= \langle e^{-T-hU} (\mathcal{H}_K U - U \mathcal{H}_K) e^{T+hU} \Phi_0, V \Phi_0 \rangle \big|_{h=0} \\
&= \langle e^{-T} (\mathcal{H}_K U - U \mathcal{H}_K) e^T \Phi_0, V \Phi_0 \rangle,
\end{aligned}$$

so using the commutativity of the cluster operators, we get

$$\langle \mathcal{A}'(t) u, v \rangle = \langle [\mathcal{H}_K(t), U] \Phi_0, V \Phi_0 \rangle \tag{4.12}$$

for all $t, u, v \in \mathbb{V}$. Expanding $U^{\dagger} V \Phi_0 \in \mathfrak{H}_K^1$ in the $\mathfrak{L}^2$-orthonormal basis $\{\Phi_\alpha\}_\alpha \subset \mathfrak{H}_K^1$,

$$U^{\dagger} V \Phi_0 = \sum_\alpha \langle U \Phi_\alpha, V \Phi_0 \rangle \Phi_\alpha,$$

we obtain using $\langle \mathcal{H}_K(t_*) \Phi_0, \Phi_\alpha \rangle = 0$ for all $\alpha \in \Xi(G)$,

$$\langle U \mathcal{H}_K(t_*) \Phi_0, V \Phi_0 \rangle = \langle \mathcal{H}_K(t_*) \Phi_0, U^{\dagger} V \Phi_0 \rangle$$
$$= \mathcal{E}_{\mathrm{CC}}(t_*) \langle U \Phi_0, V \Phi_0 \rangle + \sum_{\alpha \in \Xi(G)^c} \langle \mathcal{H}_K(t_*) \Phi_0, \Phi_\alpha \rangle \langle X_\alpha U \Phi_0, V \Phi_0 \rangle$$

for all $u, v \in \mathbb{V}$. Inserting this into (4.12) with $t = t_*$, we obtain the stated formula. □

As we noted in Section 1.1, previous analyses of the CC mapping assumed the local strong monotonicity at a zero $t_*$, *i.e.* that there is a $\delta > 0$ and a constant $C_{\mathrm{SM}}(t_*, \delta) > 0$ such that

$$\langle \mathcal{A}(t) - \mathcal{A}(s), t - s \rangle \geq C_{\mathrm{SM}}(t_*, \delta) \| t - s \|_{\mathbb{V}}^2, \quad \text{for all} \quad t, s \in B_{\mathbb{V}}(t_*, \delta). \tag{4.13}$$

The following elementary theorem makes the observations in [36] more precise.

**Theorem 4.7.** *Let $t_* \in \mathbb{V}$ be a zero of $\mathcal{A} : \mathbb{V} \to \mathbb{V}^*$.*

(i) *If $\mathcal{A}$ is strongly monotone in $B_{\mathbb{V}}(t_*, \delta)$ for some $\delta > 0$, then there exists $\delta' > 0$ such that $\mathcal{A}'(t_* + u)$ is $\mathbb{V}$-coercive for all $\|u\|_{\mathbb{V}} < \delta'$ with some constant $0 < \gamma \leq C_{\mathrm{SM}}(t_*, \delta)$, i.e.*

$$\langle \mathcal{A}'(t_* + u) v, v \rangle \geq \gamma \| v \|_{\mathbb{V}}^2 \quad \text{for all} \quad v \in \mathbb{V} \text{ and } \|u\|_{\mathbb{V}} < \delta'. \tag{4.14}$$

(ii) *Conversely, if (4.14) holds with $u = 0$, then (4.13) holds true with $\delta > 0$ chosen so that $C_{\mathrm{SM}}(t_*, \delta) := \gamma - M_\delta \delta > 0$, where*

$$M_\delta = \sup_{\|\zeta\|_{\mathbb{V}} \leq \delta} \| \mathcal{A}''(t_* + \zeta) \|_{\mathcal{L}(\mathbb{V} \times \mathbb{V}, \mathbb{V}^*)}. \tag{4.15}$$

*Proof.* To see (i), fix $\delta' > 0$ and $\|u\|_{\mathbb{V}} < \delta'$ and write for any $\|r\|_{\mathbb{V}} < \delta' < \delta$,

$$C_{\mathrm{SM}}(t_*, \delta) \| r - u \|_{\mathbb{V}}^2 \leq \langle \mathcal{A}(t_* + r) - \mathcal{A}(t_* + u), r - u \rangle \leq \langle \mathcal{A}'(t_* + u)(r - u), r - u \rangle + \frac{1}{2} M_\delta \| r - u \|_{\mathbb{V}}^3.$$

This implies

$$(C_{\mathrm{SM}}(t_*, \delta) - M_\delta \delta') \| r - u \|_{\mathbb{V}}^2 \leq \langle \mathcal{A}'(t_* + u)(r - u), r - u \rangle.$$

Any vector $v \in \mathbb{V}$ can be expressed as $v = \alpha(r - u)$ for some $\alpha > 0$ and $\|r\|_{\mathbb{V}} < \delta'$, from which $\mathbb{V}$-coercivity follows with $\gamma = C_{\mathrm{SM}}(t_*, \delta) - M_\delta \delta'$, by choosing $\delta'$ sufficiently small.

Next, to prove (ii), write the Taylor expansions of $\mathcal{A}$ at $t_*$,

$$\mathcal{A}(t_* + r) = \mathcal{A}'(t_*) r + \mathcal{R}_2(t_*; r),$$
$$\mathcal{A}(t_* + r') = \mathcal{A}'(t_*) r' + \mathcal{R}_2(t_*; r'),$$

for any $\|r\|_{\mathbb{V}}, \|r'\|_{\mathbb{V}} < \delta$ for some $\delta > 0$, from which we obtain

$$\langle \mathcal{A}(t_* + r) - \mathcal{A}(t_* + r'), r - r' \rangle = \langle \mathcal{A}'(t_*)(r - r'), r - r' \rangle + \langle \mathcal{R}_2(t_*; r) - \mathcal{R}_2(t_*; r'), r - r' \rangle$$
$$\geq \gamma \| r - r' \|_{\mathbb{V}}^2 + \langle \mathcal{R}_2(t_*; r) - \mathcal{R}_2(t_*; r'), r - r' \rangle.$$

Using the intermediate value inequality, we have

$$\| \mathcal{R}_2(t_*; r) - \mathcal{R}_2(t_*; r') \|_{\mathbb{V}^*} \leq M_{r,r'} \| r - r' \|_{\mathbb{V}},$$

where

$$M_{r,r'} = \max_{\xi \in [r,r']} \|\partial_2 \mathcal{R}_2(t_*; \xi)\| = \max_{\xi \in [r,r']} \|\mathcal{A}'(t_* + \xi) - \mathcal{A}'(t_*)\|$$

$$\leq \left( \max_{\xi \in [r,r']} \max_{\zeta \in [0,\xi]} \|\mathcal{A}''(t_* + \zeta)\|_{\mathcal{L}(\mathbb{V} \times \mathbb{V}, \mathbb{V}^*)} \right) \delta$$

$$= \left( \sup_{\|\zeta\| \leq \delta} \|\mathcal{A}''(t_* + \zeta)\|_{\mathcal{L}(\mathbb{V} \times \mathbb{V}, \mathbb{V}^*)} \right) \delta = M_\delta \delta$$

This implies $\langle \mathcal{A}(t_* + r) - \mathcal{A}(t_* + r'), r - r' \rangle \geq (\gamma - M_\delta \delta) \|r - r'\|_{\mathbb{V}}^2$ for all $\|r\|_{\mathbb{V}}, \|r'\|_{\mathbb{V}} < \delta$. Setting $t = t_* + r$ and $s = t_* + r'$ proves the claim. □

**Remark 4.8.** Let $\mathbb{V}^0 \subset \mathbb{V}$ be a subspace and consider the *projected CC mapping* $\mathcal{A}^0 : \mathbb{V}^0 \to (\mathbb{V}^0)^*$ *via*

$$\langle \mathcal{A}^0(t^0), s^0 \rangle_{\mathbb{V}^* \times \mathbb{V}} = \langle \mathcal{A}(t^0), s^0 \rangle_{\mathbb{V}^* \times \mathbb{V}} \quad \text{for all} \quad t^0, s^0 \in \mathbb{V}^0.$$

Clearly, if $\mathcal{A}$ is strongly monotone on $B_{\mathbb{V}}(t_*, \delta)$ with a constant $C_{\mathrm{SM}} > 0$ at a zero $t_*$, then $\mathcal{A}^0$ is strongly monotone on $B_{\mathbb{V}^0}(t_*^0, \sqrt{\delta^2 - \|t_*^\perp\|_{\mathbb{V}}^2})$ provided $\|t_*^\perp\|_{\mathbb{V}}^2$ is small enough, with the same constant $C_{\mathrm{SM}}$, where we have set $t_*^0 = \Pi_{\mathbb{V}^0} t_*$ and $t_*^\perp = (I - \Pi_{\mathbb{V}^0}) t_*$. Here, $\Pi_{\mathbb{V}^0} : \mathbb{V} \to \mathbb{V}$ denotes the $\ell^2$-orthogonal projector onto $\mathbb{V}^0$. Note that $t_*^0$ is *not*, in general, a zero of $\mathcal{A}^0$, hence the preceding theorem is not applicable to $\mathcal{A}^0$.

**Remark 4.9.** The quantity $M_\delta$ contains the second derivative of $\mathcal{A}$. It is easy to see that

$$\langle \mathcal{A}''(t)(u, v), w \rangle = \langle [[\mathcal{H}_K(t), U], V]\Phi_0, W\Phi_0 \rangle$$

for all $u, v, w \in \mathbb{V}$. Using (4.4), we have $\langle \mathcal{A}''(t)(u, v), w \rangle = \langle [[\mathcal{W}_K(t), U], V]\Phi_0, W\Phi_0 \rangle$, so that $\mathcal{A}''(t)$ only involves the fluctuation operator $\mathcal{W}_K$.

**Remark 4.10** (Perturbative regime)**.** Let $\mathbb{V}(G)$ be rank-regular, and consider the case when $t_* \approx 0$, which is the case considered in [36, 37]. Then roughly speaking, we have $\mathcal{H}_K(t_*) \approx \mathcal{H}_K$. Note that

$$\langle \mathcal{A}'(t_*) r, r \rangle = \langle (\mathcal{H}_K - \mathcal{E}_{\mathrm{CC}}(t_*)) R\Phi_0, R\Phi_0 \rangle + \mathcal{O}(\|t_*\|_{\mathbb{V}}).$$

Consequently, if

$$\langle (\mathcal{H}_K - \mathcal{E}_{\mathrm{CC}}(t_*)) R\Phi_0, R\Phi_0 \rangle \geq c(t_*) \|r\|_{\mathbb{V}}^2,$$

where $c(t_*) > 0$, then local strong monotonicity holds with constant $C_{\mathrm{SM}} = c(t_*) - M' \|t_*\|_{\mathbb{V}} - 2M_\delta \delta$ for $t_*$ sufficiently close to 0. In Lemma 3.5 of [36], it is shown that such a $c(t_*)$ exists under the assumption that $\mathcal{H}_K$ has a spectral gap and that $\Phi_0$ is a sufficiently good approximation of the ground state $e^{T_*}\Phi_0$ (*i.e.* that $t_*$ is sufficiently close to 0).

**Proposition 4.11.** *If $\mathcal{A}$ is locally strongly monotone at a zero $t_*$, then $t_*$ is non-degenerate.*

*Proof.* Suppose that $\ker \mathcal{A}'(t_*) \neq \{0\}$ and that $\mathcal{A}$ is locally strongly monotone near $t_*$. Then for any $0 \neq r \in \ker \mathcal{A}'(t_*)$ sufficiently close to 0, we have

$$C_{\mathrm{SM}}(\delta) \|r\|_{\mathbb{V}}^2 \leq \langle \mathcal{A}(t_* + r), r \rangle = \frac{1}{2} \langle \mathcal{A}''(t_*)(r, r), r \rangle + o(\|r\|_{\mathbb{V}}^4).$$

Rescaling $r$ by $\alpha > 0$ small, and letting $\alpha \to 0$ we obtain that $C_{\mathrm{SM}} = 0$, a contradiction. □

**Remark 4.12.** When applying topological degree theory, we will view $\mathcal{A} : \mathbb{V} \to \mathbb{V}^*$ as a mapping $\mathbb{R}^n \to \mathbb{R}^n$ by identifying $\mathbb{V}$ and $\mathbb{V}^*$ with $\mathbb{R}^n$. Following Section 1.3 of [10], we fix a basis $\{\tau_\alpha\}_{\alpha \in \Xi(G)}$ of $\mathbb{V}$ and define the linear homeomorphism $h : \mathbb{V} \to \mathbb{R}^n$ with

$$\mathbb{V} \ni t = \sum_{\alpha \in \Xi(G)} \widehat{t}_\alpha \tau_\alpha \mapsto h(t) = \sum_{\alpha \in \Xi(G)} \widehat{t}_\alpha e_\alpha \in \mathbb{R}^n,$$

where $\{e_\alpha\}_{\alpha \in \Xi(G)}$ is the standard (ordered) basis in $\mathbb{R}^n$. Also, fix a basis $\{\tau_\alpha^*\}_{\alpha \in \Xi(G)}$ of $\mathbb{V}^*$ and define the linear homeomorphism $g : \mathbb{V}^* \to \mathbb{R}^n$ analogously. Then

$$\widehat{\mathcal{A}} := g \circ \mathcal{A} \circ h^{-1} : \mathbb{R}^n \to \mathbb{R}^n$$

gives the desired mapping. Now suppose that two other bases $\{\widetilde{\tau}_\alpha\}_{\alpha \in \Xi(G)} \subset \mathbb{V}$ and $\{\widetilde{\tau}_\alpha^*\}_{\alpha \in \Xi(G)} \subset \mathbb{V}^*$ are given and let $\widetilde{h} : \mathbb{V} \to \mathbb{R}^n$ and $\widetilde{g} : \mathbb{V}^* \to \mathbb{R}^n$ be the corresponding linear homeomorphism. But then

$$g^{-1} \circ g \circ \mathcal{A} \circ h^{-1} \circ h = \mathcal{A} = \widetilde{g}^{-1} \circ \widetilde{g} \circ \mathcal{A} \circ \widetilde{h}^{-1} \circ \widetilde{h},$$

which implies $\widetilde{\mathcal{A}} := \widetilde{g} \circ \mathcal{A} \circ \widetilde{h}^{-1} = m \circ \widehat{\mathcal{A}} \circ \widetilde{m}$, where $m = \widetilde{g} \circ g^{-1} : \mathbb{R}^n \to \mathbb{R}^n$ and $\widetilde{m} = h \circ \widetilde{h}^{-1} : \mathbb{R}^n \to \mathbb{R}^n$. Using Lemma 1.3.1 of [10] ([9], Lem. 6.1.1), we obtain

$$\deg\left(\widetilde{\mathcal{A}}, \widetilde{h}(D), \widetilde{g}(0)\right) = (\operatorname{sgn} \det m)(\operatorname{sgn} \det \widetilde{m}) \deg\left(\widehat{\mathcal{A}}, h(D), g(0)\right)$$

for any open and bounded set $D \subset \mathbb{V}$ with $0 \notin \mathcal{A}(\partial D)$. We can conclude that the topological degree is independent of the choice of the basis if $\mathbb{V}$ and $\mathbb{V}^*$ are oriented the same.

Next, we determine the topological index of a zero of $\mathcal{A}$. The fact that the topological index of $t_*$ is related to its CC energy $\mathcal{E}_{\mathrm{CC}}(t_*)$ and the eigenvalues of the operator $\widehat{\mathcal{H}_K}(t_*)_{\mathfrak{V}}$ is interesting on its own right.

**Theorem 4.13** (Index formula for SRCC – non-degenerate case)**.** *Let $t_*$ be a zero of the CC mapping $\mathcal{A} : \mathbb{V} \to \mathbb{V}^*$. Then $t_*$ is non-degenerate if and only if $\mathcal{E}_{\mathrm{CC}}(t_*) \notin \sigma(\widehat{\mathcal{H}_K}(t_*)_{\mathfrak{V}})$, and in this case $t_*$ is an isolated zero and the topological index of $\mathcal{A}$ at $t_*$ is given by*

$$i(\mathcal{A}, t_*) = (-1)^\nu,$$

*where*

$$\nu = \left|\left\{j : \mathcal{E}_j\left(\widehat{\mathcal{H}_K}(t_*)_{\mathfrak{V}}\right) \in \mathbb{R}, \ \mathcal{E}_j\left(\widehat{\mathcal{H}_K}(t_*)_{\mathfrak{V}}\right) < \mathcal{E}_{\mathrm{CC}}(t_*)\right\}\right|.$$

*Proof.* It is trivial to see that if $t_*$ is non-degenerate, then it is isolated: assume that $\ker \mathcal{A}'(t_*) = \{0\}$ and write

$$\langle \mathcal{A}(t_* + r), \mathcal{A}'(t_*)r \rangle_{\mathbb{V}^*} = \|\mathcal{A}'(t_*)r\|_{\mathbb{V}^*}^2 + o\bigl(\|r\|_{\mathbb{V}}^3\bigr), \tag{4.16}$$

for all $\|r\|_{\mathbb{V}} = \varepsilon$, where $\varepsilon > 0$ is sufficiently small. This implies that $\mathcal{A}(t_* + r) \neq 0$ for all $0 < \|r\|_{\mathbb{V}} < \varepsilon$.

We can apply Theorem 3.7 with the mappings $h : \mathbb{V} \to \mathbb{R}^n$ and $g : \mathbb{V}^* \to \mathbb{R}^n$ defined in Remark 4.12. Using the notations of the said remark and (4.11), we have

$$\left\langle \widehat{\mathcal{A}}'(h^{-1}(t_*))e_\alpha, e_\beta \right\rangle_{\mathbb{R}^n} = \left\langle g\mathcal{A}'(t_*)h^{-1}(e_\alpha), e_\beta \right\rangle_{\mathbb{R}^n} = \left\langle \mathcal{A}'(t_*)h^{-1}(e_\alpha), g^\dagger(e_\beta) \right\rangle$$

$$= \left\langle \left(\widehat{\mathcal{H}_K}(t_*) - \mathcal{E}_{\mathrm{CC}}(t_*)\right)U_\alpha \Phi_0, V_\beta \Phi_0 \right\rangle,$$

where $u_\alpha = h^{-1}(e_\alpha)$ and $v_\beta = g^\dagger(e_\beta)$ and $\alpha, \beta \in \Xi(G)$. Here, $g^\dagger : \mathbb{R}^n \to \mathbb{V}$ is the adjoint of $g$. Therefore, using an appropriate basis transformation

$$i(\mathcal{A}, t_*) = \operatorname{sgn} \det \widehat{\mathcal{A}}'(h^{-1}(t_*)) = \operatorname{sgn}\left(\prod_{j \geq 0} \left(\mathcal{E}_j\left(\widehat{\mathcal{H}_K}(t_*)_{\mathfrak{V}}\right) - \mathcal{E}_{\mathrm{CC}}(t_*)\right)\right).$$

The proof is completed by noting that the elements of the matrix $\widehat{\mathcal{H}_K}(t_*)_{\mathfrak{V}}$ are real, so its complex eigenvalues come in conjugate pairs, hence only real eigenvalues contribute to the product above. $\qquad\square$

**Proposition 4.14.** *If $\mathcal{A}$ is locally strongly monotone near a zero $t_*$, then we have $i(\mathcal{A}, t_*) = 1$.*

*Proof.* We have that in particular $\widehat{\mathcal{A}} : \mathbb{R}^n \to \mathbb{R}^n$ is monotone near $h(t_*)$, so according to Theorem 3.13, $\widehat{\mathcal{A}}$ is orientation-preserving near $h(t_*)$. But then Theorem 3.12 implies that $i(\mathcal{A}, t_*) = i(\widehat{\mathcal{A}}, h(t_*)) > 0$. $\qquad\square$

Using the "unlinked" form, we can determine the topological index in the Full CC case. Recall that the eigenvalues $\mathcal{E}_n(\mathcal{H}_K)$, $n = 0, 1, \ldots$, are assumed to be increasingly ordered.

**Theorem 4.15** (Index formula for FCC – non-degenerate case)**.** *Let $\mathbb{V} = \mathbb{V}(G^{\mathrm{full}})$ and assume that $t_* \in \mathbb{V}$ is a zero of the FCC mapping $\mathcal{A} : \mathbb{V} \to \mathbb{V}^*$. Then $e^{T_*}\Phi_0 \in \mathfrak{H}^2$ is an (intermediately normalized) eigenfunction corresponding to some non-degenerate eigenvalue $\mathcal{E}_\nu(\mathcal{H}_K)$ if and only if $t_*$ is non-degenerate, and in this case $i(\mathcal{A}, t_*) = (-1)^\nu$.*

*Proof.* First, note that $\mathbb{V}(G^{\mathrm{full}})$ is rank-regular so $\widehat{\mathcal{H}_K}(t)_{\mathfrak{V}} = \mathcal{H}_K(t)_{\mathfrak{V}}$. We have $\mathcal{E}_{\mathrm{CC}}(t_*) = \mathcal{E}_\nu(\mathcal{H}_K)$ by the equivalence of FCC and FCI (see Thm. 2.2 of Part I).

According to Lemma 4.1, we have that $e^{T_*}\Phi_0$ is a non-degenerate intermediately normalized eigenfunction if and only if $\mathcal{E}_{\mathrm{CC}}(t_*) \notin \sigma(\mathcal{H}_K(t_*)_{\mathfrak{V}})$. In fact, $\mathcal{E}_{\mathrm{CC}}(t_*) \in \sigma(\mathcal{H}_K(t_*)_{\mathfrak{V}})$ if and only if there exists $R\Phi_0 \in \mathfrak{V}$ nonzero, such that

$$\langle e^{-T_*}\mathcal{H}_K e^{T_*} R\Phi_0, S\Phi_0 \rangle = \mathcal{E}_{\mathrm{CC}}(t_*)\langle R\Phi_0, S\Phi_0 \rangle,$$

for all $s \in \mathbb{V}$. Since $\mathbb{V}(G^{\mathrm{full}})$ is excitation complete, according to Lemma 4.3 the preceding equation is equivalent to

$$\langle \mathcal{H}_K R e^{T_*}\Phi_0, S\Phi_0 \rangle = \mathcal{E}_{\mathrm{CC}}(t_*)\langle R e^{T_*}\Phi_0, S\Phi_0 \rangle, \tag{4.17}$$

for all $s \in \mathbb{V}$. But this precisely means that the FCI eigenstate $e^{T_*}\Phi_0$ is degenerate, because $R e^{T_*}\Phi_0$ is another eigenvector corresponding to the same eigenvalue $\mathcal{E}_{\mathrm{CC}}(t_*)$.

We also conclude from (4.6) that $\sigma(\mathcal{H}_K(t_*)_{\mathfrak{V}}) = \sigma(\mathcal{H}_K) \smallsetminus \{\mathcal{E}_{\mathrm{CC}}(t_*)\}$. Applying Theorem 4.13, we obtain that $t_*$ is non-degenerate and $i(\mathcal{A}, t_*) = (-1)^\nu$. $\qquad\square$

It is worth noting that, in the FCC case, the zero $t_*$ representing the intermediately normalized, non-degenerate *ground state* (*i.e.* $\mathcal{E}_{\mathrm{CC}}(t_*) = \mathcal{E}_0(\mathcal{H}_K)$) has $i(\mathcal{A}, t_*) = 1$. Note that this is not necessarily true in the truncated case. While the CC method is most commonly aimed at the ground state, it can also be used to find other intermediately normalized eigenfunctions as well. Furthermore, it can also be used to obtain eigenfunctions which are orthogonal to the reference $\Phi_0$ according to the remark below.

**Remark 4.16.** The *Equation-of-Motion Coupled-Cluster* (EOM-CC) method [15] is aimed at calculating *excited* energies and states (*i.e.* $\mathcal{E}_n(\mathcal{H}_K)$ for $n > 0$, and the corresponding eigenvectors) based on a CC ground-state solution. This is done in two steps. Let $\mathbb{V} = \mathbb{V}(G^{\mathrm{full}})$. Firstly, a conventional CC calculation determines the ground state $\Psi = e^{T_*}\Phi_0$ such that $\mathcal{A}(t_*) = 0$, *i.e.* $\mathcal{H}_K\Psi = \mathcal{E}\Psi$. Secondly, the targeted excited state is of the form $\Psi_{\mathrm{ex}} = (r_0 I + R)e^{T_*}\Phi_0$, where $R$ is a cluster operator, see Lemma 4.1. We have

$$\mathcal{H}_K(t_*)_{\mathfrak{V}} R\Phi_0 = \mathcal{E}_{\mathrm{ex}} R\Phi_0. \tag{4.18}$$

In other words, we need to solve the eigenproblem of the projected similarity-transformed Hamiltonian $\mathcal{H}_K(t_*)_{\mathfrak{V}}$. Furthermore, similarly to the proof of Lemma 4.3, it is easy to see that $\mathcal{A}(t_*) = 0$ implies

$$\langle R\mathcal{H}_K(t_*)\Phi_0, S\Phi_0 \rangle = \mathcal{E}_{\mathrm{CC}}(t_*)\langle R\Phi_0, S\Phi_0 \rangle \quad \text{for all} \quad s \in \mathbb{V}.$$

Subtracting this from (4.18), we get the "commutator form" of the EOM-CC equation:

$$\langle [\mathcal{H}_K(t_*), R]\Phi_0, S\Phi_0 \rangle = \Delta\mathcal{E}\langle R\Phi_0, S\Phi_0 \rangle, \quad \text{where} \quad \Delta\mathcal{E} = \mathcal{E}_{\mathrm{ex}} - \mathcal{E}, \tag{4.19}$$

and $\mathcal{E} = \mathcal{E}_{\mathrm{CC}}(t_*)$ is the ground-state energy as given by the CC method. Let $\mathbb{V}$ be an arbitrary amplitude space. Recalling the expression (4.12) for $\mathcal{A}'(t)$, we can rephrase the EOM-CC equation (4.19) as *the weak eigenvalue problem for* $\mathcal{A}'(t_*) : \mathbb{V} \to \mathbb{V}^*$ (*cf.* [16], Sect. 13.6.3), *i.e.*

$$\langle \mathcal{A}'(t_*) r_j, s \rangle = \Delta\mathcal{E}_j \langle r_j, s \rangle, \tag{4.20}$$

for $j = 1, \ldots, J$[5] and $s \in \mathbb{V}$ is arbitrary[6]. Notice that the $\Delta\mathcal{E}_j$'s are in general complex. Using Theorem 4.13 we can obtain the following. Suppose that $\Delta\mathcal{E}_1, \ldots, \Delta\mathcal{E}_\mu$ are given by (4.20) and are all nonzero. Then

$$i(\mathcal{A}, t_*) = (-1)^\nu, \quad \nu = |\{j : \Delta\mathcal{E}_j \in \mathbb{R}, \ \Delta\mathcal{E}_j < 0\}|. \tag{4.21}$$

Due to the *nonvariational property* of truncated CC (see Sect. 2.3 of Part I), it is not *a priori* clear whether the (real) excited energies are higher than the ground-state energy, *i.e.* whether $\Delta\mathcal{E}_j > 0$. Therefore, equation (4.21) quantifies this nonvariational property through the topological index $i(\mathcal{A}, t_*)$.

Next, we draw a connection between the degeneracy of a zero $t_*$ and the Fock-splitting (2.4) of the Hamiltonian. Define $\omega_0(t_*) = \langle \mathcal{W}_K(t_*)\Phi_0, \Phi_0 \rangle$, which is the CC correction to the lowest eigenvalue $\Lambda_0$ of $\mathcal{F}$, so that the CC energy at $t_*$ is obtained as $\mathcal{E}_{\mathrm{CC}}(t_*) = \Lambda_0 + \omega_0(t_*)$.

**Proposition 4.17.** *Let* $\mathbb{V}(G)$ *be a rank-regular amplitude space and* $t_*$ *a zero of* $\mathcal{A}$*. Define the linear operator* $\mathcal{Q}(t_*) : \mathfrak{V} \to \mathfrak{V}$ *via its matrix in the Slater determinant basis as*

$$[\mathcal{Q}(t_*)]_{\alpha\beta} = \varepsilon_\alpha \delta_{\alpha\beta} + \sum_{\gamma \in \Xi(G)} t_{*,\gamma} \varepsilon_\gamma \langle X_\gamma \Phi_\beta, \Phi_\alpha \rangle \quad \text{for all} \quad \alpha, \beta \in \Xi(G).$$

*Then* $\omega_0(t_*) \notin \sigma(\mathcal{Q}(t_*) + \mathcal{W}_K(t_*)_\mathfrak{V})$ *is equivalent to* $\mathcal{E}_{\mathrm{CC}}(t_*) \notin \sigma(\mathcal{H}_K(t_*)_\mathfrak{V})$*, i.e. to the fact that* $t_*$ *is a non-degenerate zero of* $\mathcal{A}$*.*

*Proof.* We have using (4.3),

$$\begin{aligned}
\mathcal{E}_{\mathrm{CC}}(t_*) &= \langle e^{-T_*}\mathcal{F}_K e^{T_*}\Phi_0, \Phi_0 \rangle + \langle \mathcal{W}_K(t_*)\Phi_0, \Phi_0 \rangle \\
&= \langle \mathcal{F}_K\Phi_0, \Phi_0 \rangle + \langle [\mathcal{F}_K, T_*]\Phi_0, \Phi_0 \rangle + \langle \mathcal{W}_K(t_*)\Phi_0, \Phi_0 \rangle \\
&= \Lambda_0 + \langle \mathcal{W}_K(t_*)\Phi_0, \Phi_0 \rangle = \Lambda_0 + \omega_0(t_*).
\end{aligned}$$

Similarly,

$$\begin{aligned}
\langle \mathcal{H}_K(t_*)\Phi_\beta, \Phi_\alpha \rangle &= \langle \mathcal{F}_K\Phi_\beta, \Phi_\alpha \rangle + \langle [\mathcal{F}_K, T_*]\Phi_\beta, \Phi_\alpha \rangle + \langle \mathcal{W}_K(t_*)\Phi_\beta, \Phi_\alpha \rangle \\
&= (\Lambda_0 + \varepsilon_\alpha)\delta_{\alpha\beta} + \sum_{\gamma \in \Xi(G)} t_{*,\gamma} \varepsilon_\gamma \langle X_\gamma \Phi_\beta, \Phi_\alpha \rangle + \langle \mathcal{W}_K(t_*)\Phi_\beta, \Phi_\alpha \rangle.
\end{aligned}$$

Then, in the Slater determinant basis

$$\mathcal{H}_K(t_*)_\mathfrak{V} = \Lambda_0 I + \mathcal{Q}(t_*) + \mathcal{W}_K(t_*)_\mathfrak{V}.$$

Hence, $\mathcal{E}_{\mathrm{CC}}(t_*) \notin \sigma(\mathcal{H}_K(t_*)_\mathfrak{V})$ is equivalent to $\omega_0(t_*) \notin \sigma(\mathcal{Q}(t_*) + \mathcal{W}_K(t_*)_\mathfrak{V})$, which finishes the proof. $\square$

We now consider the case of a degenerate zero. Clearly, if $r \in \ker \mathcal{A}'(t_*) \neq \{0\}$, we have to consider higher-order terms of the Taylor polynomial of $\mathcal{A}$ at $t_*$,

$$\mathcal{A}(t_* + r) = \mathcal{A}'(t_*)r + \frac{1}{2}\mathcal{A}''(t_*)(r, r) + \mathcal{R}_3(t_*; r),$$

---

[5] Here, $R_j$ is not to be confused with the rank-decomposition (3.6) of Part I.

[6] A similar relation holds if $t_*$ does not represent the ground state.

where $\mathcal{A}''(t_*) : \mathbb{V} \times \mathbb{V} \to \mathbb{V}^*$ is a bounded bilinear mapping. Here, we only consider the second-order information.

Assume from now on that $\mathbb{V}$ is rank-regular, so that $\widehat{\mathcal{H}_K}(t)_{\mathfrak{V}} = \mathcal{H}_K(t)_{\mathfrak{V}}$. Suppose that $\mathcal{E}_{\mathrm{CC}}(t_*) \in \sigma(\mathcal{H}_K(t_*)_{\mathfrak{V}})$ and that $R_1\Phi_0, \ldots, R_\mu\Phi_0 \in \mathfrak{V}$ are the *right* eigenvectors of $\mathcal{H}_K(t_*)_{\mathfrak{V}}$ corresponding to $\mathcal{E}_{\mathrm{CC}}(t_*)$,

$$\langle \mathcal{H}_K(t_*) R_j \Phi_0, S\Phi_0 \rangle = \mathcal{E}_{\mathrm{CC}}(t_*) \langle R_j \Phi_0, S\Phi_0 \rangle \quad \text{for all } j = 1, \ldots, \mu \text{ and all } s \in \mathbb{V}.$$

Also, suppose that $L_1\Phi_0, \ldots, L_\mu\Phi_0 \in \mathfrak{V}$ are the *left* eigenvectors of $\mathcal{H}_K(t_*)_{\mathfrak{V}}$ corresponding to $\mathcal{E}_{\mathrm{CC}}(t_*)$,

$$\left\langle \mathcal{H}_K(t_*)^\dagger L_j \Phi_0, S\Phi_0 \right\rangle = \mathcal{E}_{\mathrm{CC}}(t_*) \langle L_j \Phi_0, S\Phi_0 \rangle \quad \text{for all } j = 1, \ldots, \mu \text{ and all } s \in \mathbb{V}.$$

The corresponding right-, and left eigenspaces are

$$W_R = \ker \mathcal{A}'(t_*) = \mathrm{Span}\{r_1, \ldots, r_\mu\},$$
$$W_L = \ker \mathcal{A}'(t_*)^\dagger = \mathrm{Span}\{\ell_1, \ldots, \ell_\mu\},$$

and let $Q : \mathbb{V} \to \mathbb{V}$ be the orthogonal projector onto $W_L$. Further, define $\widehat{Q} : \mathfrak{V} \to \mathfrak{V}$ *via* $\left\langle \widehat{Q} U\Phi_0, V\Phi_0 \right\rangle = \langle Qu, v \rangle$ for all $u, v \in \mathbb{V}$. We introduce the mapping $\mathcal{B} : \mathbb{V} \to \mathbb{V}^*$ *via*

$$\langle \mathcal{B}(t), s \rangle = \frac{1}{2}\left\langle \widehat{Q}[[\mathcal{H}_K(t_*), T], T]\Phi_0, S\Phi_0 \right\rangle = \frac{1}{2}\left\langle \widehat{Q}[[\mathcal{W}_K(t_*), T], T]\Phi_0, S\Phi_0 \right\rangle, \tag{4.22}$$

that is, the $Q$-projection of $\frac{1}{2}\mathcal{A}''(t_*)(t, t)$. In the second equality, we used (4.4). Also, note that $\mathcal{B}$ is homogeneous of degree 2, *i.e.* $\mathcal{B}(\alpha t) = \alpha^2 \mathcal{B}(t)$. The next theorem follows esentially from Leray's second reduction formula (Thm. 3.9).

**Theorem 4.18** (Index formula for SRCC – degenerate case). *Let $t_*$ be zero of the CC mapping $\mathcal{A} : \mathbb{V} \to \mathbb{V}^*$. Suppose that $\mathcal{E}_{\mathrm{CC}}(t_*) \in \sigma(\mathcal{H}_K(t_*)_{\mathfrak{V}})$ and let $W_R$, $W_L$ and $Q$ be as above. Assume that*

$$\mathcal{B}(t) \neq 0 \quad \text{for all} \quad t \in \partial B_{\mathbb{V}}(0, 1). \tag{4.23}$$

*Then, $t_*$ is an isolated zero and the topological index of $\mathcal{A}$ at $t_*$ is given by*

$$i(\mathcal{A}, t_*) = i(\mathcal{A}'(t_*) + Q, 0)\, i(\mathcal{B}|_{W_R}, 0).$$

*Proof.* First, we prove that $t_*$ is isolated. When $r \notin \ker \mathcal{A}'(t_*)$ and small, it follows that $\mathcal{A}(t_* + r) \neq 0$ similarly to (4.16). If, however $r \in \ker \mathcal{A}'(t_*)$, then we may write

$$\langle \mathcal{A}(t_* + r), \mathcal{A}''(t_*)(r, r) \rangle_{\mathbb{V}^*} = \langle \mathcal{A}'(t_*)r_0, \mathcal{A}''(t_*)(r, r) \rangle_{\mathbb{V}^*} + \frac{1}{2}\|\mathcal{A}''(t_*)(r, r)\|_{\mathbb{V}^*}^2 + \mathcal{O}(\|r\|_{\mathbb{V}}^5)$$

for all $r \in B(0, \varepsilon)$ for sufficiently small $\varepsilon > 0$. Condition (4.23) implies that $\mathcal{A}(t_* + r) \neq 0$ for all $r \in B_{\mathbb{V}}^*(0, \varepsilon)$.

Next, we apply Corollary 3.10 with the choice $D = B_{\mathbb{V}}(0, \varepsilon)$, $L = \mathcal{A}'(t_*)$ and

$$\langle \mathcal{N}(t, \lambda), s \rangle = \sum_{k=2}^{2N} \frac{\lambda^{k-2}}{k!} \langle [\mathcal{H}_K(t_*), T]_{(k)}\Phi_0, S\Phi_0 \rangle,$$

where we used (4.9). Because $\mathrm{ran}\, Q = \ker \mathcal{A}'(t_*)^\dagger$, it follows that

$$\ker Q = (\mathrm{ran}\, Q)^\perp = (\ker \mathcal{A}'(t_*)^\dagger)^\perp = \mathrm{ran}\, \mathcal{A}'(t_*) = \mathrm{ran}\, L.$$

Moreover, since $t_*$ is an isolated zero, it is possible to choose $\delta > 0$ so that the equation

$$\lambda^{-1}\mathcal{A}(t_* + \lambda t) = \mathcal{A}'(t_*)t + \lambda \mathcal{N}(t, \lambda) = 0$$

does not admit a solution $t \in \partial B_{\mathbb{V}}(0, \delta)$ for any $\lambda \in (0, 1]$.

Note that $Q\mathcal{N}(t, 0) = \mathcal{B}(t) \neq 0$ for all $t \in B_{\mathbb{V}}^*(0, \delta)$ by assumption (4.23) and the homogenity of $\mathcal{B}$. We see that conditions (i) and (ii) of Corollory 3.10 are satisfied and the result follows. $\qquad\square$

The preceding theorem reduces the computation of the index to a low-dimensional problem but the zero is still degenerate. In fact, since

$$\langle \mathcal{B}'(t)u, v \rangle = \frac{1}{2} \Big\langle \widehat{Q}([[\mathcal{H}_K(t_*), U], T] + [[\mathcal{H}_K(t_*), T], U])\Phi_0, V\Phi_0 \Big\rangle,$$

we have that $t = 0$ is a degenerate zero of $\mathcal{B}|_{W_R}$ and by assumption the only zero. Therefore, we need to apply Theorem 3.8 to determine $i(\mathcal{B}|_{W_R}, 0)$.

**Corollary 4.19.** *For a degenerate, isolated zero $t_*$ of $\mathcal{A}$, $\dim W_R = 1$, and for which (4.23) holds, we have $i(\mathcal{A}, t_*) = 0$.*

*Proof.* Let $W_R = \mathrm{Span}\{r\}$ and $W_L = \mathrm{Span}\{\ell\}$. We apply Theorem 3.8 to the mapping $\mathcal{B}|_{W_R}$ with $D = B_{\mathbb{V}}(0, \varepsilon)$, $\varepsilon > 0$ arbitrary. Fix $z' = \eta\ell \in \mathrm{ran}\, Q = W_L$ such that $0 < |\eta| < \delta$. Since $t = cr \in W_R$, the equation $\mathcal{B}(t) = z'$ in $B_{\mathbb{V}}(0, \varepsilon) \cap W_R$ is equivalent to finding $0 \neq c \in \mathbb{R}$ such that

$$\frac{c^2}{2}\big\langle (\mathcal{H}_K(t_*) - \mathcal{E}_{\mathrm{CC}}(t_*))R^2\Phi_0, L\Phi_0 \big\rangle = \eta\langle L\Phi_0, L\Phi_0 \rangle. \tag{4.24}$$

Note that the inner product on the left-hand side is nonzero by assumption (4.23). Choose $\eta$ to be of opposite sign as the inner product on the left-hand side. Then there are no solutions $c$, so $i(\mathcal{B}|_{W_R}, t_*) = 0$ and therefore $i(\mathcal{A}, t_*) = 0$ by Theorem 4.18. $\qquad\square$

**Corollary 4.20.** *Let $t_*$ be an isolated zero of the CC mapping $\mathcal{A} : \mathbb{V} \to \mathbb{V}^*$. Suppose that $\mathcal{E}_{\mathrm{CC}}(t_*) \in \sigma(\mathcal{H}_K(t_*)_{\mathfrak{V}})$ and let $W_R$ as above. Assume that $\ker \mathcal{B}'(t) = \{0\}$ for all $0 \neq t \in W_R$. If $\mu := \dim W_R$ is odd, then $i(\mathcal{A}, t_*) = 0$ and if $\mu$ is even, then $i(\mathcal{A}, t_*)$ is even. In particular, under the above hypotheses, $i(\mathcal{A}, t_*)$ cannot be $\pm 1, \pm 3, \pm 5, \ldots$*

*Proof.* Let $z' \in B_{\mathbb{V}}^*(0, \delta) \cap W_L$. Then the equation $\mathcal{B}(t) = z'$ for $0 \neq t \in W_R$ is equivalent to

$$\frac{1}{2}\langle [[\mathcal{H}_K(t_*), T], T]\Phi_0, L\Phi_0 \rangle = \langle Z'\Phi_0, L\Phi_0 \rangle,$$

for all $\ell \in W_L$. Let $\mathcal{T}$ denote the set of solutions $t$ of the preceding equation (which can be empty). By Theorem 3.8, $|\mathcal{T}| = m$ for some $m$ finite. Notice that $\mathcal{T}$ is closed under the operation $t \mapsto -t$, so $m$ is even (we also used that $0 \notin \mathcal{T}$), and let

$$\mathcal{T} = \Big\{ t_1, \ldots, t_{\frac{m}{2}}, -t_1, \ldots, -t_{\frac{m}{2}} \Big\}$$

using some appropriate indexing. From the linearity of $t \mapsto \mathcal{B}'(t)$, we get *via* Theorem 3.7,

$$\begin{aligned}
\deg(\mathcal{B}|_{W_R}, B_{\mathbb{V}}(0, r) \cap W_R, z') &= \sum_{i=1}^{\frac{m}{2}} \mathrm{sgn}\det g\mathcal{B}'(t_i)h^{-1} + \mathrm{sgn}\det g\mathcal{B}'(-t_i)h^{-1} \\
&= (1 + (-1)^{\mu}) \sum_{i=1}^{\frac{m}{2}} \mathrm{sgn}\det g\mathcal{B}'(t_i)h^{-1},
\end{aligned} \tag{4.25}$$

for some sufficiently large $r > 0$. $\qquad\square$

We close this section with two remarks.

**Remark 4.21.**

(i) Note that, according to Theorem 3.8, a zero $t_*$ of topological index 0 is "numerically unstable", because one could miss zeros altogether if the equations are solved with finite precision arithmetic, even in the FCC case. Therefore, the degenerate zeros of the SRCC mapping $\mathcal{A}$ are not robust in general. We have already seen that in the FCC case, the CC energy $\mathcal{E}_{\mathrm{CC}}(t_*)$ of a degenerate zero $t_*$ is a degenerate eigenvalue $\mathcal{E}$ of the Hamiltonian (see Rem. 4.2). Thus, we can conclude that the SRCC method is in general unsuitable for finding degenerate eigenstates and eigenvalues – an empirical fact that is well known among the practitioners of the SRCC method.

(ii) Proposition 4.14 and the preceding calculations imply that any approach that stipulates the local strong monotonicity of $\mathcal{A}$ near $t_*$ can only provide an incomplete description of the SRCC method.

## 4.3. Local properties – complex case

We now discuss what happens when *complex* amplitude spaces are considered instead. We explain the complex case in detail, because the differences from the real case are somewhat subtle. Assume that $\mathbb{V}$ is a complex amplitude space and let $\mathbb{V}^*$ denote its anti-dual. It is clear that $t \mapsto \langle \mathcal{A}(t), s \rangle$ is a (complex) polynomial for fixed $s \in \mathbb{V}$, hence with the appropriate identifications, $\mathcal{A}_{\mathbb{C}} : \mathbb{V} \to \mathbb{V}^*$ is a holomorphic mapping, where we used the subscript $\mathbb{C}$ to highlight the difference[7]. Of course, a real zero $t_* \in \mathbb{V}$ to $\mathcal{A}(t_*) = 0$ is automatically a "complex" zero: $\mathcal{A}_{\mathbb{C}}(t_*) = 0$. Further, using the fact that the Hamiltonian is real (by which we mean $\langle \mathcal{H}_K \Phi_\alpha, \Phi_\beta \rangle \in \mathbb{R}$), $\mathcal{A}_{\mathbb{C}}(t_*) = 0$ if and only if $\mathcal{A}_{\mathbb{C}}(\bar{t}_*) = 0$. Also, $\mathcal{E}_{\text{CC}}(\bar{t}) = \overline{\mathcal{E}_{\text{CC}}(t)}$.

From Theorem 3.16(i) we immediately get that $\deg(\mathcal{A}_{\mathbb{C}}, U, 0) \geq 0$ for every bounded open $U \subset \mathbb{V}$. In particular, $i(\mathcal{A}_{\mathbb{C}}, t_*) \geq 0$ for every isolated zero $t_*$. Notice that, even if $t_*$ is a zero of both $\mathcal{A}$ and $\mathcal{A}_{\mathbb{C}}$, its real and complex indices $i(\mathcal{A}, t_*)$ and $i(\mathcal{A}_{\mathbb{C}}, t_*)$ may differ; for instance we know that $i(\mathcal{A}, t_*)$ can have a sign, while $i(\mathcal{A}_{\mathbb{C}}, t_*)$ cannot. Also, Theorem 3.16(ii) implies that $i(\mathcal{A}, t_*) \geq 2$ for an isolated, degenerate zero $t_*$. Moreover, the following is true.

**Theorem 4.22.** *If $t_* \in \mathbb{V}$ is a* real *isolated zero of both $\mathcal{A}$ and $\mathcal{A}_{\mathbb{C}}$, then*

$$|i(\mathcal{A}, t_*)| \leq i(\mathcal{A}_{\mathbb{C}}, t_*), \quad i(\mathcal{A}, t_*) \equiv i(\mathcal{A}_{\mathbb{C}}, t_*) \mod 2.$$

*Proof.* The result follows from Theorem 3.18 with the choice $D = B_{\mathbb{V}}(t_*, \delta)$. □

For simplicity, we assume from now on that $\mathbb{V}$ is rank-regular (Def. 4.4), so that $\widehat{\mathcal{H}_K}(t)_{\mathfrak{V}} = \mathcal{H}_K(t)_{\mathfrak{V}}$. Adapting the proofs of Theorem 4.13 and 4.18 in the complex case, we have

**Theorem 4.23** (Index formula for SRCC – complex case)**.** *Let $t_*$ be an isolated zero of $\mathcal{A}_{\mathbb{C}} : \mathbb{V} \to \mathbb{V}^*$. Then the following hold true.*

(i) *The zero $t_*$ is non-degenerate if and only if $\mathcal{E}_{\text{CC}}(t_*) \notin \sigma(\mathcal{H}_K(t_*)_{\mathfrak{V}})$, and in this case $i(\mathcal{A}_{\mathbb{C}}, t_*) = 1$.*
(ii) *Suppose that $\mathcal{E}_{\text{CC}}(t_*) \in \sigma(\mathcal{H}_K(t_*)_{\mathfrak{V}})$ and that the second-order regularity assumption (4.23) holds true. Then*

$$i(\mathcal{A}_{\mathbb{C}}, t_*) = i(\mathcal{B}|_{W_R}, 0).$$

(iii) *Suppose that $\mathcal{E}_{\text{CC}}(t_*) \in \sigma(\mathcal{H}_K(t_*)_{\mathfrak{V}})$ and that the second-order regularity assumption (4.23) holds true. Assume further, that $\ker \mathcal{B}'(t) = \{0\}$ for all $0 \neq T\Phi_0 \in W_R$. Then $i(\mathcal{A}_{\mathbb{C}}, t_*) = m \geq 2$, where $m$ is the number of solutions $0 \neq R\Phi_0 \in W_R$ to the equation*

$$\langle (\mathcal{H}_K(t_*) - \mathcal{E}_{\text{CC}}(t_*))R^2\Phi_0, L\Phi_0 \rangle = \langle Z\Phi_0, L\Phi_0 \rangle \quad (L\Phi_0 \in W_L)$$

*for any $Z\Phi_0 \in W_L \cap B^*(0, \delta)$ for sufficiently small $\delta > 0$.*

*Proof.* For (i), it is enough to note that $\ker \mathcal{A}_{\mathbb{C}}'(t_*) \neq \{0\}$ follows *via* the same calculation as in the proof of Theorem 4.13. Part (ii) follows since $i(\mathcal{A}_{\mathbb{C}}'(t_*) + Q) = 1$. For the proof of (iii), notice that (4.25) now reads

$$\deg(\mathcal{B}|_{W_R}, B_{\mathbb{V}}(0, r) \cap W_R, z) = \sum_{i=1}^{m} \text{sgn}|\det g\mathcal{B}'(t_i)h^{-1}|^2 = m.$$

□

**Corollary 4.24.** *For a degenerate, isolated zero $t_*$ of $\mathcal{A}_{\mathbb{C}}$, $\dim W_R = 1$, and for which (4.23) holds, we have $i(\mathcal{A}_{\mathbb{C}}, t_*) = 2$.*

---

[7]Note that the realification of $\mathcal{A}_{\mathbb{C}}$, $(\mathcal{A}_{\mathbb{C}})_{\mathbb{R}}$ does *not* equal $\mathcal{A}$ (they are mappings of different type: $(\mathcal{A}_{\mathbb{C}})_{\mathbb{R}} : \mathbb{R}^{2n} \to \mathbb{R}^{2n}$).

*Proof.* The proof is analogous to that of Corollary 4.19, but (4.24) always has exactly two nonzero complex solutions $c$. ☐

We close this section with a few remarks.

**Remark 4.25.**

(i) Luckily, for a real zero $t_* \in \mathbb{V}$, the condition $\mathcal{E}_{\mathrm{CC}}(t_*) \notin \sigma(\mathcal{H}_K(t_*)_{\mathfrak{V}})$ is formally the same as in the real case, therefore a non-degenerate real zero is automatically a non-degenerate zero of $\mathcal{A}_{\mathbb{C}}$.

(ii) The degeneracy of a complex zero $t_*$ manifests itself in numerical computations as follows. Suppose the hypotheses of Theorem 4.23(iii) hold true. Combining this with Theorem 3.17, we get that the perturbed equation $\mathcal{A}(t) = z'$ has *exactly m* solutions for almost all $z' \in \mathbb{V}$ sufficiently close to zero. This is in contrast with the real case, when one might completely "lose" solutions when the index is zero (Rem. 4.21(i)). The appearance of multiple complex zeros in degenerate situations was conjectured based on numerical observations in [32, 33].

(iii) If the Hamiltonian is real (see above), then the index of the complex conjugate zero is the same: $i(\mathcal{A}_{\mathbb{C}}, t_*) = i(\mathcal{A}_{\mathbb{C}}, \bar{t}_*)$.

(iv) The classical *Bézout theorem* states that if a polynomial system

$$\left.\begin{aligned} P_1(x_1, \ldots, x_d) &= 0 \\ P_2(x_1, \ldots, x_d) &= 0 \\ &\cdots \\ P_n(x_1, \ldots, x_d) &= 0 \end{aligned}\right\}$$

has a finite number of zeros in $\mathbb{C}^d$, then the number of zeros (counting multiplicities) is at most $\Delta = \Delta_1 \cdots \Delta_d$ (called the *Bézout number*), where $\Delta_k$ denotes the degree of $P_k$. As remarked earlier, according to the Baker–Campbell–Hausdorff expansion (4.8), for the polynomials constituting the system $\mathcal{A}(t) = 0$ there holds $\Delta_1 = \ldots = \Delta_d = 4$, hence the Bézout number of the CC equations is $\Delta = 4^d$, where $d = \dim \mathbb{V}$. This is typically a huge number. However, it is known that the Bézout number often grossly overestimates the number of zeros. In fact, it was observed numerically that the number of zeros for the (truncated) CC equations is much less than the Bézout number [32].

## 4.4. Continuation of solutions

In this section we discuss how solutions of different CC methods can be "connected" in a systematic way. The idea is not new to this field (and certainly not new to nonlinear analysis, see *e.g.* [43]), and it has been a subject of both theoretical and numerical investigations in the CC literature, as we have already mentioned in the introduction.

The main theoretical tool we use to describe the aforementioned connection is a specific type of homotopy.

**Definition 4.26.** Let $\mathbb{V}^1$ be an amplitude space with direct sum decomposition $\mathbb{V}^1 = \mathbb{V}^0 \oplus \mathbb{V}^{\angle}$. Let $\mathcal{A}^j : \mathbb{V}^j \to (\mathbb{V}^j)^*$ be continuous mappings for $j = 0, 1$. A continuous map $\mathcal{K} : \mathbb{V}^1 \times [0, 1] \to (\mathbb{V}^1)^*$ is said to be an *admissible homotopy*, if

(i) $\langle \mathcal{K}(t^1, 0), s^0 \rangle = \langle \mathcal{A}^0(t^0), s^0 \rangle$ for all $t^1 \in \mathbb{V}^1$, $s^0 \in \mathbb{V}^0$, and

(ii) $\mathcal{K}(\cdot, 1) = \mathcal{A}^1$.

Furthermore, an admissible homotopy $\mathcal{K} : \mathbb{V}^1 \times [0, 1] \to (\mathbb{V}^1)^*$ is said to be *faithful*, if for every $t^0_{**} \in \mathbb{V}^0$ such that $\mathcal{A}^0(t^0_{**}) = 0$, there exists $t^{\angle}_{**} \in \mathbb{V}^{\angle}$ so that with $t^1_{**} = t^0_{**} + t^{\angle}_{**} \in \mathbb{V}^1$, there holds $\mathcal{K}(t^1_{**}, 0) = 0$.

**Example 4.27.** When $\mathcal{A}^0$ is the projection of $\mathcal{A} := \mathcal{A}^1$ onto $\mathbb{V}^0$, a simple admissible homotopy can be given by

$$\langle \mathcal{K}(t^1, \lambda), s^1 \rangle = \langle \mathcal{A}(t^0 + \lambda t^{\angle}), s^0 \rangle + \langle \mathcal{A}(t^1), s^{\angle} \rangle, \tag{4.26}$$

for all $t^1 \in \mathbb{V}^1$, $s^1 \in \mathbb{V}^1$ and $\lambda \in [0, 1]$ (*cf.* (4.33)).

Let $\mathbb{V}^1 = \mathbb{V}(G^{\text{full}})$ and $\mathcal{A}^1$ be the FCC mapping, $\mathbb{V}^0$ some rank-truncated space and $\mathcal{A}^0$ the truncated CC mapping. This case is particularly important due to the equivalence of FCC and FCI (see Thm. 4.4 of Part I), so existence of a solution to the FCI problem (essentially the Schrödinger equation) can be exploited to infer the existence of a *truncated* CC solution. Furthermore, the topological index of the CC solution can be determined by the results of Section 4.2 and the homotopy invariance of the topological degree can be used relate these quantities in certain situations (see (4.29) below).

The *zero set* of an (admissible) homotopy $\mathcal{K}$ is defined as

$$\mathcal{Z}(\mathcal{K}) = \{(t_*^1, \lambda) \in \mathbb{V}^1 \times [0,1] : \mathcal{K}(t_*^1, \lambda) = 0\}. \tag{4.27}$$

We omit $\mathcal{K}$ from the notation $\mathcal{Z}(\mathcal{K})$ whenever it is clear from the context. The $\lambda$-sections of $\mathcal{Z}$ are denoted as $\mathcal{Z}_\lambda = \{t_*^1 \in \mathbb{V}^1 : (t_*^1, \lambda) \in \mathcal{Z}\}$. Clearly, $(\mathcal{Z}(\mathcal{K}))_1 = (\mathcal{A}^1)^{-1}(0)$ for *any* admissible homotopy $\mathcal{K}$. Furthermore, $\Pi_{\mathbb{V}^0}(\mathcal{Z}(\mathcal{K}))_0 = (\mathcal{A}^0)^{-1}(0)$ for any faithful, admissible homotopy $\mathcal{K}$. We will sometimes explicitly label $t_*^1$ with the $\lambda$ to which it corresponds as $t_*^1(\lambda)$.

Recall that any topological space can be partitioned into a family of *connected components*, which are maximal (w.r.t. set inclusion), closed connected sets. We are interested in the connected components of $\mathcal{Z}(\mathcal{K})$. The Leray–Schauder continuation principle ([10], Thm. 2.1.3) immediately implies the following.

**Theorem 4.28.** *Suppose that $\mathcal{D} \subset \mathbb{V}^1 \times [0,1]$ is a bounded open set and let $\mathcal{K} : \mathbb{V} \times [0,1] \to \mathbb{V}^*$ be an admissible homotopy such that*

$$\mathcal{K}(t^1, \lambda) \neq 0 \quad \text{for all} \quad (t^1, \lambda) \in \partial\mathcal{D}. \tag{4.28}$$

*Then the following statements hold true.*

(i) $\deg(\mathcal{K}(\cdot, 0), \mathcal{D}_0, 0) = \deg(\mathcal{A}^1, \mathcal{D}_1, 0) =: d$.
(ii) *If $d \neq 0$, then there is a connected component $\mathcal{C}$ of $\mathcal{Z}$, such that*

$$\mathcal{C} \cap (\mathcal{Z}_j \times \{j\}) \neq \emptyset \quad \text{for} \quad j = 0, 1.$$

This general theorem can be used to prove existence results, which essentially consists of establishing the boundary condition (4.28). Furthermore, note that (i) implies that

$$\sum_{t_{**}^1 \in \mathcal{Z}_0 \cap \mathcal{D}_0} i(\mathcal{K}(\cdot, 0), t_{**}^1) = \sum_{t_*^1 \in \mathcal{Z}_1 \cap \mathcal{D}_1} i(\mathcal{A}^1, t_*^1) = d. \tag{4.29}$$

This is a necessary condition for zeros in $\mathcal{Z}_0$ and $\mathcal{Z}_1$ be in the same bounded connected component.

**Remark 4.29.** In the case when $\mathcal{D}$ is possibly *unbounded*, one can invoke ([10], Thm. 2.1.4) to obtain the following statement: If $\mathcal{Z}_0$ is bounded and $\deg(\mathcal{K}(\cdot, 0), D, 0) \neq 0$ for some bounded open set $D \supset \mathcal{Z}_0$, then there is a connected component $\mathcal{C}$ of $\mathcal{Z}$ intersecting $\mathcal{Z}_0 \times \{0\}$ which either intersects $\mathcal{Z}_1 \times \{1\}$ or is unbounded. The unboundedness of $\mathcal{C}$ has been observed numerically for a specific type of homotopy and its significance is discussed in the next remark.

**Remark 4.30.** In [32], a specific type of admissible homotopy (essentially (4.26), which will be discussed in Sect. 4.5 below) is used as the basis to distinguish "physical" truncated solutions from "unphysical" ones. Roughly speaking, they call a truncated solution "physical" if it shares a *bounded* connected component with an FCC solution (although they require more regularity of the connected component in question). They also found that some FCC (resp. truncated) solutions cannot be "connected" to any truncated (resp. FCC) solution. While this approach to the classification of solutions seems attractive at first, it has a serious conceptual drawback: it is unclear why one particular homotopy is preferred over the (infinitely many) others. For instance, it is conceivable that an admissible homotopy classifies a truncated solution as "unphysical" while another homotopy classifies it as "physical". Moreover, their homotopy itself does not admit an obvious physical interpretation. We will not pursue this line of thought any further in this work, and view homotopies as purely mathematical devices.

As an example of an admissible homotopy, we consider the *linear homotopy*. Let $\mathbb{V}^0 \subset \mathbb{V}^1$ be any subspace (a typical choice is based on some truncation: $\mathbb{V}^0 = \mathbb{V}(G(1,\ldots,\rho))$) and $\mathbb{V}^{\angle} := (\mathbb{V}^0)^{\perp}$, where the orthogonal complement is taken with respect to the $\mathbb{V}$-inner product. Also, to emphasize that $\mathbb{V}^{\angle}$ is actually given by the $\mathbb{V}$-orthogonal complement of $\mathbb{V}^0$, we shall write $t^1 = t^0 + t^{\perp}$ for some unique $t^0 \in \mathbb{V}^0$ and $t^{\perp} \in (\mathbb{V}^0)^{\perp}$, for any $t^1 \in \mathbb{V}^1$. Let $\mathcal{K}_{\mathrm{L}} : \mathbb{V}^1 \times [0,1] \to (\mathbb{V}^1)^*$ be given by

$$\langle \mathcal{K}_{\mathrm{L}}(t^1, \lambda), s^1 \rangle = (1-\lambda)\big(\langle \mathcal{A}^0(t^0), s^0 \rangle + \alpha |\langle t^{\perp} - u^{\perp}, s^{\perp} \rangle_{\mathbb{V}}|\big) + \lambda \langle \mathcal{A}^1(t^1), s^1 \rangle$$

for any $t^1, s^1 \in \mathbb{V}^1$ and $\lambda \in [0,1]$, and some fixed constants $\alpha > 0$ and $u^{\perp} \in (\mathbb{V}^0)^{\perp}$. Then, clearly $\mathcal{K}_{\mathrm{L}}(\cdot, 1) = \mathcal{A}^1$. Further, $\mathcal{K}_{\mathrm{L}}(t^1_{**}, 0) = 0$ is equivalent to $\mathcal{A}^0(t^0_{**}) = 0$ and $t^{\perp}_{**} = u^{\perp}$. Therefore, $\mathcal{K}_{\mathrm{L}}$ is a faithful, admissible homotopy. Also, $\mathcal{K}_{\mathrm{L}}$ has the following trivial property: if $t^1_* \in \mathbb{V}^1$ is such that $\mathcal{A}^0(t^0_*) = 0$ and $\mathcal{A}^1(t^1_*) = 0$, then $\mathcal{K}_{\mathrm{L}}(t^1_*, \lambda) = 0$ for all $\lambda \in [0,1]$.

As a simple application of the linear homotopy, we give a variant of the existence result ([36], Thm. 4.1).

**Theorem 4.31.** *Suppose that $t^1_* \in \mathbb{V}^1$ is a zero of $\mathcal{A}^1$. Let $\varkappa = \|t^{\perp}_*\|_{\mathbb{V}}$. Assume the following hold true.*

(i) $\mathcal{A}^j : \mathbb{V}^j \to \mathbb{V}^j$ *is* $C^1$.
(ii) $L^0 \varkappa = \sup_{\|u^0\|_{\mathbb{V}}=1} \langle \mathcal{A}^0(t^0_*) - \mathcal{A}^1(t^1_*), u^0 \rangle$ *for some* $L^0 < \infty$.
(iii) $\mathcal{A}^j$ *is strongly monotone in* $B_{\mathbb{V}^j}(t^j_*, \delta_j)$ *with constant* $C^j_{\mathrm{SM}} = C^j_{\mathrm{SM}}(\delta_j) > 0$ *for* $j = 0, 1$.
(iv) $\varkappa < \delta \frac{C^0_{\mathrm{SM}}}{L^0}$, *where* $\delta = \min\{\delta_0, \delta_1\}$.

*Then there exists a unique* $t^0_{**} \in B_{\mathbb{V}^0}(t^0_*, \delta)$ *such that* $\mathcal{A}^0(t^0_{**}) = 0$. *Furthermore,* $i(\mathcal{K}_{\mathrm{L}}(\cdot, 0), t^1_{**}) = i(\mathcal{A}^0, t^0_*) = 1$.

*Proof.* Set $\alpha := C^0_{\mathrm{SM}}$ and $u^{\perp} := t^{\perp}_*$ in the definition of $\mathcal{K}_{\mathrm{L}}$. We prove that the boundary condition $\mathcal{K}_{\mathrm{L}}(t^1, \lambda) \neq 0$ holds true for all $t^1 \in \partial B_{\mathbb{V}}(t^1_*, \delta)$ and $\lambda \in [0,1]$. Write $t^1 = t^1_* + r^1$, where $\|r^1\|_{\mathbb{V}} = \delta$ and

$$\langle \mathcal{K}_{\mathrm{L}}(t^1_* + r^1, \lambda), r^1 \rangle =: (1-\lambda)A_0 + \lambda A_1.$$

Here,

$$\begin{aligned}
A_0 &= \langle \mathcal{A}^0(t^0_* + r^0), r^0 \rangle + C^0_{\mathrm{SM}} \langle t^{\perp}_* + r^{\perp} - t^{\perp}_*, r^{\perp} \rangle_{\mathbb{V}} \\
&= \langle \mathcal{A}^0(t^0_* + r^0) - \mathcal{A}^0(t^0_*), r^0 \rangle + \langle \mathcal{A}^0(t^0_*) - \mathcal{A}^1(t^1_*), r^0 \rangle + C^0_{\mathrm{SM}} \|r^{\perp}\|^2_{\mathbb{V}} \\
&\geq C^0_{\mathrm{SM}} \|r^0\|^2_{\mathbb{V}} - L^0 \varkappa \|r^0\|_{\mathbb{V}} + C^0_{\mathrm{SM}} \|r^{\perp}\|^2_{\mathbb{V}} \\
&\geq (C^0_{\mathrm{SM}} \|r^1\|_{\mathbb{V}} - L^0 \varkappa) \|r^1\|_{\mathbb{V}} = (C^0_{\mathrm{SM}} \delta - L^0 \varkappa) \delta > 0,
\end{aligned}$$

where we in the last step used that $\|r^1\|^2_{\mathbb{V}} = \|r^0\|^2_{\mathbb{V}} + \|r^{\perp}\|^2_{\mathbb{V}}$. Furthermore,

$$A_1 = \langle \mathcal{A}^1(t^1_* + r^1), r^1 \rangle = \langle \mathcal{A}^1(t^1_* + r^1) - \mathcal{A}^1(t^1_*), r^1 \rangle \geq C^1_{\mathrm{SM}} \|r^1\|^2_{\mathbb{V}} > 0.$$

Using Proposition 4.14 and the uniqueness of the zero $t^1_*$ in $B_{\mathbb{V}}(t^1_*, \delta)$,

$$\deg(\mathcal{K}_{\mathrm{L}}(\cdot, 1), B_{\mathbb{V}}(t^1_*, \delta), 0) = \deg(\mathcal{A}^1, B_{\mathbb{V}}(t^1_*, \delta), 0) = 1.$$

Since we have already seen that $\mathcal{K}_{\mathrm{L}}(t^1, \lambda) \neq 0$ for all $t^1 \in \partial B_{\mathbb{V}}(t^1_*, \delta)$ and $\lambda \in [0,1]$, the homotopy invariance of the degree can be applied to get

$$\deg(\mathcal{K}_{\mathrm{L}}(\cdot, 0), B_{\mathbb{V}}(t^1_*, \delta), 0) = 1.$$

Using the existence property of the degree Corollary 3.3(ii), there exists $t^0_{**} \in \mathbb{V}^0$ such that $(t^0_{**}, t^{\perp}_*) \in B_{\mathbb{V}^1}(t^1_*, \delta)$ and $\langle \mathcal{K}_{\mathrm{L}}(t^0_{**}, 0), s^0 \rangle = \langle \mathcal{A}^0(t^0_{**}), s^0 \rangle = 0$ for all $s^0 \in \mathbb{V}^0$, which is what we wanted to prove. Uniqueness follows form the local strong monotonicity of $\mathcal{A}^0$. $\qquad\square$

In the special case when $\mathcal{A}^0$ is given as a projection of $\mathcal{A}^1$ (see Rem. 4.8), we obtain the following.

**Corollary 4.32.** *Suppose that $t^1_* \in \mathbb{V}^1$ is a zero of $\mathcal{A}^1$. Let $\varkappa = \|t^\perp_*\|_\mathbb{V}$ and suppose the following hold true.*

(i) $L^0 \varkappa = \sup_{\|u^0\|_\mathbb{V}=1} \langle \mathcal{A}^1(t^0_*) - \mathcal{A}^1(t^1_*), u^0 \rangle$ *for some $L^0 < \infty$.*

(ii) $\mathcal{A}^1$ *is strongly monotone in $B_{\mathbb{V}^1}(t^1_*, \delta)$ with constant $C_{\mathrm{SM}} > 0$.*

(iii) $\varkappa < \delta \frac{C_{\mathrm{SM}}}{C_{\mathrm{SM}}+L^0}$.

*Then there exists a unique $t^0_{**} \in B_{\mathbb{V}^0}(t^1_*, \delta - \varkappa)$ such that $\mathcal{A}^1(t^0_{**}) = 0$.*

*Proof.* It is enough to recall that $\mathcal{A}^0 = \Pi_{\mathbb{V}^0} \mathcal{A}^1|_{\mathbb{V}^0}$ is strongly monotone with constant $C_{\mathrm{SM}}$ in

$$B_{\mathbb{V}^0}\left(t^0_*, \sqrt{\delta^2 - \varkappa^2}\right) \supset B_{\mathbb{V}^0}\left(t^0_*, \delta - \varkappa\right).$$

Applying Theorem 4.31 with $C^0_{\mathrm{SM}} = C_{\mathrm{SM}}$ and $\delta_0 = \delta - \varkappa$ gives the result. $\qquad\square$

Condition (i) certainly holds, since the CC mapping is locally Lipschitz (see [36, 37]), but the constant $L^0$ here is possibly smaller than the usual Lipschitz constant (which, in turn, is bounded from below by the strong monotonicity constant).

## 4.5. Kowalski–Piecuch homotopy

In this section, we consider another homotopy that can be used to analyze the connection between the solutions to CC methods of different rank-truncation levels. The approach was pioneered by the chemists K. Kowalski and P. Piecuch [32], who conducted a comprehensive numerical study based on this idea[8]. They considered complex amplitudes and the following discussion is easily extended to that case.

**Assumption.** *Let $\mathbb{V}^1 = \mathbb{V}^0 \oplus \mathbb{V}^\angle$ be an $\ell^2$-orthogonal direct sum decomposition of a real amplitude space $\mathbb{V}^1 = \mathbb{V}(G^1)$, where $\mathbb{V}^0$ is an amplitude space corresponding to some lower truncation level. More precisely, assume that there is a rank $\rho \geq 1$, such that $\mathbb{V}^0$ contains all amplitudes with rank $\leq \rho$ and $\mathbb{V}^\angle := (\mathbb{V}^0)^{\perp_{\ell^2}}$ contains all amplitudes with rank $> \rho$. In the notations introduced in Part I,*

$$\mathbb{V}^0 = \mathbb{V}(G^0) = \mathbb{V}(G(1,\dots,\rho)) \cap \mathbb{V}^1, \text{ and } \mathbb{V}^\angle = \mathbb{V}(G^\perp) = \mathbb{V}(G(\rho+1,\dots,N)) \cap \mathbb{V}^1,$$

*where $G^0 = G(1,\dots,\rho) \cap G^1$ and $G^\perp = G(\rho+1,\dots,N) \cap G^1$. Hence, any $t^1 \in \mathbb{V}^1$ may be uniquely decomposed as $t^1 = t^0 + t^\angle$, where $t^0 \in \mathbb{V}^0$, $t^\angle \in \mathbb{V}^\angle$ and $\langle t^0, t^\angle \rangle_{\ell^2} = 0$.*

Let $\mathbb{V} = \mathbb{V}(G^{\mathrm{full}})$ be the full amplitude space and suppose that $\mathcal{A} : \mathbb{V}^1 \to (\mathbb{V}^1)^*$ is the SRCC mapping (4.1). Write

$$
\begin{aligned}
\langle \mathcal{A}(t^1), s^1 \rangle &= \left\langle e^{-T^1} \mathcal{H}_K e^{T^1} \Phi_0, S^0 \Phi_0 \right\rangle + \left\langle e^{-T^1} \mathcal{H}_K e^{T^1} \Phi_0, S^\angle \Phi_0 \right\rangle \\
&= \left\langle \left( e^{-T^0} + e^{-T^0}\left( e^{-T^\angle} - I \right) \right) \mathcal{H}_K \left( e^{T^0} + e^{T^0}\left( e^{T^\angle} - I \right) \right) \Phi_0, S^0 \Phi_0 \right\rangle + \left\langle e^{-T^1} \mathcal{H}_K e^{T^1} \Phi_0, S^\angle \Phi_0 \right\rangle \\
&= \langle \mathcal{A}^0(t^0), s^0 \rangle + \left\langle e^{-T^0}\left( e^{-T^\angle} - I \right) \mathcal{H}_K e^{T^0}(e^{T^\angle} - I)\Phi_0, S^0 \Phi_0 \right\rangle \\
&\quad + \left\langle e^{-T^0}\left( e^{-T^\angle} - I \right) \mathcal{H}_K e^{T^0} \Phi_0, S^0 \Phi_0 \right\rangle \\
&\quad + \left\langle e^{-T^0} \mathcal{H}_K e^{T^0}(e^{T^\angle} - I)\Phi_0, S^0 \Phi_0 \right\rangle + \left\langle e^{-T^1} \mathcal{H}_K e^{T^1} \Phi_0, S^\angle \Phi_0 \right\rangle \\
&= \langle \mathcal{A}(t^0), s^0 \rangle + \left\langle e^{-T^0} \mathcal{H}_K e^{T^0} \left( e^{T^\angle} - I \right) \Phi_0, S^0 \Phi_0 \right\rangle + \langle \mathcal{A}(t^1), s^\angle \rangle,
\end{aligned}
$$

---

[8]They call it the "$\beta$–nested equations", $\beta$ being the homotopy parameter.

where in the last step the second and the third terms of the penultimate expression vanish because

$$\left(e^{-T^{\angle}} - I\right)^{\dagger} S^0 \Phi_0 = -\left(T^{\angle}\right)^{\dagger} S^0 \Phi_0 - \frac{1}{2}\left(\left(T^{\angle}\right)^{\dagger}\right)^2 S^0 \Phi_0 - \ldots = 0,$$

due to the definition of the spaces $\mathbb{V}^0$ and $\mathbb{V}^{\angle}$. Note that this last relation would not hold in the case $\mathbb{V}^0 = \mathbb{V}(G(\mathrm{D})) \cap \mathbb{V}^1$, $\mathbb{V}^1 = \mathbb{V}(G^{\mathrm{full}})$, which is actually excluded by assumption.

Motivated by the preceding calculation in (4.30), we define the *Kowalski–Piecuch homotopy* $\mathcal{K}_{\mathrm{KP}} : \mathbb{V}^1 \times [0,1] \to (\mathbb{V}^1)^*$ *via* the instruction

$$\left\langle \mathcal{K}_{\mathrm{KP}}(t^1, \lambda), s^1 \right\rangle := \left\langle \mathcal{H}_K(t^0)\Phi_0, S^0\Phi_0 \right\rangle + \left\langle \mathcal{H}_K(t^1)\Phi_0, S^{\angle}\Phi_0 \right\rangle + \lambda\left\langle \mathcal{H}_K(t^0)(e^{T^{\angle}} - I)\Phi_0, S^0\Phi_0 \right\rangle \quad (4.30)$$

for all $t^1, s^1 \in \mathbb{V}^1$ and $\lambda \in [0,1]$. Here, the relation $\mathcal{K}_{\mathrm{KP}}(t^1_*, \lambda) = 0$ is the same as the system ([32], Eqs. (90), (91) and (93)).

It is obvious that $\mathcal{K}_{\mathrm{KP}}$ is an admissible homotopy. However, it is unclear whether it is faithful or not. In fact, $\mathcal{K}_{\mathrm{KP}}(t^1_{**}, 0) = 0$ is equivalent to the system

$$\left\langle \mathcal{H}_K(t^0_{**})\Phi_0, S^0\Phi_0 \right\rangle = 0, \quad (4.31)$$

$$\left\langle \mathcal{H}_K(t^0_{**} + t^{\angle}_{**})\Phi_0, S^{\angle}\Phi_0 \right\rangle = 0, \quad (4.32)$$

for all $s^1 \in \mathbb{V}^1$. In other words, the usual SRCC equation $\mathcal{A}(t^0_{**}) = 0$ (4.31), is augmented with an additional equation for $t^{\angle}_{**}$ (4.32), which, in turn, depends on $t^0_{**}$. It is not obvious at all that (4.32) has a solution $t^{\angle}_{**}$ for a given zero $t^0_{**}$. The extensive numerical evidence in [32] clearly indicates that (4.32) admits a solution in various circumstances.

Before stating our existence result we recast the KP homotopy into a more convenient form, which we already encountered in (4.26).

**Lemma 4.33.** *The following formula holds true:*

$$\left\langle \mathcal{K}_{\mathrm{KP}}(t^1, \lambda), s^1 \right\rangle = \left\langle \mathcal{A}(t^0 + \lambda t^{\angle}), s^0 \right\rangle + \left\langle \mathcal{A}(t^1), s^{\angle} \right\rangle, \quad (4.33)$$

*for all $t^1, s^1 \in \mathbb{V}^1$ and $\lambda \in [0,1]$.*

*Proof.* It is enough to prove that

$$\left\langle \mathcal{K}_{\mathrm{KP}}(t^1, \lambda), s^1 \right\rangle = \left\langle e^{-T^0}\mathcal{H}_K e^{T^0 + \lambda T^{\angle}}\Phi_0, S^0\Phi_0 \right\rangle + \left\langle e^{-T^1}\mathcal{H}_K e^{T^1}\Phi_0, S^{\angle}\Phi_0 \right\rangle, \quad (4.34)$$

because $(e^{-\lambda T^{\angle}})^{\dagger} S^0\Phi_0 = S^0\Phi_0$. To see (4.34), note that in the expansion $e^{T^0 + \lambda T^{\angle}} = e^{T^0}(I + \lambda T^{\angle} + \frac{\lambda^2}{2!}(T^{\angle})^2 + \ldots + \frac{\lambda^N}{N!}(T^{\angle})^N)$ the quadratic and higher-order terms do not contribute because the excitation rank of $\mathcal{H}_K(T^{\angle})^k\Phi_0$ exceeds $\rho$ for $k = 2, \ldots, N$, due to the fact that $\mathcal{H}_K$ is a two-body operator. $\square$

To formulate the existence result, suppose that $t^1_* \in \mathbb{V}^1$ is a non-degenerate zero of $\mathcal{A}$. Since $\mathbb{V}^1$ is assumed to be finite-dimensional, it is always possible to choose an $\alpha > 0$ so that

$$\left\langle (\mathcal{A}'(t^1_*) + \alpha I)r^1, r^1 \right\rangle \geq \gamma_\alpha \|r^1\|_{\mathbb{V}}^2 \quad \text{for all} \quad r^1 \in \mathbb{V}^1, \quad (4.35)$$

for some $\gamma_\alpha > 0$. This is also true in the complex case with a "Re" added to the left-hand side.

Define the operator $\Theta_\alpha : \mathbb{V} \to \mathbb{V}$ *via*[9]

$$\left\langle \mathcal{A}'(t^1_*)u^1, \Theta_\alpha v^1 \right\rangle = \left\langle (\mathcal{A}'(t^1_*) + \alpha I)u^1, v^1 \right\rangle \quad \text{for all} \quad u^1, v^1 \in \mathbb{V}^1. \quad (4.36)$$

Then $\Theta_\alpha$ is well-defined, as long as $\ker \mathcal{A}'(t^1_*) = \{0\}$, *i.e.* that $t^1_*$ is non-degenerate.

---

[9]The authors learned this trick from [5, 6].

**Theorem 4.34** (Existence for KP). *Let $t_*^1 \in \mathbb{V}^1$ be a non-degenerate zero of $\mathcal{A}$. Suppose the following.*

(i) *With $\theta_0 = \|\Pi_0(\Theta_\alpha - I)\Pi_0\|^2_{\mathcal{L}(\mathbb{V})}$ and $\theta_\angle = \|\Pi_0(\Theta_\alpha - I)\Pi_\angle\|^2_{\mathcal{L}(\mathbb{V})}$, there holds*

$$\eta := (1-g)\frac{\gamma_\alpha - \frac{1}{2}M_\delta\|\Theta_\alpha\|_{\mathcal{L}(\mathbb{V})}\delta}{\Delta(t_*^1) + M_\delta\delta} - \frac{1}{2}\max\{\varepsilon + 2(1+\varepsilon^{-1})\theta_0, 2(1+\varepsilon^{-1})\theta_\angle\} > \frac{1}{2},$$

*for some $\varepsilon > 0$ and $\delta > 0$, where $M_\delta$ was defined in (4.15). Here,*

$$\Delta(t_*^1) = \max_{\substack{\|u^0\|_\mathbb{V}=1 \\ \|v^\angle\|_\mathbb{V}=1}} \left|\left\langle (\mathcal{H}_K + [(T_*^1)_1^\dagger, \mathcal{H}_K])U^0\Phi_0, V^\angle\Phi_0\right\rangle\right|.$$

*Also, $\alpha > 0$ and $\gamma_\alpha > 0$ satisfy (4.35) and $0 < g < 1$ is such that $|\langle t^0, t^\angle\rangle_\mathbb{V}| \leq g\|t^0\|_\mathbb{V}\|t^\angle\|_\mathbb{V}$ for all $t^1 \in \mathbb{V}^1$.*

(ii) *With $\varkappa = \|t_*^\angle\|_\mathbb{V}$, there holds*

$$\varkappa < \frac{2\sqrt{\eta} - \sqrt{2}}{2 - \sqrt{2} + 2\sqrt{\eta}}\delta. \tag{4.37}$$

*Let $D = \{t_*^1 + r^1 \in \mathbb{V}^1 : \|r^0\|_\mathbb{V}^2 + \|r^\angle\|_\mathbb{V}^2 < \frac{1}{2}(\delta - \varkappa)^2\}$. Then, for any $\lambda \in [0, 1)$, there exists $t_{**}^1(\lambda) \in D$ such that $\mathcal{K}_{\mathrm{KP}}(t_{**}^1(\lambda), \lambda) = 0$. Furthermore, $\deg(\mathcal{K}_{\mathrm{KP}}(\cdot, \lambda), D, 0) \equiv d \neq 0$ for all $\lambda \in [0, 1]$. In particular, there exists $t_{**}^1 \in D$ such that $\mathcal{A}(t_{**}^0) = 0$.*

For the proof, see Appendix A.

**Remark 4.35.**

(i) A crucial difference between the linear homotopy and the KP homotopy is that while the decomposition used for $\mathcal{K}_{\mathrm{L}}$ is $\mathbb{V}$-orthogonal, the decomposition for $\mathcal{K}_{\mathrm{KP}}$ is $\ell^2$-orthogonal – the computation (4.30) and Lemma 4.33 exploit this heavily. Nevertheless, we used the $\mathbb{V}$-inner product in the existence result for $\mathcal{K}_{\mathrm{KP}}$. This geometric discrepancy is reflected in condition (i) above.

(ii) Using the Cauchy–Schwarz inequality it is clear that $|\langle t^0, t^\angle\rangle_\mathbb{V}| < \|t^0\|_\mathbb{V}\|t^\angle\|_\mathbb{V}$ for all $t^1 \in \mathbb{V}^1$, due to the fact that $\mathbb{V}^1 = \mathbb{V}^0 \oplus \mathbb{V}^\angle$. The maximum of the function $t^1 \mapsto \frac{\langle t^0, t^\angle\rangle_\mathbb{V}}{\|t^0\|_\mathbb{V}\|t^\angle\|_\mathbb{V}}$ is attained on $\|t^0\|_\mathbb{V} = 1$, $\|t^\angle\|_\mathbb{V} = 1$ and this maximum may be taken as the $0 < g < 1$ of condition (i).

(iii) In the coercive case, *i.e.* when (4.35) holds with $\alpha = 0$, we have $\Theta_0 = I$, so that $\theta_0 = \theta_\angle = 0$. Letting $\varepsilon \to 0$, condition (i) simplifies to

$$\eta = \frac{(1-g)(\gamma_0 - \frac{1}{2}M_\delta\delta)}{\Delta(t_*^1) + M_\delta\delta} > \frac{1}{2}.$$

This last condition in turn reduces to $\gamma_0 > M_\delta\delta$ as $g$ and $\Delta(t_*)$ approaches zero.

(iv) It is interesting to note that while the linear homotopy $\mathcal{K}_{\mathrm{L}}$ involves the "targeted" solution $t_*^1$, the KP homotopy $\mathcal{K}_{\mathrm{KP}}$ does not. In this sense, $\mathcal{K}_{\mathrm{KP}}$ is "universal".

(v) The result is straightforward to extend to the complex case.

(vi) A careful inspection of the proof shows that the result can be generalized to the case when $\mathbb{V}^1 = \mathbb{V}^0 \oplus \mathbb{V}^\angle$ is an arbitrary $\ell^2$-orthogonal direct sum decomposition as long as the homotopy $\mathcal{K}_{\mathrm{KP}}$ is *defined via* (4.33) (*cf.* (4.26)). In this more general case, $\Delta(t_*^1)$ is given by a more complicated expression.

The constant $\Delta(t_*^1)$ also deserves some explanation. Roughly speaking, the "defect" $\Delta(t_*^1)$ measures how much the subspace $\mathfrak{V}^0 \subset \mathfrak{V}$ deviates from being an invariant subspace of the operator $(\mathcal{H}_K + [(T_*^1)_1^\dagger, \mathcal{H}_K])_\mathfrak{V} : \mathfrak{V} \to \mathfrak{V}$. In addition, we can invoke (4.3) to write

$$\Delta(t_*^1) = \max_{\substack{\|u^0\|_\mathbb{V}=1 \\ \|v^\angle\|_\mathbb{V}=1}} \left|\left\langle\left(\mathcal{W}_K + \left[(T_*^1)_1^\dagger, \mathcal{W}_K\right]\right)U^0\Phi_0, V^\angle\Phi_0\right\rangle\right|.$$

Said differently, $\Delta\left(t_*^1\right) = \|(\mathcal{W}_K + [(T_*^1)_1^\dagger, \mathcal{W}_K])_\mathfrak{V}\|_{\mathcal{L}(\mathfrak{V}^0, \mathfrak{V}^\perp)}$. Note that the term $[(T_*^1)_1^\dagger, \mathcal{W}_K]$ can be eliminated *via* orbital rotations (since $(t_*^1)_1 = 0$ can be achieved according to the Thouless theorem), as it involves single excitations only, in which case the amplitude dependence of $\Delta$ is removed. Hence, $\Delta$ quantifies how much the operator $(\mathcal{W}_K)_\mathfrak{V}$ leaves $\mathfrak{V}^0$ invariant.

We summarize the above existence result in the corollary below that holds under the following structural assumptions.

**Assumption** (KPA). $\Delta\left(t_*^1\right)$ *can be made sufficiently small by an appropriate choice of the orbital basis and truncation level* $1 \leq \rho < N$.

**Assumption** (KPB). *There is a* $\delta_0 > 0$ *such that* $M_{\delta_0}$ *is sufficiently small.*

**Corollary 4.36.** *Suppose that* $t_*^1 \in \mathbb{V}^1$ *is a non-degenerate zero of* $\mathcal{A}$ *and that Assumptions (KPA) and (KPB) hold. If* $\|t_*^\angle\|_\mathbb{V}$ *is sufficiently small, then there exists* $t_{**}^1 \in \mathbb{V}^1$ *in a neigborhood of* $t_*^1$ *such that* $\mathcal{K}_{\mathrm{KP}}(t_{**}^1, 0) = 0$.

**Remark 4.37.** Recall that $M_\delta$ only involves the second derivative $\mathcal{A}''$ near $t_*^1$ (see (4.15)), so that (KPB) can be viewed as a "perturbative" assumption. Further, as we noted in Remark 4.9, $M_\delta$ involves the mapping $\zeta \mapsto [[\mathcal{W}_K(t_*^1 + \zeta), \cdot], \cdot]\Phi_0$. Therefore, (KPB) may be viewed as the higher-order generalization of assumption on the smallness of the local Lipschitz constant of the mapping $\zeta \mapsto \mathcal{W}_K(t_*^1 + \zeta)\Phi_0$ in Assumption BII of [36] (see also Rem. 4.41(v)).

In the last step of the *proof* of Theorem 4.34, we could have invoked Theorem 4.28 instead to obtain the existence of a connected component $\mathcal{C}$ of the zero set $\mathcal{Z}(\mathcal{K}_{\mathrm{KP}})$, such that $\mathcal{C} \cap (\mathcal{Z}(\mathcal{K}_{\mathrm{KP}}))_j \neq \emptyset$ for $j = 0, 1$. Under the above assumptions, this provides a theoretical basis for the "solution trajectories" observed in [32]. Further, combining Theorem 4.34 with (4.29), we get for $t_{**}^1 = t_{**}^1(0)$,

$$\sum_{t_{**}^1 \in \mathcal{C} \cap (\mathcal{Z}(\mathcal{K}_{\mathrm{KP}}))_0} i(\mathcal{K}_{\mathrm{KP}}(\cdot, 0), t_{**}^1) = i(\mathcal{A}, t_*^1).$$

Since we do not have uniqueness for $t_{**}^1$ in this case, it is possible that the left-hand side may contain multiple terms which sum up to $i(\mathcal{A}, t_*^1)$. In particular, $i(\mathcal{K}_{\mathrm{KP}}(\cdot, 0), t_{**}^1)$ does not need to be $i(\mathcal{A}, t_*^1)$.

To close this section, we calculate the topological index for a non-degenerate zero at the $\lambda = 0$ endpoint of the KP homotopy. Note that the index at $\lambda = 1$ is simply given by Theorems 4.13 and 4.18.

Fix $t^1 \in \mathbb{V}^1$ and define the operator $\widehat{\mathcal{H}_K}(t^1)$ (not to be confused with (4.10) in a different context),

$$\widehat{\mathcal{H}_K}(t^1) = \mathcal{H}_K(t^1) - \sum_{\alpha \in \Xi(G^0)} \langle \mathcal{H}_K(t^1)\Phi_0, \Phi_\alpha \rangle X_\alpha. \tag{4.38}$$

Notice that $\left\langle \widehat{\mathcal{H}_K}(t^1)\Phi_0, S^0\Phi_0 \right\rangle = 0$ for all $s^0 \in \mathbb{V}^0$. Define the linear mapping $\widehat{\mathcal{H}_K}(t^1)_{\mathfrak{V}^0, \mathfrak{V}^\angle} : \mathfrak{V}^0 \to \mathfrak{V}^\angle$ via $\left\langle \widehat{\mathcal{H}_K}(t^1)_{\mathfrak{V}^0, \mathfrak{V}^\angle}\Psi, \Psi' \right\rangle = \left\langle \widehat{\mathcal{H}_K}(t^1)\Psi, \Psi' \right\rangle$ for all $\Psi \in \mathfrak{V}^0$ and $\Psi' \in \mathfrak{V}^\angle$. We first calculate the derivative of $\mathcal{K}_{\mathrm{KP}}(\cdot, 0)$.

**Lemma 4.38.** *The derivative* $\partial_1\mathcal{K}_{\mathrm{KP}}(t_{**}^1, 0) : \mathbb{V}^1 \to (\mathbb{V}^1)^*$ *of* $\mathcal{K}_{\mathrm{KP}}(\cdot, 0)$ *at a zero* $t_{**}^1$ *is given by*

$$\langle \partial_1\mathcal{K}_{\mathrm{KP}}(t_*^1, 0)u^1, v^1 \rangle = \left\langle \begin{pmatrix} \mathcal{H}_K(t_{**}^0)_{\mathfrak{V}^0} - \mathcal{E}_{\mathrm{CC}}(t_{**}^0) & 0 \\ \widehat{\mathcal{H}_K}(t_{**}^1)_{\mathfrak{V}^0, \mathfrak{V}^\angle} & \widehat{\mathcal{H}_K}(t_{**}^1)_{\mathfrak{V}^\angle} - \mathcal{E}_{\mathrm{CC}}(t_{**}^1) \end{pmatrix} \begin{pmatrix} U^0\Phi_0 \\ U^\angle\Phi_0 \end{pmatrix}, \begin{pmatrix} V^0\Phi_0 \\ V^\angle\Phi_0 \end{pmatrix} \right\rangle$$

*for all* $u^1, v^1 \in \mathbb{V}^1$.

*Proof.* A calculation analogous to the one in the proof of Lemma 4.6 shows that $D(t^1) := \partial_1 \mathcal{K}_{\mathrm{KP}}(t^1, 0) : \mathbb{V}^1 \to (\mathbb{V}^1)^*$ is given by

$$\begin{aligned}
\langle D(t^1) u^1, v^1 \rangle &= \langle [\mathcal{H}_K(t^0), U^0] \Phi_0, V^0 \Phi_0 \rangle + \langle [\mathcal{H}_K(t^1), U^1] \Phi_0, V^{\angle} \Phi_0 \rangle \\
&=: D_1(t^1) + D_2(t^1)
\end{aligned} \tag{4.39}$$

for all $t^1, u^1, v^1 \in \mathbb{V}^1$. We set $t^1 = t^1_{**}$ and evaluate $D_1$ and $D_2$. Write

$$(U^0)^\dagger V^0 \Phi_0 = \sum_{\alpha \in \Xi(G^1) \cup \{0\}} \langle U^0 \Phi_\alpha, V^0 \Phi_0 \rangle \Phi_\alpha,$$

from which,

$$\begin{aligned}
\langle U^0 \mathcal{H}_K(t^0_{**}) \Phi_0, V^0 \Phi_0 \rangle &= \langle \mathcal{H}_K(t^0_{**}) \Phi_0, (U^0)^\dagger V^0 \Phi_0 \rangle \\
&= \left( \sum_{\alpha \in \Xi(G^0)} + \sum_{\alpha \in \Xi(G^0)^c} \right) \langle \mathcal{H}_K(t^0_{**}) \Phi_0, \Phi_\alpha \rangle \langle U^0 \Phi_\alpha, V^0 \Phi_0 \rangle + \mathcal{E}_{\mathrm{CC}}(t^0_{**}) \langle U^0 \Phi_0, V^0 \Phi_0 \rangle.
\end{aligned}$$

Here, the first sum vanishes because of (4.31) and the second by the orthogonality of $\mathbb{V}^0$ and $\mathbb{V}^{\angle}$. Consequently,

$$D_1(t^1_{**}) = \langle (\mathcal{H}_K(t^0_{**}) - \mathcal{E}_{\mathrm{CC}}(t^0_{**})) U^0 \Phi_0, V^0 \Phi_0 \rangle.$$

Analogously, (4.32) implies that

$$\langle U^{\angle} \mathcal{H}_K(t^1_{**}) \Phi_0, V^{\angle} \Phi_0 \rangle = \sum_{\alpha \in \Xi(G^0)} \langle \mathcal{H}_K(t^1_{**}) \Phi_0, \Phi_\alpha \rangle \langle X_\alpha U^{\angle} \Phi_0, V^{\angle} \Phi_0 \rangle + \mathcal{E}_{\mathrm{CC}}(t^1_{**}) \langle U^{\angle} \Phi_0, V^{\angle} \Phi_0 \rangle,$$

and

$$\langle U^0 \mathcal{H}_K(t^1_{**}) \Phi_0, V^{\angle} \Phi_0 \rangle = \sum_{\alpha \in \Xi(G^0)} \langle \mathcal{H}_K(t^1_{**}) \Phi_0, \Phi_\alpha \rangle \langle X_\alpha U^0 \Phi_0, V^{\angle} \Phi_0 \rangle.$$

Thus,

$$D_2(t^1_{**}) = \left\langle \left( \widehat{\mathcal{H}_K}(t^1_{**}) - \mathcal{E}_{\mathrm{CC}}(t^1_{**}) \right) U^{\angle} \Phi_0, V^{\angle} \Phi_0 \right\rangle + \left\langle \widehat{\mathcal{H}_K}(t^1_{**}) U^0 \Phi_0, V^{\angle} \Phi_0 \right\rangle,$$

and the stated expression now follows. $\qquad \square$

The preceding lemma implies the index formula for the KP homotopy.

**Theorem 4.39** (Index formula for KP – non-degenerate case). *Suppose that $t^1_{**} \in \mathbb{V}^1$ is a zero of $\mathcal{K}_{\mathrm{KP}}(\cdot, 0)$. Then $t^1_{**}$ is non-degenerate if and only if the conditions*

(I) $\mathcal{E}_{\mathrm{CC}}(t^0_{**}) \notin \sigma(\mathcal{H}_K(t^0_{**})_{\mathfrak{V}^0})$,
(II) $\mathcal{E}_{\mathrm{CC}}(t^1_{**}) \notin \sigma(\widehat{\mathcal{H}_K}(t^1_{**})_{\mathfrak{V}^{\angle}})$

*both hold true, and in this case the topological index of $\mathcal{K}_{\mathrm{KP}}(\cdot, 0)$ at $t^1_{**}$ is given by $i(\mathcal{K}_{\mathrm{KP}}(\cdot, 0), t^1_{**}) = (-1)^{\nu^0 + \nu^{\angle}}$, where*

$$\begin{aligned}
\nu^0 &= \left| \left\{ j : \mathcal{E}_j(\mathcal{H}_K(t^0_{**})_{\mathfrak{V}^0}) \in \mathbb{R}, \ \mathcal{E}_j(\mathcal{H}_K(t^0_{**})_{\mathfrak{V}^0}) < \mathcal{E}_{\mathrm{CC}}(t^0_{**}) \right\} \right|, \\
\nu^{\angle} &= \left| \left\{ j : \mathcal{E}_j\left( \widehat{\mathcal{H}_K}(t^1_{**})_{\mathfrak{V}^{\angle}} \right) \in \mathbb{R}, \ \mathcal{E}_j\left( \widehat{\mathcal{H}_K}(t^1_{**})_{\mathfrak{V}^{\angle}} \right) < \mathcal{E}_{\mathrm{CC}}(t^1_{**}) \right\} \right|.
\end{aligned}$$

*Proof.* The proof follows from Lemma 4.38 along similar lines as Theorem 4.13,

$$i\left( \mathcal{K}_{\mathrm{KP}}(\cdot, 0), t^1_{**} \right) = \mathrm{sgn} \prod_{j \geq 0} \left( \mathcal{E}_j(\mathcal{H}_K(t^0_{**})_{\mathfrak{V}^0}) - \mathcal{E}_{\mathrm{CC}}(t^0_{**}) \right) \mathrm{sgn} \prod_{j \geq 0} \left( \mathcal{E}_j\left( \widehat{\mathcal{H}_K}(t^1_{**})_{\mathfrak{V}^{\angle}} \right) - \mathcal{E}_{\mathrm{CC}}(t^1_{**}) \right).$$

$\square$

## 4.6. An energy error estimate

In this this section we derive an energy error estimate for general eigenstates for the KP homotopy using the results of Appendix B.

**Theorem 4.40** (Energy error estimate)**.** *Let* $\mathbb{V}^1 = \mathbb{V}(G^{\mathrm{full}})$ *and suppose that* $t_*^1 \in \mathbb{V}^1$ *is a zero of* $\mathcal{A}$*, and that* $t_{**}^1 \in \mathbb{V}^1$ *is a zero of* $\mathcal{K}_{\mathrm{KP}}(\cdot, 0) = 0$*. If the nonorthogonality condition* $\left\langle e^{T_{**}^0} \Phi_0, e^{T_*^1} \Phi_0 \right\rangle \neq 0$ *holds true, then*

$$\left| \mathcal{E}_{\mathrm{CC}}(t_{**}^1) - \mathcal{E}_{\mathrm{CC}}(t_*^1) \right| \leq C(t_{**}^1, t_*^1) \| t_{**}^{\angle} \|_{\mathbb{V}}, \tag{4.40}$$

*where*

$$C(t_{**}^1, t_*^1) = (C^2 + M(t_{**}^1)) \frac{\left\| \Pi_{\mathfrak{Y}^{\angle}} \left( e^{T_{**}^0} \right)^{\dagger} \Pi_{\mathfrak{Y}^{\angle}} e^{T_*^1} \Phi_0 \right\|_{\mathfrak{H}^1}}{\left| \left\langle e^{T_{**}^0} \Phi_0, e^{T_*^1} \Phi_0 \right\rangle \right|}$$

*is bounded as* $\| t_{**}^{\angle} \|_{\mathbb{V}} \to 0$*. Here,* $C$ *is the norm equivalence constant from Remark 2.6 and* $M(t_{**}^1) = \max_{\xi \in [t_{**}^0, t_{**}^1]} \| u^1 \mapsto \Pi_{\mathfrak{Y}^{\angle}} [\mathcal{W}_K(\xi), U^1] \|_{\mathcal{L}(\mathbb{V}, \mathfrak{H}^{-1})}$*.*

*Proof.* Setting $\Psi = e^{T_*^1} \Phi_0$ and $\lambda = 0$ in Theorem B.1(II), we have

$$\left| \mathcal{E}_{\mathrm{CC}}(t_{**}^1) - \mathcal{E}_{\mathrm{CC}}(t_*^1) \right| = \frac{\left| \left\langle (\mathcal{H}_K(t_{**}^1) - \mathcal{H}_K(t_{**}^0)) \Phi_0, \Pi_{\mathfrak{Y}^{\angle}} \left( e^{T_{**}^0} \right)^{\dagger} \Pi_{\mathfrak{Y}^{\angle}} e^{T_*^1} \Phi_0 \right\rangle \right|}{\left| \left\langle e^{T_{**}^0} \Phi_0, e^{T_*^1} \Phi_0 \right\rangle \right|}. \tag{4.41}$$

Note that using (4.3) we may write

$$\begin{aligned} (\mathcal{H}_K(t_{**}^1) - \mathcal{H}_K(t_{**}^0)) \Phi_0 &= [\mathcal{F}_K, T_{**}^{\angle}] \Phi_0 + (\mathcal{W}_K(t_{**}^1) - \mathcal{W}_K(t_{**}^0)) \Phi_0 \\ &= \sum_{\gamma \in \Xi(G^{\angle})} \varepsilon_{\gamma}(t_{**}^{\angle})_{\gamma} \Phi_{\gamma} + (\mathcal{W}_K(t_{**}^1) - \mathcal{W}_K(t_{**}^0)) \Phi_0. \end{aligned}$$

Hence, we can bound the numerator of the right-hand side of (4.41) as

$$\left| \sum_{\gamma \in \Xi(G^{\angle})} \varepsilon_{\gamma}(t_{**}^{\angle})_{\gamma} \left\langle \Phi_{\gamma}, \Pi_{\mathfrak{Y}^{\angle}} \left( e^{T_{**}^0} \right)^{\dagger} \Pi_{\mathfrak{Y}^{\angle}} e^{T_*^1} \Phi_0 \right\rangle \right| + \left| \left\langle (\mathcal{W}_K(t_{**}^1) - \mathcal{W}_K(t_{**}^0)) \Phi_0, \Pi_{\mathfrak{Y}^{\angle}} \left( e^{T_{**}^0} \right)^{\dagger} \Pi_{\mathfrak{Y}^{\angle}} e^{T_*^1} \Phi_0 \right\rangle \right|.$$

Using the Cauchy–Schwarz inequality, the first term may be further bounded as

$$\sum_{\gamma \in \Xi(G^{\angle})} \varepsilon_{\gamma}(t_{**}^{\angle})_{\gamma} \left\langle \Phi_{\gamma}, \Pi_{\mathfrak{Y}^{\angle}} \left( e^{T_{**}^0} \right)^{\dagger} \Pi_{\mathfrak{Y}^{\angle}} e^{T_*^1} \Phi_0 \right\rangle \leq \|\! \| t_{**}^{\angle} \|\! \| \left\| \Pi_{\mathfrak{Y}^{\angle}} \left( e^{T_{**}^0} \right)^{\dagger} \Pi_{\mathfrak{Y}^{\angle}} e^{T_*^1} \Phi_0 \right\|,$$

where the Fock norm $\|\! \| \cdot \|\! \|$ was defined in Remark 2.6. For the second term, we use the intermediate value inequality to obtain the bound $M(t_{**}^1) \| t_{**}^{\angle} \|_{\mathbb{V}} \| \Pi_{\mathfrak{Y}^{\angle}} (e^{T_{**}^0})^{\dagger} \Pi_{\mathfrak{Y}^{\angle}} e^{T_*^1} \Phi_0 \|_{\mathfrak{H}^1}$. $\square$

**Remark 4.41.**

(i) If the nonorthogonality condition holds and $t_{**}^{\angle} = 0$, then according to (4.31) and (4.32), $t_{**}^0$ is an FCC solution such that $\left\langle e^{T_{**}^0} \Phi_0, e^{T_*^1} \Phi_0 \right\rangle \neq 0$. In this case, (4.40) implies that the energy error is zero: $\mathcal{E}_{\mathrm{CC}}(t_{**}^0) = \mathcal{E}_{\mathrm{CC}}(t_*^1)$.

(ii) In the practically relevant case $\rho \geq 2$, using (B.1), we have that $\mathcal{E}_{\mathrm{CC}}(t_{**}^1) = \mathcal{E}_{\mathrm{CC}}(t_{**}^0)$ so the left-hand side of (4.40) does not involve $t_{**}^{\angle}$ at all.

(iii) The appearance of the quantity $\|t_{**}^{\angle}\|_{\mathbb{V}}$ allows us to view the auxiliary equation (4.32) as providing an *a posteriori* error estimate.

(iv) If the nonorthogonality condition does not hold, *i.e.* $\left\langle e^{T_{**}^0}\Phi_0, e^{T_*^1}\Phi_0 \right\rangle = 0$, then $e^{T_{**}^0}\Phi_0$ and $e^{T_*^1}\Phi_0$ represent different eigenstates if $t_*^1$ is assumed to be non-degenerate. While $e^{T_{**}^0}\Phi_0$ itself does not satisfy the Schrödinger equation, it must be viewed as an approximation to an eigenstate *different* from $e^{T_*^1}\Phi_0$.

(v) Note that a local Lipschitz assumption (on a ball including $t_{**}^1$ and $t_{**}^0$) with constant $L$ on the mapping $t \mapsto \mathcal{W}_K(t)\Phi_0$ can be used to obtain the result of the theorem with $M$ replaced by $L$. Such an assumption is akin to Assumption B.II in [36] where it is used for guaranteeing the local strong monotonicity of the CC mapping.

## 5. Conclusions and further work

In this second part of a series of two articles, we analyzed the SRCC method. We provided background material for the setting of the quantum-mechanical problems the SRCC method is aimed at in Section 2. Then, in Section 3, we gave a brief summary of topological degree theory since this served as our main tool for the analysis.

The main discussion is contained in Section 4, where we began with the definition of the SRCC mapping and some elementary considerations in Section 4.1. We then considered the local properties of the SRCC mapping in Section 4.2. It turned out that the topological index of the CC mapping (Thm. 4.13) is connected with the nonvariational property of the CC method (Rem. 4.16), and the eigenvalues of the Fock operator (Prop. 4.17). In the degenerate case, the classic Leray reduction formula provided the topological index (Thm. 4.18). We also discussed the case when the cluster amplitudes are allowed to be complex in Section 4.3.

In Section 4.4, we discussed how certain homotopies can be used to analyze the CC method, in particular to prove the existence of a truncated CC solution through the use of topological degree theory. This was done using an idea well-known both in nonlinear analysis and in quantum chemistry: that an appropriate homotopy "connects" the truncated problem with the exact problem (essentially the Schrödinger equation), therefore one is able to infer (homotopy-invariant) information regarding the former problem from the latter. As an introductory example, we considered the linear homotopy (Thm. 4.31).

Next, in Section 4.5, motivated by the works of the chemists Kowalski and Piecuch, we considered a homotopy that connects CC mappings corresponding to different truncation levels. Using this, we proved an existence result for the said homotopy (Thm. 4.34), which also implies the existence of a truncated CC solution under certain assumptions. The index formula for the KP homotopy was also derived in the non-degenerate case (Thm. 4.39). Using a known result about the KP homotopy (Appendix B), we also derived an energy error estimate in Section 4.6.

Finally, let us discuss some possible directions of research. Clearly, it would be interesting to extend our analysis to the infinite-dimensional case. An obvious next step would be the analysis of the JM-MRCC method (see Sect. 4.2 of Part I). It would be also interesting to look at the Extended CC (ECC) [1, 24] and the Unitary CC (UCC) methods [3].

## Appendix A. Proof of Theorem 4.34

We first prove that $\mathcal{K}_{\mathrm{KP}}(\cdot, \lambda) \neq 0$ on $\partial D$ for all $\lambda \in [0, 1]$. Set $t^1 = t_*^1 + r^1$ and $s^1 = \Theta_\alpha r^1$ in (4.33) where $r^1 \in \partial D$, and write

$$\begin{aligned}
\left\langle \mathcal{K}_{\mathrm{KP}}(t_*^1 + r^1, \lambda), \Theta_\alpha r^1 \right\rangle &= \left\langle \mathcal{A}(t_*^0 + r^0 + \lambda t_*^{\angle} + \lambda r^{\angle}), \Pi_0 \Theta_\alpha r^1 \right\rangle + \left\langle \mathcal{A}(t_*^1 + r^1), \Pi_{\angle}\Theta_\alpha r^1 \right\rangle \\
&= \left\langle \mathcal{A}(t_*^0 + \lambda t_*^{\angle} + r^0 + \lambda r^{\angle}) - \mathcal{A}(t_*^1 + r^1), \Pi_0\Theta_\alpha r^1 \right\rangle + \left\langle \mathcal{A}(t_*^1 + r^1), \Theta_\alpha r^1 \right\rangle \\
&=: (\mathrm{I}) + (\mathrm{II}).
\end{aligned}$$

For (I), Taylor expansion around $t_*^1$ leads to

$$(\mathrm{I}) = (\lambda - 1)\langle \mathcal{A}'(t_*^1)(t_*^\angle + r^\angle), \Pi_0 \Theta_\alpha r^1 \rangle + \langle \mathcal{R}_2(t_*^1, (\lambda-1)t_*^\angle + r^0 + \lambda r^\angle) - \mathcal{R}_2(t_*^1, r^1), \Pi_0 \Theta_\alpha r^1 \rangle$$
$$\geq (\lambda - 1)\langle \mathcal{A}'(t_*^1)(t_*^\angle + r^\angle), \Pi_0 \Theta_\alpha r^1 \rangle - M\|t_*^\angle + r^\angle\|\|\Pi_0 \Theta_\alpha r^1\|$$
$$\geq -(\Delta(t_*^1) + M)\|t_*^\angle + r^\angle\|_{\mathbb{V}}\|\Pi_0 \Theta_\alpha r^1\|_{\mathbb{V}},$$

where we used the intermediate value inequality. Here, letting $\Psi^0 \in \mathfrak{V}^0$ correspond to the amplitude $\Pi_0 \Theta_\alpha r^1$,

$$\langle \mathcal{A}'(t_*^1)(t_*^\angle + r^\angle), \Pi_0 \Theta_\alpha r^1 \rangle = \langle (\mathcal{H}_K(t_*^1) - \mathcal{E}_{\mathrm{CC}}(t_*^1))(T_*^\angle + R^\angle)\Phi_0, \Psi^0 \rangle = \langle \mathcal{H}_K(t_*^1)(T_*^\angle + R^\angle)\Phi_0, \Psi^0 \rangle$$
$$= \langle (\mathcal{H}_K + [\mathcal{H}_K, T_*^1] + \ldots)(T_*^\angle + R^\angle)\Phi_0, \Psi^0 \rangle \leq \Delta(t_*^1)\|t_*^\angle + r^\angle\|_{\mathbb{V}}\|\Pi_0 \Theta_\alpha r^1\|_{\mathbb{V}},$$

where the first equality is (4.11), and in the last inequality we exploited that $\mathcal{H}_K$ decreases the excitation rank at most by 2 (so that the single amplitudes $(t_*^1)_1$ of $t_*^1$ only contribute). Also, the following estimate holds,

$$M = \max_{\xi \in [(\lambda-1)t_*^\angle + r^0 + \lambda r^\angle, r^1]} \|\partial_2 \mathcal{R}_2(t_*^1, \xi)\|_{\mathcal{L}(\mathbb{V}, \mathbb{V}^*)}$$
$$= \max_{\xi \in [(\lambda-1)t_*^\angle + r^0 + \lambda r^\angle, r^1]} \|\mathcal{A}'(t_* + \xi) - \mathcal{A}'(t_*)\|_{\mathcal{L}(\mathbb{V}, \mathbb{V}^*)}$$
$$\leq \max_{\xi \in [(\lambda-1)t_*^\angle + r^0 + \lambda r^\angle, r^1]} \|\xi\|_{\mathbb{V}} \max_{\zeta \in [0, \xi]} \|\mathcal{A}''(t_*^1 + \zeta)\|_{\mathcal{L}(\mathbb{V} \times \mathbb{V}, \mathbb{V}^*)}$$
$$\leq M_\delta \max_{\xi \in [(\lambda-1)t_*^\angle + r^0 + \lambda r^\angle, r^1]} \|\xi\|_{\mathbb{V}} \leq M_\delta(\|r^0\|_{\mathbb{V}} + \|r^\angle\|_{\mathbb{V}} + \varkappa) \leq M_\delta \delta.$$

For (II), we have using (4.35), (4.36) and Taylor's theorem,

$$(\mathrm{II}) = \langle \mathcal{A}'(t_*^1)r^1, \Theta_\alpha r^1 \rangle - \langle \mathcal{R}_2(t_*^1, r^1), \Theta_\alpha r^1 \rangle \geq \gamma_\alpha \|r^1\|_{\mathbb{V}}^2 - \tfrac{1}{2}M_\delta\|r^1\|_{\mathbb{V}}^2\|\Theta_\alpha r^1\|_{\mathbb{V}}$$
$$\geq (\gamma_\alpha - \tfrac{1}{2}M_\delta\|\Theta_\alpha\|_{\mathcal{L}(\mathbb{V})}\|r^1\|_{\mathbb{V}})\|r^1\|_{\mathbb{V}}^2 \geq (\gamma_\alpha - \tfrac{1}{2}M_\delta\|\Theta_\alpha\|_{\mathcal{L}(\mathbb{V})}\delta)\|r^1\|_{\mathbb{V}}^2.$$

In summary, using the definitions of $\theta_0$ and $\theta_\angle$, and setting $\widetilde{\gamma} = \gamma_\alpha - \tfrac{1}{2}M_\delta\|\Theta_\alpha\|_{\mathcal{L}(\mathbb{V})}\delta$,

$$\langle \mathcal{K}_{\mathrm{KP}}(t_*^1 + r^1, \lambda), \Theta_\alpha r^1 \rangle$$
$$\geq \widetilde{\gamma}\|r^1\|_{\mathbb{V}}^2 - (\Delta(t_*^1) + M_\delta\delta)\|t_*^\angle + r^\angle\|_{\mathbb{V}}\|\Pi_0\Theta_\alpha r^1\|_{\mathbb{V}}$$
$$\geq \widetilde{\gamma}\|r^1\|_{\mathbb{V}}^2 - \frac{\Delta(t_*^1) + M_\delta\delta}{2}(\|t_*^\angle + r^\angle\|_{\mathbb{V}}^2 + \|\Pi_0\Theta_\alpha r^1\|_{\mathbb{V}}^2)$$
$$\geq (1-g)\widetilde{\gamma}(\|r^0\|_{\mathbb{V}}^2 + \|r^\angle\|_{\mathbb{V}}^2) - \frac{\Delta(t_*^1) + M_\delta\delta}{2}(\varkappa^2 + 2\varkappa\|r^\angle\|_{\mathbb{V}} + \|r^\angle\|_{\mathbb{V}}^2 + \|\Pi_0\Theta_\alpha r^1\|_{\mathbb{V}}^2)$$
$$\geq \left((1-g)\widetilde{\gamma} - \frac{\Delta(t_*^1) + M_\delta\delta}{2}(1 + \max\{\varepsilon + 2(1+\varepsilon^{-1})\theta_0, 2(1+\varepsilon^{-1})\theta_\angle\})\right)$$
$$\times (\|r^0\|_{\mathbb{V}}^2 + \|r^\angle\|_{\mathbb{V}}^2) - \frac{\Delta(t_*^1) + M_\delta\delta}{2}(\varkappa^2 + \sqrt{2}\varkappa(\delta - \varkappa))$$
$$\geq \left((1-g)\widetilde{\gamma} - \frac{\Delta(t_*^1) + M_\delta\delta}{2}\max\{\varepsilon + 2(1+\varepsilon^{-1})\theta_0, 2(1+\varepsilon^{-1})\theta_\angle\}\right)\frac{(\delta - \varkappa)^2}{2}$$
$$- \frac{\Delta(t_*^1) + M_\delta\delta}{2}\left(\varkappa + \frac{\sqrt{2}}{2}(\delta - \varkappa)\right)^2.$$

The positivity of the last expression follows from (4.37). We also used the bound

$$\|r^\angle\|_{\mathbb{V}}^2 + \|\Pi_0\Theta_\alpha r^1\|_{\mathbb{V}}^2 \leq (1 + \varepsilon^{-1})\|\Pi_0(\Theta_\alpha - I)r^1\|_{\mathbb{V}}^2 + (1 + \varepsilon)\|r^0\|_{\mathbb{V}}^2 + \|r^\angle\|_{\mathbb{V}}^2$$

$$\leq 2\big(1+\varepsilon^{-1}\big)(\|\Pi_0(\Theta_\alpha-I)\Pi_0\|_{\mathcal{L}(\mathbb{V})}^2\|r^0\|_{\mathbb{V}}^2+\|\Pi_0(\Theta_\alpha-I)\Pi_\angle\|_{\mathcal{L}(\mathbb{V})}^2\|r^\angle\|_{\mathbb{V}}^2)+(1+\varepsilon)\|r^0\|_{\mathbb{V}}^2+\|r^\angle\|_{\mathbb{V}}^2$$
$$\leq (1+\max\{\varepsilon+2\big(1+\varepsilon^{-1}\big)\theta_0, 2\big(1+\varepsilon^{-1}\big)\theta_\angle\})(\|r^0\|_{\mathbb{V}}^2+\|r^\angle\|_{\mathbb{V}}^2).$$

We can now apply the homotopy invariance of the degree to get $\deg(\mathcal{K}_{\mathrm{KP}}(\cdot,\lambda),D,0)\equiv d\neq0$ for all $\lambda\in[0,1]$, with $\delta$ decreased if necessary.

## APPENDIX B. A SHORT PROOF OF THE KOWALSKI–PIECUCH THEOREM

To close this section, we present the main theorem[10] of Kowalski and Piecuch [32] in a somewhat different form. Define the *KP energy* as

$$\mathcal{E}_{\mathrm{KP}}(t^1,\lambda)=\big\langle\mathcal{H}_K(t^0+\lambda t^\angle)\Phi_0,\Phi_0\big\rangle=\big\langle\mathcal{H}_K e^{T^0+\lambda T^\angle}\Phi_0,\Phi_0\big\rangle,$$

so that $\mathcal{E}_{\mathrm{KP}}(t^1,0)=\mathcal{E}_{\mathrm{CC}}(t^0)$ and $\mathcal{E}_{\mathrm{KP}}(t^1,1)=\mathcal{E}_{\mathrm{CC}}(t^1)$. If $\rho\geq2$, then according to (2.10) of Part I, we simply have

$$\mathcal{E}_{\mathrm{KP}}(t^1,\lambda)=\big\langle\mathcal{H}_K(t^0)\Phi_0,\Phi_0\big\rangle=\big\langle\mathcal{H}_K e^{T^0}\Phi_0,\Phi_0\big\rangle,\tag{B.1}$$

although we will not exploit this property in the proof.

**Theorem B.1** (Kowalski–Piecuch)**.**

(I) *Suppose that* $\Psi=(c_0I+C^0)\Phi_0\in\mathfrak{H}_K^1$, *with* $c_0\in\mathbb{R}$ *and* $c^0\in\mathbb{V}^0$, *satisfies the (weak) Schrödinger equation*

$$\langle\mathcal{H}_K\Psi,\Phi\rangle=\mathcal{E}\langle\Psi,\Phi\rangle\quad\text{for all}\quad\Phi\in\mathfrak{H}_K^1\tag{B.2}$$

*for some* $\mathcal{E}\in\mathbb{R}$ *and that*

$$\big\langle e^{T_{**}^0(\lambda)+\lambda T_{**}^\angle(\lambda)}\Phi_0,\Psi\big\rangle\neq0,\tag{B.3}$$

*where* $t_{**}^1(\lambda)\in\mathbb{V}^1$ *is a zero of* $\mathcal{K}_{\mathrm{KP}}(\cdot,\lambda)$ *for all* $\lambda\in[0,1]$. *Then*

$$\mathcal{E}_{\mathrm{KP}}(t_{**}^1(\lambda),\lambda)\equiv\mathcal{E}\quad\text{for all}\quad\lambda\in[0,1].$$

(II) *Suppose that* $\Psi=(c_0I+C^1)\Phi_0$, *with* $c_0\in\mathbb{R}$ *and* $c^1\in\mathbb{V}^1$, *satisfies* (B.2) *for some* $\mathcal{E}\in\mathbb{R}$. *If* (B.3) *holds true for* some $\lambda\in[0,1]$ *and* $t_{**}^1=t_{**}^1(\lambda)\in\mathbb{V}^1$ *with* $\mathcal{K}_{\mathrm{KP}}(t_{**}^1,\lambda)=0$, *then*

$$\mathcal{E}_{\mathrm{KP}}\big(t_{**}^1(\lambda),\lambda\big)-\mathcal{E}=\frac{\big\langle\big(\mathcal{H}_K\big(t_{**}^1\big)-\mathcal{H}_K\big(t_{**}^0+\lambda t_{**}^\angle\big)\big)\Phi_0,\Pi_{\mathfrak{Y}^\angle}\big(e^{T_{**}^0+\lambda T_{**}^\angle}\big)^\dagger C^\angle\Phi_0\big\rangle}{\big\langle e^{T_{**}^0+\lambda T_{**}^\angle}\Phi_0,\Psi\big\rangle}$$
$$=(1-\lambda)\frac{\big\langle\Gamma\big(t_{**}^1,\lambda\big)\Phi_0,\Pi_{\mathfrak{Y}^\angle}\big(e^{T_{**}^0+\lambda T_{**}^\angle}\big)^\dagger C^\angle\Phi_0\big\rangle}{\big\langle e^{T_{**}^0+\lambda T_{**}^\angle}\Phi_0,\Psi\big\rangle},$$

*where* $\Gamma(t^1,\lambda)$ *is given by* (B.6) *below.*

Furthermore, in case the energy blows up, we have the following.

**Theorem B.2.** *Suppose that* $\Psi=(c_0I+C^1)\Phi_0$, *with* $c_0\in\mathbb{R}$ *and* $c^1\in\mathbb{V}^1$, *satisfies* (B.2) *for some* $\mathcal{E}\in\mathbb{R}$. *Furthermore, assume that* $t_{**}^1(\lambda)\in\mathbb{V}^1$ *is a zero of* $\mathcal{K}_{\mathrm{KP}}(\cdot,\lambda)$ *for all* $\lambda$ *in a neighborhood of some* $\lambda_0\in[0,1]$. *If* $|\mathcal{E}_{\mathrm{KP}}(t_{**}^1(\lambda),\lambda)|\to\infty$ *as* $\lambda\to\lambda_0$ *and*

$$\big\langle\Gamma\big(t_{**}^1(\lambda),\lambda\big)\Phi_0,\Pi_{\mathfrak{Y}^\angle}\big(e^{T_{**}^0+\lambda T_{**}^\angle}\big)^\dagger C^\angle\Phi_0\big\rangle=\mathcal{O}\Big(\frac{1}{1-\lambda}\Big)\quad(\lambda\to\lambda_0),\tag{B.4}$$

*then*

$$\big\langle e^{T_{**}^0(\lambda)+\lambda T_{**}^\angle(\lambda)}\Phi_0,\Psi\big\rangle\to0\quad(\lambda\to\lambda_0).$$

---

[10]They call it the "Fundamental Theorem of $\beta$-Nested Equation Formalism".

Part (I) of Theorem B.1 says that if one can solve the Schrödinger equation *exactly* on $\mathbb{V}^0$ for an eigenvalue $\mathcal{E}$ and (B.3) holds true for a zero $t_{**}^1(\lambda) \in \mathbb{V}^1$ of $\mathcal{K}_{\mathrm{KP}}(\cdot, \lambda)$ for all $\lambda$, then the KP energy $\mathcal{E}_{\mathrm{KP}}(t_{**}^1(\lambda), \lambda)$ is identically $\mathcal{E}$. Notice that $t_{**}^1(1)$ is not required to represent $\Psi$, *i.e.* $e^{T_{**}^1(1)}\Phi_0 \neq \Psi$ is allowed. Also, no regularity of the trajectory $\lambda \mapsto t_{**}^1(\lambda)$ is demanded.

Part (II) stipulates that the Schrödinger equation can be solved on $\mathbb{V}^1$ with an eigenvalue $\mathcal{E}$ and that the nonorthogonality condition (B.3) holds true for some $t_{**}^1 \in \mathbb{V}^1$ zero of $\mathcal{K}_{\mathrm{KP}}(\cdot, \lambda)$ for some $\lambda$. Then, the error in the energy can be expressed by the stated formula. If one assumes the hypothesis for all $\lambda \in [0,1]$ in a neighborhood of 1, then we can conclude that the KP energy $\mathcal{E}_{\mathrm{KP}}(t_{**}^1(\lambda), \lambda)$ tends to $\mathcal{E}$ smoothly, as $\lambda \to 1$. Again, no regularity of $\lambda \mapsto t_{**}^1(\lambda)$ is needed.

Finally, Theorem B.2 considers the case when the KP energy diverges as $\lambda \to \lambda_0$. Assuming the growth condition (B.4), we can conclude that the KP solution $e^{T_{**}^0(\lambda) + \lambda T_{**}^{\angle}(\lambda)}\Phi_0$ becomes orthogonal to the eigenstate $\Psi$.

For the proof, we need the following lemma which recasts the KP equations in an "unlinked" form.

**Lemma B.3.** *Suppose that $t_{**}^1 = t_{**}^1(\lambda) \in \mathbb{V}^1$ is such that $\mathcal{K}_{\mathrm{KP}}(t_{**}^1, \lambda) = 0$ for some $\lambda \in [0,1]$. Then,*

$$\Big\langle \big(\mathcal{H}_K - \mathcal{E}_{\mathrm{KP}}(t_{**}^1, \lambda)\big)e^{T_{**}^0 + \lambda T_{**}^{\angle}}\Phi_0, S^1\Phi_0 \Big\rangle = -\Big\langle \mathcal{G}(t_{**}^1, \lambda)\Phi_0, \Pi_{\mathfrak{V}^{\angle}}\big(e^{T_{**}^0 + \lambda T_{**}^{\angle}}\big)^{\dagger}S^{\angle}\Phi_0 \Big\rangle, \tag{B.5}$$

*where*

$$\mathcal{G}(t_{**}^1, \lambda) = \mathcal{H}_K(t_{**}^1) - \mathcal{H}_K(t_{**}^0 + \lambda t_{**}^{\angle}).$$

*Moreover, $\mathcal{G}(t_{**}^1, \lambda) = (1 - \lambda)\Gamma(t_{**}^1, \lambda)$, where*

$$\Gamma(t_{**}^1, \lambda) = \sum_{k=1}^{2N} \frac{(1 - \lambda)^{k-1}}{k!} e^{-(T_{**}^0 + \lambda T_{**}^{\angle})}\big[\mathcal{H}_K, T_{**}^{\angle}\big]_{(k)} e^{T_{**}^0 + \lambda T_{**}^{\angle}}. \tag{B.6}$$

*Proof.* Assume that $\mathcal{K}_{\mathrm{KP}}(t_{**}^1, \lambda) = 0$, so using (4.33) we have

$$\big\langle \mathcal{H}_K(t_{**}^0 + \lambda t_{**}^{\angle})\Phi_0, S^0\Phi_0 \big\rangle + \big\langle \mathcal{H}_K(t_{**}^1)\Phi_0, S^{\angle}\Phi_0 \big\rangle = 0 \quad \text{for all} \quad s^1 \in \mathbb{V}^1.$$

This can be rewritten as

$$\big\langle \mathcal{H}_K(t_{**}^0 + \lambda t_{**}^{\angle})\Phi_0, S^1\Phi_0 \big\rangle = -\big\langle \big(\mathcal{H}_K(t_{**}^1) - \mathcal{H}_K(t_{**}^0 + \lambda t_{**}^{\angle})\big)\Phi_0, S^{\angle}\Phi_0 \big\rangle \quad \text{for all} \quad s^1 \in \mathbb{V}^1,$$

or,

$$\big\langle \mathcal{H}_K(t_{**}^0 + \lambda t_{**}^{\angle})\Phi_0, S^1\Phi_0 \big\rangle = -\big\langle \mathcal{G}(t_{**}^1, \lambda)\Phi_0, S^{\angle}\Phi_0 \big\rangle \quad \text{for all} \quad s^1 \in \mathbb{V}^1.$$

Then we can write

$$\Big\langle \big(\mathcal{H}_K - \mathcal{E}_{\mathrm{KP}}(t_{**}^1, \lambda)\big)e^{T_{**}^0 + \lambda T_{**}^{\angle}}\Phi_0, S^1\Phi_0 \Big\rangle = \Big\langle e^{-(T_{**}^0 + \lambda T_{**}^{\angle})}\big(\mathcal{H}_K - \mathcal{E}_{\mathrm{KP}}(t_{**}^1, \lambda)\big)e^{T_{**}^0 + \lambda T_{**}^{\angle}}\Phi_0, \big(e^{T_{**}^0 + \lambda T_{**}^{\angle}}\big)^{\dagger}S^1\Phi_0 \Big\rangle$$

$$= \Big\langle e^{-(T_{**}^0 + \lambda T_{**}^{\angle})}\mathcal{H}_K e^{T_{**}^0 + \lambda T_{**}^{\angle}}\Phi_0, \Pi_{\mathfrak{V}^1}\big(e^{T_{**}^0 + \lambda T_{**}^{\angle}}\big)^{\dagger}S^1\Phi_0 \Big\rangle$$

$$= -\Big\langle \mathcal{G}(t_{**}^1, \lambda)\Phi_0, \Pi_{\mathfrak{V}^{\angle}}\big(e^{T_{**}^0 + \lambda T_{**}^{\angle}}\big)^{\dagger}S^1\Phi_0 \Big\rangle$$

$$= -\Big\langle \mathcal{G}(t_{**}^1, \lambda)\Phi_0, \Pi_{\mathfrak{V}^{\angle}}\big(e^{T_{**}^0 + \lambda T_{**}^{\angle}}\big)^{\dagger}S^{\angle}\Phi_0 \Big\rangle$$

for all $s^1 \in \mathbb{V}^1$. In the last step we used that $\Pi_{\mathfrak{V}^{\angle}}(e^{T_{**}^0 + \lambda T_{**}^{\angle}})^{\dagger}$ maps $\mathfrak{V}^0$ to zero. The "moreover" part is a simple expansion using (4.9),

$$\mathcal{G}(t_{**}^1, \lambda) = e^{-(T_{**}^0 + \lambda T_{**}^{\angle})}\Big(e^{-(1-\lambda)T_{**}^{\angle}}\mathcal{H}_K e^{(1-\lambda)T_{**}^{\angle}} - \mathcal{H}_K\Big)e^{T_{**}^0 + \lambda T_{**}^{\angle}}$$

$$= \sum_{k=1}^{2N} \frac{(1-\lambda)^k}{k!} e^{-(T_{**}^0 + \lambda T_{**}^{\angle})}\big[\mathcal{H}_K, T_{**}^{\angle}\big]_{(k)} e^{T_{**}^0 + \lambda T_{**}^{\angle}}.$$

$\square$

*Proof of Theorem* (B.1)*.* We have using Lemma B.3,

$$
\begin{aligned}
0 &= \left\langle \left(\mathcal{H}_K - \mathcal{E}_{\mathrm{KP}}\big(t^1_{**}(\lambda), \lambda\big)\right) e^{T^0_{**}(\lambda) + \lambda T^{\angle}_{**}(\lambda)} \Phi_0, \big(c_0 I + C^0\big) \Phi_0 \right\rangle \\
&= \big(\mathcal{E} - \mathcal{E}_{\mathrm{KP}}\big(t^1_{**}(\lambda), \lambda\big)\big) \left\langle e^{T^0_{**}(\lambda) + \lambda T^{\angle}_{**}(\lambda)} \Phi_0, \Psi \right\rangle.
\end{aligned}
$$

Similarly, for part (II)

$$
\begin{aligned}
&\left\langle \left(\mathcal{H}_K\big(t^1_*\big) - \mathcal{H}_K\big(t^0_{**} + \lambda t^{\angle}_{**}\big)\right) \Phi_0, \Pi_{\mathfrak{V}^{\angle}} \left(e^{T^0_{**} + \lambda T^{\angle}_{**}}\right)^{\dagger} C^{\angle} \Phi_0 \right\rangle \\
&= \left\langle \left(\mathcal{H}_K - \mathcal{E}_{\mathrm{KP}}\big(t^1_{**}(\lambda), \lambda\big)\right) e^{T^0_{**} + \lambda T^{\angle}_{**}} \Phi_0, \big(c_0 I + C^0 + C^{\angle}\big) \Phi_0 \right\rangle \\
&= \big(\mathcal{E} - \mathcal{E}_{\mathrm{KP}}\big(t^1_{**}(\lambda), \lambda\big)\big) \left\langle e^{T^0_{**} + \lambda T^{\angle}_{**}} \Phi_0, \Psi \right\rangle.
\end{aligned}
$$

Theorem B.2 also follows from the previous equality. $\square$

## References

[1] J.S. Arponen, Variational principles and linked-cluster exp S expansions for static and dynamic many-body problems. *Ann. Phys.* **151** (1983) 311–382.

[2] V. Bach, E.H. Lieb, M. Loss and J.P. Solovej, There are no unfilled shells in unrestricted Hartree–Fock theory, in The Stability of Matter: From Atoms to Stars. Springer (1997) 309–311.

[3] R.J. Bartlett, S.A. Kucharski and J. Noga, Alternative coupled-cluster ansätze II. The unitary coupled-cluster method. *Chem. Phys. Lett.* **155** (1989) 133–140.

[4] N. Benedikter, Hartree–Fock theory. Online lecture notes (2017).

[5] A. Buffa, Remarks on the discretization of some noncoercive operator with applications to heterogeneous Maxwell equations. *SIAM J. Numer. Anal.* **43** (2005) 1–18.

[6] A. Buffa, R. Hiptmair, T. von Petersdorff and C. Schwab, Boundary element methods for Maxwell transmission problems in lipschitz domains. *Numer. Math.* **95** (2003) 459–485.

[7] E. Cances, M. Defranceschi, W. Kutzelnigg, C. Le Bris and Y. Maday, Computational quantum chemistry: a primer. *Handb. Numer. Anal.* **10** (2003) 3–270.

[8] J. Cronin, Fixed Points and Topological Degree in Nonlinear Analysis. Vol. 11. American Mathematical Society (1995).

[9] G. Dinca and J. Mawhin, Brouwer degree and applications. Preprint (2009).

[10] G. Dinca and J. Mawhin, Brouwer Degree (The Core of Nonlinear Analysis). Birkhäuser Basel (2021).

[11] P. Drábek and J. Milota, Methods of Nonlinear Analysis: Applications to Differential Equations. Springer Science & Business Media (2007).

[12] F.M. Faulstich, A. Laestadius, O. Legeza, R. Schneider and S. Kvaal, Analysis of the tailored coupled-cluster method in quantum chemistry. *SIAM J. Numer. Anal.* **57** (2019) 2579–2607.

[13] F.M. Faulstich, M. Máté, A. Laestadius, M.A. Csirik, L. Veis, A. Antalik, J. Brabec, R. Schneider, J. Pittner, S. Kvaal and O. Legeza, Numerical and theoretical aspects of the DMRG-TCC method exemplified by the nitrogen dimer. *J. Chem. Theory Comput.* **15** (2019) 2206–2220.

[14] G. Friesecke, The multiconfiguration equations for atoms and molecules: charge quantization and existence of solutions. *Arch. Ration. Mech. Anal.* **169** (2003) 35–71.

[15] J. Geertsen, M. Rittby and R.J. Bartlett, The equation-of-motion coupled-cluster method: excitation energies of Be and CO. *Chem. Phys. Lett.* **164** (1989) 57–62.

[16] T. Helgaker, P. Jorgensen and J. Olsen, Molecular Electronic-Structure Theory. John Wiley & Sons (2014).

[17] K. Jankowski and K. Kowalski, Physical and mathematical content of coupled-cluster equations. II. On the origin of irregular solutions and their elimination via symmetry adaptation. *J. Chem. Phys.* **110** (1999) 9345–9352.

[18] K. Jankowski and K. Kowalski, Physical and mathematical content of coupled-cluster equations. IV. Impact of approximations to the cluster operator on the structure of solutions. *J. Chem. Phys.* **111** (1999) 2952–2959.

[19] K. Jankowski, K. Kowalski and P. Jankowski, Multiple solutions of the single-reference coupled-cluster equations. I. H4 model revisited. *Int. J. Quantum Chem.* **50** (1994) 353–367.

[20] K. Jankowski, K. Kowalski, I. Grabowski and H. Monkhorst, Correspondence between physical states and solutions to the coupled-cluster equations. *Int. J. Quantum Chem.* **75** (1999) 483–496.

[21] F. Kossoski, A. Marie, A. Scemama, M. Caffarel and P.-F. Loos, Excited states from state specific orbital optimized pair coupled cluster. *J. Chem. Theory Comput.* **17** (2021) 4756–4768.

[22] K. Kowalski and K. Jankowski, Full solution to the coupled-cluster equations: the H4 model. *Chem. Phys. Lett.* **290** (1998) 180–188.

[23] A. Laestadius and F.M. Faulstich, The coupled-cluster formalism – a mathematical perspective. *Mol. Phys.* **117** (2019) 2362–2373.

[24] A. Laestadius and S. Kvaal, Analysis of the extended coupled-cluster method in quantum chemistry. *SIAM J. Numer. Anal.* **56** (2018) 660–683.

[25] M. Lewin, Geometric methods for nonlinear many-body quantum systems. *J. Funct. Anal.* **260** (2011) 3535–3595.

[26] M. Lewin, Existence of Hartree–Fock excited states for atoms and molecules. *Lett. Math. Phys.* **108** (2018) 985–1006.

[27] E.H. Lieb and B. Simon, On solutions to the Hartree–Fock problem for atoms and molecules. *Journal Chem. Phys.* **61** (1974) 735–736.

[28] P.-L. Lions, Solutions of Hartree–Fock equations for Coulomb systems. *Commun. Math. Phys.* **109** (1987) 33–97.

[29] D. O'Regan, Y.J. Cho and Y.-Q. Chen, Topological Degree Theory and Applications. CRC Press (2006).

[30] J. Paldus, M. Takahashi and B. Cho, Degeneracy and coupled-cluster approaches. *Int. J. Quantum Chem.* **26** (1984) 237–244.

[31] W.V. Petryshyn, Approximation-Solvability of Nonlinear Functional and Differential Equations. Vol. 171. CRC Press (1992).

[32] P. Piecuch and K. Kowalski, In search of the relationship between multiple solutions characterizing coupled-cluster theories, in Computational Chemistry: Reviews of Current Trends. World Scientific (2000) 1–104.

[33] P. Piecuch, S. Zarrabian, J. Paldus and J. Čížek, Coupled-cluster approaches with an approximate account of triexcitations and the optimized-inner-projection technique. II. Coupled-cluster results for cyclic-polyene model systems. *Phys. Rev. B* **42** (1990) 3351.

[34] V.V. Prasolov, Problems and Theorems in Linear Algebra. Vol. 134. American Mathematical Society (1994).

[35] T. Rohwedder, The continuous Coupled Cluster formulation for the electronic Schrödinger equation. *ESAIM: Math. Modell. Numer. Anal.-Modél. Math. Anal. Numér.* **47** (2013) 421–447.

[36] T. Rohwedder and R. Schneider, Error estimates for the coupled cluster method. *ESAIM: Math. Modell. Numer. Anal.-Modél. Math. Anal. Numér.* **47** (2013) 1553–1582.

[37] R. Schneider, Analysis of the projected coupled cluster method in electronic structure calculation. *Numer. Math.* **113** (2009) 433–471.

[38] I.V. Skrypnik, Methods for Analysis of Nonlinear Elliptic Boundary Value Problems. Vol. 139. American Mathematical Society (1994).

[39] J.P. Solovej, Many body quantum mechanics. Lecture Notes (2007).

[40] L.T. Watson, S.C. Billups and A.P. Morgan, Algorithm 652: Hompack: a suite of codes for globally convergent homotopy algorithms. *ACM Trans. Math. Softw. (TOMS)* **13** (1987) 281–310.

[41] H. Yserentant, Regularity and Approximability of Electronic Wave Functions. Springer (2010).

[42] P. Zabrejko, Rotation of vector fields: definition, basic properties, and calculation, in Topological Nonlinear Analysis II. Springer (1997) 445–601.

[43] E. Zeidler and P.R. Wadsack, Nonlinear Functional Analysis and its Applications: Fixed-Point Theorems/Transl. by Peter R. Wadsack. Springer-Verlag (1993).

[44] T.P. Živković, Existence and reality of solutions of the coupled-cluster equations. *Int. J. Quantum Chem.* **12** (1977) 413–420.

[45] T.P. Živković and H.J. Monkhorst, Analytic connection between configuration–interaction and coupled-cluster solutions. *J. Math. Phys.* **19** (1978) 1007–1022.