

COUPLED-CLUSTER THEORY REVISITED. PART I: DISCRETIZATION*

MIHÁLY A. CSIRIK¹ AND ANDRE LAESTADIUS^{2,1}

Abstract. In a series of two articles, we propose a comprehensive mathematical framework for Coupled-Cluster-type methods. These methods aim at accurately solving the many-body Schrödinger equation. In this first part, we rigorously describe the discretization schemes involved in Coupled-Cluster methods using graph-based concepts. This allows us to discuss different methods in a unified and more transparent manner, including multireference methods. Moreover, we derive the single-reference and the Jeziorski–Monkhorst multireference Coupled-Cluster equations in a unified and rigorous manner.

1991 Mathematics Subject Classification. 81V55, 81-08, 81-10.

The dates will be set by the publisher.

1. INTRODUCTION

The Coupled-Cluster (CC) method is one of the most popular methods in computational quantum chemistry among Hartree–Fock (HF) and Density-Functional Theory (DFT). In its full generality, the quantum many-body problem is intractable, and it is one of the greatest challenges of quantum mechanics to devise practically useful methods to approximate the solutions of the many-body Schrödinger equation. Although the stationary Schrödinger equation itself is a linear eigenvalue problem, it is extremely high-dimensional even for a few particles and a low-dimensional one-particle space.¹ Here, we focus on those fermionic systems which are described by the so-called *molecular Hamilton operator*—on which most electronic-structure models are based in quantum chemistry. The Galerkin method applied to the Schrödinger equation (sometimes combined with an initial HF “guess”) is branded Configuration Interaction (CI) in computational quantum chemistry; unfortunately, its applicability is limited due to the aforementioned high-dimensionality issue. The HF method is perhaps conceptually the simplest, whereby the ground state is approximated by minimizing the energy of the system over Slater determinants; the resulting Euler–Lagrange equations constitute a nonlinear eigenvalue problem that yields the HF ground state. HF theory has attracted much interest in the mathematical physics community, see e.g. [2, 3, 7, 8, 12, 19, 25, 26, 36]. The spiritual successor to the statistical mechanics-motivated Thomas–Fermi

Keywords and phrases: quantum mechanics, many-body problem, quantum chemistry, electronic structure, coupled-cluster theory

* This work has received funding from the Norwegian Research Council through Grant Nos. 287906 (CCerror) and 262695 (CoE Hylleraas Center for Quantum Molecular Sciences).

¹ Hylleraas Centre for Quantum Molecular Sciences, Department of Chemistry, University of Oslo, P.O. Box 1033 Blindern, N-0315 Oslo, Norway (e-mail: m.a.csirik@kjemi.uio.no)

² Department of Computer Science, Oslo Metropolitan University, P.O. Box 4 St. Olavs plass, NO-0130 Oslo, Norway (e-mail: andre.laestadius@oslomet.no)

¹The dimension is $\binom{K}{N}$, where N is the number of particles, and K is dimension of the one-particle Hilbert space.

© EDP Sciences, SMAI 1999

theory—DFT—is the single most used method in quantum chemistry, and some of its mathematical aspects are also highly non-trivial [10, 20–22].

CC theory is a vast and highly active subfield of quantum chemistry, consisting of many variants and refinements. However, among the aforementioned methods, the CC approach has arguably received the least attention in the mathematics community.

1.1. Previous work

It is beyond the scope of this paper to give a historical review of the CC method and its vast number of variants. The interested reader is pointed to [4, 5, 14, 18, 35]. The survey article [27] is somewhat more mathematically-oriented and also proposes a rather general framework.

The analysis of the single-reference CC method by R. Schneider and T. Rohwedder [31–33] serves as a starting point of our description of the CC discretization. In Part II, we will briefly summarize [33] and its follow-up works, where the context is more appropriate.

1.2. Outline

It is our intention to present both known and new results in a self-contained manner and primarily with a mathematical audience in mind. In Section 2, we describe the setting of the quantum-mechanical problems the CC theory is aimed at. Next, Section 2.3 gives a rough sketch of the most basic CI and CC methods.

We begin our discussion in Section 3.1 with the definition of a partial order relation which will be used to encode the relevant transitions of the system, called *excitations*. This partial order, and the induced lattice operations will be used in Section 3.2 to define the *excitation graph*, which fully describes the CC discretization scheme. We give a few examples of the generality of our concepts and also extend the definition of the excitation graph to the multireference (MR) case. After this, the corresponding *excitation operators* (Section 3.3), *cluster operators* (Section 3.4) and *cluster amplitude spaces* (Section 3.5) are constructed, which are the essential building blocks for the formulation of any CC-like method.

In Section 4, we give short derivations of the conventional single-reference (SRCC) and Jeziorski–Monkhorst multireference (JM-MRCC) methods. We do so by generalizing the known procedure which is based on perturbation theory.

In Appendix A we calculate various graph-theoretic properties of the excitation graph. In Appendix B we propose a method based on linear programming to select reference determinants for the multi-reference setting in an optimal way.

2. BACKGROUND

In this section we collect the concepts and results that are necessary for the forthcoming discussion. For proofs and more about the mathematical foundations of quantum mechanics, see e.g. [13, 15, 24, 28–30, 38].

The spectrum of a linear operator A is written $\sigma(A)$, the elements of its discrete spectrum as $\mathcal{E}_n(A)$, where $n = 0, 1, 2, \dots$, if A is bounded from below. We use the usual notation $[A, B] = AB - BA$ for the commutator. The 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 $V \rightarrow W$ endowed with the operator norm $\|\cdot\|_{\mathcal{L}(V, W)}$. Furthermore, V^* denotes the (continuous) dual space.

2.1. Function spaces

In the context of many-body quantum mechanics, the Lebesgue-, and Sobolev spaces $L^2(\mathbb{R}^3)$ and $H^1(\mathbb{R}^3)$ are viewed as “one-particle spaces”. We ignore spin for simplicity and consider $L^2(\mathbb{R}^3)$ and $H^1(\mathbb{R}^3)$ as *real* Hilbert spaces, again for clarity. These choices are justified for our model Hamiltonian (see Section 2.2 below), but we remark that all the forthcoming considerations can be trivially extended to the more general setting. The

one-particle spaces are then used to define the N -particle *fermionic* spaces (see e.g. [23])

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

endowed with the inner products

$$\langle \Psi, \Phi \rangle = \int_{\mathbb{R}^{3N}} \Psi(\mathbf{X}) \Phi(\mathbf{X}) \, d\mathbf{X}$$

and

$$\langle \Psi, \Phi \rangle_{\mathfrak{H}^1} = \langle \Psi, \Phi \rangle + \sum_{k=1}^N \int_{\mathbb{R}^{3N}} \nabla_{\mathbf{x}_k} \Psi(\mathbf{X}) \cdot \nabla_{\mathbf{x}_k} \Phi(\mathbf{X}) \, d\mathbf{X},$$

respectively. Here, $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N) \in \mathbb{R}^{3N}$ and $\mathbf{z} \cdot \mathbf{w}$ denotes the Euclidean inner product. Also, $\nabla_{\mathbf{x}_k} = (\partial_{x_k^1}, \partial_{x_k^2}, \partial_{x_k^3})$ is the distributional gradient operator acting on the k th triple of the arguments. The norms corresponding to $\langle \cdot, \cdot \rangle$ and $\langle \cdot, \cdot \rangle_{\mathfrak{H}^1}$ are denoted as $\| \cdot \|$ and $\| \cdot \|_{\mathfrak{H}^1}$, respectively. We define the second order Sobolev space as $\mathfrak{H}^2 = \mathfrak{L}^2 \cap H^2(\mathbb{R}^{3N})$.

Let $K \geq N$ or $K = \infty$ and assume that an L^2 -orthonormal (*spin-*)*orbital set* $\mathcal{B}_K = \{\varphi_p\}_{p=1}^K \subset H^1(\mathbb{R}^3)$ is given. We define the subspace $H_K^1(\mathbb{R}^3) = \text{Span } \mathcal{B}_K \subset H^1(\mathbb{R}^3)$. Corresponding to \mathcal{B}_K we can construct the set of *Slater determinants*

$$\mathfrak{B}_K = \{\Phi_\alpha \in \mathfrak{H}^1 : 1 \leq \alpha_1 < \dots < \alpha_N \leq K, \Phi_\alpha(\mathbf{X}) = N!^{-1/2} \det(\varphi_{\alpha_i}(\mathbf{x}_j))_{1 \leq i, j \leq N}\}.$$

Then \mathfrak{B}_K is \mathfrak{L}^2 -orthonormal. Set

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

The negative exponent Sobolev space \mathfrak{H}^{-1} will also be used in the sequel, which is given by the continuous dual space $(\mathfrak{H}^1)^*$. We will exploit that the dense continuous embeddings $\mathfrak{H}^1 \subset \mathfrak{L}^2 \subset \mathfrak{H}^{-1}$ hold true (see e.g. [1]), i.e. the spaces in question form a *Gelfand triple*.

Remark 2.1. An important remark is in order. Recall that $V \subset H \subset V^*$ are said to form a Gelfand triple if V is a real reflexive Banach space, H is real separable Hilbert space and the embedding $V \subset H$ is continuous and V is dense in H (see e.g. [39, Section 23.4]). It follows that for any $\Psi \in H$ there is a $\widehat{\Psi} \in V^*$ such that $\langle \widehat{\Psi}, \Phi \rangle_{V^* \times V} = \langle \Psi, \Phi \rangle_H$ for all $\Phi \in V$, and the mapping $H \rightarrow V^*$ given by $\Psi \mapsto \widehat{\Psi}$ is linear, injective and continuous. Hence, the embedding $H \subset V^*$ is also continuous (and dense). Henceforth, we will write Ψ in place of $\widehat{\Psi}$ for brevity.

Convention: We will drop the subscript from $\langle \cdot, \cdot \rangle$, as it is *either* obvious from its arguments that the duality pairing $\langle \cdot, \cdot \rangle_{V^* \times V}$ has to be used, *or* both $\langle \cdot, \cdot \rangle_{V^* \times V}$ and $\langle \cdot, \cdot \rangle_H$ are acceptable due to the Gelfand triple setting as discussed above. In particular, we apply this convention to $\mathfrak{H}^1 \subset \mathfrak{L}^2 \subset \mathfrak{H}^{-1}$, to the cluster amplitude spaces discussed in Section 3.5 and also in the general framework of Section 4.

2.2. Schrödinger Hamiltonian

In this section, we introduce the model Hamiltonian for concreteness. Let $V, w : \mathbb{R}^3 \rightarrow \mathbb{R}$ be Kato class potentials: $V, w \in L^{3/2}(\mathbb{R}^3) + L_\varepsilon^\infty(\mathbb{R}^3)^2$ with w even and define the quadratic form \mathcal{E} on \mathfrak{H}^1 as

$$\mathcal{E}(\Psi) = \frac{1}{2} \|\nabla \Psi\|^2 + \int_{\mathbb{R}^{3N}} \left(\sum_{1 \leq i \leq N} V(\mathbf{x}_i) + \sum_{1 \leq i < j \leq N} w(\mathbf{x}_i - \mathbf{x}_j) \right) |\Psi(\mathbf{X})|^2 \, d\mathbf{X}$$

²By definition $f \in L^{3/2}(\mathbb{R}^3) + L_\varepsilon^\infty(\mathbb{R}^3)$, if for every $\varepsilon > 0$ there is an $f_1 \in L^{3/2}(\mathbb{R}^3)$ and $f_2 \in L^\infty(\mathbb{R}^3)$ with $\|f_2\|_\infty < \varepsilon$ so that $f = f_1 + f_2$.

for any $\Psi \in \mathfrak{H}^1$. For every $\varepsilon > 0$, there is a $C_\varepsilon > 0$ so that Kato’s inequality (see e.g. [12] for a detailed proof),

$$\frac{1-\varepsilon}{2} \|\nabla\Psi\|^2 - C_\varepsilon \|\Psi\|^2 \leq \mathcal{E}(\Psi) \leq \frac{1+\varepsilon}{2} \|\nabla\Psi\|^2 + C_\varepsilon \|\Psi\|^2 \quad \text{for all } \Psi \in \mathfrak{H}^1,$$

holds true. This implies that the quadratic form induced by V and w is infinitesimally form bounded with respect to $-\Delta$. The KLMN theorem (see e.g. [29]) implies that there exists a unique self-adjoint operator $\mathcal{H} : D(\mathcal{H}) \rightarrow \mathfrak{L}^2$ associated to \mathcal{E} , having form domain $Q(\mathcal{H}) = Q(\mathcal{E}) = \mathfrak{H}^1$ and being lower semibounded. This \mathcal{H} is given by

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

for all $\Psi \in D(\mathcal{H})$ and $\mathbf{X} \in \mathbb{R}^{3N}$. Kato’s inequality implies that there is a constant $M > 0$, such that

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

for all $\Psi, \Phi \in \mathfrak{H}^1$. Thus, \mathcal{H} can be extended to a bounded mapping $\mathfrak{H}^1 \rightarrow \mathfrak{H}^{-1}$, 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 $\Phi \in \mathfrak{H}^1$.

As far as the finite-dimensional case $K < \infty$ is concerned, we simply consider the Galerkin projection of the weak Schrödinger equation. More precisely, let $\mathfrak{H}_K^1 \subset \mathfrak{H}^1$ be as defined in Section 2.1. Then $\Psi \in \mathfrak{H}_K^1$ and $\mathcal{E} \in \mathbb{R}$ are said to satisfy the *projected Schrödinger equation* if $\langle \mathcal{H}\Psi, \Phi \rangle = \mathcal{E} \langle \Psi, \Phi \rangle$ for all $\Phi \in \mathfrak{H}_K^1$.

The so-called (electronic) molecular Hamiltonian \mathcal{H} corresponds to the special case

$$V(\mathbf{x}) = - \sum_{1 \leq j \leq M} \frac{Z_j}{|\mathbf{x} - \mathbf{r}_j|} \quad \text{and} \quad w(\mathbf{x}) = \frac{1}{|\mathbf{x}|},$$

where $Z_j \in \mathbb{N}$ ($j = 1, \dots, M$) and $\mathbf{r}_1, \dots, \mathbf{r}_M \in \mathbb{R}^3$ denote the charges and the positions of the $M \in \mathbb{N}$ nuclei.

2.3. The CI and the CC method

We now give a very rough description of the single-reference CI and CC methods. For the rigorous derivations, we refer to Section 4.

In a preliminary step—typically using the Hartree–Fock method—the *reference determinant*

$$\Phi_0 = N!^{-1/2} \det(\varphi_i(\mathbf{x}_j))_{1 \leq i, j \leq N}$$

is constructed and normalized so that $\|\Phi_0\| = 1$. We restrict our discussion here to the case when relevant the function spaces are real. The *occupied orbitals* $\mathcal{B}_{\text{occ}} = \{\varphi_p\}_{p=1}^N \subset H^1(\mathbb{R}^3)$ are extended to a basis $\mathcal{B}_K = \{\varphi_p\}_{p=1}^K \subset H_K^1(\mathbb{R}^3)$ by adding $K - N$ *virtual orbitals* $\mathcal{B}_{K, \text{virt}} = \{\varphi_p\}_{p=N+1}^K$, so that $\mathcal{B}_K = \mathcal{B}_{\text{occ}} \cup \mathcal{B}_{K, \text{virt}}$. Here, $K = \infty$ is allowed. The orthonormal set \mathcal{B}_K generates the Slater determinants \mathfrak{B}_K and the subspace $\mathfrak{H}_K^1 \subset \mathfrak{H}^1$ (see Section 2.1). For later convenience, we introduce the space $\mathfrak{H}^{1, \perp}$ as the \mathfrak{L}^2 -orthogonal complement of $\text{Span}\{\Phi_0\}$ in \mathfrak{H}^1 . Further, we also set $\mathfrak{H}_K^{1, \perp} = \mathfrak{H}^{1, \perp} \cap \mathfrak{H}_K^1$. Further, we define $\mathfrak{L}^{2, \perp}$ as the \mathfrak{L}^2 -orthogonal complement of $\text{Span}\{\Phi_0\}$ in \mathfrak{L}^2 .

In both the CI and the CC method, the Schrödinger equation $\mathcal{H}\Psi = \mathcal{E}\Psi$ is solved based on the reference wavefunction Φ_0 . For simplicity,³ we consider the case when Ψ is sought after in the form $\Psi = \Phi_0 + \underline{\Psi}$, where $\langle \underline{\Psi}, \Phi_0 \rangle = 0$. In other words, Ψ is calculated via a *correction* $\underline{\Psi}$ to Φ_0 . Note that $\langle \Psi, \Phi_0 \rangle = 1$, which is called the *intermediate normalization* condition. If the “targeted” wavefunction Ψ happens to be orthogonal to the reference determinant Φ_0 , then the Ansatz $\Psi = \Phi_0 + \underline{\Psi}$ cannot yield a solution.

³Although the CI method is more general.

The *Full Configuration Interaction* (FCI) method can be summarized as follows: find $\underline{\Psi} \in \mathfrak{H}_K^{1,\perp}$ such that

$$\langle \mathcal{H}(\Phi_0 + \underline{\Psi}), \Phi \rangle = \mathcal{E}_{\text{CI}} \langle \Phi_0 + \underline{\Psi}, \Phi \rangle \quad \text{for all } \Phi \in \mathfrak{H}_K^{1,\perp}. \quad (2.2)$$

Here, $\mathcal{E}_{\text{CI}} = \|\Psi\|^{-2} \langle \mathcal{H}\Psi, \Psi \rangle$ is called the *CI*-, or *variational energy*. The *projected CI* method is simply the Galerkin projection of the previous problem to some finite dimensional subspace $\mathfrak{V}_d \subset \mathfrak{H}_K^{1,\perp}$, i.e. to find $\underline{\Psi}_d \in \mathfrak{V}_d$ such that

$$\langle \mathcal{H}(\Phi_0 + \underline{\Psi}_d), \Phi_d \rangle = \mathcal{E}_{d,\text{CI}} \langle \Phi_0 + \underline{\Psi}_d, \Phi_d \rangle \quad \text{for all } \Phi_d \in \mathfrak{V}_d. \quad (2.3)$$

The choice of the Galerkin subspace \mathfrak{V}_d is typically based on so-called truncation rank, for instance $\mathfrak{V}_d = \mathfrak{V}_{\text{SD}}$, is the span of singly-, and doubly excited determinants in \mathfrak{B}_K . The corresponding (projected) CI method in this case is designated as ‘‘CISD’’.

The CI equations are more commonly expressed using *cluster operators*. A cluster operator $C : \mathfrak{L}^2 \rightarrow \mathfrak{L}^2$ is a bounded linear operator that is a linear combination of special products of fermionic creation and annihilation operators a_i^\dagger and a_i (see Part II for a definition), so that the action of each such product is to replace some occupied orbitals \mathcal{B}_{occ} with the same number of virtual orbitals $\mathcal{B}_{K,\text{virt}}$ in a Slater determinant (see Remark 3.18). A cluster operator C can therefore be parametrized with the said linear-combination coefficients, which are denoted by the lower case c and are called *cluster amplitudes*. The vector space of all cluster amplitudes is denoted by \mathbb{V} . There is a one-to-one correspondence between functions in $\mathfrak{L}^{2,\perp}$ (resp. $\mathfrak{H}^{1,\perp}$) and functions of the form $C\Phi_0$, where $C : \mathfrak{L}^2 \rightarrow \mathfrak{L}^2$ (resp. $C : \mathfrak{H}^1 \rightarrow \mathfrak{H}^1$) is a cluster operator (see [31]). Therefore, (2.2) can be expressed as follows: find a cluster operator C (or, equivalently cluster amplitudes c), such that

$$\langle \mathcal{H}(I + C)\Phi_0, S\Phi_0 \rangle = \mathcal{E}_{\text{CI}}(c) \langle (I + C)\Phi_0, S\Phi_0 \rangle \quad \text{for all cluster operators } S. \quad (2.4)$$

Here, $\mathcal{E}_{\text{CI}}(c) = \|(I + C)\Phi_0\|^{-2} \langle \mathcal{H}(I + C)\Phi_0, (I + C)\Phi_0 \rangle$. Although this might seem an unnecessary complication at first, cluster operators are essential for the formulation of the CC method.

In the CC method, the ‘‘exponential Ansatz’’ is assumed for the intermediately normalized wavefunction Ψ . Substituting $\Psi = e^T \Phi_0$ into the Schrödinger equation, where T is a cluster operator, we get

$$\mathcal{H}e^T \Phi_0 = \mathcal{E}_{\text{CC}} e^T \Phi_0, \quad (2.5)$$

for some $\mathcal{E}_{\text{CC}} \in \mathbb{R}$. First, to determine \mathcal{E}_{CC} we premultiply (2.5) by e^{-T} (e^T is always invertible), and take the inner product with Φ_0 to obtain the *CC energy*

$$\mathcal{E}_{\text{CC}} := \mathcal{E}_{\text{CC}}(t) = \langle e^{-T} \mathcal{H}e^T \Phi_0, \Phi_0 \rangle, \quad (2.6)$$

where we used the normalization $\|\Phi_0\| = 1$. Second, by premultiplying (2.5) by e^{-T} again, but now testing against functions in $\mathfrak{H}_K^{1,\perp}$, we get the *Full CC* (FCC) method: find cluster amplitudes $t_* \in \mathbb{V}$ such that

$$\langle e^{-T_*} \mathcal{H}e^{T_*} \Phi_0, S\Phi_0 \rangle = 0, \quad \text{for all } s \in \mathbb{V}. \quad (2.7)$$

The *projected CC* method is the Galerkin projection of the FCC problem with respect to some subspace $\mathbb{V}_d \subset \mathbb{V}$. More precisely, the task is to find $t_{d,*} \in \mathbb{V}_d$ such that

$$\langle e^{-T_{d,*}} \mathcal{H}e^{T_{d,*}} \Phi_0, S_d \Phi_0 \rangle = 0 \quad \text{for all } s_d \in \mathbb{V}_d. \quad (2.8)$$

For the moment, we denote the corresponding CC energy by $\mathcal{E}_{d,\text{CC}}$. Again, \mathbb{V}_d is based on some truncation, such as SD, in which case the corresponding method is called ‘‘CCSD’’.

We now discuss the relation between CI and CC. It was shown that the FCI (2.2) and the FCC (2.7) equations are equivalent, see [31, Theorem 5.3].

Theorem 2.2 (Equivalence of FCI and FCC). *The problems (2.4) and (2.7) are equivalent, i.e. a full CC solution t_* is also full CI solution $I + C_* = e^{T_*}$, and vice versa. Moreover, $\mathcal{E}_{\text{CC}}(t_*) = \mathcal{E}_{\text{CI}}(c_*)$ holds true.*

However, the corresponding Galerkin-projected problems are *not* equivalent. Further, while $\mathcal{E}_{\text{CI}} \leq \mathcal{E}_{d,\text{CI}}$ due to the Rayleigh–Ritz variational principle, the same is not true for the CC method and numerical experience undoubtedly shows that *there is no obvious relation* in general between $\mathcal{E}_{\text{CC}} = \mathcal{E}_{\text{CI}}$ and $\mathcal{E}_{d,\text{CC}}$; this last phenomenon is called the *nonvariational property* of CC theory. Note that according to Theorem 2.2, FCC is variational.

Despite this, the gains of CC over CI are significant. First, by construction, the CC method is size-consistent, even when truncated [33, Theorem 4.10]. This property is crucial for the precise determination of various chemical properties. Second, the evaluation of expressions involving the *similarity-transformed Hamilton operator* $e^{-T}\mathcal{H}e^T$ is greatly eased by the formula

$$e^{-T}\mathcal{H}e^T = \sum_{j=0}^4 \frac{1}{j!} [\mathcal{H}, T]_{(j)}, \quad (2.9)$$

see [33, Theorem A.1],⁴ where the iterated commutators are given by $[\mathcal{H}, T]_{(0)} = \mathcal{H}$ and $[\mathcal{H}, T]_{(j)} = [[\mathcal{H}, T]_{(j-1)}, T]$ for $j \geq 1$. Equation (2.9) may be referred to as the terminating Baker–Campbell–Hausdorff series, and it makes the computer implementation of CC methods feasible even for moderately sized systems. In particular, the Slater–Condon rules imply that the CC energy can be computed as⁵

$$\mathcal{E}_{\text{CC}}(t) = \langle \mathcal{H}(I + T_1 + T_2 + \frac{1}{2}T_1^2)\Phi_0, \Phi_0 \rangle. \quad (2.10)$$

Furthermore, (2.9) also implies that the polynomial system (2.7) (and hence its Galerkin projection (2.8)) is quartic in terms of the cluster amplitudes t . Despite their apparent simplicity, the CC equations usually involve many complicated terms and even their assembly is a nontrivial task. In summary, the CC method approximates an extremely high-dimensional linear problem (2.2) by a low-dimensional nonlinear problem (2.8).

3. COUPLED-CLUSTER DISCRETIZATION

Using an appropriate string of creation and annihilation operators, any fermionic state can be changed to any other one (see Part II, or [14, 37]). In our context, a set of N occupied orbitals is given; its complement is called the set of virtual orbitals. The action of an excitation operator on a Slater determinant consists of annihilating a few occupied orbitals and creating the same number of virtual orbitals (hence the particle number N is conserved). A de-excitation operator amounts to the reverse action: annihilating some virtual orbitals and creating the same number of occupied ones. Obviously, any N -particle Slater determinant can be obtained by acting with an appropriate excitation operator on the “reference state”, which is the N -particle Slater determinant of all the occupied orbitals. However, it might also be possible to arrive at the same Slater determinant from another state through successive excitations. The concrete relationships are nontrivial and this section is devoted to their description.

3.1. Excitation order

Let Λ be a countable set called the *orbital set* and let 2^Λ denote the power set of Λ . In concrete examples, we will often use the numbers $\Lambda = \{1, 2, 3, \dots\}$ to label the elements of Λ for the sake of simplicity, and set $K = |\Lambda|$. Let $N \geq 1$ denote the number of particles, and set $S = \{\alpha \in 2^\Lambda : |\alpha| = N\}$, the elements of which are called (N -particle) *states*. The particle number N is assumed to be fixed throughout. Fix $M \geq 1$ *reference states*

$$\Omega = \{0_1, \dots, 0_M\} \subset S.$$

⁴Their proof is straightforward to adapt to the more general Hamilton operator defined in Section 2.2.

⁵Actually, the term $\langle \mathcal{H}T_1\Phi_0, \Phi_0 \rangle$ vanishes if Φ_0 is the Hartree–Fock solution (Brillouin theorem).

For every $m = 1, \dots, M$ define

$$L_m = S \setminus (\Omega \setminus \{0_m\})$$

and on it, the partial order relation

$$\alpha \preceq_m \beta \iff \underline{\beta}_m \subset \underline{\alpha}_m \quad \text{and} \quad \bar{\alpha}^m \subset \bar{\beta}^m$$

for any $\alpha, \beta \in L_m$, where

$$\underline{\alpha}_m = \alpha \cap 0_m \quad \text{and} \quad \bar{\alpha}^m = \alpha \cap (0_m)^c,$$

and the complement is to be understood relative to Λ . According to commonly used nomenclature, we call $\underline{\alpha}_m$ the *occupied part of α w.r.t. 0_m* and $\bar{\alpha}^m$ the *virtual part of α w.r.t. 0_m* . This partial order relation is a generalization of [31, Definition 4.2]. By definition, $L_m = \{\alpha \in S : 0_m \preceq_m \alpha\}$ and for the sake of convenience, we introduce the notations $\bar{S} = S \setminus \Omega$ and $\bar{L}_m = L_m \setminus \{0_m\}$. Note that the reference states are defined *not* to be comparable with respect to \preceq_m with each other.

The partial order \preceq_m generates the *join* and *meet* lattice operations

$$\begin{aligned} \alpha \vee_m \beta &= (\underline{\alpha}_m \cap \underline{\beta}_m) \cup (\bar{\alpha}^m \cup \bar{\beta}^m), \\ \alpha \wedge_m \beta &= (\underline{\alpha}_m \cup \underline{\beta}_m) \cup (\bar{\alpha}^m \cap \bar{\beta}^m), \end{aligned}$$

for all $\alpha, \beta \in L_m$. Furthermore, we introduce the orthocomplementation $\alpha^\perp = \Lambda \setminus \alpha$.

For the so-called *single-reference* (SR) case, $M = 1$ and we will make the convention that all the m indices are dropped from the notation. For the next result, we extend \preceq , \vee and \wedge to the whole 2^Λ .

Proposition 3.1. *The structure $B = (2^\Lambda, \vee, \wedge, 0, 1, \perp)$ is a Boolean algebra, that is, a distributive, bounded lattice in which the de Morgan laws hold true. Here, we set $1 := \Lambda$, the identity for \wedge .*

A similar statement holds true in the *multi-reference* (MR) case, for the individual structures $B_m = (2^\Lambda, \vee_m, \wedge_m, 0_m, 1, \perp)$. Even though the algebraic structure on B is nice, the subset S loses this structure. In fact, S is *not* a sublattice of B , since for example $\alpha \vee_m \beta, \alpha \wedge_m \beta \notin S$ for distinct α and β with $\underline{\alpha} = \underline{\beta} = \emptyset$. The reason why we stated Proposition 3.1, however, is because we will exploit the operational rules for \vee , \wedge and \perp on a few occasions; for instance, in the following trivial result.

Lemma 3.2. *Let $\gamma, \beta \in 2^\Lambda$ be such that $\beta \preceq_m \gamma$. Then, $\alpha \vee_m \beta = \gamma$ if and only if $\alpha = \beta^\perp \wedge_m \gamma$.*

Proof. We have

$$\alpha \vee_m \beta = (\gamma \wedge_m \beta^\perp) \vee_m \beta = (\gamma \vee_m \beta) \wedge_m (\beta^\perp \vee_m \beta) = (\gamma \vee_m \beta) \wedge_m 1 = \gamma \vee_m \beta = \gamma,$$

where in the last step we used $\beta \preceq_m \gamma$. Further, if $\alpha' \vee_m \beta = \gamma$ as well, then $\alpha' \vee_m \beta = \alpha \vee_m \beta$. By joining β^\perp to both sides, we get $\alpha' = \alpha$. \square

The poset (L_m, \preceq_m) also admits a rank function which makes it a *graded poset*. Being a graded poset means that there is a *rank function* $\text{rk}_m : L_m \rightarrow \mathbb{N}$ satisfies $\text{rk}_m(\alpha) < \text{rk}_m(\beta)$ whenever $\alpha \prec_m \beta$, and $\text{rk}_m(\beta) = \text{rk}_m(\alpha) + 1$ if there is no element γ such that $\alpha \prec_m \gamma \prec_m \beta$. The choice $\text{rk}_m(\alpha) = |\bar{\alpha}^m|$ is easily seen to satisfy the requirements. Obviously, the maximum value that $\text{rk}_m(\alpha)$ can take is N . For a geometric description of the rank function, see Appendix B.

3.2. Excitation graphs

As we remarked in the previous section, L_m fails to be a sublattice of the Boolean algebra B_m . Therefore, let us consider pairs $(\alpha, \beta) \in L_m \times L_m$ for which $\alpha \vee_m \beta \in L_m$. In other words, pairs $(\alpha, \beta) \in L_m \times L_m$ for which $|\underline{\alpha}_m \cap \underline{\beta}_m| + |\bar{\alpha}^m \cup \bar{\beta}^m| = N$, or, using the inclusion-exclusion principle $|A \cup B| = |A| + |B| - |A \cap B|$, we can equivalently write

$$|\underline{\alpha}_m \cup \underline{\beta}_m| + |\bar{\alpha}^m \cap \bar{\beta}^m| = N, \tag{3.1}$$

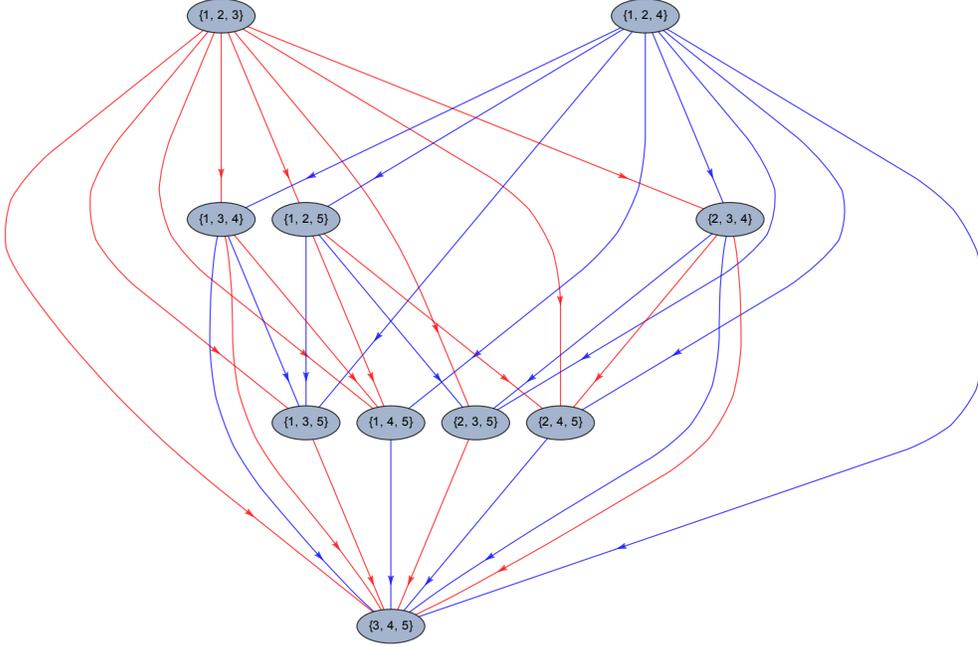


FIGURE 1. Full multi-reference excitation multigraph for $\Lambda = \{1, \dots, 5\}$ and $O_1 = \{1, 2, 3\}$, $O_2 = \{1, 2, 4\}$. The edges corresponding to O_1 and O_2 are shown in red and blue, respectively.

since $|\underline{\alpha}_m| + |\bar{\alpha}^m| = N$ and $|\underline{\beta}_m| + |\bar{\beta}^m| = N$ by hypothesis. While still $\alpha \vee_m \beta \in L_m$ even in the case $\bar{\alpha}^m \cap \bar{\beta}^m \neq \emptyset$, we wish to avoid that possibility on physical grounds. Namely, such an operation would introduce a repeated virtual orbital in the state $\alpha \vee_m \beta$ and that is not allowed on the account of the Pauli exclusion principle. We note in passing that $\alpha \vee_m \beta \in L_m$ is equivalent to $\alpha \wedge_m \beta \in L_m$. In conclusion, we restrict our attention to the set

$$\mathcal{L}_m = \{(\alpha, \beta) \in L_m \times L_m : |\underline{\alpha}_m \cup \underline{\beta}_m| = N \text{ and } |\bar{\alpha}^m \cap \bar{\beta}^m| = 0\}. \quad (3.2)$$

Hence, if $(\alpha, \beta) \in \mathcal{L}_m$, then we have $\alpha \vee_m \beta \in L_m$. The set \mathcal{L}_m is symmetric to the diagonal (which it does not contain), i.e. $(\alpha, \beta) \in \mathcal{L}_m$ iff $(\beta, \alpha) \in \mathcal{L}_m$ and $(\alpha, \alpha) \notin \mathcal{L}_m$. Also, $(0_m, \alpha), (\alpha, 0_m) \in \mathcal{L}_m$ for any $\alpha \in L_m$. Furthermore, the rank function rk_m is additive on \mathcal{L}_m in the sense that

$$\text{rk}_m(\alpha \vee_m \beta) = \text{rk}_m(\alpha) + \text{rk}_m(\beta),$$

for any $(\alpha, \beta) \in \mathcal{L}_m$. This property may also be seen to be a reason why we want to exclude the case $\bar{\alpha}^m \cap \bar{\beta}^m \neq \emptyset$. Indeed, it could also be taken as the *definition* of \mathcal{L}_m .

Proposition 3.3. *The set \mathcal{L}_m can be written as*

$$\mathcal{L}_m = \{(\alpha, \beta) \in L_m \times L_m : \alpha \vee_m \beta \in L_m \text{ and } \text{rk}_m(\alpha \vee_m \beta) = \text{rk}_m(\alpha) + \text{rk}_m(\beta)\}.$$

Proof. Let \mathcal{L}'_m denote the set on the right hand side of the preceding equation. Then, it is clear from the above that $\mathcal{L}_m \subset \mathcal{L}'_m$. Conversely, suppose that $(\alpha, \beta) \in \mathcal{L}'_m$. Then, $|\bar{\alpha}^m \cup \bar{\beta}^m| = |\bar{\alpha}^m| + |\bar{\beta}^m|$, from which $|\bar{\alpha}^m \cap \bar{\beta}^m| = 0$ using the inclusion-exclusion principle. Since $\alpha \vee_m \beta \in L_m$, (3.1) holds true, and we have that $|\underline{\alpha}_m \cup \underline{\beta}_m| = N$, so $(\alpha, \beta) \in \mathcal{L}_m$. Hence, $\mathcal{L}_m \supset \mathcal{L}'_m$. \square

The set \mathcal{L}_m is used for our main definition.

Definition 3.4. The digraph $G_m^{\text{full}} = (L_m, E_m^{\text{full}})$ is called the *full (SR) excitation graph w.r.t. 0_m* , where

$$E_m^{\text{full}} = \{(\beta, \alpha \vee_m \beta) \in L_m \times L_m : (\alpha, \beta) \in \mathcal{L}_m, \alpha \neq 0_m\}.$$

A subgraph $G_m = (L_m, E_m)$, $E_m \subset E_m^{\text{full}}$ is said to be an *(SR) excitation (sub)graph w.r.t. 0_m* .

Notice that we excluded $\alpha = 0_m$ to omit loop edges. Lemma 3.2 can be refined in the following manner.

Lemma 3.5. *Let $(\beta, \gamma) \in E_m^{\text{full}}$. Then $\alpha = \beta^\perp \wedge_m \gamma \in L_m$ is the unique α such that $\alpha \vee_m \beta = \gamma$.*

Proof. Using Lemma 3.2, we can uniquely solve the equation $\alpha \vee_m \beta = \gamma$ for α to obtain $\alpha = \beta^\perp \wedge_m \gamma \in 2^\Lambda$. Therefore, $(\beta, \alpha \vee_m \beta) \in E_m^{\text{full}}$, which implies that $\alpha \in L_m$ using the definition of E_m^{full} . \square

Corollary 3.6. *The digraph G_m^{full} does not contain parallel edges.*

Various graph-theoretic quantities of the single-reference excitation graph are calculated in Appendix A.

A digraph $G = (V, E)$ is said to be *transitive* if $(u, v) \in E$ and $(v, w) \in E$ imply $(u, w) \in E$. It follows by induction that, if G is transitive, and whenever G contains a directed path $((v_0, v_1), (v_1, v_2), \dots, (v_{n-1}, v_n))$, then $(v_0, v_n) \in E$.

Proposition 3.7. *The digraph G_m^{full} is transitive.*

Proof. Suppose that $(\gamma_0, \gamma_1) \in E_m^{\text{full}}$ and $(\gamma_1, \gamma_2) \in E_m^{\text{full}}$. Then there exists α and β such that $\gamma_1 = \alpha \vee_m \gamma_0$ and $\gamma_2 = \beta \vee_m \gamma_1$. Since $\gamma_2 = (\alpha \vee_m \beta) \vee \gamma_0$ by the associativity of \vee_m , it follows easily from Proposition 3.3 that $(\gamma_0, \alpha \vee_m \beta) \in \mathcal{L}_m$. Therefore,

$$(\gamma_0, \gamma_2) = (\gamma_0, (\alpha \vee_m \beta) \vee_m \gamma_0) \in E_m^{\text{full}},$$

which is what we wanted to show. \square

Transitivity of certain subgraphs, and of G_m^{full} itself will come up later, since vaguely speaking this property will imply the algebraic closedness of the set of excitation operators that we attach to the edges (see Section 3.3 and Section 3.4).

We label the edges of G_m^{full} with their corresponding α . Thus, to every directed edge $(\beta, \alpha \vee_m \beta) \in E_m^{\text{full}}$ there corresponds a map $x_{m,\alpha} : L_m \rightarrow L_m$ defined with the instruction $x_{m,\alpha}(\beta) = \alpha \vee_m \beta$. This way, the digraph G_m^{full} may be interpreted as a commutative diagram (cf. Section 3.3). Note that a label $x_{m,\alpha}$ may appear on multiple edges.

Furthermore, for any subgraph $G_m = (L_m, E_m)$, we introduce the *set of excitations* $\Xi(G_m) \subset \bar{L}_m$ of G_m via

$$\Xi(G_m) = \{\alpha \in \bar{L}_m : (\beta, \alpha \vee_m \beta) \in E_m \text{ for some } \beta \in L_m\}. \quad (3.3)$$

Note that the excitations are indexed with the same set L_m as the states themselves, but in general $\Xi(G_m) \neq \bar{L}_m$. Nonetheless, for the full excitation graph G_m^{full} , we have in fact $\Xi(G_m^{\text{full}}) = \bar{L}_m$.

The reason why explicitly stated that we are considering the “full” excitation graphs is that, in practice, one is forced to ignore the “degree of freedom” (called “cluster amplitudes”, see Section 3.5) corresponding to some edges.⁶ This is done by considering certain subsets of the full edge set E_m^{full} .

Definition 3.8. An excitation subgraph $G_m = (L_m, E_m)$ is said to be a *consistent subgraph (of G_m^{full})* if $E_m \subset E_m^{\text{full}}$, and whenever $(\beta, \alpha \vee_m \beta) \in E_m$ for some $\beta \in L_m$ and $\alpha \in L_m$, then $(\beta', \alpha \vee_m \beta') \in E_m$ for all $\beta' \in L_m$.

The consistency criterion can be rephrased as follows: for a fixed $\alpha \in \Xi(G_m)$, either E_m contains the whole “orbit” $\{(\beta, \alpha \vee_m \beta) \in E_m : \beta \in L_m\}$ or it does not contain it at all. Note that the set $\Xi(G_m)$ can equally well be used to define a consistent subgraph.

⁶Note that the vertex set is still the “full” vertex set L_m —some vertices might become isolated.

Definition 3.9. For a given $r = 1, \dots, N$, define $G_m(r) = (L_m, E_m(r))$, where

$$E_m(r) = \{(\beta, \alpha \vee_m \beta) \in E_m^{\text{full}} : \beta \in L_m, \alpha \in \bar{L}_m \text{ such that } \text{rk}_m(\alpha) = r\}.$$

The subgraph $G_m(r_1, \dots, r_\rho) = (L_m, E_m(r_1, \dots, r_\rho))$ is called a *rank-truncated excitation subgraph* if

$$E_m(r_1, \dots, r_\rho) = E_m(r_1) \cup \dots \cup E_m(r_\rho) \quad \text{for } r_1, \dots, r_\rho \in \{1, \dots, N\}.$$

We refer to $G_m(1)$, $G_m(1, 2)$, $G_m(1, 2, 3)$, etc. more colloquially as $G_m(\text{S})$, $G_m(\text{SD})$, $G_m(\text{SDT})$, etc.

Rank-truncation does not introduce isolated vertices in $G_m(r_1, \dots, r_\rho)$ as long as one of the r_j 's is 1. However, in the doubles (D) case, $G(\text{D})$ *does* in fact produce isolated vertices so that vertices of odd rank cannot be reached. Also, note that these truncated subgraphs like $G_m(\text{S})$ and $G_m(\text{SD})$ are *not* transitive in general.

We shall summarize these observations in the next theorem. Recall that a digraph is said to be *weakly connected* if every pair of vertices has an undirected path between them.

Theorem 3.10. *Let $G_m = G_m(r_1, \dots, r_\rho)$ be a rank-truncated excitation subgraph. Then the following is true.*

- (i) G_m is a consistent subgraph.
- (ii) G_m is weakly connected if one of the r_j 's is 1.

Proof. Obvious from the definition. □

Next, we briefly consider two rather “exotic” CC-like methods to demonstrate the generality of the excitation graph concept.

Example 3.11. The excitation graph corresponding to the Tailored CC method (see e.g. [11]) can be described as follows. In this SR method ($M = 1$), the orbital set Λ is partitioned according to $\Lambda_{\text{CAS}} = \{1, \dots, N, N + 1, \dots, k\}$ and $\Lambda_{\text{ext}} = \Lambda \setminus \Lambda_{\text{CAS}}$ for some $k = N, \dots, |\Lambda|$. This induces a splitting $L = L(\text{CAS}) \dot{\cup} L(\text{ext})$, where

$$L(\text{CAS}) = \{\alpha \in L : \alpha \subset \Lambda_{\text{CAS}}\}, \quad L(\text{ext}) = L \setminus L(\text{CAS}).$$

Furthermore, the edge set E^{full} may also be split accordingly

$$E(\text{CAS}) = \{(\beta, \alpha \vee \beta) \in E^{\text{full}} : \alpha \subset \Lambda_{\text{CAS}}, \beta \in L\}, \quad \text{and} \quad E(\text{ext}) = E^{\text{full}} \setminus E(\text{CAS}).$$

In other words, $E(\text{CAS})$ contains excitations which change CAS occupied orbitals to CAS virtual ones, and as such, no edge in $E(\text{CAS})$ leaves $L(\text{CAS})$ that starts from $L(\text{CAS})$. It is easy to see that both $G(\text{CAS}) = (L, E(\text{CAS}))$ and $G(\text{ext}) = (L, E(\text{ext}))$ are transitive and consistent subgraphs.

Example 3.12. A generalization of $E(\text{CAS})$ is the “CAS-type subalgebra” (denoted as “ $\mathfrak{g}^{(N)}(R, S)$ ” in [17]), which is constructed from two given subsets $\Lambda_R \subset \{1, \dots, N\}$ and $\Lambda_S \subset \{N + 1, \dots\}$. Define $\Lambda_{\text{int}} = \Lambda_R \dot{\cup} \Lambda_S$ and $\Lambda_{\text{ext}} = \Lambda \setminus \Lambda_{\text{int}}$. This induces a splitting $L = L(\text{int}) \dot{\cup} L(\text{ext})$, where

$$L(\text{int}) = \{\alpha \in L : \alpha \subset \Lambda_{\text{int}}\}, \quad \text{and} \quad L(\text{ext}) = L \setminus L(\text{int}).$$

The edge set E^{full} decomposes as

$$E(\text{int}) = \{(\beta, \alpha \vee \beta) \in E^{\text{full}} : \alpha \subset \Lambda_{\text{int}}, \beta \in L\}, \quad \text{and} \quad E(\text{ext}) = E^{\text{full}} \setminus E(\text{int}).$$

In other words, $E(\text{int})$ contains excitations that replace some orbitals in Λ_R with ones in Λ_S . Then $G(\text{int}) = (L, E(\text{int}))$ and $G(\text{ext}) = (L, E(\text{ext}))$ are transitive and consistent subgraphs. Clearly, Example 3.11 can be recovered with the choice $\Lambda_R = \{1, \dots, N\}$, $\Lambda_S = \{N + 1, \dots, k\}$.

Finally, we define excitation graph in the multireference case, which is a natural extension of the above concepts.

Remark 3.13. An important warning is in order. In general, $\alpha \vee_m \beta$ may or may not be equal to $\alpha \vee_\ell \beta$ for $m \neq \ell$. In fact, take $\Lambda = \{1, 2, \dots, 7\}$ and $0_1 = \{1, 2, 3\}$, $0_2 = \{1, 2, 4\}$. Then, with $\alpha = \{1, 3, 5\}$ and $\beta = \{2, 6, 7\}$, we have $\alpha \vee_1 \beta = \alpha \vee_2 \beta = \{5, 6, 7\}$. On the other hand, with $\alpha = \{2, 3, 4\}$ and $\beta = \{1, 2, 5\}$, we have $\alpha \vee_1 \beta = \{2, 4, 5\}$, but $\alpha \vee_2 \beta = \{2, 3, 5\}$. Note that in the first case, we actually have $(\alpha, \alpha \vee_1 \beta) \in E_1^{\text{full}}$ and $(\alpha, \alpha \vee_2 \beta) \in E_2^{\text{full}}$, i.e. a double edge.

Definition 3.14. The *full MR excitation multigraph w.r.t. Ω* , $G^{\text{full}} = (L, E^{\text{full}})$ is defined as the union of the individual full SR excitation graphs $G_m^{\text{full}} = (L_m, E_m^{\text{full}})$ for all $m = 1, \dots, M$, i.e.

$$L = \bigcup_{m=1}^M L_m, \quad E^{\text{full}} = \biguplus_{m=1}^M E_m^{\text{full}},$$

where \biguplus denotes multiset union.

Note that as opposed to the SR graph G_m^{full} , the MR graph G^{full} might have parallel edges (called “redundant” excitations), this justifies that G^{full} was introduced as a multigraph. Notice that other references cannot be “reached” from a given one (see Section 3.1). An algorithm for choosing the set of reference states $\Omega = \{0_m\}_{m=1}^M$ in an optimal way, adhering to some given criteria is described in Appendix B.

3.3. Excitation operators

Recall that $\Omega = \{0_m\}_{m=1}^M$ denotes the set of references, and that L_m does *not* contain the other reference states $\Omega \setminus \{0_m\}$. The construction described below is to be repeated for every $m = 1, \dots, M$ separately.

First, we fix an ordering of the indices in $\alpha \in S$. Then, for every element $\alpha = \{\alpha_1, \dots, \alpha_N\} \in S$ we assign the lexicographically ordered N -tuple

$$\alpha^< = (\alpha_1^<, \dots, \alpha_N^<) \in \Lambda^N, \quad \alpha_1^< < \dots < \alpha_N^<, \quad \text{where } \alpha_j^< \in \alpha.$$

Without loss of generality, we can assume that the orbital indices contained in 0_m are strictly less than the virtual indices $\Lambda \setminus 0_m$.

As in Section 2.1, fix an orthonormal set $\mathcal{B}_K = \{\varphi_p\}_{p \in \Lambda} \subset H^1(\mathbb{R}^3)$ and the corresponding Slater determinants

$$\mathfrak{B}_K = \{\Phi_\alpha \in \mathfrak{H}^1 : \alpha \in S, \Phi_\alpha(\mathbf{X}) = N!^{-1/2} \det(\varphi_{\alpha_i^<}(\mathbf{x}_j))_{1 \leq i, j \leq N}\}. \quad (3.4)$$

Recall the notation $\mathfrak{H}_K^1 \subset \mathfrak{H}^1$ for the subspace spanned by \mathfrak{B}_K ; which is allowed to be finite-, or infinite-dimensional depending on $K = |\Lambda|$.

Definition 3.15. Let $G_m = (L_m, E_m)$ be a subgraph of G_m^{full} . The family of linear operators $X_\alpha^{(m)} := X_\alpha(G_m) : \mathfrak{H}_K^1 \rightarrow \mathfrak{H}_K^1$ given by

$$X_\alpha(G_m)\Phi_\beta = \begin{cases} \sigma(\alpha, \beta)\Phi_{\alpha \vee_m \beta} & (\beta, \alpha \vee_m \beta) \in E_m \\ 0 & (\beta, \alpha \vee_m \beta) \notin E_m \end{cases}$$

for each $\alpha \in \Xi(G_m)$ and $\beta \in S$, and extended boundedly and linearly to the whole space \mathfrak{H}_K^1 (see [31]) is called the family of *excitation operators on G_m* . Here, $\sigma(\alpha, \beta)$ is the sign of the permutation $\pi(\alpha, \beta)$ that puts the N -tuple $((\bar{\beta}^m)^<, (\bar{\alpha}^m)^<)$ in lexicographical order.

Assuming $\Xi(G_m) \neq \emptyset$, by the definition of $\Xi(G_m)$ (see (3.3)) for every $\alpha \in \Xi(G_m)$ there is some $\beta \in L_m$ such that $(\beta, \alpha \vee_m \beta) \in E_m$ and therefore $X_\alpha(G_m) \neq 0$. Recalling $\text{rk}_m(\alpha \vee_m \beta) = \text{rk}_m(\alpha) + \text{rk}_m(\beta)$ (see Proposition 3.3), we can roughly say that an excitation operator $X_\alpha(G_m)$ increases the rank by $\text{rk}_m(\alpha)$.

Since $G_m = (L_m, E_m)$ is a subgraph of $G_m^{\text{full}} = (L_m, E_m^{\text{full}})$, some excitations might be missing, i.e. $\Xi(G_m) \subset \Xi(G_m^{\text{full}})$. The next result shows that the excitation operators constructed for a *consistent subgraph* G_m (see Definition 3.8) are precisely the same as the ones constructed for G_m^{full} , with some of the excitation operators possibly missing. This explains the use of the word “consistent”.

Theorem 3.16. *Let G_m be a consistent subgraph of G_m^{full} . Then,*

$$X_\alpha(G_m) \equiv X_\alpha(G_m^{\text{full}}) \quad \text{for all } \alpha \in \Xi(G_m).$$

Proof. Fix $\alpha \in \Xi(G_m)$, then by (3.3) and Definition 3.8, $(\beta, \alpha \vee_m \beta) \in E_m$ for all $\beta \in L_m$. Consequently, $X_\alpha(G_m)\Phi_\beta = X_\alpha(G_m^{\text{full}})\Phi_\beta$ for all $\beta \in S$. \square

Based on this result, if G_m is consistent, it is safe to drop the “ G_m ” from the notation $X_\alpha(G_m)$ and simply denote the excitation operators by $X_\alpha^{(m)}$, or by X_α in the SR case. However, it is important to note that for a given α , $X_\alpha^{(m)} \neq X_\alpha^{(\ell)}$ in general for differing reference states $m \neq \ell$, see Remark 3.13.

The excitation operators enjoy nice algebraic properties which we summarize in the next theorem (cf. [31, Lemma 2.5]).

Theorem 3.17. *Let $G_m = (L_m, E_m)$ be a consistent subgraph of G_m^{full} and let $\{X_\alpha^{(m)}\}_{\alpha \in \Xi(G_m)}$ denote the set of excitation operators on G_m . Then the following properties hold true.*

(i) (commutativity) *For all $\alpha, \beta \in \Xi(G_m)$, there holds $X_\alpha^{(m)}X_\beta^{(m)} = X_\beta^{(m)}X_\alpha^{(m)}$. In detail, for any $\gamma \in S$,*

$$X_\alpha^{(m)}X_\beta^{(m)}\Phi_\gamma = \begin{cases} \sigma(\alpha, \beta \vee_m \gamma)\sigma(\beta, \gamma)\Phi_{\alpha \vee_m \beta \vee_m \gamma} & (\beta \vee_m \gamma, \alpha \vee_m \beta \vee_m \gamma), \\ & (\gamma, \beta \vee_m \gamma) \in E_m \\ 0 & \text{otherwise} \end{cases}$$

(ii) *If G_m is transitive, then $\{0\} \cup \{\pm X_\alpha^{(m)}\}_{\alpha \in \Xi(G_m)}$ is multiplicatively closed. In particular, $\{0\} \cup \{\pm X_\alpha^{(m)}\}_{\alpha \in \Xi(G_m^{\text{full}})}$ is multiplicatively closed.*

(iii) (nilpotency) *For all $\alpha \in \Xi(G_m)$, $(X_\alpha^{(m)})^2 = 0$.*

Proof. To see (i), first observe that if $(\beta \vee_m \gamma, \alpha \vee_m (\beta \vee_m \gamma)), (\gamma, \beta \vee_m \gamma) \in E_m$, then $(\alpha \vee_m \gamma, \beta \vee_m (\alpha \vee_m \gamma)), (\gamma, \alpha \vee_m \gamma) \in E_m$ due to the consistent subgraph property of G_m . It is obvious that $\Phi_{\alpha \vee_m \beta \vee_m \gamma} = \Phi_{\beta \vee_m \alpha \vee_m \gamma}$ from the commutativity of \vee_m . It remains to prove $\sigma(\alpha, \beta \vee_m \gamma)\sigma(\beta, \gamma) = \sigma(\beta, \alpha \vee_m \gamma)\sigma(\alpha, \gamma)$. Let π_1, π_2 and τ_1, τ_2 be the permutations that put $((\bar{\beta} \cup \bar{\gamma})^<, \bar{\alpha}^<)$, $(\bar{\gamma}^<, \bar{\beta}^<)$ and $((\bar{\alpha} \cup \bar{\gamma})^<, \bar{\beta}^<)$, $(\bar{\gamma}^<, \bar{\alpha}^<)$, respectively, in lexicographic order. Then $\pi_1 \circ \pi_2 = \tau_1 \circ \tau_2 = \sigma$, where σ is the permutation that puts $(\bar{\alpha}, \bar{\beta}, \bar{\gamma})$ in lexicographic order. The claim follows from the multiplicativity of the sgn function on permutations.

For (ii), suppose that G_m is transitive and that $\alpha, \beta \in \Xi(G_m)$. Using (i), either $X_\alpha^{(m)}X_\beta^{(m)} = \pm X_{\alpha \vee_m \beta}^{(m)}$ or $X_\alpha^{(m)}X_\beta^{(m)} = 0$. In the former case, $(\beta \vee_m \gamma, \alpha \vee_m \beta \vee_m \gamma), (\gamma, \beta \vee_m \gamma) \in E_m$ implies that $(\gamma, \alpha \vee_m \beta \vee_m \gamma) \in E_m$ by the transitivity of G_m , so $\alpha \vee_m \beta \in \Xi(G_m)$.

For (iii), it is enough to notice that $(\alpha \vee_m \alpha \vee_m \gamma, \alpha \vee_m \gamma) = (\alpha \vee_m \gamma, \alpha \vee_m \gamma) \notin E_m^{\text{full}}$, because G_m^{full} does not contain loop edges by definition. \square

It is important to note that in general *excitation operators corresponding to different reference states do not commute*: $X_\alpha^{(m)}X_\beta^{(\ell)} \neq X_\beta^{(m)}X_\alpha^{(\ell)}$ for $m \neq \ell$, again, because of Remark 3.13.

Remark 3.18. The excitation operators are traditionally expressed using the language of second quantization. Let a_p^\dagger and a_p denote the fermionic creation and annihilation operators. Then $\Phi_\beta = a_{\beta_1}^\dagger \cdots a_{\beta_N}^\dagger |\text{vac}\rangle$, where $\beta = \{\beta_1 < \dots < \beta_N\}$, and

$$X_\alpha = a_{p_1}^\dagger a_{q_1} \cdots a_{p_n}^\dagger a_{q_n}.$$

Here, $|\text{vac}\rangle$ is the Fock vacuum state, $\{q_1, \dots, q_n\} = 0 \setminus \underline{\alpha}$ and $\{p_1, \dots, p_n\} = \bar{\alpha}$ with $q_1 < \dots < q_n$ and $p_1 < \dots < p_n$. In other words, X_α changes the orbitals $0 \setminus \underline{\alpha}$ to $\bar{\alpha}$, as expected. Although the excitation operators commute with each other, they *do not* commute in general with the Hamiltonian.

We now define a family of operators which “reverse” the action of $X_\alpha^{(m)}$.

Definition 3.19. Let $G_m = (L_m, E_m)$ be a subgraph of G_m^{full} . For all $\alpha \in \Xi(G_m)$, the linear operators $(X_\alpha^{(m)})^\dagger : \mathfrak{H}_K^1 \rightarrow \mathfrak{H}_K^1$ defined via

$$(X_\alpha^{(m)})^\dagger \Phi_\beta = \begin{cases} \sigma(\alpha, \alpha^\perp \wedge_m \beta) \Phi_{\alpha^\perp \wedge_m \beta} & (\alpha^\perp \wedge_m \beta, \beta) \in E_m \\ 0 & (\alpha^\perp \wedge_m \beta, \beta) \notin E_m \end{cases}$$

for any $\beta \in S$, and extended boundedly and linearly to the whole space \mathfrak{H}_K^1 , are called *de-excitation operators* on G_m .

It is easy to see using Lemma 3.5 and Proposition 3.3 that

$$\text{rk}_m(\alpha^\perp \wedge_m \beta) = \text{rk}_m(\beta) - \text{rk}_m(\alpha), \quad (3.5)$$

whenever $(\alpha^\perp \wedge_m \beta, \beta) \in E_m$. Therefore, we may roughly say that the de-excitation operator $(X_\alpha^{(m)})^\dagger$ decreases the rank by $\text{rk}_m(\alpha)$. Of course, the notation \dagger is not coincidental, and $(X_\alpha^{(m)})^\dagger$ is in fact the \mathfrak{L}^2 -adjoint of $X_\alpha^{(m)}$.

Theorem 3.20. *Suppose that $\{X_\alpha^{(m)}\}$ and $\{(X_\alpha^{(m)})^\dagger\}$ are the set of excitation and de-excitation operators corresponding to the excitation graph G_m . Then*

$$\langle (X_\alpha^{(m)})^\dagger \Phi, \Psi \rangle = \langle \Phi, X_\alpha^{(m)} \Psi \rangle \quad \text{for all } \Phi, \Psi \in \mathfrak{H}_K^1 \text{ and } \alpha \in \Xi(G_m).$$

Proof. It is enough to prove the relation for $\Phi = \Phi_\gamma$ and $\Psi = \Phi_\beta$, as the general statement follows by linearity. Suppose that $(\alpha^\perp \wedge_m \gamma, \gamma) \in E_m$, then

$$\langle (X_\alpha^{(m)})^\dagger \Phi_\gamma, \Phi_\beta \rangle = \sigma(\alpha, \alpha^\perp \wedge_m \gamma) \langle \Phi_{\alpha^\perp \wedge_m \gamma}, \Phi_\beta \rangle = \sigma(\alpha, \beta) \langle \Phi_\gamma, \Phi_{\alpha \vee_m \beta} \rangle = \langle \Phi_\gamma, X_\alpha^{(m)} \Phi_\beta \rangle,$$

where we used that $\alpha^\perp \wedge_m \gamma = \beta \in L_m$ if and only if $\alpha \vee_m \beta = \gamma \in L_m$ (Lemma 3.5). \square

Theorem 3.21. *Let $G_m = (L_m, E_m)$ be a consistent subgraph of G_m^{full} and let $\{X_\alpha^{(m)}\}_{\alpha \in \Xi(G_m)}$ and $\{(X_\alpha^{(m)})^\dagger\}_{\alpha \in \Xi(G_m)}$ denote the set of excitation-, and de-excitation operators on G_m . Then the following properties hold true.*

(i) *(commutativity) For all $\alpha, \beta \in \Xi(G_m)$, there holds*

$$(X_\alpha^{(m)})^\dagger (X_\beta^{(m)})^\dagger = (X_\beta^{(m)})^\dagger (X_\alpha^{(m)})^\dagger.$$

(ii) *For any $\alpha, \beta \in \Xi(G_m)$ and $\gamma \in S$, the following formula holds true:*

$$(X_\alpha^{(m)})^\dagger X_\beta^{(m)} \Phi_\gamma = \sigma(\alpha, \alpha^\perp \wedge_m (\beta \vee_m \gamma)) \sigma(\beta, \gamma) \Phi_{\alpha^\perp \wedge (\beta \vee \gamma)}$$

if $(\gamma, \beta \vee_m \gamma) \in E_m$ and $(\alpha^\perp \wedge_m (\beta \vee_m \gamma), \beta \vee_m \gamma) \in E_m$ both hold true. Otherwise, $(X_\alpha^{(m)})^\dagger X_\beta^{(m)} \Phi_\gamma = 0$.

In particular, $(X_\alpha^{(m)})^\dagger \Phi_\alpha = \Phi_{0_m}$.

(iii) *$(X_\alpha^{(m)})^\dagger \Phi_{0_\ell} = 0$ for any $m, \ell = 1, \dots, M$ and $\alpha \in \Xi(G_m)$.*

(iv) *(nilpotency) $((X_\alpha^{(m)})^\dagger)^2 = 0$ for any $\alpha \in \Xi(G_m)$.*

Proof. Part (i) follows from Theorem 3.20 combined with Theorem 3.17 (i). Part (ii) follows directly from the definitions. Part (iii) comes from the fact that there are no edges between different 0_m 's. Part (iv) follows from Theorem 3.17 (iii). \square

It is highly important to stress that in general excitation-, and de-excitation operators do not commute with each other:

$$X_\alpha^{(m)} (X_\alpha^{(m)})^\dagger \neq (X_\alpha^{(m)})^\dagger X_\alpha^{(m)},$$

in other words, the $X_\alpha^{(m)}$'s are *nonnormal* operators. Also, $[(X_\alpha^{(m)})^\dagger, X_\beta^{(m)}] \neq 0$ in general. This fact is the source of many technical obstacles in the analysis of the CC method, primarily because it implies that the similarity-transformed Hamilton operator (2.9) is nonnormal.

3.4. Cluster operators

From now on, we omit the reference index m from the notations, with the understanding that the considerations hold true for every reference independently. Suppose that we constructed the set of excitation operators $\{X_\alpha\}_{\alpha \in \Xi(G)}$ for a given consistent subgraph $G = (L, E)$. The completion of their linear hull

$$\mathfrak{v}(G) = \overline{\text{Span}\{X_\alpha\}_{\alpha \in \Xi(G)}}^{\|\cdot\|_{\mathcal{L}(\mathfrak{H}^1, \mathfrak{H}^1)}}$$

is called the *space of cluster operators on G* endowed with operator norm $\|\cdot\|_{\mathcal{L}(\mathfrak{H}^1, \mathfrak{H}^1)}$. As mentioned earlier, if G is not the full excitation graph G^{full} , then certain excitation operators will be absent and therefore, they will be missing from $\mathfrak{v}(G)$ as well.

Proposition 3.22. *For any $T \in \mathfrak{v}(G)$, we have $T^{N+1} = 0$.*

Proof. It is enough to prove that an arbitrary product of $N+1$ excitation operators is zero. In fact, by definition every excitation operator either increases the rank of a Slater determinant by at least 1 or maps it to zero. But the rank cannot increase above N , so the product must be zero. \square

It is well-known that the vector space $\mathfrak{v}(G^{\text{full}})$ constructed on the full excitation graph G^{full} forms a *commutative algebra* (see e.g. [33, Lemma 4.2]) with the usual multiplication (a subalgebra of the algebra of bounded linear operators $\mathcal{L}(\mathfrak{H}_K^1, \mathfrak{H}_K^1)$). According to Proposition 3.22, it is also *nilpotent*. More generally, we have

Theorem 3.23. *$\mathfrak{v}(G)$ is a nilpotent, commutative algebra for any transitive excitation graph G .*

Proof. Follows from Theorem 3.17 (ii). \square

If, however, G is not transitive, then $\mathfrak{v}(G)$ is *not* an algebra in general—for instance in $\mathfrak{v}(G(\text{SD}))$ there are no excitation operators of rank 3 and above, but the rank of the products of excitation operators can be arbitrary ($\leq N$).

Example 3.24. We observed in Example 3.11 that the CAS-subgraph $G(\text{CAS})$ corresponding to the TCC method is transitive and consistent, hence $\mathfrak{v}(G(\text{CAS}))$ forms a subalgebra of $\mathfrak{v}(G^{\text{full}})$ (cf. [17]). Similarly, for $G(\text{int})$ in Example 3.12, $\mathfrak{v}(G(\text{int}))$ also forms a subalgebra. However, in a truncated setting, where only certain low-rank edges of $E(\text{CAS})$ (or $E(\text{int})$) are retained, transitivity, hence the subalgebra property is lost.

Let now the excitation graph $G = (L, E)$ be arbitrary. A cluster operator $C \in \mathfrak{v}(G)$ may be decomposed according to the excitation ranks of its constituent excitations as

$$C = \sum_{r=1}^N C_r, \quad \text{where} \quad C_r = \sum_{\text{rk}(\alpha)=r} c_\alpha X_\alpha. \quad (3.6)$$

We say that C is of rank r if it contains excitation operators of rank at most r . Note that the graded structure of G is compatible with this decomposition in the sense that if C and D are of ranks r and s , respectively, then CD is of rank at most $r+s$.

Remark 3.25. In the SR case, the cluster operators can be used to express any wavefunction in \mathfrak{H}_K^1 if the full excitation graph G^{full} is used for their construction. In fact, in this case, $X_\alpha \Phi_0 = \Phi_\alpha$ for every $\alpha \in \bar{L}$, so we

may express any function in \mathfrak{H}_K^1 through a linear combination of the excitation operators *and* the identity I . More precisely, if

$$\Psi = \sum_{\alpha \in L} c_\alpha \Phi_\alpha = c_0 \Phi_0 + \sum_{\alpha \in \bar{L}} c_\alpha \Phi_\alpha, \quad \text{then} \quad \Psi = \left[c_0 I + \sum_{\alpha \in \bar{L}} c_\alpha X_\alpha \right] \Phi_0,$$

for some scalars $\{c_\alpha\}_{\alpha \in L}$. Recall that in Section 2.3 we assumed the intermediate normalization condition $\langle \Psi, \Phi_0 \rangle = 1$, which implies $c_0 = 1$. There is a one-to-one correspondence between functions $\Psi \in \mathfrak{H}_K^{1,\perp}$ and the cluster operators C_Ψ defined as

$$C_\Psi = \sum_{\alpha \in \bar{L}} c_\alpha X_\alpha, \quad \text{where} \quad c_\alpha = \langle \Psi, \Phi_\alpha \rangle. \quad (3.7)$$

It is not clear, however, that $C_\Psi \in \mathcal{L}(\mathfrak{H}_K^1, \mathfrak{H}_K^1)$. See Theorem 3.26 below for the precise statement of this nontrivial fact. Also, if the excitation graph does not contain every edge of the form $(0, \alpha)$ —which is typically the case if some truncation is used—then it is *not* possible to assign a cluster operator (3.7) to every $\Psi \in \mathfrak{H}_K^{1,\perp}$.

The following important result makes the aforementioned correspondence between functions and cluster operators precise.

Theorem 3.26. [31, Theorem 4.1 and Lemma 5.1] *Fix $\Psi \in \mathfrak{H}^{1,\perp}$. Then, the following hold true.*

(1) *The cluster operator C_Ψ (3.7) satisfies $C_\Psi \in \mathcal{L}(\mathfrak{H}^1, \mathfrak{H}^1)$. Furthermore, there is a constant $b > 0$ independent of Ψ such that*

$$\|\Psi\|_{\mathfrak{H}^1} \leq \|C_\Psi\|_{\mathcal{L}(\mathfrak{H}^1, \mathfrak{H}^1)} \leq b \|\Psi\|_{\mathfrak{H}^1}.$$

(2) *$C_\Psi^\dagger \in \mathcal{L}(\mathfrak{H}^1, \mathfrak{H}^1)$, and there is a constant $b' > 0$ independent of Ψ 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 $\|\Psi\|_{\mathfrak{H}^1}$.

(3) *C_Ψ can be extended to $\mathcal{L}(\mathfrak{H}^{-1}, \mathfrak{H}^{-1})$.*

Next, we consider the so-called *exponential Ansatz*, which is the representation

$$I + C = e^T, \quad \text{where} \quad T = \sum_{\alpha \in \Xi(G^{\text{full}})} t_\alpha X_\alpha \in \mathfrak{v}(G^{\text{full}}),$$

and $C \in \mathfrak{v}(G^{\text{full}})$. Here, e^T is simply a finite sum due to the nilpotency of T , i.e.

$$e^T = I + T + \frac{1}{2!}T^2 + \dots + \frac{1}{N!}T^N.$$

The inverse of the exponential should be the logarithm, as one would expect.

Theorem 3.27. [31, Lemma 5.2] *For any cluster operator $C \in \mathfrak{v}(G^{\text{full}})$ there exists a unique cluster operator $T \in \mathfrak{v}(G^{\text{full}})$, such that $e^T = I + C$. Furthermore,*

$$T = \log(I + C) = C - \frac{1}{2}C^2 + \frac{1}{3}C^3 - \dots + \frac{(-1)^{N-1}}{N}C^N.$$

Moreover, the exponential map is a bijection between

$$\mathcal{S} = \left\{ S \in \mathcal{L}(\mathfrak{H}^1, \mathfrak{H}^1) : S = \sum_{\alpha \in \bar{L}} s_\alpha X_\alpha \right\} \quad \text{and} \quad I + \mathcal{S}.$$

Furthermore, the result also holds true if $\mathcal{L}(\mathfrak{H}^1, \mathfrak{H}^1)$ is replaced with $\mathcal{L}(\mathfrak{H}^{-1}, \mathfrak{H}^{-1})$.

It is important to note that if some proper excitation subgraph $G = (L, E)$ is considered instead of G^{full} , the previous result does *not* hold. For instance, if $G(\text{SD})$ is considered, then it might not be possible to represent $I + C$ as e^T , where $C \in \mathfrak{v}(G^{\text{full}})$ and $T \in \mathfrak{v}(G)$. This in particular implies that wavefunctions of the form $e^T \Phi_0$ where $T \in \mathfrak{v}(G)$ is *not* the totality of intermediately normalized wavefunctions.

In the multireference (MR) case, the analogue of the exponential Ansatz is called the *Jeziorski–Monkhorst (JM) Ansatz*, see Section 4.2 below. In the JM-MRCC method, M wavefunctions, say Ψ_1, \dots, Ψ_M are “targeted”, and the expansion

$$\Psi_j = \sum_{m=1}^M a_j^{(m)} e^{T^{(m)}} \Phi_{0_m}, \quad \text{where } a_j^{(m)} \in \mathbb{R}, \quad (3.8)$$

is utilized. In the untruncated case, suppose that $\Psi_j = (I + C^{(j)}) \Phi_{0_j} = e^{T^{(j)}} \Phi_{0_j}$, as above, for all $j = 1, \dots, M$. Then the JM expansion coefficients $a_j^{(m)}$ of Ψ_j are simply δ_{jm} .

3.5. Cluster amplitude spaces

The linear combination coefficients of the excitation operators making up a cluster operator are called *cluster amplitudes*. Let $\ell^2(G)$ denote Hilbert space of square summable real-, or complex-valued sequences indexed by the edge labels of the excitation graph G , i.e.

$$\ell^2(G) = \{t = (t_\alpha)_{\alpha \in \Xi(G)} : \|t\|_{\ell^2} < \infty\}.$$

The (real or complex) Hilbert space

$$\mathbb{V}(G) = \{t \in \ell^2(G) : \|T\Phi_0\|_{\mathfrak{H}^1} < \infty\},$$

endowed with the \mathfrak{H}^1 -inner product $\langle t, s \rangle_{\mathbb{V}} = \langle T\Phi_0, S\Phi_0 \rangle_{\mathfrak{H}^1}$ is called the (*cluster*) *amplitude space corresponding to G* . Nevertheless, from now on we use the convention that the unmarked $\langle t, s \rangle = \langle T\Phi_0, S\Phi_0 \rangle_{\mathfrak{L}^2}$ and $\|\cdot\|$ refers to the ℓ^2 -inner product and ℓ^2 -norm. Clearly, $\|t\| \leq \|t\|_{\mathbb{V}}$.

Remark 3.28. Similarly to $\mathfrak{H}^1 \hookrightarrow \mathfrak{L}^2 \hookrightarrow \mathfrak{H}^{-1}$, the spaces $\mathbb{V}(G) \hookrightarrow \ell^2(G) \hookrightarrow \mathbb{V}(G)^*$ also form a Gelfand triple.

It is clear that the space of cluster operators $\mathfrak{v}(G)$ is canonically isomorphic to $\mathbb{V}(G)$ via

$$\mathfrak{v}(G) \ni \sum_{\alpha \in \Xi(G)} c_\alpha X_\alpha = C \mapsto c = (c_\alpha)_{\alpha \in \Xi(G)} \in \mathbb{V}(G).$$

As customary in CC theory, we will never explicitly denote this isomorphism, and instead use capital letters S, T, U, V, W , etc. to denote the cluster operators and small letters s, t, u, v, w , etc. to denote their corresponding cluster amplitudes.

Furthermore, to every amplitude space $\mathbb{V}(G)$ there corresponds a *functional amplitude space* $\mathfrak{A}(G) \subset \mathfrak{H}^{1,\perp}$ through the (ℓ^2, \mathfrak{L}^2) -isometric isomorphism $\mathbb{V}(G) \rightarrow \mathfrak{A}(G)$ given by

$$\mathbb{V}(G) \ni c \mapsto C\Phi_0 = \sum_{\alpha \in \Xi(G)} c_\alpha \Phi_\alpha \in \mathfrak{A}(G).$$

Clearly, an appropriate subset of the Slater determinant basis \mathfrak{B}_K (see (3.4)) forms a basis of the functional amplitude space $\mathfrak{A}(G)$.

Given a closed subspace $\mathfrak{U} \subset \mathfrak{A}(G)$, we will sometimes use the orthogonal projector $\Pi_{\mathfrak{U}} : \mathfrak{L}^2 \rightarrow \mathfrak{U} \subset \mathfrak{L}^2$ onto \mathfrak{U} , defined as

$$\langle \Pi_{\mathfrak{U}} \Psi, \Phi \rangle = \langle \Psi, \Phi \rangle, \quad \text{for all } \Psi \in \mathfrak{L}^2, \Phi \in \mathfrak{U}.$$

Hence, the inclusion map $I_{\mathfrak{U}} : \mathfrak{U} \rightarrow \mathfrak{L}^2$, given by $I_{\mathfrak{U}} \Phi = \Phi$ for all $\Phi \in \mathfrak{U}$ satisfies $I_{\mathfrak{U}}^\dagger = \Pi_{\mathfrak{U}}$.

We continue by recalling an important notion due to [33].

Definition 3.29. The excitation graph G is said to be *excitation complete*, if $\alpha^\perp \wedge \beta \in \Xi(G)$ for all $\alpha, \beta \in \Xi(G)$ with $(\alpha^\perp \wedge \beta, \beta) \in E$ and $\alpha \neq \beta$.

It is easy to see using (3.5), that commonly used rank-truncated graphs such as $G(1, 2, \dots, \rho)$ and $G(D)$ are excitation complete.

Proposition 3.30. [33, Lemma 5.5] Suppose that G is excitation complete, let $\mathfrak{V} = \mathfrak{V}(G)$ and $\mathfrak{V}_0 = \text{Span}\{\Phi_0\} \oplus \mathfrak{V}$. Fix $t \in \mathbb{V}$.

- (i) The linear mappings $e^{\pm T^\dagger} I_{\mathfrak{V}_0} : \mathfrak{V}_0 \rightarrow \mathfrak{V}_0$ are bijective.
- (ii) The linear mappings $\Pi_{\mathfrak{V}} e^{\pm T^\dagger} I_{\mathfrak{V}} : \mathfrak{V} \rightarrow \mathfrak{V}$ are surjective.

The result follows easily from the next lemma.

Lemma 3.31. [33, Lemma 5.4] Suppose that G is excitation complete. Then, for every $\alpha, \beta \in \Xi(G)$ we have $X_\alpha^\dagger \Phi_\beta \in \mathfrak{V}(G) \cup \{\Phi_0\}$.

Proof. From Theorem 3.21 (ii), we have

$$X_\alpha^\dagger \Phi_\beta = \sigma(\alpha, \alpha^\perp \wedge \beta) \Phi_{\alpha^\perp \wedge \beta},$$

if $(\alpha^\perp \wedge \beta, \beta) \in E$. If $\alpha \neq \beta$, then right-hand side is in $\mathfrak{V}(G)$, since G is excitation complete. If $\alpha = \beta$, then the right-hand side is simply Φ_0 . \square

Proof of Proposition 3.30. By linearity, Lemma 3.31 implies that the mapping $T^\dagger : \mathfrak{V}_0(G) \rightarrow \mathfrak{V}_0(G)$ and so $e^{\pm T^\dagger} : \mathfrak{V}_0(G) \rightarrow \mathfrak{V}_0(G)$ as well. But $(e^{T^\dagger})^{-1} = e^{-T^\dagger}$, which proves (i). Part (ii) follows easily from this. \square

4. DERIVATION OF THE COUPLED-CLUSTER EQUATIONS

In this section, we give derivations of the SRCC-, and a variant of the MRCC equations. The approach presented here is based on [27]. We would like to stress that the discussion only applies to the *full* (that is, untruncated) CC methods.

The essence of the following theorem seems to be well-known in the physics and quantum chemistry literature, and the method itself is generally attributed to C. Bloch [6], who devised it in the context of perturbation theory.

Theorem 4.1. Let \mathfrak{H} and \mathfrak{L} be (real or complex) Hilbert spaces so that they form a Gelfand triple: $\mathfrak{H} \subset \mathfrak{L} \subset \mathfrak{H}^*$. Let $\mathcal{H} : \mathfrak{H} \rightarrow \mathfrak{H}^*$ be a bounded operator. Let $\mathfrak{M}, \mathfrak{N} \subset \mathfrak{H}$ be any pair of closed subspaces so that the following complementarity condition holds:

$$\mathfrak{M} \oplus \mathfrak{N}^\perp = \mathfrak{H}. \quad (4.1)$$

Then the following are equivalent.

- (i) $\mathfrak{M} \subset \mathfrak{H}$ is weakly \mathcal{H} -invariant: for every $\Phi \in \mathfrak{M}$ there exists $\tilde{\Phi} \in \mathfrak{M}$ such that $\langle \mathcal{H}\Phi, \Phi' \rangle = \langle \tilde{\Phi}, \Phi' \rangle$ for all $\Phi' \in \mathfrak{H}$.
- (ii) (weak Bloch equation) There holds

$$\langle \mathcal{H}\Xi\Phi, (I - \Xi^\dagger)\Phi' \rangle = 0 \quad \text{for all } \Phi \in \mathfrak{N}, \Phi' \in \mathfrak{N}^\perp, \quad (4.2)$$

where $\Xi : \mathfrak{H} \rightarrow \mathfrak{H}$ denotes the (oblique) projector onto \mathfrak{M} along \mathfrak{N}^\perp , i.e. $\text{ran } \Xi = \mathfrak{M}$ and $\ker \Xi = \mathfrak{N}^\perp$.

Furthermore, if

$$\mathfrak{M} = \text{Span}\{\Psi_j \in \mathfrak{H} : j = 1, \dots, J\}, \quad \text{where } \langle \mathcal{H}\Psi_j, \bar{\Phi} \rangle = \mathcal{E}_j \langle \Psi_j, \bar{\Phi} \rangle \quad (\bar{\Phi} \in \mathfrak{H}) \quad (4.3)$$

for some $\mathcal{E}_j \in \mathbb{C}$, then with the effective Hamiltonian $\mathcal{H}^{\text{eff}} : \mathfrak{N} \rightarrow \mathfrak{N}$, given by $\langle \mathcal{H}^{\text{eff}}\Phi, \Phi' \rangle = \langle \mathcal{H}\Xi\Phi, \Phi' \rangle$ for all $\Phi, \Phi' \in \mathfrak{N}$, we have

$$\langle \mathcal{H}^{\text{eff}}\Pi\Psi_j, \Phi \rangle = \mathcal{E}_j \langle \Pi\Psi_j, \Phi \rangle \quad \text{for all } \Phi \in \mathfrak{N}, \quad (4.4)$$

where $\Pi : \mathfrak{H} \rightarrow \mathfrak{H}$ denotes the \mathfrak{L} -orthogonal projector onto \mathfrak{N} , i.e. $\text{ran } \Pi = \mathfrak{N}$ and $\ker \Pi = \mathfrak{N}^\perp$.

Proof. For (i) \implies (ii), note that using $\ker \Xi = \mathfrak{N}^\perp$ and $\text{ran } \Xi = \mathfrak{M}$, it follows from (i) that for every $\Phi \in \mathfrak{N}$ there exists $\tilde{\Phi} \in \mathfrak{M}$ such that $\langle \mathcal{H}\Xi\Phi, \bar{\Phi} \rangle = \langle \tilde{\Phi}, \bar{\Phi} \rangle$ for all $\bar{\Phi} \in \mathfrak{H}$. Put $\bar{\Phi} = (I - \Xi^\dagger)\Phi'$ to obtain

$$\langle \mathcal{H}\Xi\Phi, (I - \Xi^\dagger)\Phi' \rangle = \langle \tilde{\Phi}, (I - \Xi^\dagger)\Phi' \rangle = 0 \quad \text{for all } \Phi \in \mathfrak{N}, \Phi' \in \mathfrak{H},$$

where we used that $\tilde{\Phi} \in \mathfrak{M}$ and $\text{ran}(I - \Xi^\dagger) = \mathfrak{M}^\perp$. From this, (4.2) follows.

To see (ii) \implies (i), fix $\Phi \in \mathfrak{M}$ and note that (4.2) implies $F_\Phi(\Phi') = 0$ for all $\Phi' \in \mathfrak{M}^\perp$, where $F_\Phi(\Phi') := \langle \mathcal{H}\Phi, \Phi' \rangle$ for all $\Phi' \in \mathfrak{H}$. Here, $F_\Phi(\cdot)$ is a bounded linear functional on the dense subspace $\mathfrak{H} \subset \mathfrak{L}$. Extend F_Φ to a bounded linear functional \widehat{F}_Φ on \mathfrak{L} . The Riesz representation theorem implies that there is a $\tilde{\Phi} \in \mathfrak{L}$ such that $\widehat{F}_\Phi(\Phi') = \langle \tilde{\Phi}, \Phi' \rangle$ for all $\Phi' \in \mathfrak{L}$. But $0 = F_\Phi(\Phi') = \widehat{F}_\Phi(\Phi') = \langle \tilde{\Phi}, \Phi' \rangle$ for all $\Phi' \in \mathfrak{M}^\perp$, so $\tilde{\Phi} \in \mathfrak{M}^{\perp\perp} = \mathfrak{M}$. Therefore, we constructed a $\tilde{\Phi} \in \mathfrak{M}$ such that $\langle \mathcal{H}\Phi, \Phi' \rangle = \langle \tilde{\Phi}, \Phi' \rangle$ for all $\Phi' \in \mathfrak{H}$, which is what we wanted to prove.

To prove the ‘‘furthermore’’ part, first note that \mathfrak{M} is weakly \mathcal{H} -invariant. We now claim that $\Xi\Pi = \Xi$. In fact, $\text{ran}(I - \Pi) = \ker \Pi = \ker \Xi$, so $\Xi(I - \Pi) = 0$. Continuing the proof, note that the second relation of (4.3) is equivalent to

$$\langle \mathcal{H}\Xi\Pi\Psi_j, \bar{\Phi} \rangle = \mathcal{E}_j \langle \Xi\Pi\Psi_j, \bar{\Phi} \rangle \quad \text{for all } \bar{\Phi} \in \mathfrak{H}.$$

Using (4.2), this can be further written as

$$\langle \mathcal{H}\Xi\Pi\Psi_j, \Xi^\dagger\bar{\Phi} \rangle = \mathcal{E}_j \langle \Pi\Psi_j, \Xi^\dagger\bar{\Phi} \rangle \quad \text{for all } \bar{\Phi} \in \mathfrak{H}.$$

The desired result follows by noting that $\text{ran } \Xi^\dagger = \mathfrak{N}$. \square

In practice, \mathfrak{M} (called the ‘‘exact model space’’) is unknown and \mathfrak{N} (called the ‘‘model space’’) is chosen in a way that it provides a ‘‘reasonable approximation’’ to \mathfrak{M} , i.e. that (4.1) holds. In particular, $\mathfrak{M} \subset \mathfrak{N}^\perp$ is not permitted. Then, the unknown ‘‘wave operator’’ Ξ (hence \mathfrak{M}) can be determined by solving the weak Bloch equation (4.2). Next, the eigenvalue problem for \mathcal{H}^{eff} is solved to obtain the energies $\mathcal{E}_1, \dots, \mathcal{E}_M$ and (some of the) eigenvectors.

Remark 4.2.

- (i) It is important to note that solving the Bloch equation only provides a weakly \mathcal{H} -invariant subspace \mathfrak{M} and it might *not* be a direct sum of (weak) eigenspaces in general. In other words, \mathfrak{M} might be spanned by an incomplete set of eigenvectors. Clearly, in such a situation some of the eigenvectors cannot be recovered through solving the eigenproblem for the effective Hamiltonian \mathcal{H}^{eff} .
- (ii) The Bloch equation (4.2) is more commonly given in the ‘‘strong’’ form ‘‘ $\Xi\mathcal{H}\Xi = \mathcal{H}\Xi$ ’’.

The situation is greatly simplified, when one considers one-dimensional subspaces \mathfrak{N} and \mathfrak{M} , because a one-dimensional invariant subspace is always an eigenspace.

Corollary 4.3. *Let $\dim \mathfrak{N} = \dim \mathfrak{M} = 1$, and set $\mathfrak{N} = \text{Span}\{\Phi_0\}$ for some $\Phi_0 \in \mathfrak{H}$. Further, let $\mathfrak{M} = \text{Span}\{\Psi\}$ for some $\Psi \in \mathfrak{H}$, and suppose that $\langle \Psi, \Phi_0 \rangle = 1$. Then, the following are equivalent.*

- (i) $\langle \mathcal{H}\Psi, \bar{\Phi} \rangle = \mathcal{E} \langle \Psi, \bar{\Phi} \rangle$ for all $\bar{\Phi} \in \mathfrak{H}$ and some scalar \mathcal{E} .
- (ii) $\langle \mathcal{H}\Xi\Phi_0, (I - \Xi^\dagger)\Phi' \rangle = 0$ for all $\Phi' \in \mathfrak{N}^\perp$.

Furthermore, $\mathcal{E} = \langle \mathcal{H}\Xi\Phi_0, \Phi_0 \rangle$.

4.1. The SRCC method

The single-reference Coupled-Cluster method easily follows from Corollary 4.3 through the exponential parametrization of the wave operator. In the following theorem, we re-establish [31, Theorem 5.3] (see Theorem 2.2).

Theorem 4.4. *Let $\mathcal{H} : \mathfrak{H}_K^1 \rightarrow (\mathfrak{H}_K^1)^*$ be a bounded operator. Fix $\Phi_0 \in \mathfrak{H}_K^1$ with $\|\Phi_0\| = 1$ and suppose that $\Psi \in \mathfrak{H}_K^1$ is such that $\langle \Psi, \Phi_0 \rangle = 1$. Then the following are equivalent.*

- (i) $\langle \mathcal{H}\Psi, \Phi \rangle = \mathcal{E}\langle \Psi, \Phi \rangle$ for all $\Phi \in \mathfrak{H}_K^1$ for some scalar \mathcal{E} .
(ii) (Full CC) $\Psi = e^{T_*}\Phi_0$ for some $t_* \in \mathbb{V}(G^{\text{full}})$ such that

$$\langle e^{-T_*}\mathcal{H}e^{T_*}\Phi_0, S\Phi_0 \rangle = 0 \quad \text{for all } s \in \mathbb{V}(G^{\text{full}}). \quad (4.5)$$

Furthermore, $\mathcal{E} = \langle e^{-T_*}\mathcal{H}e^{T_*}\Phi_0, \Phi_0 \rangle$.

- (iii) (Full CI) $\Psi = (I + C_*)\Phi_0$ for some $c_* \in \mathbb{V}(G^{\text{full}})$ such that

$$\langle \mathcal{H}(I + C_*)\Phi_0, S\Phi_0 \rangle = \mathcal{E}_{\text{CI}}\langle (I + C_*)\Phi_0, S\Phi_0 \rangle \quad \text{for all } s \in \mathbb{V}(G^{\text{full}}), \quad (4.6)$$

where $\mathcal{E}_{\text{CI}} = \langle \mathcal{H}(I + C_*)\Phi_0, \Phi_0 \rangle$. Furthermore, $\mathcal{E} = \mathcal{E}_{\text{CI}}$.

Proof. Let $\mathfrak{H} = \mathfrak{H}_K^1$ and $\mathfrak{L} = \mathfrak{L}^2$. First, we prove (i) \iff (ii). We apply Corollary 4.3 with the SRCC wave operator

$$\Xi = e^{T_*}\Pi_{\Phi_0},$$

where T_* is some cluster operator and Π_{Φ_0} is the orthogonal projector onto $\mathfrak{N} = \text{Span}\{\Phi_0\}$. Note that $\mathfrak{N}^\perp = \mathfrak{V}(G^{\text{full}})$. It is easy to see that Ξ is idempotent, and that $\ker \Xi = \mathfrak{N}^\perp$. By an appropriate choice of T_* , $\text{ran } \Xi = \mathfrak{M}$ using $\langle \Psi, \Phi_0 \rangle = 1$ and Theorem 3.27. Furthermore, $\text{Span}\{e^{T_*}\Phi_0\} = \text{ran } \Xi \subset \mathfrak{H}$ due to Theorem 3.26. Applying Corollary 4.3, (i) holds if and only if $\Psi = e^{T_*}\Phi_0$ and T_* satisfies the weak Bloch equation

$$\langle \mathcal{H}e^{T_*}\Phi_0, (I - \Pi_{\Phi_0}e^{T_*^\dagger})S'\Phi_0 \rangle = 0 \quad \text{for all } s' \in \mathbb{V}(G^{\text{full}}).$$

Recalling Proposition 3.30 (ii), and using the change of variables $S' = e^{-T_*^\dagger}S$,

$$\langle e^{-T_*}\mathcal{H}e^{T_*}\Phi_0, S\Phi_0 \rangle = 0 \quad \text{for all } s \in \mathbb{V}(G^{\text{full}}).$$

Here we used that e^{-T_*} can be extended to a bounded $\mathfrak{H}^{-1} \rightarrow \mathfrak{H}^{-1}$ operator (Theorem 3.26).⁷ Note that \mathcal{H}^{eff} is now a one-dimensional linear map (i.e. a multiplication by a scalar), so $\sigma(\mathcal{H}^{\text{eff}}) = \langle e^{-T_*}\mathcal{H}e^{T_*}\Phi_0, \Phi_0 \rangle = \mathcal{E}$.

Next, we prove (i) \iff (iii). We now apply Corollary 4.3 with the SRCI wave operator

$$\Xi = (I + C_*)\Pi_{\Phi_0},$$

where C_* is some cluster operator and the claim follows from a straightforward calculation. Further, now $\sigma(\mathcal{H}^{\text{eff}}) = \langle \mathcal{H}(I + C_*)\Phi_0, \Phi_0 \rangle = \mathcal{E}$. \square

4.2. The Jeziorski–Monkhorst MRCC method

In MRCC methods the “model space” \mathfrak{N} is chosen to be the space spanned by M orthonormal reference determinants,

$$\mathfrak{N} = \text{Span}\{\Phi_{0_m} : m = 1, \dots, M\}.$$

The *Jeziorski–Monkhorst method* [16] uses the following Ansatz for the wave operator:

$$\Xi = \sum_{m=1}^M e^{T^{(m)}}\Pi_{\Phi_{0_m}}, \quad (4.7)$$

which corresponds to (3.8).

Theorem 4.5. *Let \mathfrak{N} be defined as above and set $\mathfrak{M} = \text{Span}\{\Psi_m : m = 1, \dots, M\}$, where $\{\Psi_m\}_{m=1}^M \subset \mathfrak{H}_K^1$ is \mathfrak{L}^2 -orthogonal. Suppose that for every $m = 1, \dots, M$, $\langle \Psi_m, \Phi_{0_n} \rangle \neq 0$ for at least one $n = 1, \dots, M$. Then, the following are equivalent.*

⁷We refer the reader to the proof of [31, Theorem 5.3] for more details.

- (i) \mathfrak{M} is weakly \mathcal{H} -invariant: for every Ψ_m ($m = 1, \dots, M$) there exists $\tilde{\Psi}_m \in \mathfrak{M}$ such that $\langle \mathcal{H}\Psi_m, \Phi' \rangle = \langle \tilde{\Psi}_m, \Phi' \rangle$ for all $\Phi' \in \mathfrak{H}_K^1$.
- (ii) (Full JM-MRCC) $\mathfrak{M} = \text{Span}\{e^{T_*^{(m)}}\Phi_{0_m} : m = 1, \dots, M\}$, where $t_*^{(m)} \in \mathbb{V}(G_m^{\text{full}})$ satisfies

$$\langle e^{-T_*^{(m)}}\mathcal{H}e^{T_*^{(m)}}\Phi_{0_m}, S^{(m)}\Phi_{0_m} \rangle = \sum_{n=1}^M \mathcal{H}_{mn}^{\text{eff}} \langle e^{-T_*^{(m)}}e^{T_*^{(n)}}\Phi_{0_n}, S^{(m)}\Phi_{0_m} \rangle, \quad (4.8)$$

for all $s^{(m)} \in \mathbb{V}(G_m^{\text{full}})$ and $m = 1, \dots, M$, where the matrix elements of the effective Hamiltonian are given by $\mathcal{H}_{mn}^{\text{eff}} = \langle e^{-T_*^{(m)}}\mathcal{H}e^{T_*^{(n)}}\Phi_{0_n}, \Phi_{0_m} \rangle$.

- (iii) (Full MRCI) $\mathfrak{M} = \text{Span}\{(I + C_*^{(m)})\Phi_{0_m} : m = 1, \dots, M\}$, where $c_*^{(m)} \in \mathbb{V}(G_m^{\text{full}})$ satisfies

$$\langle \mathcal{H}(I + C_*^{(m)})\Phi_{0_m}, S^{(m)}\Phi_{0_m} \rangle = \sum_{n=1}^M \hat{\mathcal{H}}_{mn}^{\text{eff}} \langle (I + C_*^{(n)})\Phi_{0_n}, S^{(m)}\Phi_{0_m} \rangle, \quad (4.9)$$

for all $s^{(m)} \in \mathbb{V}(G_m^{\text{full}})$ and $m = 1, \dots, M$, where the matrix elements of the effective Hamiltonian are given by $\hat{\mathcal{H}}_{mn}^{\text{eff}} = \langle \mathcal{H}(I + C_*^{(n)})\Phi_{0_n}, \Phi_{0_m} \rangle$.

Furthermore, suppose that $\langle \mathcal{H}\Psi_m, \bar{\Phi} \rangle = \mathcal{E}_m \langle \Psi_m, \bar{\Phi} \rangle$ for all $\bar{\Phi} \in \mathfrak{H}_K^1$ and $m = 1, \dots, M$. Then the following hold true.

- (a) Suppose \mathfrak{M} is given as in (ii). Then the coefficients $a_j^{(m)}$ in the expansion $\Psi_j = \sum_{n=1}^M a_j^{(n)} e^{T_*^{(n)}}\Phi_{0_n}$ are given as the solution to the eigenvalue problem

$$\sum_{n=1}^M \mathcal{H}_{nm}^{\text{eff}} a_j^{(n)} = \mathcal{E}_j a_j^{(m)} \quad \text{where } m = 1, \dots, M.$$

- (b) Suppose \mathfrak{M} is given as in (iii). Then the coefficients $\hat{a}_j^{(m)}$ in the expansion $\Psi_j = \sum_{n=1}^M \hat{a}_j^{(n)} (I + C_*^{(n)})\Phi_{0_n}$ are given as the solution to the eigenvalue problem

$$\sum_{n=1}^M \hat{\mathcal{H}}_{nm}^{\text{eff}} \hat{a}_j^{(n)} = \mathcal{E}_j \hat{a}_j^{(m)} \quad \text{where } m = 1, \dots, M.$$

Proof. Let $\mathfrak{H} = \mathfrak{H}_K^1$. First, we prove (i) \iff (ii) by applying Theorem 4.4. Clearly, for the JM wave operator (4.7) we have $\Xi^2 = \Xi$ and $\ker \Xi = \mathfrak{N}^\perp$ and

$$\text{ran } \Xi = \text{Span}\{e^{T_*^{(m)}}\Phi_{0_m} : m = 1, \dots, M\}.$$

The weak Bloch equation (4.2) is equivalent to

$$\langle \mathcal{H}e^{T_*^{(m)}}\Phi_{0_m}, \Phi' \rangle = \sum_{n=1}^M \langle \mathcal{H}e^{T_*^{(m)}}\Phi_{0_m}, \Pi_{\Phi_{0_n}} e^{(T_*^{(n)})^\dagger} \Phi' \rangle$$

for all $\Phi' \in \mathfrak{N}^\perp$ and $m = 1, \dots, M$. Setting $\Phi' = S^{(m)}\Phi_{0_m}$, we obtain

$$\begin{aligned} \langle \mathcal{H}e^{T_*^{(m)}}\Phi_{0_m}, S^{(m)}\Phi_{0_m} \rangle &= \sum_{n=1}^M \langle \mathcal{H}e^{T_*^{(m)}}\Phi_{0_m}, \Pi_{\Phi_{0_n}} e^{(T_*^{(n)})^\dagger} S^{(m)}\Phi_{0_m} \rangle \\ &= \sum_{n=1}^M \langle \mathcal{H}e^{T_*^{(m)}}\Phi_{0_m}, \Phi_{0_n} \rangle \langle e^{(T_*^{(n)})^\dagger} S^{(m)}\Phi_{0_m}, \Phi_{0_n} \rangle \\ &= \sum_{n=1}^M \langle e^{-T_*^{(m)}} \mathcal{H}e^{T_*^{(m)}}\Phi_{0_m}, \Phi_{0_n} \rangle \langle e^{T_*^{(n)}}\Phi_{0_n}, S^{(m)}\Phi_{0_m} \rangle \end{aligned}$$

for all $s^{(m)} \in \mathbb{V}(G_m^{\text{full}})$. Here, we used that $(T^{(m)})^\dagger\Phi_{0_n} = 0$, see Theorem 3.21 (iii). The proof of (4.8) is finished by invoking Proposition 3.30 (ii) and replacing $S^{(m)}$ by $(e^{-T_*^{(m)}})^\dagger S^{(m)}$.

Next, we prove (i) \iff (iii). The MRCI wave operator reads

$$\Xi = \sum_{m=1}^M (I + C_*^{(m)}) \Pi_{\Phi_{0_m}}.$$

With this choice (4.2) is equivalent to

$$\langle \mathcal{H}(I + C_*^{(m)})\Phi_{0_m}, \Phi' \rangle = \sum_{n=1}^M \langle \mathcal{H}(I + C_*^{(m)})\Phi_{0_m}, \Pi_{\Phi_{0_n}} (I + C_*^{(n)})^\dagger \Phi' \rangle$$

for all $\Phi' \in \mathfrak{N}^\perp$ and $m = 1, \dots, M$. Setting $\Phi' = S^{(m)}\Phi_{0_m}$, this can be written as

$$\begin{aligned} \langle \mathcal{H}(I + C_*^{(m)})\Phi_{0_m}, S^{(m)}\Phi_{0_m} \rangle &= \sum_{n=1}^M \langle \mathcal{H}(I + C_*^{(m)})\Phi_{0_m}, \Pi_{\Phi_{0_n}} (I + C_*^{(n)})^\dagger S^{(m)}\Phi_{0_m} \rangle \\ &= \sum_{n=1}^M \langle \mathcal{H}(I + C_*^{(m)})\Phi_{0_m}, \Phi_{0_n} \rangle \langle (I + C_*^{(n)})^\dagger S^{(m)}\Phi_{0_m}, \Phi_{0_n} \rangle \\ &= \sum_{n=1}^M \langle \mathcal{H}(I + C_*^{(m)})\Phi_{0_m}, \Phi_{0_n} \rangle \langle (I + C_*^{(n)})\Phi_{0_n}, S^{(m)}\Phi_{0_m} \rangle, \end{aligned}$$

which is what we wanted to prove.

For the ‘‘furthermore’’ part of (a), expanding Ψ_j as $\Psi_j = \sum_{n=1}^M a_j^{(n)} e^{T_*^{(n)}}\Phi_{0_n}$, for some scalars $a_j^{(n)}$, we find that $a_j^{(m)} = \langle \Psi_j, \Phi_{0_m} \rangle$. It is easy to see that (4.4) now reads

$$\sum_{n=1}^M \langle \mathcal{H}e^{T_*^{(n)}}\Phi_{0_n}, \Phi_{0_m} \rangle a_j^{(n)} = \mathcal{E}_j a_j^{(m)}$$

for all $j = 1, \dots, M$. The proof of the ‘‘furthermore’’ part of (b) is similar. \square

5. CONCLUSIONS AND FURTHER WORK

In this first part of a series of two articles, we proposed a framework to describe the discretization scheme involved in CC-like methods. At the core of the description is the concept of the excitation graph (Definition 3.4),

which completely determines all necessary building blocks such as excitation operators (Section 3.3), cluster operators (Section 3.4) and cluster amplitude spaces (Section 3.5). The excitation graph concept admits a straightforward extension to the multireference case (Definition 3.14). Another advantage of our approach is that it avoids the use of second-quantized formalism and hence allowed us to prove the basic results (such as Theorem 3.17 and Theorem 3.21) in a more transparent manner. Besides these, we also pointed out a number of structural properties of the excitation graph in Section 3.2. It is important to note that some of these graph-theoretic properties are reflected in the algebraic structure of the excitation operators (Theorem 3.16 and Theorem 3.23). Some relevant combinatorial quantities have been calculated in Appendix A. Furthermore, we proposed an algorithm to determine the reference states in an optimal fashion for the multireference case in Appendix B.

In Section 4, we provided unified and rigorous derivations of both the single-reference- (Section 4.1), and a multireference (Section 4.2) CC method. The derivations used a general theorem (Theorem 4.1) motivated by a known method based on perturbation theory.

6. ACKNOWLEDGEMENTS

The authors would like to thank Fabian M. Faulstich and Simen Kvaal for helpful discussions and comments on the manuscript. The useful suggestions of anonymous reviewer are gratefully acknowledged.

APPENDIX A. PROPERTIES OF THE EXCITATION GRAPH

Here, we restrict ourselves to the single-reference case ($M = 1$) and drop the subscript m 's from the notation. Recall that K denotes the cardinality of the orbital set Λ . Given $\gamma \in L$, we introduce the set of paths of length n from 0 to γ in G ,

$$\mathbb{P}^n(\gamma) = \{\boldsymbol{\alpha} \in L \times \dots \times L : \text{there is a path } 0 \rightarrow \gamma \text{ in } G \text{ having edges } \boldsymbol{\alpha}\}.$$

The following theorem sheds light on the combinatorial structure of the excitation graph.

Theorem A.1. *Let $G^{\text{full}} = (L, E^{\text{full}})$ be the full SR excitation graph with K orbitals and $2N \leq K$ particles. Then the following properties hold true.*

- (i) *The number of vertices in G is given by $|L| = \binom{K}{N}$.*
- (ii) *The number of vertices of rank r is $|L(r)| = \binom{N}{r} \binom{K-N}{r}$.*
- (iii) *There are no edges in E^{full} entirely inside $L(r)$, and the number of edges from $L(r)$ to $L(r+s)$ is given by*

$$|E(r, r+s)| = \binom{K-N}{r} \binom{K-N-r}{s} \binom{N}{s+r} \binom{s+r}{r},$$

for all $r = 0, 1, \dots, N$ and $s = 0, \dots, N-r$, and $|E(r, r+s)| = 0$ if $s = N-r+1, \dots, N$. Furthermore, the symmetry property $|E(r, r+s)| = |E(s, r+s)|$ holds true.

- (iv) *The total number of edges is given by*

$$|E^{\text{full}}| = \sum_{r=1}^N \binom{N}{r} \binom{K-N}{r} \binom{K-2r}{N-r}.$$

- (v) *The number of directed paths of length $n \leq r = \text{rk}(\gamma)$ from 0 to γ is given by $|\mathbb{P}^n(\gamma)| = p(r, n)$, where*

$$p(r, n) = \sum_{\substack{r_1 + \dots + r_n = r \\ r_1, \dots, r_n \geq 1}} \left(\frac{r!}{r_1! \dots r_n!} \right)^2. \quad (\text{A.1})$$

Proof. (i) is trivial, so is (ii). As for (iii), we enumerate the pairs (α, β) in E^{full} as follows. Fix α with $\text{rk}(\alpha) = r$, then β must satisfy $r + s \leq N$, where $\text{rk}(\beta) = s$, so that $|\underline{\alpha} \cup \underline{\beta}| = N$ is possible. In $\underline{\beta}$, we must choose the missing internal letters from $\underline{\alpha}$ and there are r of them. For the remaining $N - s - r$ elements, we may choose freely: there are $\binom{N-r}{N-s-r}$ possibilities to do this. Next, $\overline{\beta}$ must be disjoint from $\overline{\alpha}$, so there are $M - N - r$ letters to choose from, giving $\binom{M-N-r}{s}$ possibilities. Multiplying these independent choices by the number of ways α can be chosen for fixed r , we get

$$\binom{N}{r} \binom{M-N}{r} \binom{N-r}{N-s-r} \binom{M-N-r}{s} \quad (\text{A.2})$$

for $s = 1, \dots, N - r$. This can be rewritten using the formula $\binom{n}{h} \binom{n-h}{k} = \binom{n}{k} \binom{n-k}{h}$ as

$$\binom{M-N}{r} \binom{M-N-r}{s} \binom{N}{s+r} \binom{s+r}{r}.$$

Using the aforementioned formula for the first two factors, we also get the desired symmetry property.

Next, to derive (iv) we sum up (A.2),

$$|E^{\text{full}}| = \sum_{r=0}^N \sum_{s=1}^{N-r} \binom{N}{r} \binom{M-N}{r} \binom{N-r}{N-s-r} \binom{M-N-r}{s}.$$

Using Vandermonde's identity,

$$\sum_{s=1}^{N-r} \binom{N-r}{N-s-r} \binom{M-N-r}{s} = \binom{M-2r}{N-r} - 1,$$

we get

$$|E^{\text{full}}| = \sum_{r=1}^N \binom{N}{r} \binom{M-N}{r} \binom{M-2r}{N-r},$$

where we used Vandermonde's identity once more.

Next, we prove (v). We need to change 0 into γ in n steps (edges). Suppose that the rank-increment of each step is r_1, \dots, r_n , and are such that $r_1 + \dots + r_n = r$. In the k th step we replace letters $(\alpha_1, \dots, \alpha_{r_k})$ with $(\beta_1, \dots, \beta_{r_k})$. These choices can be done independently, so there are $r!^2$ possibilities. However, the order of the α 's and β 's is irrelevant in each step so we have to divide by $(r_1! \cdots r_n!)^2$. Summing over all r_1, \dots, r_n gives the stated formula. \square

Remark A.2.

(i) It follows that the vertex density per rank is hypergeometric,

$$\nu_r = \frac{\binom{N}{r} \binom{K-N}{M-N-r}}{\binom{K}{N}}, \quad \text{where } r = 0, 1, \dots, N. \quad (\text{A.3})$$

Therefore, its mean is $\frac{N}{K}(K-N)$ and its variance is $\frac{(K-N)^2 N^2}{(K-1)K^2}$.

(ii) The formula (A.1) implies that $|\mathbb{P}^n(\gamma)|$ is independent of N and M and is constant for all γ of fixed rank r .

(iii) If S truncation is in effect, we have $p_S(r, n) = r!^2$ if $r = n$ and 0 otherwise.

- (iv) For the SD truncation, note that the number of (r_1, \dots, r_n) tuples with $r_j \in \{1, 2\}$, $r_1 + \dots + r_n = r$ and $|\{j : r_j = 2\}| = k$ is given by $\binom{n}{k}$ if $r = n + k$ and 0 otherwise. Therefore,

$$p_{\text{SD}}(r, n) = \frac{r!^2}{4^{r-n}} \binom{n}{r-n}.$$

- (v) According to the proof of [31, Lemma 4.4.],

$$|\{\beta \in L : \beta \preceq \alpha\}| = \sum_{s=1}^{r-1} \binom{r}{s} \binom{r-1}{r-s},$$

where $r = \text{rk}(\alpha)$.

APPENDIX B. OPTIMAL CHOICE OF MULTIREFERENCE DETERMINANTS

In this appendix, we describe an algorithm that can be used to automatically determine an optimal set of multireference determinants. Let $J \in \mathbb{N}$ and let

$$\{\gamma_1, \dots, \gamma_J\} \subset S$$

be a fixed set of determinants. Also, fix an excitation rank truncation, e.g. S, SD, SDT, etc. We want to select a *minimal* set of reference elements $\Omega = \{0_1, \dots, 0_M\}$, so that each γ_j is reachable through a *direct* S, SD, SDT, etc. excitation from Ω , this is called “first-order interaction space” in MRCC theory.

Recall that each $\alpha \in 2^\Lambda$ can be represented as a binary characteristic vector $\vec{\alpha} \in \{0, 1\}^K$ such that

$$\vec{\alpha}^t = \begin{cases} 1 & t \in \alpha \\ 0 & t \notin \alpha \end{cases}$$

The set $\{0, 1\}^K$ endowed with the Hamming metric

$$d_{\text{H}}(\vec{\alpha}, \vec{\beta}) = |\{t : \vec{\alpha}^t \neq \vec{\beta}^t, t = 1, \dots, K\}|$$

is a complete metric space, called the Hamming space. The closed balls and the spheres in this space are denoted as $B_{\text{H}}(\vec{\alpha}, R)$ and $S_{\text{H}}(\vec{\alpha}, R)$. Using this language, S is simply $S_{\text{H}}(\vec{0}, N)$, where $\vec{0} = (0, \dots, 0)$.⁸ Further,

$$\text{rk}_m(\alpha) = \frac{1}{2} d_{\text{H}}(\vec{0}_m, \vec{\alpha})$$

for any $m = 1, \dots, M$. Notice that $d_{\text{H}}(\vec{\alpha}, \vec{\beta}) \geq 2$ for distinct $\vec{\alpha}, \vec{\beta} \in S_{\text{H}}(\vec{0}, N)$.

This way, our optimization problem may be formulated as a covering problem in Hamming space. Let ρ denote the excitation rank truncation, e.g. $\rho = 1, 2, 3, \dots$ for S, SD, SDT, etc. Fix $J \in \mathbb{N}$ and $\Gamma = \{\vec{\gamma}_1, \dots, \vec{\gamma}_J\} \subset S_{\text{H}}(\vec{0}, N)$. We need to find a minimal set of Hamming balls $\{B_{\text{H}}(\vec{0}_m, 2\rho) : m = 1, \dots, M\}$ with $\vec{0}_m \in S_{\text{H}}(\vec{0}, N)$ such that

$$\Gamma \subset \bigcup_{m=1}^M B_{\text{H}}(\vec{0}_m, 2\rho) \cap S_{\text{H}}(\vec{0}, N).$$

⁸We warn the reader that the notation $\vec{0}_m$ for the vector representation of 0_m is slightly colliding with $\vec{0}$, the actual zero vector for the Hamming space.

Obviously, $\vec{0}_m \in \Gamma_{2\rho}$, where

$$\Gamma_{2\rho} = \bigcup_{j=1}^J B_{\mathbb{H}}(\vec{\gamma}_j, 2\rho) \cap S_{\mathbb{H}}(\vec{0}, N).$$

In other words, it is sufficient to look for the $\vec{0}_m$'s in the much smaller set $\Gamma_{2\rho}$. Let $n = |\Gamma_{2\rho}|$, and introduce some indexing in $\Gamma_{2\rho}$, say $\Gamma_{2\rho} = \{\vec{\alpha}_1, \dots, \vec{\alpha}_n\}$. The geometric form of the covering problem may be rephrased as a binary integer linear program (BILP) [34],

$$\left. \begin{array}{l} \sum_{\nu=1}^n \mathbf{c}_{\nu} \mathbf{x}_{\nu} \rightarrow \min! \\ \sum_{\substack{\vec{\gamma}_j \in B_{\mathbb{H}}(\vec{\alpha}_{\nu}, 2\rho) \\ 1 \leq \nu \leq n}} \mathbf{x}_{\nu} \geq 1, \quad j = 1, \dots, J \\ \mathbf{x} \in \{0, 1\}^n \end{array} \right\}$$

where $\mathbf{c} \in \mathbb{Q}^n$ is a given rational cost vector.

Remark B.1. If $\mathbf{c}_{\nu} = 0$ for some ν , then we will automatically have $\mathbf{x}_{\nu} = 1$ in the solution, even if $B_{\mathbb{H}}(\vec{\alpha}_{\nu}, 2\rho)$ does not cover. On the other hand, assigning a larger (resp. infinite) cost \mathbf{c}_{ν} will likely (resp. surely) end up $\mathbf{x}_{\nu} = 0$ in the solution.

The above problem is called a ‘‘multidimensional knapsack problem’’ in the optimization community, which seems to be extensively studied. However, we just naively solve the BILP using general ILP methods available in *Mathematica*. In our experience, the BILP can be built up and solved in a small amount of time for practically relevant parameters N , K and J , even on an older machine.

The reason for the apparent efficiency might be that the number of variables n in the the BILP above is significantly less than $|S| = \binom{K}{N}$. In fact, using the binary entropy function $H(x) = -x \log_2 x - (1-x) \log_2 (1-x)$, we have the rough estimate

$$\frac{n}{|S|} \leq \frac{J |B_{\mathbb{H}}(\vec{0}, 2\rho)|}{|B_{\mathbb{H}}(\vec{0}, N)|} \leq J \sqrt{8K\lambda'(1-\lambda')} 2^{-K(H(\lambda')-H(\lambda))},$$

where $\lambda = 2\rho/K$ and $\lambda' = N/K$ valid for $0 < \lambda, \lambda' < \frac{1}{2}$ [9, Lemma 2.4.4]. Notice that $H(\lambda') \geq H(\lambda)$, so $n/|S| \rightarrow 0$ as $K \rightarrow \infty$.

REFERENCES

- [1] R. A. Adams and J. J. Elsevier. *Sobolev spaces*. Elsevier, 2003.
- [2] V. Bach. Error bound for the Hartree-Fock energy of atoms and molecules. *Communications in mathematical physics*, 147(3):527–548, 1992.
- [3] 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*, pages 309–311. Springer, 1997.
- [4] R. J. Bartlett and M. Musiał. Coupled-cluster theory in quantum chemistry. *Reviews of Modern Physics*, 79(1):291, 2007.
- [5] R. Bishop. An overview of coupled cluster theory and its applications in physics. *Theoretica chimica acta*, 80(2-3):95–148, 1991.
- [6] C. Bloch. Sur la th orie des perturbations des  tats li s. *Nuclear Physics*, 6:329–347, 1958.
- [7] E. Canc es, M. Defranceschi, W. Kutzelnigg, C. Le Bris, and Y. Maday. Computational quantum chemistry: a primer. *Handbook of numerical analysis*, 10:3–270, 2003.
- [8] E. Canc es and C. Le Bris. On the convergence of SCF algorithms for the Hartree-Fock equations. *ESAIM: Mathematical Modelling and Numerical Analysis*, 34(4):749–774, 2000.
- [9] G. Cohen, I. Honkala, S. Litsyn, and A. Lobstein. *Covering codes*. Elsevier, 1997.
- [10] H. Eschrig. *The fundamentals of density functional theory*, volume 32. Springer, 1996.

- [11] F. M. Faulstich, A. Laestadius, O. Legeza, R. Schneider, and S. Kvaal. Analysis of the tailored coupled-cluster method in quantum chemistry. *SIAM Journal on Numerical Analysis*, 57(6):2579–2607, 2019.
- [12] G. Friesecke. The multiconfiguration equations for atoms and molecules: charge quantization and existence of solutions. *Archive for rational mechanics and analysis*, 169(1):35–71, 2003.
- [13] S. J. Gustafson and I. M. Sigal. *Mathematical concepts of quantum mechanics*. Springer Science & Business Media, 2011.
- [14] T. Helgaker, P. Jorgensen, and J. Olsen. *Molecular electronic-structure theory*. John Wiley & Sons, 2014.
- [15] P. D. Hislop and I. M. Sigal. *Introduction to spectral theory: With applications to Schrödinger operators*, volume 113. Springer Science & Business Media, 2012.
- [16] B. Jeziorski and H. J. Monkhorst. Coupled-cluster method for multideterminantal reference states. *Physical Review A*, 24(4):1668, 1981.
- [17] K. Kowalski. Properties of coupled-cluster equations originating in excitation sub-algebras. *The Journal of Chemical Physics*, 148(9):094104, 2018.
- [18] H. Kümmel, K. H. Lührmann, and J. G. Zabolitzky. Many-fermion theory in exps-(or coupled cluster) form. *Physics Reports*, 36(1):1–63, 1978.
- [19] M. Lewin. Existence of Hartree–Fock excited states for atoms and molecules. *Letters in Mathematical Physics*, 108(4):985–1006, 2018.
- [20] M. Lewin. Semi-classical limit of the Levy–Lieb functional in density functional theory. *Comptes Rendus Mathématique*, 356(4):449–455, 2018.
- [21] M. Lewin, E. H. Lieb, and R. Seiringer. The local density approximation in density functional theory. *Pure and Applied Analysis*, 2(1):35–73, 2020.
- [22] E. H. Lieb. Density functionals for Coulomb systems. *International Journal of Quantum Chemistry*, 24(3):243–277, 1983.
- [23] E. H. Lieb and M. Loss. Analysis. In *Amer. Math. Soc*, 2001.
- [24] E. H. Lieb and R. Seiringer. *The stability of matter in quantum mechanics*. Cambridge University Press, 2010.
- [25] E. H. Lieb and B. Simon. On solutions to the Hartree-Fock problem for atoms and molecules. *The Journal of Chemical Physics*, 61(2):735–736, 1974.
- [26] P.-L. Lions. Solutions of Hartree-Fock equations for Coulomb systems. *Communications in Mathematical Physics*, 109(1):33–97, 1987.
- [27] J. Paldus. Coupled cluster theory. In S. Wilson and G. H. Diercksen, editors, *Methods in computational molecular physics*, volume 293, pages 99–184. Springer Science & Business Media, 1991.
- [28] M. Reed and B. Simon. *Methods of modern mathematical physics I: Functional analysis*, volume 1. Elsevier, 1972.
- [29] M. Reed and B. Simon. *Methods of modern mathematical physics II: Fourier Analysis, Self-Adjointness*, volume 2. Elsevier, 1975.
- [30] M. Reed and B. Simon. *Methods of modern mathematical physics IV: Analysis of Operators*, volume 4. Elsevier, 1978.
- [31] T. Rohwedder. The continuous Coupled Cluster formulation for the electronic Schrödinger equation. *ESAIM: Mathematical Modelling and Numerical Analysis-Modélisation Mathématique et Analyse Numérique*, 47(2):421–447, 2013.
- [32] T. Rohwedder and R. Schneider. Error estimates for the coupled cluster method. *ESAIM: Mathematical Modelling and Numerical Analysis-Modélisation Mathématique et Analyse Numérique*, 47(6):1553–1582, 2013.
- [33] R. Schneider. Analysis of the projected coupled cluster method in electronic structure calculation. *Numerische Mathematik*, 113(3):433–471, 2009.
- [34] A. Schrijver. *Theory of linear and integer programming*. John Wiley & Sons, 1998.
- [35] I. Shavitt and R. J. Bartlett. *Many-body methods in chemistry and physics: MBPT and coupled-cluster theory*. Cambridge university press, 2009.
- [36] J. P. Solovej. The ionization conjecture in Hartree-Fock theory. *Annals of mathematics*, pages 509–576, 2003.
- [37] J. P. Solovej. Many body quantum mechanics. *Lecture Notes.*, 2007.
- [38] H. Yserentant. *Regularity and approximability of electronic wave functions*. Springer, 2010.
- [39] E. Zeidler. *Nonlinear Functional Analysis and Its Applications: Part 2A. Linear monotone operators*. Springer-Verlag, 1985.