

ANALYSIS OF THE TAILORED COUPLED-CLUSTER METHOD IN QUANTUM CHEMISTRY*

FABIAN M. FAULSTICH[†], ANDRE LAESTADIUS[†], ÖRS LEGEZA[‡], REINHOLD SCHNEIDER[§], AND SIMEN KVAAL[†]

Abstract. In quantum chemistry, one of the most important challenges is the static correlation problem when solving the electronic Schrödinger equation for molecules in the Born–Oppenheimer approximation. In this article, we analyze the tailored coupled-cluster method (TCC), one particular and promising method for treating molecular electronic-structure problems with static correlation. The TCC method combines the single-reference coupled-cluster (CC) approach with an approximate reference calculation in a subspace [complete active space (CAS)] of the considered Hilbert space that covers the static correlation. A one-particle spectral gap assumption is introduced, separating the CAS from the remaining Hilbert space. This replaces the nonexisting or nearly nonexisting gap between the highest occupied molecular orbital and the lowest unoccupied molecular orbital usually encountered in standard single-reference quantum chemistry. The analysis covers, in particular, CC methods tailored by tensor-network states (TNS-TCC methods). The problem is formulated in a nonlinear functional analysis framework, and, under certain conditions such as the aforementioned gap, local uniqueness and existence are proved using Zarantonello’s lemma. From the Aubin–Nitscheduality method, a quadratic error bound valid for TNS-TCC methods is derived, e.g., for linear-tensor-network TCC schemes using the density matrix renormalization group method.

Key words. Multi-reference Coupled-Cluster Method, Tailored Coupled-Cluster Method, Density Matrix Renormalization Group Method, Tensor Network States, Error Estimates, Existence and Uniqueness

AMS subject classifications. 65Z05, 81-08, 81V55, 81P40

1. Introduction. In this article, we present an analysis of the coupled-cluster (CC) method tailored by tensor-network states (TNS) for statically correlated electronic systems in quantum chemistry, thereby providing one of the first mathematically rigorous analyses of a multireference CC (MRCC) method.

The CC method is today the most popular wavefunction-based computational method in quantum chemistry [4]. The CCSD(T) scheme, the CC approach with single, double and perturbative triple excitations [33, 5], is referred to as the *gold standard of quantum chemistry*, as it yields computational results within error bars of practical experiments for small- and medium-sized molecules at a reasonable cost [23]. However, a severe disadvantage of conventional CC theory is that it fails dramatically for multireference systems, that is, systems whose wavefunction cannot be well approximated by a single Slater determinant reference function [13]. Such systems are said to be *statically correlated*, opposed to systems that are well approximated by a single Slater determinant, which are said to be *dynamically correlated* only.

*Submitted to the editors October 15, 2019.

Funding: This work has received funding from the Research Council of Norway (RCN) under CoE Grant No. 262695 (Hylleraas Centre for Quantum Molecular Sciences), from ERC-STG-2014 under grant No. 639508, from the Hungarian National Research, Development and Innovation Office (NKFIH) through Grant Nos. K120569, NN110360 and from the Hungarian Quantum Technology National Excellence Program (Project No. 2017-1.2.1-NKP-2017-00001). Ö. L. also acknowledges financial support from the Alexander von Humboldt foundation.

[†]Hylleraas Centre for Quantum Molecular Sciences, Department of Chemistry, University of Oslo, P.O. Box 1033 Blindern, N-0315 Oslo, Norway (f.m.faulstich@kjemi.uio.no).

[‡]Strongly Correlated Systems "Lendület" Research Group, Wigner Research Centre for Physics, H-1525, Budapest, Hungary.

[§]Modeling, Simulation and Optimization in Science, Department of Mathematics, Technische Universität Berlin, Sekretariat MA 5-3, Straße des 17. Juni 136, 10623 Berlin, Germany

Even if most molecules are single-reference systems in their equilibrium configuration, multireference character arises even in the simplest of chemical reactions, e.g., dissociation of N_2 . Yet, the static correlation problem is a long-lasting challenge in quantum chemistry. Many different MRCC approaches have been formulated to deal with the problem of static correlation. However, aside from formal difficulties and implementational complications, none of these methods have become a widely applicable tool. A review of different MRCC approaches is beyond the scope of this article, and we refer to Lyakh *et al.* [26] for a detailed description of the different benefits and disadvantages.

We are here concerned with an MRCC method that is based on the single-reference methodology (also called an externally corrected ansatz): The tailored CC (TCC) method extends a precomputed solution for a chosen subsystem of the full system by including further electron correlations via CC theory. We refer to the subsystem as the complete active space (CAS) and to the remaining system as the external space. Given the single-reference CC method’s major drawback, this subsystem needs to contain the static correlations. Consequently, the TCC method can be seen as a special type of an embedding method. Mathematically this corresponds to a division of excitation operators in two disjoint sub-algebras [19]. Nevertheless, in comparison with other “genuine” MRCC schemes, the TCC method suffers from the drawback that it is based on a single-reference theory and therewith introduces a certain bias towards a particular reference determinant. A possible remedy for this drawback is a large CAS covering the static correlations. The exponential scaling of the CAS makes an efficient approximation scheme for statically correlated systems indispensable for a TCC implementation of practical significance. To that end, the TCC method was recently combined with the density matrix renormalization group (DMRG) method. The DMRG method [45] is a high accuracy tool for statically correlated systems [7], nonetheless, for dynamically correlated problems it requires high computational resources making a wide-ranging application—at this time—intractable. Hence, it is the symbiosis of the DMRG and the CC method that creates a high efficiency scheme suitable for multireference systems [42, 43, 44, 10, 2, 22]. Granted that the DMRG-TCCSD method is the major motivation for the following analysis, we highlight that the applicability of this article’s results exceeds the DMRG-TCCSD method and, more generally, the TNS-TCC method.

This paper is organized as follows. We start by giving a short mathematical introduction to quantum chemistry. In Section 3, we introduce the TCC method with its major caveat: the CAS choice. Our main results—presented in Section 4—rest on certain assumptions that are connected to the structure of the one-particle basis from which the N -electron wavefunctions are constructed. Generalizing the concept of a HOMO-LUMO gap (see Section 4.1), we introduce a gap between the CAS and the external space (Assumption (A)). This allows us to derive various norm equivalences that can be used to establish continuity of the considered cluster operators with respect to different topologies. Moreover, a more technical constraint (Assumption (B)) enters our analysis when we assume that the fluctuation potential, i.e., an operator that models a part of the electron–electron interaction, cannot be too large when restricted to the external space. This manifests the importance of a well-chosen CAS as mentioned above. Also, as far as the multireference character of systems included in our treatment is concerned, we have to assume that those determinants that contribute the most in the N -electron CAS have energies very similar to the reference determinant. Other determinants can contribute too, but their weight must become smaller the larger the energy difference with respect to the reference determinant be-

comes. We then use Zarantonello’s lemma to derive local existence and uniqueness of TCC solutions under Assumption (A) and (B). In Section 4.2, we perform an energy error analysis and present major differences to the single-reference CC method. Via the Aubin–Nitsche-duality method we are able to derive a quadratic energy error bound valid for TCC schemes like the TNS-TCC method.

2. The Electronic Schrödinger Equation. In general, a Hamilton operator is an elliptic differential operator, formally defined by

$$(1) \quad H\psi = -\frac{1}{2}\Delta\psi + V\psi .$$

The function $V : \mathbb{R}^n \rightarrow \mathbb{R}$ is called the potential of the operator. Such differential operators are in general well studied [9, 11, 12, 34]. However, the numerical treatment of physical systems, especially electronic systems, is still challenging. In the spirit of mathematical rigor, we summarize the weak formulation of the Hamilton operator in Eq. (1):

The Hamilton operator induces a bilinear form $\mathcal{A}_V : \mathcal{C}_c^\infty(\mathbb{R}^n) \times \mathcal{C}_c^\infty(\mathbb{R}^n)$ by

$$(2) \quad \mathcal{A}_V(\tilde{\psi}, \psi) = \frac{1}{2}\langle \nabla\tilde{\psi}, \nabla\psi \rangle_{(L^2(\mathbb{R}^n))^n} + \langle \tilde{\psi}, V\psi \rangle_{L^2(\mathbb{R}^n)} ,$$

where $\mathcal{C}_c^\infty(\mathbb{R}^n)$ is the space of smooth functions on \mathbb{R}^n with finite support. Assuming boundedness of $V(x)(\cdot) : \mathcal{C}_c^\infty(\mathbb{R}^n) \rightarrow L^2(\mathbb{R}^n)$, the Cauchy-Schwarz inequality yields $\mathcal{A}_V(\tilde{\psi}, \psi) \leq C\|\tilde{\psi}\|_{H^1(\mathbb{R}^n)}\|\psi\|_{H^1(\mathbb{R}^n)}$, for all $\tilde{\psi}, \psi \in \mathcal{C}_c^\infty(\mathbb{R}^n)$. Since $\mathcal{C}_c^\infty(\mathbb{R}^n)$ is dense in $H^1(\mathbb{R}^n)$, we can extend \mathcal{A}_V to a bounded and symmetric bilinear form on $H^1(\mathbb{R}^n) \times H^1(\mathbb{R}^n)$.

Subsequently we omit the domain of the function space whenever it is clear from context. In this article, we assume that H satisfies Gårding’s inequality [34], i.e., there exist $c, e \in \mathbb{R}$ with $c > 0$ such that

$$(3) \quad \mathcal{A}_V(\psi, \psi) + e\langle \psi, \psi \rangle_{L^2} \geq c\|\psi\|_{H^1}^2 .$$

We furthermore define the Rayleigh–Ritz quotient $\mathcal{R}_V(\psi) = \mathcal{A}_V(\psi, \psi) / \langle \psi, \psi \rangle_{L^2}$ for all $\psi \in H^1 \setminus \{0\}$. Then $E_0 = \inf_{\psi \in H^1 \setminus \{0\}} \mathcal{R}_V(\psi)$ is well defined even though the infimum need not be attained. However, if such a minimizer exists it is called a ground state. Under the assumption that H attains a ground state $\psi_0 \in H^1$ we can recast the Schrödinger equation $\mathcal{A}_V(\tilde{\psi}, \psi_0) = E_0\langle \tilde{\psi}, \psi_0 \rangle_{L^2}$ for all $\tilde{\psi} \in H^1$ (i.e. $H\psi_0 = E_0\psi_0$) by means of the Rayleigh–Ritz variational principle:

$$(4) \quad E_0 = \min_{\psi \in H^1 \setminus \{0\}} \mathcal{R}_V(\psi) .$$

Note, whenever $\gamma = \inf\{\mathcal{R}_V(\psi) : \psi \in H^1, \psi \neq 0, \langle \psi_0, \psi \rangle_{L^2} = 0\} - E_0 > 0$, ψ_0 is (up to a phase) the unique ground state of H and γ is called the spectral gap.

This article focuses on the electronic Schrödinger equation obtained from the Born–Oppenheimer approximation [40, 6]. In Hartree atomic units, the Hamilton operator of a Coulomb system that consists of N electrons and N_{nuc} nuclei reads

$$H\psi(x) = -\sum_{i=1}^N \frac{1}{2}\Delta_i\psi(x) + \underbrace{\left(\frac{1}{2} \sum_{i=1}^N \sum_{j \neq i}^N \frac{1}{|r_i - r_j|} - \sum_{i=1}^N \sum_{j=1}^{N_{\text{nuc}}} \frac{Z_j}{|r_i - R_j|} \right)}_{=V_C} \psi(x) ,$$

with V_C the Coulomb potential. Here $\psi(x) = \psi(x_1, \dots, x_N)$, where the argument $x_i = (r_i, s_i)$ for $i \in \{1, \dots, N\}$ is associated with the position of the i -th electron $r_i \in \mathbb{R}^3$ and its spin $s \in \{\pm 1/2\}$. As a result of the Born–Oppenheimer approximation, the nuclei positions $R_j \in \mathbb{R}^3$ and charges $Z_j > 0$, $j \in \{1, \dots, M\}$, enter as fixed parameters. This general formulation is so far independent of spin as an explicit variable. Moreover, solutions to the above Hamiltonian do not naturally fulfill Pauli’s principle, i.e., fermionic state functions need to be antisymmetric with respect to permutations of the coordinates x_i . Considering these further constraints, the set of admissible wavefunctions is given by

$$(5) \quad \mathcal{H} = H^1 \left(\left(\mathbb{R}^3 \times \left\{ \pm \frac{1}{2} \right\} \right)^N \right) \cap \bigwedge_{i=1}^N L^2 \left(\mathbb{R}^3 \times \left\{ \pm \frac{1}{2} \right\} \right),$$

where \wedge is the antisymmetric tensor product that guarantees Pauli’s principle. We conclude, the minimization problem Eq. 4 corresponding to electronic structure calculations is given by

$$(6) \quad E_0 = \min_{\psi \in \mathcal{H} \setminus \{0\}} \mathcal{R}_{V_C}(\psi).$$

Remark 1. The Hamilton operator is here a map $H : H^1 \supseteq \mathcal{H} \rightarrow H^{-1}$, where H^{-1} is the dual space of H^1 . In particular, this means that instead of the L^2 -inner product we need to consider the dual pairing $\langle \cdot, \cdot \rangle_{H^1, H^{-1}}$. To justify the use of the inner product we recall that H^1 is continuously embedded in L^2 and that H^1 is dense in L^2 , i.e. H^1 is densely embedded in L^2 and we write $H^1 \xrightarrow{d} L^2$. For such a Hilbert space structure, we define the Gelfand triple $H^1 \xrightarrow{d} L^2 \xrightarrow{d} H^{-1}$ (also called rigged Hilbert space), identifying $L^2 \simeq (L^2)'$. Note that as a consequence we are no longer allowed to identify $H^1 \simeq H^{-1}$. One advantage of the Gelfand triple is that the use of the L^2 inner product instead of the dual pairing $\langle \cdot, \cdot \rangle_{H^1, H^{-1}}$ becomes meaningful [46]: Given the Gelfand triple $H^1 \xrightarrow{d} L^2 \xrightarrow{d} H^{-1}$ and the scalar product $\langle \cdot, \cdot \rangle_{L^2}$ on $L^2 \times L^2$, we find $\langle x, y \rangle_{L^2} = \langle x, y \rangle_{H^1 \times H^{-1}}$ for all $x \in H^1$ and $y \in L^2$ since $H^1 \subseteq L^2$ and $L^2 \subseteq H^{-1}$. By Hahn–Banach we can therefore continuously extend $\langle x, \cdot \rangle_{L^2}$ from L^2 to H^{-1} for arbitrary but fixed $x \in H^1$.

Remark 1 becomes important when considering quantum molecular systems on the infinite dimensional Hilbert space \mathcal{H} . We subsequently make use of the inner product notation, emphasizing that the reader should keep this detail in mind. Moreover, henceforth we use the short notation $\langle \cdot, \cdot \rangle$ rather than $\langle \cdot, \cdot \rangle_{L^2}$ or $\langle \cdot, \cdot \rangle_{l^2}$ whenever the meaning is clear from context.

3. Approximate Solutions of the Schrödinger Equation. The high dimensionality of Eq. (6) makes a direct minimization in general intractable. The variety of possible approximations, depending on the chemical problem and required accuracy, is rich [13, 27, 14]. However, most wavefunction based schemes rely on an antisymmetrized product ansatz. The factors of this exterior product are called *spin-orbitals* and the functions spanning the solution space are denoted *Slater determinants*. Subsequently, we denote the spin-orbitals by χ and Slater determinants by ϕ . For an N -electron problem, let $N < K$ and $\mathcal{B} = \{\chi_1, \dots, \chi_K\} \subseteq H^1(\mathbb{R}^3 \times \{\pm \frac{1}{2}\})$ denote an $L^2(\mathbb{R}^3 \times \{\pm \frac{1}{2}\})$ -orthonormal set of functions, called spin-orbitals. An N -particle wavefunction fulfilling Pauli’s exclusion principle is obtained by forming the exterior

product of N spin-orbitals $\{\chi_{\mu_1}, \dots, \chi_{\mu_N}\}$

$$(7) \quad \phi[\mu_1, \dots, \mu_N](x_1, \dots, x_N) = \frac{1}{\sqrt{N!}} \bigwedge_{i=1}^N \chi_{\mu_i}(x_1, \dots, x_N) = \frac{1}{\sqrt{N!}} \det(\chi_{\mu_i}(x_j))_{i,j=1}^N,$$

where the indices $\mu_1, \dots, \mu_N \in \{1, \dots, K\}$ are in canonical order, i.e., $\mu_1 < \dots < \mu_N$. We see immediately that Slater determinants inherit L^2 -orthonormality from the spin-orbital basis. The corresponding Galerkin space \mathcal{H}_K is then spanned by all possible exterior products of length N in \mathcal{B} . This construction yields a combinatorial scaling of \mathcal{H}_K —also called the full configuration-interaction (FCI) space. An L^2 -orthonormal basis \mathcal{B}_K of \mathcal{H}_K is obtained by imposing a canonical ordering of the spin-orbitals in the exterior products, i.e.,

$$\mathcal{B}_K = \{\phi[\mu_1, \dots, \mu_N] : \mu_i \in \{1, \dots, K\}, \mu_1 < \dots < \mu_N\}.$$

Subsequently we use the notation $\phi_\mu = \phi[\mu_1, \dots, \mu_N]$ and without loss of generality define the reference determinant $\phi_0 = \phi[1, \dots, N]$. Furthermore, we use the standard terminology of quantum chemistry and call spin-orbitals defining ϕ_0 occupied and the remaining ones virtual. Indices I, J, K, \dots are assumed to be occupied (i.e. smaller or equal than N) while A, B, C, \dots are assumed to be virtual (i.e. greater than N).

Essential to the CC theory is the L^2 -bounded commutative algebra of cluster operators \mathcal{C}_K , defined via single-excitation operators. We define a single-excitation operator X_I^A as follows: $X_I^A \phi_\mu$ replaces χ_I by χ_A for any ϕ_μ if $\mu_i = I$ for some i and $\mu_j \neq A$ for all j , otherwise $X_I^A \phi_\mu = 0$. Since Slater determinants are normalized, this defines X_I^A as an L^2 -bounded operator. Higher order excitation operators are then defined as product of single-excitation operators, e.g. , the double excitation operator $X_{IJ}^{AB} = X_I^A X_J^B$. The fermionic commutation relations, i.e., $[a_i, a_j^\dagger]_+ = \delta_{ij}$ and $[a_i^\dagger, a_j^\dagger]_+ = [a_i, a_j]_+ = 0$, yield that excitation operators commute. The set of excitation operators is then trivially an L^2 -bounded and commutative algebra. Furthermore we define the rank of an excitation operator as the length of the product, when written as product of single-excitation operators. Note that by antisymmetry of Slater determinants, the product $X_{I_1 \dots I_n}^{A_1 \dots A_n}$ is antisymmetric under permutations of $\{I_1, \dots, I_n\}$ and $\{A_1, \dots, A_n\}$, respectively. Similar to \mathcal{H}_K , a basis of \mathcal{C}_K is obtained by imposing a canonical ordering of the product of single-excitation operators with respect to the orbital indices.

PROPOSITION 2. *We can induce a norm on \mathcal{C}_K via $\|X_\mu\|_{\mathcal{C}_K} = \|X_\mu \phi_0\|_{H^1}$. Then \mathcal{C}_K is isometrically isomorphic to $\text{span}\{\phi_0\}^\perp$, where \perp denotes L^2 -orthogonal complement in \mathcal{H}_K .*

Proof. For any $\phi_\mu \in \mathcal{B}_K$, there exists a unique excitation operator such that $\phi_\mu = X_\mu \phi_0$ up to a sign factor, i.e., ϕ_μ is generated from ϕ_0 by repeated substitution of occupied spin-orbitals. Conversely, for any excitation operator X_μ there is a unique $\phi_\mu \in \mathcal{B}_K$ such that $\phi_\mu = X_\mu \phi_0$ up to a sign factor. Hence, we can define a bijective homomorphism between \mathcal{C}_K and $\text{span}\{\phi_0\}^\perp$, where we impose canonical vector-space operations on the respective spaces, i.e., vector addition and scalar multiplication. By construction this map is trivially an isometry, which proves the claim. \square

Subsequently, we refer to the basis index μ as an excitation index and switch to the more common multi-index notation, i.e., $\mu = \binom{A_1, \dots, A_r}{I_1, \dots, I_r}$ with occupied indices $\{I_1, \dots, I_r\}$ and virtual indices $\{A_1, \dots, A_r\}$. The set of all possible excitation indices

is denoted \mathcal{J} , where we dropped the dependence on K and the reference state due to notational simplicity. Using the canonical ordering, the number of possible excitation indices up to a certain excitation rank $n \leq N$ is given by

$$|\mathcal{J}| = \sum_{k=1}^n \binom{N}{k} \binom{K-N}{k}.$$

In practice the spin-orbitals in \mathcal{B} and thus the reference wavefunction ϕ_0 come from a preliminary Hartree–Fock calculation [13, 24, 25]: In a nutshell, starting with an initial spin-orbital basis $\{\chi_i^{(0)}\}_{i=1}^K$ we minimize Eq. 6 with a mean-field potential. This yields a nonlinear K -dimensional eigenvalue problem $\bar{F}(\chi_1, \dots, \chi_N)\chi_i = \lambda_i\chi_i$, for $i = 1, \dots, N$, where the Fock matrix \bar{F} depends on the N occupied spin-orbitals. The Fock matrix is symmetric, implying that the N eigenvectors can be completed with $K - N$ additional eigenvectors. It is these eigenfunctions that form \mathcal{B} .

We observe that the Hartree–Fock calculation depends on the dimension K in a manner which is not entirely controlled: In general, it is unclear whether the $\{\chi_i\}_{i=1}^K$ form a global minimum of the Rayleigh–Ritz minimization problem and whether the solution converges as $K \rightarrow \infty$. Such questions are beyond the scope of the present article, but is relevant in context of the $K \rightarrow \infty$ limit of the TCC method, see Remark 16. The Hartree–Fock calculation induces a splitting of the Hamilton operator $H = F + W$ with $F = \sum_{i=1}^N \bar{F}(i)$, where $\bar{F}(i) = I \otimes \dots \otimes I \otimes \bar{F}(i) \otimes I \otimes \dots \otimes I$ indicating by $\bar{F}(i)$ that \bar{F} appears on the i -th position in the Kronecker product. Subsequently, we will refer to F as the Fock operator and to W as the fluctuation potential.

We define for any multi-index μ of excitation rank $n \leq N$ the number

$$\varepsilon_\mu = \sum_{j=1}^n (\lambda_{A_j} - \lambda_{I_j}),$$

i.e., the sum of the single-particle Hartree–Fock energy differences of the occupied and virtual spin-orbitals in μ . Defining $\Lambda_0 = \sum_{i=1}^N \lambda_i$ —the sum over the N first single-particle Hartree–Fock energies—we see that the Slater determinants \mathcal{B}_K , formed by the single-particle Hartree–Fock eigenfunctions, are the N -particle Hartree–Fock eigenfunctions with $F\phi_\mu = (\Lambda_0 + \varepsilon_\mu)\phi_\mu$.

Returning to the Schrödinger equation, the L^2 -normalization constraint on $\psi \in \mathcal{H}_K$ is subsequently replaced by the intermediate normalization, i.e., $\langle \phi_0, \psi \rangle = 1$. Hence, $\psi = (I + S)\phi_0$ holds for an operator $S = \sum_{\mu \in \mathcal{J}} s_\mu X_\mu \in \mathcal{C}_K$ and we denote the basis coefficients $(s_\mu)_{\mu \in \mathcal{J}} = (\langle \phi_\mu, \psi \rangle)_{\mu \in \mathcal{J}}$ *excitation amplitudes*. Inserting this parameterization of wavefunctions into the Schrödinger equation, we find that Eq. (6) is equivalent to the linear problem

$$(8) \quad \begin{cases} E_0^{(\text{FCI})} = \langle \phi_0, H\psi_0^{(\text{FCI})} \rangle, \\ 0 = \langle \phi_\mu, (H - E_0^{(\text{FCI})})\psi_0^{(\text{FCI})} \rangle, \quad \forall \mu \in \mathcal{J}, \end{cases}$$

which is known as the FCI scheme. For a derivation of the corresponding amplitude equations we refer the reader to [13].

3.1. Projected Single-Reference Coupled-Cluster Method. The previously described FCI approach suffers from the *curse of dimensionality* since \mathcal{H}_K grows exponential with the number of particles, i.e., $\dim(\mathcal{H}_K) \in \mathcal{O}(K^N)$. Furthermore, truncating the operator $S \in \mathcal{C}_K$ at rank- n excitations reduces the computational cost but

yields CI methods that are no longer energy size-extensive nor size-consistent [13], which are quantum chemical concepts relating to the correct energy behavior with respect to the system's size and dissociation [39]. Alternatively to the linear manifold used in Eq. (8), an exponential parameterization of wavefunctions can be used [15, 16]: Let $\psi \in \mathcal{H}_K$ be intermediately normalized, i.e., $\psi = (I + S)\phi_0$ for some $S \in \mathcal{C}_K$. Then there exists a unique $T \in \mathcal{B}(H^1, H^{-1})$ with $\psi = e^T \phi_0$ [38] (for the result in the limit $K \rightarrow \infty$ see [36]), where

$$(9) \quad T = \sum_{\mu \in \mathcal{J}} t_\mu X_\mu \quad \text{and} \quad T = \log(I + S) .$$

This exponential parameterization has the benefit that it is multiplicatively separable with respect to subsystems that are separated by distance, thereby regaining size-extensivity and consistency under mild assumptions on the reference determinant [39].

To solve the Schrödinger equation, it remains to determine the *cluster amplitudes* $(t_\mu)_{\mu \in \mathcal{J}}$. This is the pursuit of the CC method. The linked CC equations describing the cluster amplitudes are given by [13]:

$$(10) \quad \begin{cases} E_0^{(CC)} = \langle \phi_0, e^{-T} H e^T \phi_0 \rangle , \\ 0 = \langle \phi_\mu, e^{-T} H e^T \phi_0 \rangle , \quad \forall \mu \in \mathcal{J} . \end{cases}$$

The equivalence to the Schrödinger equation (6), is straightforwardly established [13]: Given an intermediately normalized minimizer of Eq. (6) $\psi = e^T \phi_0$, we obtain

$$E_0 e^T \phi_0 = H e^T \phi_0 \Rightarrow E_0 \phi_0 = e^{-T} H e^T \phi_0 \Rightarrow \begin{cases} E_0 = \langle \phi_0, e^{-T} H e^T \phi_0 \rangle , \\ 0 = \langle \phi_\mu, e^{-T} H e^T \phi_0 \rangle , \quad \forall \mu \in \mathcal{J} . \end{cases}$$

Conversely, given a solution $\psi = e^T \phi_0$ fulfilling Eqs. (10), we find

$$\begin{aligned} H e^T \phi_0 &= e^T e^{-T} H e^T \phi_0 = \sum_{\mu \in \mathcal{J}} e^T \phi_\mu \langle \phi_\mu, e^{-T} H e^T \phi_0 \rangle + e^T \phi_0 \langle \phi_0, e^{-T} H e^T \phi_0 \rangle \\ &= E_0^{(CC)} e^T \phi_0 . \end{aligned}$$

Note that this equivalence does in general not hold true under truncations of T , e.g., considering only single- and double-excitations in T (the CCSD method). In this case, the CC method is no longer variational. For a more detailed discussion on this topic see [20].

We emphasize that there exists a one-to-one relation between cluster amplitudes $(t_\mu)_{\mu \in \mathcal{J}}$ and the therewith defined cluster operators $T = \sum_{\mu \in \mathcal{J}} t_\mu X_\mu$ [36]. Therefore, we shall denote cluster amplitudes with small letters and the corresponding cluster operators with the respective capital letter. Let $\mathcal{V}_K^{(CC)} = \{t \in \mathbb{R}^{|\mathcal{J}|} : \|t\|_{\mathcal{V}_K^{(CC)}} < +\infty\}$ be the (Hilbert) space of cluster amplitudes, where

$$\|\cdot\|_{\mathcal{V}_K^{(CC)}} : \mathbb{R}^{|\mathcal{J}|} \rightarrow [0, +\infty]; \quad t \mapsto \|t\|_{\mathcal{V}_K^{(CC)}} = \sqrt{\sum_{\mu \in \mathcal{J}} \varepsilon_\mu |t_\mu|^2} .$$

We see that $\|\cdot\|_{\mathcal{V}_K^{(CC)}}$ is a norm if $\varepsilon_\mu > 0$ for all $\mu \in \mathcal{J}$. We then refer to $\|\cdot\|_{\mathcal{V}_K^{(CC)}}$ as the cluster amplitude norm. This is guaranteed by assuming a HOMO-LUMO gap, i.e., $\varepsilon_0 = \lambda_{N+1} - \lambda_N > 0$.

Although a HOMO-LUMO gap is very common in electronic structure analysis, it limits the results to a subset of systems. For statically correlated systems the Fock operator usually has a degenerate or almost degenerate spectrum, i.e., there exists no HOMO-LUMO gap or it is negligibly small. In either case, this yields divergence of the used quasi-Newton method since the HOMO-LUMO gap enters inversely in the approximate Jacobian.

Formally, the linked CC equations can be defined using the CC function

$$f_{\text{CC}} : \mathcal{V}_K^{(\text{CC})} \rightarrow (\mathcal{V}_K^{(\text{CC})})'; \quad t \mapsto (\langle \phi_\mu, e^{-T} H e^T \phi_0 \rangle)_{\mu \in \mathcal{J}}$$

with the energy functional $\mathcal{E}_{\text{CC}} : \mathcal{V}_K^{(\text{CC})} \rightarrow \mathbb{R}; \quad t \mapsto \langle \phi_0, e^{-T} H e^T \phi_0 \rangle$. Consequently, we can write Eqs. (10) as

$$\begin{cases} E_0^{(\text{CC})} = \mathcal{E}_{\text{CC}}(t), \\ 0 = \langle v, f_{\text{CC}}(t) \rangle, \quad \text{for all } v \in \mathcal{V}_K. \end{cases}$$

This shows that the projected CC method is a nonlinear Galerkin scheme. A corresponding analysis can be found in [38].

3.2. The Tailored Coupled-Cluster Method. A major drawback of the projected CC theory is the intractability of statically correlated systems. Many attempts have been taken to remedy this impediment but so far no panacea has been found [4]. The TCC method, as an externally corrected CC method, is not based on the Jeziorski–Monkhorst ansatz [4, 26, 18], however, it is still able to compute statically correlated systems with comparable accuracy [42, 43, 44, 10, 2]. Using a basis splitting approach [32, 31, 1, 30] it is possible to combine the single-reference CC method with CAS computations [17]. To that end, the wavefunction is split into two parts: a fixed part imported from a prior CAS calculation and an external part, which is adjusted in the presence of that fixed CAS part. We use the following basis splitting.

DEFINITION 3. Let $\{\chi_1, \dots, \chi_K\} \subseteq H^1$ be a set of L^2 -orthonormal spin-orbitals with $K > N$ and ϕ_0 the considered reference Slater determinant. We define

$$\mathcal{B}_{\text{CAS}} = \underbrace{\{\chi_1, \dots, \chi_N\}}_{\text{occupied}}, \underbrace{\{\chi_{N+1}, \dots, \chi_k\}}_{\text{unoccupied}}, \quad \mathcal{B}_{\text{ext}} = \underbrace{\{\chi_{k+1}, \dots, \chi_K\}}_{\text{external}}$$

and furthermore $\mathcal{B}_{\text{CAS}} = \{\phi[\mu_1, \dots, \mu_N] : \mu_i \in \{1, \dots, k\}, \mu_1 < \dots < \mu_N\}$. The corresponding FCI space \mathcal{H}_{CAS} is then defined as the span of \mathcal{B}_{CAS} . We define \mathcal{H}_{ext} to be the L^2 -orthogonal space of \mathcal{H}_{CAS} , i.e., $\mathcal{H}_K = \mathcal{H}_{\text{CAS}} \oplus \mathcal{H}_{\text{ext}}$. Analogously, we split the set of excitation-indices \mathcal{J} describing the set of possible excitations, i.e., $\mathcal{J}_{\text{CAS}} = \{\mu \in \mathcal{J} : X_\mu \phi_0 \in \mathcal{H}_{\text{CAS}}\}$ and $\mathcal{J}_{\text{ext}} = \{\mu \in \mathcal{J} : X_\mu \phi_0 \notin \mathcal{H}_{\text{CAS}}\}$.

Remark 4. We note that \mathcal{J}_{ext} does not only contain excitations into states purely excited in \mathcal{B}_{ext} but also into mixed states, i.e., for $\mu = \binom{A_1, \dots, A_n}{I_1, \dots, I_n}$ there exists at least one $l \in \{1, \dots, n\}$ such that $A_l \in \{k+1, \dots, K\}$.

We highlight that the basis splitting in practice cannot be arbitrary. For the correctness of the TCC method it is of utmost importance that \mathcal{B}_{CAS} covers all statically correlated spin-orbitals. Moreover, \mathcal{B}_{ext} should only consist of spin-orbitals with dynamic electron correlation. A well-chosen basis splitting can be obtained using concepts of quantum information theory as has been introduced in [41]. This caveat

will be further discussed in Section 3.3. We also refer to [10] for a case study on the N_2 molecule illustrating the TCC method’s sensitivity to the CAS choice.

Given an intermediately normalized approximate CAS-solution ϕ_{CAS} , we can write $\phi_{\text{CAS}} = e^{T^{\text{CAS}}} \phi_0 \approx \psi_{\text{CAS}}^{(\text{FCI})}$. The TCC solution is then given by $\psi_*^{(\text{TCC})} = e^{T^{\text{ext}}} e^{T^{\text{CAS}}} \phi_0$, where T^{ext} is obtained by solving the linked TCC equations:

$$(11) \quad \begin{cases} E_0^{(\text{TCC})} = \langle \phi_0, e^{-T^{\text{CAS}}} e^{-T^{\text{ext}}} H e^{T^{\text{CAS}}} e^{T^{\text{ext}}} \phi_0 \rangle, \\ 0 = \langle \phi_\mu, e^{-T^{\text{CAS}}} e^{-T^{\text{ext}}} H e^{T^{\text{CAS}}} e^{T^{\text{ext}}} \phi_0 \rangle, \quad \mu \notin \mathcal{J}_{\text{CAS}}. \end{cases}$$

We emphasize that for the TCC method, the CAS-solution ϕ_{CAS} and therewith T^{CAS} is fixed. Similar to the analysis in [38], a useful measure for the dynamical correction is a weighted l^2 -norm of the external cluster amplitudes. Let

$$\mathcal{V}_{\text{ext}} = \{t \in \mathbb{R}^{|\mathcal{J}_{\text{ext}}|} : \|t\|_{\mathcal{V}_{\text{ext}}} < +\infty\}$$

be the space of external cluster amplitudes, where

$$\|\cdot\|_{\mathcal{V}_{\text{ext}}} : \mathbb{R}^{|\mathcal{J}_{\text{ext}}|} \rightarrow [0, +\infty]; \quad t \mapsto \|t\|_{\mathcal{V}_{\text{ext}}} = \sqrt{\sum_{\mu \in \mathcal{J}_{\text{ext}}} \varepsilon_\mu |t_\mu|^2}.$$

The map $\|\cdot\|_{\mathcal{V}_{\text{ext}}}$ is a norm if $\varepsilon_\mu > 0$ for all $\mu \in \mathcal{J}_{\text{ext}}$. Assumptions on the considered systems to ensure such structure will be elaborated in Section 4.1. Using this framework we can define the N -electron TCC function as follows.

DEFINITION 5. *Let $K, N \in \mathbb{N}$ with $K > N$ be fixed, $\mathcal{B} = \{\chi_1, \dots, \chi_K\} \subseteq H^1$ a set of L^2 -orthonormal spin-orbitals and $\phi_0 \in \mathcal{H}_K$ the considered reference state. Further, assume the splitting $\mathcal{B} = \mathcal{B}_{\text{CAS}} \dot{\cup} \mathcal{B}_{\text{ext}}$ of \mathcal{B} and the CAS-solution $\phi_{\text{CAS}} = e^{T^{\text{CAS}}} \phi_0$, with corresponding amplitudes $t^{\text{CAS}} = (t_\mu^{\text{CAS}})_{\mathcal{J}_{\text{CAS}}}$. We define the TCC function*

$$f(\cdot; t^{\text{CAS}}) : \mathcal{V}_{\text{ext}} \rightarrow (\mathcal{V}_{\text{ext}})'; \quad t \mapsto f(t; t^{\text{CAS}}),$$

where $(f(t; t^{\text{CAS}}))_\mu = \langle \phi_\mu, e^{-T^{\text{CAS}}} e^{-T} H e^T e^{T^{\text{CAS}}} \phi_0 \rangle$ for $\mu \in \mathcal{J}_{\text{ext}}$. In addition, let the TCC-energy functional be given by

$$\mathcal{E}(t; t^{\text{CAS}}) = \langle \phi_0, e^{-T^{\text{CAS}}} e^{-T} H e^T e^{T^{\text{CAS}}} \phi_0 \rangle.$$

Using the TCC function, the linked TCC equations (11) become

$$\begin{cases} E_0^{(\text{TCC})} = \mathcal{E}(t; t^{\text{CAS}}), \\ 0 = \langle v, f(t; t^{\text{CAS}}) \rangle, \quad \forall v \in \mathcal{V}_{\text{ext}}. \end{cases}$$

This formulation resembles the single-reference CC method. Indeed, $f(t; t^{\text{CAS}}) = P_{\mathcal{V}_{\text{ext}}} f_{\text{CC}}(t \oplus t^{\text{CAS}})$ with the orthogonal projection $P_{\mathcal{V}_{\text{ext}}}$ onto \mathcal{V}_{ext} , relates the TCC function to the classical CC function in Eq. (10). Note that the CAS-part of the cluster amplitudes is still fixed. Despite this close connection to the CC method, we shall see that the TCC scheme differs heavily from the single-reference CC method in its computational performance and analysis.

3.3. Entropy based CAS choice. We start this section by noting that any Slater determinant can be uniquely described by an occupation tensor $\mathbf{e}^{m_1} \otimes \dots \otimes \mathbf{e}^{m_K}$, where $\mathbf{e}^0 = (1, 0)^T, \mathbf{e}^1 = (0, 1)^T \in \mathbb{R}^2$. This identification is part of the second

quantization [13] and is in fact an isometric isomorphism (see the Jordan–Wigner transformation [29]). Consequently, we can interpret any real wavefunction as an element in the 2^K dimensional linear space $\mathcal{W}_K = \bigotimes_{i=1}^K \mathbb{R}^2$ with given basis $\{\phi_{\mathbf{m}} = \mathbf{e}^{\mathbf{m}_1} \otimes \dots \otimes \mathbf{e}^{\mathbf{m}_K} : \mathbf{m}_i \in \{0, 1\}\}$. Given a low-rank DMRG solution ψ_{DMRG} on \mathcal{W}_K , i.e., $\psi_{\text{DMRG}} = \sum_{\mathbf{m}=1}^K c_{\mathbf{m}} \phi_{\mathbf{m}}$, we introduce the quantum information theory concepts used to choose a CAS. We start by considering the i -mode matricization $\mathbf{U}[i] \in \mathbb{R}^{2^{K-1} \times 2}$ of the solution tensor ψ_{DMRG} , i.e., the matrix obtained from ψ_{DMRG} by transforming the basis elements $\phi_{\mathbf{m}}$ by taking \mathbf{m}_i as row index and all remaining indices as compound column index. We introduce the elementwise notation $U[i]_{(\mathbf{m}_1, \dots, \mathbf{m}_i, \dots, \mathbf{m}_K), (\mathbf{m}_i)}$, where \mathbf{m}_i means that \mathbf{m}_i is removed from the binary string \mathbf{m} and all remaining indices are combined to one compound index. We then compute the single-orbital entropy for the i -mode matricization denoted $s(i)$, i.e., $s(i) = -\text{Tr}(\mathbf{D}[i] \ln \mathbf{D}[i]) \in [0, \ln(2)]$, where $\mathbf{D}[i] = \mathbf{U}[i]^T \mathbf{U}[i] \in \mathbb{R}^{2 \times 2}$ is the single-orbital density matrix. Based on Szalay *et al.* [41], the single-orbital entropy can be used to describe the degree of electron correlation, i.e., a large value of $s(i)$ indicates static correlations. However, since the electron correlation is a two particle effect, we need to measure the information flow for all possible electron pairs. This is done via the mutual information: We start by computing the two-orbital entropy $s(i, j)$. Similarly to the single-orbital entropy $s(i)$, the two-orbital entropy $s(i, j) = -\text{Tr}(\mathbf{D}[i, j] \ln \mathbf{D}[i, j]) \in [0, \ln(4)]$ where $\mathbf{D}[i, j] \in \mathbb{R}^{4 \times 4}$ is the two-orbital density matrix obtained from $U[i, j]_{(\mathbf{m}_i, \mathbf{m}_j), (\mathbf{m}_1, \dots, \mathbf{m}_i, \dots, \mathbf{m}_j, \dots, \mathbf{m}_K)}$. Given the single- and two-orbital entropies, we can compute the mutual information, $I(i, j) = s(i) + s(j) - s(i, j)$ for $i, j \in 1, \dots, K$. This quantifies the electron correlations between orbital i and j as they are embedded in the whole system [35]. The large values of $I(i, j)$ describe static correlations while the small matrix elements stand for the dynamic correlation. In certain cases, the decreasingly ordered values of $I(i, j)$ show a jump, which clearly distinguishes a set of statically correlated orbitals, and suggests a basis splitting at this jump. However, general mutual information profiles do not need to show such behavior. Then the *a priori* thresholds \underline{s} and \underline{n} are introduced to identify orbitals with $s(i) > \underline{s}$ and $I(i, j) > \underline{n}$. It is these orbitals that are then used to define \mathcal{B}_{CAS} and therewith the basis splitting. In practice, \underline{s} and \underline{n} are systematically lowered until convergence of the DMRG-TCC method is reached. This approach is heuristic but provides an efficient tool for obtaining well-chosen \mathcal{B}_{CAS} and \mathcal{B}_{ext} , which is essential for the TCC method's success. We highlight that the above procedure is feasible for larger systems since the used quantities are qualitatively very robust with respect to the bond-dimension, i.e., a CAS choice can be obtained from a low rank calculation on \mathcal{H}_K [10]. For more details and numerical investigations on the CAS choice we refer the reader to [10].

4. Analysis of the TCC Method. We focus here on the mathematical analysis of the TCC method for a finite spin-orbital set, i.e., $K < \infty$. Several caveats of the limit process $K \rightarrow \infty$ are subsequently addressed, but a full investigation is relegated to future work.

First, we show the consistency of the TCC method, in the sense that exact solutions of the Schrödinger equation are reproduced. We denote $\psi_* = e^{T_*^{\text{FCI}}} \phi_0$ the exact solution on \mathcal{H}_K . We split the amplitudes such that $T_*^{\text{FCI}} = T_*^{\text{CAS}} + T_*^{\text{ext}}$ with $t_*^{\text{CAS}} \in \mathcal{V}_{\text{CAS}}$ and $t_*^{\text{ext}} \in \mathcal{V}_{\text{ext}}$.

THEOREM 6. *Let E be any eigenvalue of H and assume ψ_* satisfies the Schrödinger equation. Then $f(t_*^{\text{ext}}; t_*^{\text{CAS}}) = 0$ and $E = \mathcal{E}(t_*^{\text{ext}}; t_*^{\text{CAS}})$.*

Proof. Let $\mu \in \mathcal{J}_{\text{ext}}$ and choose $\psi' = e^{-(T_*^{\text{ext}})^\dagger} e^{-(T_*^{\text{CAS}\dagger})} \phi_\mu \in \mathcal{H}_K$. By assumption

$$0 = \langle \psi', (H - E)\psi_* \rangle = \langle \phi_\mu, e^{-T_*^{\text{CAS}}} e^{-T_*^{\text{ext}}} (H - E) e^{T_*^{\text{ext}}} e^{T_*^{\text{CAS}}} \phi_0 \rangle = (f(t_*^{\text{ext}}; t_*^{\text{CAS}}))_\mu .$$

Inserting instead $\psi' = e^{-(T_*^{\text{ext}})^\dagger} e^{-(T_*^{\text{CAS}\dagger})} \phi_0 \in \mathcal{H}_K$ gives $E = \mathcal{E}(t_*^{\text{ext}}; t_*^{\text{CAS}})$. \square

Remark 7. An important observation is that tailoring the CC method with a FCI solution on the CAS, i.e., a solution that corresponds to $t_{\text{FCI}}^{\text{CAS}}$, does not necessarily reproduce the FCI solution on \mathcal{H}_K . More precisely, let $f(t_{\text{FCI}}^{\text{ext}}; t_{\text{FCI}}^{\text{CAS}}) = 0$ then $\psi_*^{(\text{TCC})} = e^{T_{\text{FCI}}^{\text{ext}}} e^{T_{\text{FCI}}^{\text{CAS}}} \phi_0$ is not necessarily a minimizer of Eq. (6) and does therewith in general not fulfill the Schrödinger equation. However, Theorem 6 shows that $f(t^{\text{ext}}; t^{\text{CAS}}) = 0$ is a necessary condition for $\psi = e^{T^{\text{ext}}} e^{T^{\text{CAS}}} \phi_0$ to solve the Schrödinger equation on \mathcal{H}_K . In the continuous formulation of the traditional CC theory, equivalence has been proven in [36] see Theorem 5.3. Equivalence for the projected CC method has been shown in [20] see Section 2.2., using [28].

We emphasize that the CAS part T_*^{CAS} of the exact cluster operator T_*^{FCI} is not equal to the cluster operator that corresponds to the FCI solution on \mathcal{H}_{CAS} . The CAS amplitudes $((t_*^{\text{CAS}})_\mu)_{\mu \in \mathcal{J}_{\text{CAS}}}$ on \mathcal{H}_K are solutions of equations that depend on the external amplitudes. The FCI solution $\psi_{\text{CAS}}^{(\text{FCI})} = e^{T_{\text{FCI}}^{\text{CAS}}} \phi_0$ on \mathcal{H}_{CAS} , however, depends on the Hamilton operator projected onto the CAS space. Hence, in general $T_{\text{FCI}}^{\text{CAS}} \neq T_*^{\text{CAS}}$.

Remark 8. Theorem 6 does not imply *local uniqueness* of $t_* \in \mathcal{V}_{\text{ext}}$, even if t_*^{FCI} is locally unique.

Throughout Subsection 4.1 we consider a fixed and sufficiently good CAS solution, i.e., $\phi_{\text{CAS}} \approx \psi_{\text{CAS}}^{(\text{FCI})} \approx P_{\mathcal{V}_{\text{CAS}}} \psi_*$. As a consequence we will simplify the notation by neglecting the parametric dependency of f and \mathcal{E} on t^{CAS} . We also highlight that the following analysis holds for any TCC scheme, but in particular for TNS-TCC schemes like the DMRG-TCCSD method.

4.1. Local Uniqueness and Residual Bounds. The single-reference CC method, as well as the considered TCC method, are formulated as nonlinear Galerkin schemes. This suggests the use of Zarantonello’s lemma [48] to characterize local uniqueness and residual bounds. This is in line with previous studies on single-reference CC methods [38, 37, 21]. We state without proof:

LEMMA 9 (Local Version of Zarantonello’s lemma [48]). *Let $g : X \rightarrow X'$ be a map between a Hilbert space $(X, \langle \cdot, \cdot \rangle, \|\cdot\|)$ and its dual X' , and let $x_* \in B_\delta$ be a root, $g(x_*) = 0$, where B_δ is an open ball of radius δ around x_* . Assume that g is Lipschitz continuous and locally strongly monotone in B_δ with constants $L > 0$ and $\gamma > 0$, respectively.*

Then the root x_ is unique in B_δ . Indeed, there is a ball $C_\varepsilon \subset X'$ with $0 \in C_\varepsilon$ such that the solution map $g^{-1} : C_\varepsilon \rightarrow X$ exists and is Lipschitz continuous, implying that the equation $g(x_* + x) = y$ has a unique solution $x = g^{-1}(y) - x_*$, depending continuously on y , with norm $\|x\| \leq \delta$. Moreover, let $X_d \subset X$ be a closed subspace such that x_* can be approximated sufficiently well, i.e., the distance $d(x_*, X_d)$ is sufficiently small. Then, the projected problem $g_d(x_d) = 0$ has a unique solution $x_d \in X_d \cap B_\delta$ and*

$$\|x_* - x_d\| \leq \frac{L}{\gamma} d(x_*, X_d) .$$

We emphasize that the above theorem depends strongly on the topology of the considered Hilbert space. We already made the particular choice of $\|\cdot\|_{\mathcal{V}_{\text{ext}}}$ to measure the dynamical correction. This is motivated by the fact that $(\varepsilon_\mu)_{\mu \in \mathcal{J}_{\text{ext}}}$ is computationally accessible. A major difference between the presented analysis and the single-reference CC case [38, 36, 37] is that the assumption of a HOMO-LUMO gap is no longer reasonable. In the context of the TCC method it is assumed that \mathcal{B}_{CAS} and \mathcal{B}_{ext} are chosen such that $\lambda_{k+1} - \lambda_k > 0$. We therefore introduce the CAS-ext gap between λ_k and λ_{k+1} . In analogy to previous literature on analysis of the CC theory, we denote the CAS-ext gap by $\varepsilon_0 = \lambda_{k+1} - \lambda_k$. The assumption of a CAS-ext gap is reasonable under the assumption that \mathcal{H}_{CAS} captures all strong correlation such that the (one-particle) Fock operator's degenerate eigenstates are in the CAS.

Besides the single-particle spectral gap condition, we note that the Fock operator F corresponds to a Hamilton operator with a particular potential V_F in Eq. (1). Consequently, with $V = V_F$ in Eq. (3) we assume

$$(12) \quad \langle \psi, (F + e)\psi \rangle \geq c \|\psi\|_{H^1}^2, \quad \forall \psi \in H^1.$$

For a further discussion on spectral gap and Gårding inequalities in CC theories we refer to [20]. Moreover, in agreement with Section 2, we suppose

$$(13) \quad |\langle \tilde{\psi}, F\psi \rangle| \leq C \|\tilde{\psi}\|_{H^1} \|\psi\|_{H^1}, \quad \forall \psi, \tilde{\psi} \in H^1.$$

One of the main assumption of this article can then be summarized:

Assumption (A). *For the Fock operator F , Eqs. (12) and (13) hold and there exists a CAS-ext gap $\varepsilon_0 = \lambda_{k+1} - \lambda_k > 0$.*

Remark 10. Note that a gap assumption between λ_N and λ_{k+1} is also possible, i.e., $\tilde{\varepsilon}_0 = \lambda_{k+1} - \lambda_N$. We shall refer to this as the extended CAS-ext gap. The difference to ε_0 is that $\tilde{\varepsilon}_0$ is directly proportional to the size of the CAS, i.e., choosing a large CAS yields a large λ_{k+1} and therewith a large value of $\tilde{\varepsilon}_0$. Consequently, this connects the following norm estimates with the CAS. We point out that every following statement holds true for either gap condition, however, the constants involved may differ.

The main argument for considering ε_0 (or $\tilde{\varepsilon}_0$) is that the following analysis holds not only for ground-state approximation schemes but also for excited state approximations, which is a major difference to the previous analyses of single-reference CC methods [38, 36, 37, 21]. In the TCC scheme, the single-reference CC method is used to add a dynamical correction to $\phi_{\text{CAS}} \in \mathcal{H}_{\text{CAS}}$ on the external space \mathcal{H}_{ext} , i.e., it captures dynamical correlations between orbitals in \mathcal{H}_{ext} as well as dynamical correlations between orbitals in \mathcal{H}_{CAS} and \mathcal{H}_{ext} . This correction can be done for any wavefunction $\phi_{\text{CAS}} \in \mathcal{H}_{\text{CAS}}$, in particular also for approximations of excited states in \mathcal{H}_{CAS} . We emphasize that correlations between orbitals in \mathcal{B}_{ext} and \mathcal{B}_{CAS} are not considered when computing ϕ_{CAS} , which introduces a methodological error to the method [10].

Note that Assumption (A) is an assumption on the single-particle spectrum. This allows us to establish $\varepsilon_\mu > \varepsilon_0$ for all $\mu \in \mathcal{J}_{\text{ext}}$, however, it does not necessarily imply $\varepsilon_\sigma \leq \varepsilon_\mu$ for $\sigma \in \mathcal{J}_{\text{CAS}}$ and $\mu \in \mathcal{J}_{\text{ext}}$. Thus, under Assumption (A) we might not have a spectral gap in the N -particle space.

Next, we introduce the Fock norm on \mathcal{H}_{ext} .

DEFINITION 11. *The map $\|\cdot\|_F : \mathcal{H}_{\text{ext}} \rightarrow \mathbb{R}_+$ is given by $\phi \mapsto \sqrt{\langle \phi, (F - \Lambda_0)\phi \rangle}$.*

LEMMA 12. Suppose Assumption (A), then $\|\phi\|_F = \sqrt{\langle \phi, (F - \Lambda_0)\phi \rangle}$ is a norm on \mathcal{H}_{ext} and

$$(14) \quad \langle T\phi_0, (F - \Lambda_0)T\phi_0 \rangle \geq \eta \|T\phi_0\|_{H^1}^2, \quad \forall t \in \mathcal{V}_{\text{ext}},$$

where $\eta > 0$ is defined in the proof. Moreover, $\|\cdot\|_F$ is equivalent to $\|\cdot\|_{H^1}$ on \mathcal{H}_{ext} .

Proof. The assumption of a Gårding inequality of the Fock operator and a spectral gap (Eq. (13)) imply (14). The derivation is given by Lemma 11 in [21] and is here included to highlight the importance of a CAS-ext gap. Before starting the proof, we note that Eq. 12 implies $e \geq \Lambda_0$, since Λ_0 is the smallest eigenvalue of F in \mathcal{H}_K . Then we set $q = \varepsilon_0/(\varepsilon_0 + \Lambda_0 + e) > 0$ and $\eta = qc$, where e, c are the constants from the Gårding inequality (12). Assumption (A) yields $\langle T\phi_0, (F - \Lambda_0)T\phi_0 \rangle \geq \varepsilon_0 \|T\phi_0\|_{L^2}^2$, for $t \in \mathcal{V}_{\text{ext}}$. The Gårding inequality (12) implies

$$\begin{aligned} & \langle T\phi_0, (F - \Lambda_0)T\phi_0 \rangle \\ &= q \langle T\phi_0, (F + e)T\phi_0 \rangle - q \langle T\phi_0, (\Lambda_0 + e)T\phi_0 \rangle + (1 - q) \langle T\phi_0, (F - \Lambda_0)T\phi_0 \rangle \\ &\geq qc \|T\phi_0\|_{H^1}^2 + ((1 - q)\varepsilon_0 - q(\Lambda_0 + e)) \|T\phi_0\|_{L^2}^2 = \eta \|T\phi_0\|_{H^1}^2. \end{aligned}$$

Therefore $\|\phi\|_F = 0$ if and only if $\phi = 0$. The self-adjointness of F gives the triangle inequality and the homogeneity follows immediately. Hence, $\|\phi\|_F$ is a norm. The proof is completed by noting that Eq. (14) and the boundedness of F (Eq. (13)) yield the equivalence of $\|\cdot\|_F$ and $\|\cdot\|_{H^1}$ on \mathcal{H}_{ext} . \square

PROPOSITION 13. For $t \in \mathcal{V}_{\text{ext}}$ $\|t\|_{\mathcal{V}_{\text{ext}}} = \|T\phi_0\|_F$, and in particular $\|t\|_{\mathcal{V}_{\text{ext}}} \sim \|T\phi_0\|_{H^1}$.

Remark 14. Note that the spectral (CAS-ext) gap assumption of F gives

$$\|T\phi_0\|_F^2 = \langle T\phi_0, (F - \Lambda_0)T\phi_0 \rangle \geq \varepsilon_0 \|T\phi_0\|_{L^2}^2,$$

which is the same as the direct estimate $\|t\|_{\mathcal{V}_{\text{ext}}}^2 = \sum_{\mu \in \mathcal{I}_{\text{ext}}} \varepsilon_\mu t_\mu^2 \geq \varepsilon_0 \|t\|_2^2$. This makes the Fock norm natural in the following analysis.

Two useful facts regarding the Fock operator and excitation operators are stated in the following lemma (for a proof see [13]).

LEMMA 15. Let F be the Fock operator, $\mu = (A_1, \dots, A_{|\mu|})_{I_1, \dots, I_{|\mu|}}$ and $T = \sum_{\mu \in \mathcal{J}} t_\mu X_\mu$. Then

$$[F, X_\mu] = \sum_{j=1}^{|\mu|} (\lambda_{A_j} - \lambda_{I_j}) X_\mu = \varepsilon_\mu X_\mu \quad \text{and} \quad e^{-T} F e^T = F + [F, T].$$

Proof of Proposition 13. Let $t \in \mathcal{V}_{\text{ext}}$, we find by means of Lemma 15

$$\begin{aligned} \|T\phi_0\|_F^2 &= \langle T\phi_0, (F - \Lambda_0)T\phi_0 \rangle = \sum_{\mu, \nu \in \mathcal{J}_{\text{ext}}} t_\mu t_\nu \langle \phi_\mu, (F - \Lambda_0)\phi_\nu \rangle \\ &= \sum_{\mu, \nu \in \mathcal{J}_{\text{ext}}} t_\mu t_\nu \langle \phi_\mu, [F, X_\nu]\phi_0 \rangle = \sum_{\mu \in \mathcal{J}_{\text{ext}}} t_\mu^2 \varepsilon_\mu = \|t\|_{\mathcal{V}_{\text{ext}}}^2. \end{aligned}$$

\square

Remark 16. The first formula of Lemma 15 uses the fact that F is diagonal, i.e., a finite K . However, the fact that $[F, X_\mu]$ is a cluster operator can be proven using only the F -orthogonality of an occupied χ_I and an unoccupied χ_A . Thus, while in the infinite-dimensional case the first statement certainly fails due to the continuous spectrum, it is reasonable to expect that the second statement still stands.

THEOREM 17. *Under Assumption (A) the norm equivalence $\|T\|_{\mathcal{B}(H^1)} \sim \|t\|_{\mathcal{V}_{\text{ext}}}$ holds for $t \in \mathcal{V}_{\text{ext}}$.*

To show this we first prove the following lemma.

LEMMA 18. *Let $\nu \in \mathcal{J}_{\text{ext}}$ and $\alpha, \mu \in \mathcal{J}$ with $|\alpha|, |\mu| \leq |\nu|$ and $\langle \phi_\nu, X_\alpha \phi_\mu \rangle \neq 0$. Then there exists a constant $C \geq 0$ such that*

$$i) \quad \frac{\varepsilon_\nu}{\varepsilon_\mu} \leq C\varepsilon_\alpha, \text{ if } \alpha, \mu \in \mathcal{J}_{\text{ext}} \quad ii) \quad \varepsilon_\nu \leq C\varepsilon_\alpha, \text{ if } \alpha \in \mathcal{J}_{\text{ext}} \text{ and } \mu \in \mathcal{J}_{\text{CAS}} .$$

Proof. Set $\delta = (\lambda_{k+1} + \lambda_k)/2$ and define $\bar{\lambda}_\nu = \max\{\lambda_{A_j} : j = 1, \dots, |\nu|\} - \delta$, which is well-defined since K is finite. We first demonstrate, following Lemma 4.14 in [38], for all $\nu \in \mathcal{J}_{\text{ext}}$ there exists a $C > 0$ such that

$$(15) \quad C^{-1}\varepsilon_\nu \leq \bar{\lambda}_\nu \leq \varepsilon_\nu .$$

Let $\nu \in \mathcal{J}_{\text{ext}}$. It is immediate that $\varepsilon_0^{-1} \geq \varepsilon_\nu^{-1}$. From the definition of $\bar{\lambda}_\nu$, we conclude

$$\varepsilon_\nu = \sum_{j=1}^{|\nu|} (\lambda_{A_j} - \lambda_{I_j}) \leq N(\bar{\lambda}_\nu - (\lambda_1 - \delta)) .$$

Since $\bar{\lambda}_\nu \geq \lambda_{k+1} - \delta$ it follows $\bar{\lambda}_\nu \geq \varepsilon_0/2$, which is equivalent to $(2\bar{\lambda}_\nu)^{-1} \leq \varepsilon_0^{-1}$. This implies $|\lambda_1 - \delta| \leq 2|\lambda_1 - \delta|\bar{\lambda}_\nu/\varepsilon_0$. Thus,

$$N^{-1}\varepsilon_\nu \leq \bar{\lambda}_\nu + |\lambda_1 - \delta| \leq (1 + 2|\lambda_1 - \delta|/\varepsilon_0)\bar{\lambda}_\nu ,$$

which proves the first inequality of Eq. (15). For the second inequality we define $\lambda_{A_{j^*}} = \max\{\lambda_{A_j} : j = 1, \dots, |\nu|\}$ and note that $\varepsilon_\nu \geq \lambda_{A_{j^*}} - \lambda_{I_{j^*}} \geq \lambda_{A_{j^*}} - \delta = \bar{\lambda}_\nu$. We now prove the lemma considering three cases:

i) Let $\alpha, \mu \in \mathcal{J}_{\text{ext}}$ and $\bar{\lambda}_\alpha \geq \bar{\lambda}_\mu$. Then $\bar{\lambda}_\alpha = \bar{\lambda}_\nu$ and we estimate

$$\frac{\varepsilon_\nu}{\varepsilon_\mu} \leq \frac{C\bar{\lambda}_\nu}{\varepsilon_0} = \frac{C}{\varepsilon_0}\bar{\lambda}_\alpha \leq \frac{C}{\varepsilon_0}\varepsilon_\alpha .$$

ii) Let $\alpha, \mu \in \mathcal{J}_{\text{ext}}$ and $\bar{\lambda}_\alpha \leq \bar{\lambda}_\mu$. Then $\bar{\lambda}_\mu = \bar{\lambda}_\nu$ and using $(2\bar{\lambda}_\alpha)^{-1} \leq \varepsilon_0^{-1}$ we obtain

$$\frac{\varepsilon_\nu}{\varepsilon_\mu} \leq \frac{C\bar{\lambda}_\nu}{\bar{\lambda}_\mu} = \frac{2\bar{\lambda}_\alpha}{2\bar{\lambda}_\alpha}C \leq \frac{2C}{\varepsilon_0}\varepsilon_\alpha .$$

iii) Let $\alpha \in \mathcal{J}_{\text{ext}}$ and $\mu \in \mathcal{J}_{\text{CAS}}$. Then $\bar{\lambda}_\alpha = \bar{\lambda}_\nu$ and $\varepsilon_\nu \leq C\bar{\lambda}_\nu = C\bar{\lambda}_\alpha \leq C\varepsilon_\alpha$. \square

Proof of Theorem 17. Proposition 13 implies the inequality $\|t\|_{\mathcal{V}_{\text{ext}}} \lesssim \|T\phi_0\|_{H^1} \leq \|T\|_{\mathcal{B}(H^1)}\|\phi_0\|_{H^1}$. Consequently, it remains to show that $\|T\psi\|_{H^1} \leq C\|t\|_{\mathcal{V}_{\text{ext}}}\|\psi\|_{H^1}$ for $\psi \in \text{span}\{\phi_0\}^\perp$ (in the L^2 -sense). Let $\psi = \sum_{\mu \in \mathcal{J}} s_\mu \phi_\mu = S\phi_0 \in \mathcal{H}_K$, $T = \sum_{\alpha \in \mathcal{J}_{\text{ext}}} t_\alpha X_\alpha$ and $s = (s_\mu)_{\mu \in \mathcal{J}}$, where we assume without loss of generality that

$(s_\mu)_{\mu \in \mathcal{J}} = ((s_\mu)_{\mu \in \mathcal{J}_{\text{ext}}}, (s_\mu)_{\mu \in \mathcal{J}_{\text{CAS}}})$. Note that the product TS is an excitation operator with cluster amplitudes in \mathcal{V}_{ext} . Hence, Proposition 13 yields

$$(16) \quad \begin{aligned} \|T\psi\|_{H^1}^2 &= \|TS\phi_0\|_{H^1}^2 \sim \|(\langle \phi_\nu, TS\phi_0 \rangle)_{\nu \in \mathcal{J}_{\text{ext}}}\|_{\mathcal{V}_{\text{ext}}}^2 = \|(\langle \phi_\nu, T\psi \rangle)_{\nu \in \mathcal{J}_{\text{ext}}}\|_{\mathcal{V}_{\text{ext}}}^2 \\ &= \sum_{\nu \in \mathcal{J}_{\text{ext}}} \left(\varepsilon_\nu^{1/2} \left| \sum_{\alpha \in \mathcal{J}_{\text{ext}}} \sum_{\mu \in \mathcal{J}} t_\alpha s_\mu \langle \phi_\nu, X_\alpha \phi_\mu \rangle \right| \right)^2. \end{aligned}$$

We now define $A = (\langle \phi_\nu, T\phi_\mu \rangle)_{\nu \in \mathcal{J}_{\text{ext}}, \mu \in \mathcal{J}}$, $D = \text{diag}(\varepsilon_\nu^{1/2})_{\nu \in \mathcal{J}_{\text{ext}}}$ and $\tilde{D} = \text{diag}(D, I)$. The operator inequality $\|TS\phi_0\|_{H^1}^2 \leq \|S\|_{\mathcal{B}(H^1)}^2 \|T\phi_0\|_{H^1}^2$ yields with Eq. (16) that $\|t\|_{\mathcal{V}}^2 \sim \|DA\tilde{D}^{-1}\tilde{D}s\|_2^2$. We estimate $\|DA\tilde{D}^{-1}\|_2$ by means of Lemma 18:

i) Let $\mu \in \mathcal{J}_{\text{ext}}$. Then

$$\tilde{a}_{\nu, \mu} = \left(\frac{\varepsilon_\nu}{\varepsilon_\mu} \right)^{1/2} \sum_{\alpha \in \mathcal{J}_{\text{ext}}} t_\alpha \langle \phi_\nu, X_\alpha \phi_\mu \rangle \lesssim \sum_{\alpha \in \mathcal{J}_{\text{ext}}} t_\alpha \varepsilon_\alpha^{1/2} \langle \phi_\nu, X_\alpha \phi_\mu \rangle.$$

ii) Let $\mu \in \mathcal{J}_{\text{CAS}}$. Then

$$\tilde{a}_{\nu, \mu} = \varepsilon_\nu^{1/2} \sum_{\alpha \in \mathcal{J}_{\text{ext}}} t_\alpha \langle \phi_\nu, X_\alpha \phi_\mu \rangle \lesssim \sum_{\alpha \in \mathcal{J}_{\text{ext}}} t_\alpha \varepsilon_\alpha^{1/2} \langle \phi_\nu, X_\alpha \phi_\mu \rangle.$$

Hence, $\|DA\tilde{D}^{-1}\|_2^2 \leq C \sum_{\alpha \in \mathcal{J}_{\text{ext}}} t_\alpha^2 \varepsilon_\alpha = C \|t\|_{\mathcal{V}_{\text{ext}}}^2$ and $\|T\|_{\mathcal{B}(H^1)} \leq C \|t\|_{\mathcal{V}_{\text{ext}}}$. The norm equivalence follows since $\|t\|_{\mathcal{V}_{\text{ext}}} \sim \|T\psi\|_{H^1} \sim \|T\|_{\mathcal{B}(H^1)}$. \square

We show the applicability of Lemma 9 by establishing Lipschitz continuity of the TCC function.

THEOREM 19. *The function $f : \mathcal{V}_{\text{ext}} \rightarrow \mathcal{V}'_{\text{ext}}$, given in Definition 5, is differentiable at $t \in \mathcal{V}_{\text{ext}}$. Furthermore, the derivative is Lipschitz continuous as well as all higher derivatives. In particular, for any ball $B_r(t_*) \subseteq \mathcal{V}_{\text{ext}}$ there exists a Lipschitz constant L depending on r and t_* such that*

$$(17) \quad \|f(t_1) - f(t_2)\|_{\mathcal{V}'_{\text{ext}}} \leq L \|t_1 - t_2\|_{\mathcal{V}_{\text{ext}}}$$

for $t_1, t_2 \in B_r(t_*)$.

Proof. For the derivative of f we find

$$Df(t) : \mathcal{V}_{\text{ext}} \rightarrow \mathcal{V}'_{\text{ext}} ; s \mapsto \langle \phi_\mu, e^{-T} [e^{-T^{\text{CAS}}} H e^{T^{\text{CAS}}}, S] e^T \phi_0 \rangle.$$

Note that Theorem 17 yields $T^\dagger \in \mathcal{B}(H^{-1})$ for any cluster amplitude vector $t \in \mathcal{V}_{\text{ext}}$. Then, using $H : H^1 \rightarrow H^{-1}$ we obtain $|\langle Df(t)s, u \rangle| \leq C \|s\|_{\mathcal{V}_{\text{ext}}} \|u\|_{\mathcal{V}_{\text{ext}}}$ for given $s, u \in \mathcal{V}_{\text{ext}}$. This shows the boundedness of $f'(t) : \mathcal{V}_{\text{ext}} \rightarrow \mathcal{V}'_{\text{ext}}$, hence, $f : \mathcal{V}_{\text{ext}} \rightarrow \mathcal{V}'_{\text{ext}}$ is differentiable at $t \in \mathcal{V}_{\text{ext}}$. The continuity of the Coulomb potential [47] and the fluctuation potential $W = H - F$ [24] further implies the continuity of $t \mapsto f'(t)$. Hence f is local Lipschitz continuous on $B_r(t_*)$. Higher order derivatives are treated in the same way. \square

To prove that f is locally strongly monotone, we use the decomposition

$$(18) \quad H = F + PWP + (W - PWP),$$

where W is the fluctuation operator and P is the orthogonal projection onto the CAS. The decomposition is motivated from a perturbation theory point of view as follows:

Suppose $\lambda = \|W - PWP\|_{\mathcal{B}(H^1, H^{-1})} = 0$. Then it is straightforward to see that \mathcal{H}_{CAS} is an invariant subspace for H , and hence the CAS FCI problem is exact. Therefore, $t_* = 0$ is a solution in this case, as can easily be checked. Also, the CAS-ext gap at least intuitively indicates that the TCC function f is locally strongly monotone at $t_* = 0$. (This can also be checked.) Now, suppose $\lambda = \|W - PWP\|_{\mathcal{B}(H^1, H^{-1})}$ is finite and sufficiently small. It is reasonable to expect that $t_*(\lambda)$ is correspondingly small, i.e., a small perturbation of the case $W - PWP = 0$, staying within the domain of strong monotonicity. In conclusion, we expect that under some smallness assumption on $W - PWP$ it is achievable to demonstrate local strong monotonicity of the TCC function f . We also note that by enlarging the CAS, $W - PWP$ becomes smaller, so that tuning the CAS can be an important tool to achieve proper smallness in practice.

For a fixed T^{CAS} , we define the map

$$O : \mathcal{V}_{\text{ext}} \rightarrow H^{-1} ; t \mapsto (e^{-T}(W_{\text{CAS}} - PW_{\text{CAS}}P)e^T - (W_{\text{CAS}} - PW_{\text{CAS}}P))\phi_0 ,$$

where $W_{\text{CAS}} = \exp(-T^{\text{CAS}})W \exp(T^{\text{CAS}})$. Similarly to Theorem 19, we find that $O(\cdot)$ is differentiable with

$$DO(s) : \mathcal{V}_{\text{ext}} \rightarrow H^{-1} ; t \mapsto [e^{-S}(W_{\text{CAS}} - PW_{\text{CAS}}P)e^S, T]\phi_0 ,$$

which implies locally Lipschitz continuity. For technical reasons, we will make use of a Lipschitz condition with respect to the l^2 -norm, which is no restriction since all norms are equivalent in finite dimensions.

Assumption (B). *There exists a ball $B_\delta(t_*) \subset \mathcal{V}_{\text{ext}}$ such that for $t_1, t_2 \in B_\delta(t_*)$ we have*

$$\|O(t_1) - O(t_2)\|_{L^2} \leq L_* \|t_1 - t_2\|_2 ,$$

where the Lipschitz constant $L_* > 0$ fulfills

$$(19) \quad \varepsilon_0 - \omega_0 - \Omega_{\text{CAS}} > L_* ,$$

with $\Omega_{\text{CAS}} = \sum_{\sigma \in \mathcal{J}_{\text{CAS}}} |t_\sigma^{\text{CAS}} \varepsilon_\sigma|$, $\omega_0 = \langle \phi_0, W_{\text{CAS}} \phi_0 \rangle$ and ε_0 the previously defined CAS-ext gap.

Remark 20. We note that the assumption of Lipschitz continuity of $O(\cdot)$ in Assumption (B) is more than actually needed. The crucial requirement is

$$|\langle (T_1 - T_2)\phi_0, O(t_1) - O(t_2) \rangle| \leq C_* \|t_1 - t_2\|_2^2 ,$$

for some relatively small C_* . However, this constant C_* can be bounded from above in terms of the Lipschitz constant $L_* > 0$ of $O(\cdot)$ since $C_* \leq CL_*$, where by Proposition 13 a constant C exists fulfilling $\|t\|_{\mathcal{V}_{\text{ext}}} \leq C \|T\phi_0\|_{H^1}$. Furthermore,

$$(20) \quad \begin{aligned} \|DO(s)\|_{\mathcal{B}(L^2)} &\sim \delta_{W_{\text{CAS}}} := \|W_{\text{CAS}} - PW_{\text{CAS}}P\|_{\mathcal{B}(L^2)} \\ &\leq \sum_k \frac{1}{k!} \|[W - PWP, T^{\text{CAS}}]_{(k)}\|_{\mathcal{B}(L^2)} , \end{aligned}$$

such that $L_* \sim \delta_{W_{\text{CAS}}}$ and C_* fulfills Eq. (19) under the assumption that $W - PWP$ is sufficiently small related to T^{CAS} as displayed in the rhs. of Eq. (20). The latter aligns with a perturbational viewpoint of the TCC method as outlined above. Note that we *do not* impose a norm restriction on W itself but an ideal CAS, meaning that the multireference character is captured within the CAS, i.e., PWP . The norm restriction on $W - PWP$ then becomes a natural consequence of the optimal CAS choice.

Remark 21. Note that since $\phi_{\text{CAS}} = e^{T^{\text{CAS}}} \phi_0$ is an approximate solution on the CAS, ω_0 accounts for the non-trivial energy correction (vis-a-vis ϕ_0) and thus is negative for quantum-molecular systems. Typically then, $\omega_0 < 0$ and the CAS-ext gap ε_0 together with $|\omega_0|$ have to be large enough such that $\varepsilon_0 + |\omega_0| > \Omega_{\text{CAS}}$. Furthermore, Assumption (B) allows t_σ^{CAS} to be relatively large for $\sigma \in \mathcal{J}_{\text{CAS}}$ with ε_σ small. A not too big Ω_{CAS} can be guaranteed if $\{\lambda_j\}_{j=1}^k$ is densely confined because $|\varepsilon_\sigma| \leq N(\lambda_k - \lambda_1)$.

We are now able to prove that f is locally strongly monotone.

THEOREM 22. *Under Assumption (A) and (B), the TCC function f is locally strongly monotone on $B_\delta(t_*)$ for some $\delta > 0$.*

Proof. Let $t_1, t_2 \in B_\delta(t_*) \subseteq \mathcal{V}_{\text{ext}}$ and write the Hamiltonian as in Eq. (18). With the notation $\delta_f = \langle f(t_1) - f(t_2), t_1 - t_2 \rangle$, $\delta_T = T_1 - T_2$ and $H_{t_i} = e^{-T_i} H e^{T_i}$, the definition of the TCC function f and Lemma 15 yield

$$\begin{aligned} \delta_f &= \langle \delta_T \phi_0, e^{-T^{\text{CAS}}} (H_{t_1} - H_{t_2}) e^{T^{\text{CAS}}} \phi_0 \rangle \\ &= \langle \delta_T \phi_0, e^{-T^{\text{CAS}}} [F, \delta_T] e^{T^{\text{CAS}}} \phi_0 \rangle + \langle \delta_T \phi_0, (e^{-T_1} P W_{\text{CAS}} P - e^{-T_2} P W_{\text{CAS}} P) \phi_0 \rangle \\ &\quad + \langle \delta_T \phi_0, O(t_1) - O(t_2) \rangle \\ &= \delta_1 + \delta_2 + \delta_3, \end{aligned}$$

where the last equality defines δ_1, δ_2 and δ_3 .

To bound δ_1 from below, we first note that Lemma 15 implies

$$[F, e^{T^{\text{CAS}}}] = \sum_{n=1}^N \frac{1}{n!} \sum_{\mu \in \mathcal{J}_{\text{CAS}}} (t_{\text{CAS}}^{(n)})_\mu X_\mu = S.$$

Since S commutes with $e^{\pm T^{\text{CAS}}}$ and δ_T , we obtain

$$e^{-T^{\text{CAS}}} [F, \delta_T] e^{T^{\text{CAS}}} = e^{-T^{\text{CAS}}} ((S + e^{T^{\text{CAS}}} F) \delta_T - \delta_T (S + e^{T^{\text{CAS}}} F)) = F \delta_T - \delta_T F,$$

and consequently $\delta_1 = \langle \delta_T \phi_0, (F - \Lambda_0) \delta_T \phi_0 \rangle = \sum_{\mu \in \mathcal{J}_{\text{ext}}} \varepsilon_\mu (t_1 - t_2)_\mu^2$.

Next we find

$$\begin{aligned} \delta_2 &= \langle \delta_T \phi_0, (e^{-T_1} P W_{\text{CAS}} - e^{-T_2} P W_{\text{CAS}}) \phi_0 \rangle \\ &= -\langle \delta_T \phi_0, \delta_T P W_{\text{CAS}} \phi_0 \rangle + \sum_{k=2}^{\infty} \frac{(-1)^k}{k!} \langle \delta_T \phi_0, (T_2^k - T_1^k) P W_{\text{CAS}} \phi_0 \rangle \\ &= -\sum_{\mu \in \mathcal{J}_{\text{ext}}} (t_1 - t_2)_\mu^2 \langle \phi_0, P W_{\text{CAS}} \phi_0 \rangle \\ (21) \quad &\quad - \sum_{\substack{\mu \neq \nu \in \mathcal{J}_{\text{ext}} \\ \mu \ominus \nu \in \text{CAS}}} (t_1 - t_2)_\mu (t_1 - t_2)_\nu \langle \phi_{\mu \ominus \nu}, P W_{\text{CAS}} \phi_0 \rangle \\ &\quad + \sum_{k=2}^{\infty} \frac{(-1)^k}{k!} \langle \delta_T \phi_0, (T_2^k - T_1^k) P W_{\text{CAS}} \phi_0 \rangle. \end{aligned}$$

We now define $\delta\Psi = \phi_{\text{CAS}} - \psi_{\text{CAS}}^{(\text{FCI})}$ with $\phi_{\text{CAS}} = \exp(T^{\text{CAS}}) \phi_0 \approx \psi_{\text{CAS}}^{(\text{FCI})}$, where

$PHP\Psi_{\text{CAS}}^* = E_{\text{CAS}}^{(\text{FCI})}\Psi_{\text{CAS}}^*$. We know that

$$\begin{aligned} PW_{\text{CAS}}P\phi_0 &= Pe^{-T^{\text{CAS}}}H\phi_{\text{CAS}} - Pe^{-T^{\text{CAS}}}Fe^{T^{\text{CAS}}}\phi_0 \\ &= E_{\text{CAS}}^{(\text{FCI})}Pe^{-T^{\text{CAS}}}\Psi_{\text{CAS}}^* + Pe^{-T^{\text{CAS}}}H\delta\Psi - P(F + [F, T^{\text{CAS}}])\phi_0 \\ &= E_{\text{CAS}}^{(\text{FCI})}\phi_0 + Pe^{-T^{\text{CAS}}}(H - E_{\text{CAS}}^{(\text{FCI})})\delta\Psi - P(F + [F, T^{\text{CAS}}])\phi_0. \end{aligned}$$

Since we are merely interested in the projections onto ϕ_σ with $\sigma \in \mathcal{J}_{\text{CAS}}$, we set $\mathcal{R} = \langle \phi_\sigma, Pe^{-T^{\text{CAS}}}(H - E_{\text{CAS}}^{(\text{FCI})})\delta\Psi \rangle$ and obtain

$$\begin{aligned} (22) \quad \langle \phi_\sigma, PW_{\text{CAS}}\phi_0 \rangle &= \langle \phi_\sigma, Pe^{-T^{\text{CAS}}}(H - E_{\text{CAS}}^{(\text{FCI})})\delta\Psi \rangle - \langle \phi_\sigma, [F, T^{\text{CAS}}]\phi_0 \rangle \\ &= \mathcal{R} + \langle \phi_\sigma, T^{\text{CAS}}F\phi_0 \rangle - \langle \phi_\sigma, FT^{\text{CAS}}\phi_0 \rangle \\ &= \mathcal{R} + \sum_{\mu} t_{\mu}^{\text{CAS}}\Lambda_0 \langle \phi_\sigma, \phi_{\mu} \rangle - \sum_{\mu} t_{\mu}(\Lambda_0 + \varepsilon_{\sigma}) \langle \phi_\sigma, \phi_{\mu} \rangle \\ &= \mathcal{R} + t_{\sigma}\Lambda_0 - t_{\sigma}(\Lambda_0 + \varepsilon_{\sigma}) = \mathcal{R} - t_{\sigma}\varepsilon_{\sigma}. \end{aligned}$$

The quantity $\mathcal{R} \sim \|\delta\Psi\|_{L^2}$ is directly steerable by the used CAS method. Hence, assuming $\phi_{\text{CAS}} \approx \psi_{\text{CAS}}^{(\text{FCI})}$ to be a sufficiently good approximation eliminates the above \mathcal{R} dependence. Inserting Eq. (22) in the second term of Eq. (21), we find

$$\begin{aligned} (23) \quad &\sum_{\substack{\mu \neq \nu \in \mathcal{J}_{\text{ext}} \\ \mu \ominus \nu \in \text{CAS}}} |(t_1 - t_2)_{\mu}(t_1 - t_2)_{\nu}| t_{\mu \ominus \nu} |\varepsilon_{\mu \ominus \nu}| \\ &\leq \left(\sum_{\substack{\mu \neq \nu \in \mathcal{J}_{\text{ext}} \\ \mu \ominus \nu \in \text{CAS}}} (t_1 - t_2)_{\mu}^2 |t_{\mu \ominus \nu} \varepsilon_{\mu \ominus \nu}| \right)^{\frac{1}{2}} \times \left(\sum_{\substack{\mu \neq \nu \in \mathcal{J}_{\text{ext}} \\ \mu \ominus \nu \in \text{CAS}}} (t_1 - t_2)_{\nu}^2 |t_{\mu \ominus \nu} \varepsilon_{\mu \ominus \nu}| \right)^{\frac{1}{2}} \\ &\leq \Omega_{\text{CAS}} \|t_1 - t_2\|_2^2, \end{aligned}$$

where we recall that $\Omega_{\text{CAS}} = \sum_{\sigma \in \mathcal{J}_{\text{CAS}}} |t_{\sigma} \varepsilon_{\sigma}|$ as defined in Assumption (B). Since $\omega_0 = \langle \phi_0, W_{\text{CAS}}\phi_0 \rangle$ and $\|T_1^k - T_2^k\|_{L^2} \in \mathcal{O}(\|t_1 - t_2\|_2^k)$, we conclude with Proposition 13 that

$$(24) \quad \delta_2 \geq -(\omega_0 + \Omega_{\text{CAS}}) \|t_1 - t_2\|_2^2 + \mathcal{O}(\|t_1 - t_2\|_{\mathcal{V}_{\text{ext}}}^3).$$

For the last term, Assumption (B) implies that

$$\delta_3 \geq -\|\delta_T\phi_0\|_{L^2} \|O(t_1) - O(t_2)\|_{L^2} \geq -L_* \|t_1 - t_2\|_2^2.$$

Combining the different bounds above and assuming that ε_0 , ω_0 , Ω_{CAS} , and L_* fulfill Eq. (19), we conclude the existence of a $\lambda \in (0, 1)$ such that

$$\begin{aligned} \delta_f &\geq \lambda \sum_{\mu} \varepsilon_{\mu} (t_1 - t_2)_{\mu}^2 + \sum_{\mu} \left[(1 - \lambda) \varepsilon_{\mu} - \omega_0 - (L_* + \Omega_{\text{CAS}}) \right] (t_1 - t_2)_{\mu}^2 \\ &\quad + \mathcal{O}(\|t_1 - t_2\|_{\mathcal{V}_{\text{ext}}}^3) \\ &\geq \lambda \|t_1 - t_2\|_{\mathcal{V}_{\text{ext}}}^2 + \sum_{\mu} \left[(1 - \lambda) \varepsilon_0 - \omega_0 - (L_* + \Omega_{\text{CAS}}) \right] (t_1 - t_2)_{\mu}^2 \\ &\quad + \mathcal{O}(\|t_1 - t_2\|_{\mathcal{V}_{\text{ext}}}^3) \\ &\geq \lambda \|t_1 - t_2\|_{\mathcal{V}_{\text{ext}}}^2 + \mathcal{O}(\|t_1 - t_2\|_{\mathcal{V}_{\text{ext}}}^3) \geq \gamma \|t_1 - t_2\|_{\mathcal{V}_{\text{ext}}}^2 \sim \|\delta_T\phi_0\|_{H^1}^2. \end{aligned}$$

In the last step we have assumed δ to be sufficiently small such that $\mathcal{O}(\|t_1 - t_2\|_{\mathcal{V}_{\text{ext}}}^3)$ can be absorbed. \square

Remark 23. We note that Eq. (23) is a pessimistic estimation, since we neglect the conditions of the excitation indices, i.e., $\mu \neq \nu$ such that $\mu \ominus \nu \in \mathcal{J}_{\text{CAS}}$. This restriction means that the excitation rank of the CC method dictates which CAS amplitudes are considered. In particular, considering the DMRG-TCCSD method we find

$$\sum_{\substack{\nu \in \mathcal{J}_{\text{ext}} \\ \nu \neq \mu \\ \mu \ominus \nu \in \text{CAS}}} |t_{\mu \ominus \nu} \varepsilon_{\mu \ominus \nu}| \leq \sum_{\sigma \in \mathcal{J}_{\text{CAS}}^{(1)}} |t_{\sigma} \varepsilon_{\sigma}|,$$

where the superscripted $\mathcal{J}_{\text{CAS}}^{(1)}$ means that only single-excitations on the CAS are considered—which correspond to orbital rotations. Moreover, without loss of generality one can assume Brueckner type orbitals, implying that this term vanishes.

By Theorem 19 and 22, we can apply Lemma 9 to the TCC function f ensuring a locally unique and quasi-optimal approximate solutions. Next, we will show quadratic convergence of tailored coupled-cluster methods which aligns the non-variational TCC approach with any variational method in terms of convergence speed.

4.2. Error Estimate. In this section we present an estimate for the energy error introduced by truncating the TCC method, e.g. DMRG-TCCSD. In comparison to the single-reference CC method, the error is divided into different parts as a consequence of the basis splitting. The TCC function is typically parameterized by an approximation T^{CAS} of the FCI solution $T_{\text{FCI}}^{\text{CAS}}$ on \mathcal{H}_{CAS} . We emphasize that $T_{\text{FCI}}^{\text{CAS}}$ is in itself an approximation of the inaccessible T_*^{CAS} (cf. Theorem 6 and the following discussion). This of course influences the error and is here accounted for. On top of that, the truncation error of the CC method applied to $\phi_{\text{CAS}} = e^{T^{\text{CAS}}} \phi_0$ enters. For this part of the error we follow the analysis of the single-reference CC methods and use the Aubin-Nitsche-duality method for nonlinear Galerkin schemes, see [37]. We consider d -dimensional approximation spaces $\mathcal{V}_{\text{ext}}^{(d)}$, $d \leq |\mathcal{J}|$, of the external amplitude space \mathcal{V}_{ext} . For a given T^{CAS} we denote $t_d \in \mathcal{V}_{\text{ext}}^{(d)}$ the solution of $P_d f(\cdot; t^{\text{CAS}})|_{\mathcal{V}_{\text{ext}}^{(d)}} = 0$, where P_d is the l^2 -orthogonal projection onto $(\mathcal{V}_{\text{ext}}^{(d)})'$. Thus, t_d is an approximation of the full solution $t_* \in \mathcal{V}_{\text{ext}}$, where t_* solves $f(\cdot; t^{\text{CAS}}) = 0$ on \mathcal{V}_{ext} .

Remark 24. In practice, the space $\mathcal{V}_{\text{ext}}^{(d)}$ is constructed by restricting the CC amplitudes to a particular subspace, e.g., allowing excitations from the reference ϕ_0 into the external space of rank less than a fixed number, say, including up to singles and doubles. This choice is practical (the dimension d is fairly low), however, the alternative truncation that allows excitations from *any* CAS determinant ϕ_α into the external space of rank less than a fixed number yields what is called the first-order interaction space [26]. While the dimension can be much higher than the previous choice, it gives external correlation energies guaranteed to be correct through second order in $H_1 = W - PWP$. In other words, all CAS determinants are treated on equal footing, which is essential for an optimal multireference treatment. The first truncation scheme puts special significance to the reference ϕ_0 .

We will here derive a general error estimate valid for every choice of method used on \mathcal{H}_{CAS} potentially introducing an additional error on the CAS denoted δE_{CAS} . In notational consistency with the introduction of Section 4, let $\psi_* = e^{T_*^{\text{ext}}} e^{T_*^{\text{CAS}}} \phi_0$ be the exponential parameterization of the FCI solution on \mathcal{H}_K . Then, the energy error

is subsequently split as follows

$$\begin{aligned}
(25) \quad \delta E &= |\mathcal{E}(t_d; t^{\text{CAS}}) - \mathcal{E}(t_*^{\text{ext}}; t_*^{\text{CAS}})| \\
&\leq |\mathcal{E}(t_d; t^{\text{CAS}}) - \mathcal{E}(t_*; t^{\text{CAS}})| + |\mathcal{E}(t_*; t^{\text{CAS}}) - \mathcal{E}(t_*; t_{\text{FCI}}^{\text{CAS}})| \\
&\quad + |\mathcal{E}(t_*; t_{\text{FCI}}^{\text{CAS}}) - \mathcal{E}(t_*^{\text{ext}}; t_*^{\text{CAS}})| \\
&=: \delta\varepsilon + \delta\varepsilon_{\text{CAS}} + \delta\varepsilon_{\text{CAS}}^* ,
\end{aligned}$$

where the last equality defines the different error terms.

The quantity $\delta\varepsilon$ describes the error produced by truncating the TCC method parameterized by $\phi_{\text{CAS}} = e^{T^{\text{CAS}}} \phi_0$. The second term $\delta\varepsilon_{\text{CAS}}$ is connected to the usage of an approximate solution $\psi_{\text{CAS}} = e^{T^{\text{CAS}}} \phi_0$ on \mathcal{H}_{CAS} instead of the FCI solution $\phi_{\text{CAS}}^{(\text{FCI})} = e^{T_{\text{FCI}}^{\text{CAS}}} \phi_0$. We introduce $\tilde{t}_* \in \mathcal{V}_{\text{ext}}$ that solves $f(\tilde{t}_*; t_{\text{FCI}}^{\text{CAS}}) = 0$. Note that the pair $(\tilde{t}_*, t_{\text{FCI}}^{\text{CAS}}) \in \mathcal{V}_{\text{CAS}} \times \mathcal{V}_{\text{ext}}$ is the best solution possible using a given basis splitting. We emphasize, in comparison, that $t_* = (t_*^{\text{CAS}}, t_*^{\text{ext}})$ is a theoretical construct where the basis splitting has been done after computing t_* .

The main result of this section is given below in Theorem 25. The idea is to bound δE by means of the splitting above. We introduce the error δE_{CAS} in the following way: The wavefunction $e^{T_{\text{FCI}}^{\text{CAS}}} \phi_0$ is in general not an eigenfunction of H , however, it is an eigenfunction of PHP where P is the orthogonal projection on \mathcal{H}_{CAS} . We then define

$$(26) \quad \delta E_{\text{CAS}} = |\langle \phi_0, (e^{-T^{\text{CAS}}} PHP e^{T^{\text{CAS}}} - e^{-T_{\text{FCI}}^{\text{CAS}}} PHP e^{T_{\text{FCI}}^{\text{CAS}}}) \phi_0 \rangle| .$$

The energy difference δE_{CAS} describes the error induced by an approximation to the FCI solution on \mathcal{H}_{CAS} . We emphasize that this error depends on the approximation method used. Using the DMRG, which is a variational method, yields a quadratic error bound.

The error $\delta\varepsilon$ is estimated using similar techniques as described in Ref. [37]. To that end, we define the following Euler-Lagrange systems. For notational simplicity we drop again the explicit parameterization by t^{CAS} . We consider the functionals

$$\langle f(t), \cdot \rangle : \mathcal{V}_{\text{ext}} \rightarrow \mathbb{R}; \quad u \mapsto \langle U \phi_0, e^{-T^{\text{CAS}}} e^{-T} H e^T e^{T^{\text{CAS}}} \phi_0 \rangle$$

and

$$\mathcal{E}(\cdot) : \mathcal{V}_{\text{ext}} \rightarrow \mathbb{R}; \quad u \mapsto \langle \phi_0, e^{-T^{\text{CAS}}} e^{-U} H e^U e^{T^{\text{CAS}}} \phi_0 \rangle .$$

We note that $\langle f(t), \cdot \rangle$ is a real-valued linear form whereas $\mathcal{E}(\cdot)$ is a nonlinear functional. The corresponding variational problem

$$(27) \quad \langle f(t), u \rangle = 0 \quad , \forall u \in \mathcal{V}_{\text{ext}}$$

describes the cluster equations. The associated Galerkin approximation on $\mathcal{V}_{\text{ext}}^{(d)} \subseteq \mathcal{V}_{\text{ext}}$ determines $t_d \in \mathcal{V}_{\text{ext}}^{(d)}$ such that

$$(28) \quad \langle f(t_d), u_d \rangle = 0 \quad , \forall u_d \in \mathcal{V}_{\text{ext}}^{(d)} .$$

We use the Euler-Lagrange method to estimate the error $\mathcal{E}(t) - \mathcal{E}(t_d)$. Introducing the dual variable $z \in \mathcal{V}_{\text{ext}}$, we define the Lagrangian

$$(29) \quad \mathcal{L} : \mathcal{V}_{\text{ext}} \times \mathcal{V}_{\text{ext}} \rightarrow \mathbb{R}; \quad (t, z) \mapsto \mathcal{E}(t) - \langle f(t), z \rangle ,$$

and seek for stationary points $(t_*, z_*) \in \mathcal{V}_{\text{ext}} \times \mathcal{V}_{\text{ext}}$ of $\mathcal{L}(\cdot, \cdot)$, i.e.,

$$(30) \quad \mathcal{L}'(t_*, z_*)(u, v) = \begin{cases} \mathcal{E}'(t_*)u - \langle f'(t_*)u, z_* \rangle \\ - \langle f(t_*), v \rangle \end{cases} = 0 ,$$

for all $(u, v) \in \mathcal{V}_{\text{ext}} \times \mathcal{V}_{\text{ext}}$. The Galerkin approximations $(t_d, z_d) \in \mathcal{V}_{\text{ext}}^{(d)} \times \mathcal{V}_{\text{ext}}^{(d)}$ are defined by the discrete Euler-Lagrange system

$$(31) \quad \mathcal{L}'(t_d, z_d)(u_d, v_d) = \begin{cases} \mathcal{E}'(t_d)u_d - \langle f'(t_d)u_d, z_d \rangle \\ - \langle f(t_d), v_d \rangle \end{cases} = 0 ,$$

for all $(u_d, v_d) \in \mathcal{V}_{\text{ext}}^{(d)} \times \mathcal{V}_{\text{ext}}^{(d)}$. We remark that in both situations (30) and (31), the t - respectively the t_d -component of any stationary point is a solution of the cluster equations and the discrete cluster equations, respectively.

The main results of this section now reads:

THEOREM 25. *Let $\mathcal{B} = \{\chi_1, \dots, \chi_K\} \subseteq H^1$ be a set of L^2 -orthonormal spin-orbitals that are split into \mathcal{B}_{CAS} and \mathcal{B}_{ext} . We denote \mathcal{H}_K and \mathcal{H}_{CAS} the FCI space corresponding to \mathcal{B} resp. \mathcal{B}_{CAS} . Let further $t_*^{\text{CAS}} \in \mathcal{V}_{\text{CAS}}$ be the projection of the FCI amplitudes on \mathcal{H}_K onto \mathcal{H}_{CAS} , $t_{\text{FCI}}^{\text{CAS}} \in \mathcal{V}_{\text{CAS}}$ the FCI amplitudes on \mathcal{H}_{CAS} , and $t^{\text{CAS}} \in \mathcal{V}_{\text{CAS}}$ an approximation to $t_{\text{FCI}}^{\text{CAS}}$. Let $\mathcal{V}_{\text{ext}}^{(d)} \subset \mathcal{V}_{\text{ext}}$ be a subspace fulfilling*

$$(32) \quad d(t_*, \mathcal{V}_{\text{ext}}^{(d)}) \leq \frac{\gamma \delta}{\gamma + L} ,$$

where $\gamma, L > 0$ are the monotonicity and Lipschitz constants of $f(\cdot; t^{\text{CAS}})$ on $B_\delta(t_*)$. Then there is a unique solution $t_d \in \mathcal{V}_{\text{ext}}^{(d)}$ of $P_d f(\cdot; t^{\text{CAS}})|_{\mathcal{V}_{\text{ext}}^{(d)}} = 0$ that approximates the solution $t_* \in \mathcal{V}_{\text{ext}}$ of $f(\cdot; t^{\text{CAS}}) = 0$ on \mathcal{V}_{ext} . Let $(z_d, z_*) \in \mathcal{V}_{\text{ext}}^{(d)} \times \mathcal{V}_{\text{ext}}$ be the corresponding dual solutions of $(t_d, t_*) \in \mathcal{V}_{\text{ext}}^{(d)} \times \mathcal{V}_{\text{ext}}$. Further, set $\tilde{t}_* \in \mathcal{V}_{\text{ext}}$ the solution of $f(\cdot; t_{\text{FCI}}^{\text{CAS}}) = 0$ on \mathcal{V}_{ext} and $t_*^{\text{ext}} \in \mathcal{V}_{\text{ext}}$ the projection of the FCI amplitudes on \mathcal{H}_K onto $\mathcal{H}_{\text{CAS}}^\perp$. It then follows that the energy error can be bounded as

$$\begin{aligned} \delta E \lesssim & \|t_d - t_*\|_{\mathcal{V}_{\text{ext}}} (\|t_d - t_*\|_{\mathcal{V}_{\text{ext}}} + \|z_d - z_*\|_{\mathcal{V}_{\text{ext}}}) + \|t_* - t_*^{\text{ext}}\|_{\mathcal{V}_{\text{ext}}}^2 + \|t_* - \tilde{t}_*\|_{\mathcal{V}_{\text{ext}}}^2 \\ & + \|t_{\text{FCI}}^{\text{CAS}} - t_*^{\text{CAS}}\|_2^2 + \|t^{\text{CAS}} - t_{\text{FCI}}^{\text{CAS}}\|_2^2 + \sum_{\substack{\mu \in \mathcal{J}_{\text{ext}} \\ |\mu|=1}} \varepsilon_\mu (\tilde{t}_*)_\mu^2 + \delta E_{\text{CAS}} . \end{aligned}$$

Remark 26. The energy error estimate in Theorem 25 holds for any basis splitting fulfilling the presented conditions. However, in the extremal cases of a minimal or maximal basis splitting, i.e., $k = N$ and $k = K$, the TCC method collapses to the CC and CAS method, respectively.

Remark 27. Since we do not have an equivalence of Theorem 17 for sequences over \mathcal{J}_{CAS} (ε_μ are not guaranteed to be strictly greater than zero for $\mu \in \mathcal{J}_{\text{CAS}}$), we instead bound the sequences over \mathcal{J}_{CAS} using the unweighted l^2 -norm.

We will prove Theorem 25 by first establishing a series of lemmas that relates to the r.h.s. of Eq. (25). We start with the term $\delta \varepsilon_{\text{CAS}}^* = |\mathcal{E}(t_*; t_{\text{FCI}}^{\text{CAS}}) - \mathcal{E}(t_*^{\text{ext}}; t_*^{\text{CAS}})|$.

LEMMA 28. *Under the assumptions of Theorem 25 the following bound holds*

$$\delta \varepsilon_{\text{CAS}}^* \lesssim \|t_* - t_*^{\text{ext}}\|_{\mathcal{V}_{\text{ext}}}^2 + \|t_{\text{FCI}}^{\text{CAS}} - t_*^{\text{CAS}}\|_2^2 .$$

Proof. Recall that $\psi_* = e^{T_*^{\text{ext}}} e^{T_*^{\text{CAS}}} \phi_0$ corresponds to the FCI solution on \mathcal{H}_K and consequently $D\mathcal{E}(t_*^{\text{ext}}; t_*^{\text{CAS}}) = 0$. Taylor expanding $\mathcal{E}(t_*; t_{\text{FCI}}^{\text{CAS}})$ around $(t_*^{\text{CAS}}, t_*^{\text{ext}})$ yields

$$\mathcal{E}(t_*; t_{\text{FCI}}^{\text{CAS}}) - \mathcal{E}(t_*^{\text{ext}}; t_*^{\text{CAS}}) = \frac{1}{2} D^2 \mathcal{E}(t_*^{\text{ext}}; t_*^{\text{CAS}})((e, \tilde{e}), (e, \tilde{e})) + \mathcal{R}^{(3)},$$

where $\tilde{e} = t_* - t_*^{\text{ext}}$, $e = t_{\text{FCI}}^{\text{CAS}} - t_*^{\text{CAS}}$ and $\mathcal{R}^{(3)}$ describes the third order error term. For $H_{t_1+t_2} = e^{-T_1} e^{-T_2} H e^{T_2} e^{T_1}$ with amplitudes $t_1 \in \mathcal{V}_{\text{ext}}$ and $t_2 \in \mathcal{V}_{\text{CAS}}$ we compute

$$(D^2 \mathcal{E}(t_1; t_2))_{\mu, \nu} = \langle \phi_0, [[H_{t_1+t_2}, X_\nu], X_\mu] \phi_0 \rangle = \langle \phi_0, H_{t_1+t_2} X_\nu X_\mu \phi_0 \rangle.$$

Thus, with $H_* = H_{t_*^{\text{ext}} + t_*^{\text{CAS}}}$ and

$$\delta_{\tilde{T}} = \sum_{\mu \in \mathcal{J}_{\text{ext}}} (t_* - t_*^{\text{ext}})_\mu X_\mu, \quad \delta_T = \sum_{\mu \in \mathcal{J}_{\text{CAS}}} (t_{\text{FCI}}^{\text{CAS}} - t_*^{\text{CAS}})_\mu X_\mu$$

we have

$$\begin{aligned} D^2 \mathcal{E}(t_*^{\text{ext}}; t_*^{\text{CAS}})((e, \tilde{e}), (e, \tilde{e})) &= \langle \phi_0, H_*(\delta_{\tilde{T}} + \delta_T)^2 \phi_0 \rangle \\ &\leq 2 \langle \phi_0, H_* \delta_{\tilde{T}}^2 \phi_0 \rangle + 2 \langle \phi_0, H_* (\delta_T)^2 \phi_0 \rangle. \end{aligned}$$

Using Theorem 17, as well as the boundedness of H , we obtain

$$\langle \phi_0, H_* \delta_{\tilde{T}}^2 \phi_0 \rangle \leq C \|\phi_0\|_{H^1}^2 \|\delta_{\tilde{T}}\|_{\mathcal{B}(H^1)}^2 \leq C \|t_* - t_*^{\text{ext}}\|_{\mathcal{V}_{\text{ext}}}^2.$$

By direct computation, we bound the term $\langle \phi_0, H_*(\delta_T)^2 \phi_0 \rangle$ using the $l^2(\mathcal{J}_{\text{CAS}})$ norm

$$\begin{aligned} \langle \phi_0, H_*(\delta_T)^2 \phi_0 \rangle &\leq C \|\delta_T\|_{\mathcal{B}(H^1)}^2 = C \left\| \sum_{\mu \in \mathcal{J}_{\text{CAS}}} (t_{\text{FCI}}^{\text{CAS}} - t_*^{\text{CAS}})_\mu X_\mu \right\|_{\mathcal{B}(H^1)}^2 \\ &\leq C \sum_{\mu \in \mathcal{J}_{\text{CAS}}} (t_{\text{FCI}}^{\text{CAS}} - t_*^{\text{CAS}})_\mu^2 \|X_\mu\|_{\mathcal{B}(H^1)}^2 \leq C \|t_{\text{FCI}}^{\text{CAS}} - t_*^{\text{CAS}}\|_2^2. \quad \square \end{aligned}$$

Next, we analyze the energy difference $\delta \varepsilon_{\text{CAS}} = |\mathcal{E}(t_*; t_{\text{FCI}}^{\text{CAS}}) - \mathcal{E}(t_*; t_*^{\text{CAS}})|$.

LEMMA 29. *Under the assumptions of Theorem 25 the following bound holds*

$$\delta \varepsilon_{\text{CAS}} \lesssim \delta E_{\text{CAS}} + \|t_* - \tilde{t}_*\|_{\mathcal{V}_{\text{ext}}}^2 + \|(T^{\text{CAS}} - T_{\text{FCI}}^{\text{CAS}}) \phi_0\|_{H^1}^2 + \sum_{|\mu|=1} \varepsilon_\mu (\tilde{t}_*)_\mu^2.$$

Proof. Starting from the definition of $\delta \varepsilon_{\text{CAS}}$, we obtain straightforwardly

$$\delta \varepsilon_{\text{CAS}} \leq |\langle \phi_0, (e^{-T^{\text{CAS}}} H e^{T^{\text{CAS}}} - e^{-T_{\text{FCI}}^{\text{CAS}}} H e^{T_{\text{FCI}}^{\text{CAS}}}) \phi_0 \rangle| + \mathcal{R},$$

where $\mathcal{R} = |\langle \phi_0, [(e^{-T^{\text{CAS}}} H e^{T^{\text{CAS}}} - e^{-T_{\text{FCI}}^{\text{CAS}}} H e^{T_{\text{FCI}}^{\text{CAS}}}), e^{T_*}] \phi_0 \rangle|$. Since $\phi_0, e^{T_{\text{FCI}}^{\text{CAS}}} \phi_0$ and $e^{T^{\text{CAS}}} \phi_0$ are elements of \mathcal{H}_{CAS} , we find

$$\begin{aligned} \delta \varepsilon_{\text{CAS}} - \mathcal{R} &\leq |\langle \phi_0, (e^{-T^{\text{CAS}}} H e^{T^{\text{CAS}}} - e^{-T_{\text{FCI}}^{\text{CAS}}} H e^{T_{\text{FCI}}^{\text{CAS}}}) \phi_0 \rangle| \\ &\leq |\langle \phi_0, (e^{-T^{\text{CAS}}} P H P e^{T^{\text{CAS}}} - e^{-T_{\text{FCI}}^{\text{CAS}}} P H P e^{T_{\text{FCI}}^{\text{CAS}}}) \phi_0 \rangle| \\ &\quad + |\langle \phi_0, ([T^{\text{CAS}}, P] H P e^{T^{\text{CAS}}} - [T_{\text{FCI}}^{\text{CAS}}, P] H P e^{T_{\text{FCI}}^{\text{CAS}}}) \phi_0 \rangle|. \end{aligned}$$

For any excitation operator $X = \sum_{\mu \in \mathcal{J}_{\text{CAS}}} c_{\mu} X_{\mu}$, we remark that $XP\psi \in \mathcal{H}_{\text{CAS}}$ for all $\psi \in \mathcal{H}_K$. By definition of \mathcal{H}_{CAS} we also find $XQ\psi \in \mathcal{H}_{\text{ext}}$ for all $\psi \in \mathcal{H}_K$, where $Q = I - P$. Therefore $X = (P + Q)X(P + Q) = PXP + QXQ$ and consequently $[X, P] = [PXP, P] = 0$. Hence, $[T^{\text{CAS}}, P] = [T_{\text{FCI}}^{\text{CAS}}, P] = 0$. In particular,

$$\delta\varepsilon_{\text{CAS}} \leq |\langle \phi_0, (e^{-T^{\text{CAS}}} PHPe^{T^{\text{CAS}}} - e^{-T_{\text{FCI}}^{\text{CAS}}} PHPe^{T_{\text{FCI}}^{\text{CAS}}}) \phi_0 \rangle| + \mathcal{R} = \delta E_{\text{CAS}} + \mathcal{R},$$

where δE_{CAS} is defined by Eq. (26). To estimate \mathcal{R} we consider the splitting of the Hamilton operator $H = F + W$. Note that $[T^{\text{CAS}}, T_*] = [T_{\text{FCI}}^{\text{CAS}}, T_*] = 0$ which implies together with Lemma 15 that the F -dependent terms in \mathcal{R} vanish. The Baker–Campbell–Hausdorff expansion and the fact that $((T_*)^m)^{\dagger} \phi_0 = 0$ for all $m \geq 1$ then yields

$$\mathcal{R} = |\langle \phi_0, \left(\sum_{m=1} \frac{1}{m!} [W, e^{T^{\text{CAS}}}]_m - \sum_{m=1} \frac{1}{m!} [W, e^{T_{\text{FCI}}^{\text{CAS}}}]_m \right) \sum_{m=1} \frac{1}{m!} (T_*)^m \phi_0 \rangle|.$$

Since W is a two-particle operator, the Slater–Condon rules imply that the non-zero contributions in the above expansion are given for $m = 1$ and only by the single-excitation parts of the respective operators. It then follows with $((T^{\text{CAS}})_1)^{\dagger} \phi_0 = ((T_{\text{FCI}}^{\text{CAS}})_1)^{\dagger} \phi_0 = 0$ that

$$\mathcal{R} = |\langle \phi_0, W(T^{\text{CAS}} - T_{\text{FCI}}^{\text{CAS}})_1 (T_*)_1 \phi_0 \rangle|,$$

where $(\cdot)_1$ denotes the single-excitation part of the respective operator. We then estimate

$$\begin{aligned} \mathcal{R} &\leq |\langle \phi_0, W(T^{\text{CAS}} - T_{\text{FCI}}^{\text{CAS}})_1 (T_* - \tilde{T}_*)_1 \phi_0 \rangle| + |\langle \phi_0, W(T^{\text{CAS}} - T_{\text{FCI}}^{\text{CAS}})_1 (\tilde{T}_*)_1 \phi_0 \rangle| \\ &\leq \left(C_1 \|T_* - \tilde{T}_*\|_{\mathcal{B}(H^1)} + C_2 \|(\tilde{T}_*)_1\|_{\mathcal{B}(H^1)} \right) \|(T^{\text{CAS}} - T_{\text{FCI}}^{\text{CAS}}) \phi_0\|_{H^1} \\ &\leq \frac{C_1}{2} \|T_* - \tilde{T}_*\|_{\mathcal{B}(H^1)}^2 + \frac{C_2}{2} \|(\tilde{T}_*)_1\|_{\mathcal{B}(H^1)}^2 + \frac{C_1 + C_2}{2} \|(T^{\text{CAS}} - T_{\text{FCI}}^{\text{CAS}}) \phi_0\|_{H^1}^2. \end{aligned}$$

Hence, $\mathcal{R} \leq D_1 \|t_* - \tilde{t}_*\|_{\mathcal{V}_{\text{ext}}}^2 + D_2 \|(T^{\text{CAS}} - T_{\text{FCI}}^{\text{CAS}}) \phi_0\|_{H^1}^2 + D_3 \sum_{|\mu|=1} \varepsilon_{\mu} (\tilde{t}_*)_{\mu}^2$. \square

For the remaining error $\delta\varepsilon$ we use techniques that have been developed by Bangerth and Rannacher for a general functional analytic framework [3]. Hence, under the assumption that f is locally strongly monotone the following analysis holds also in the $K \rightarrow \infty$ limit. Nevertheless, before passing on to the error estimate of $\delta\varepsilon$ we characterize the approximation space $\mathcal{V}_{\text{ext}}^{(d)}$. Let $\{b_1, \dots, b_D\}$ be a basis of \mathcal{V}_{ext} , and without loss of generality, $\{b_1, \dots, b_d\}$ be the corresponding subbasis of $\mathcal{V}_{\text{ext}}^{(d)}$ with $d < D$. A key aspect for the analysis is $\mathcal{V}_{\text{ext}}^{(d)}$ being a sufficiently good approximation of \mathcal{V}_{ext} . Subsequently, we elaborate a sufficient condition for this to hold. Let $\delta > 0$ be chosen according to Assumption (B) such that Theorem 19 and 22 imply f being strongly monotone and Lipschitz continuous on $B_{\delta}(t_*)$ with constants γ and L . Further, we define

$$\kappa_d = d(t_*, \mathcal{V}_{\text{ext}}^{(d)}) = \min_{t_d \in \mathcal{V}_{\text{ext}}^{(d)}} \|t_d - t_*\|_{\mathcal{V}_{\text{ext}}}.$$

Eq. (32) in Theorem 25 yields the assumption $\kappa_d \leq \gamma\delta/(\gamma + L)$. Then, the truncated cluster equation $f|_{\mathcal{V}_{\text{ext}}^{(d)}} = 0$ has a locally unique solution on $\mathcal{V}_{\text{ext}}^{(d)} \cap B_{\delta}(t_*)$. We adapt the proof of Theorem 4.1 in [37], which rests on the following consequence of Brouwer's fixed point theorem [8]:

THEOREM 30 (Brouwer, 1965). *Equip \mathbb{R}^d with any norm $\|\cdot\|_d$. Let B_R be the closed ball of radius R centered at $x = 0$ and $h : B_R \rightarrow \mathbb{R}^d$ be continuous. If $\langle h(x), x \rangle \geq 0$ on ∂B_R then $h(x) = 0$ for some $x \in B_R$.*

Let $t_{\text{opt}} \in \mathcal{V}_{\text{ext}}^{(d)}$ with $\kappa_d = \|t_{\text{opt}} - t_*\|_{\mathcal{V}_{\text{ext}}}$, we define the continuous function $h_d : \mathbb{R}^d \rightarrow \mathbb{R}^d$; $x \mapsto (y_j)_{j=1}^d$, where $y_j = \langle f(t_{\text{opt}} + v), b_j \rangle$ and $v = \sum_{j=1}^d x_j b_j$. We chose $\|x\|_d = \|v\|_{\mathcal{V}_{\text{ext}}^{(d)}}$ as a norm on \mathbb{R}^d . Then, $h_d(t) = 0$ if and only if $f(t)|_{\mathcal{V}_{\text{ext}}^{(d)}} = 0$. By assumption $\delta - \kappa_d \geq \delta L / (\gamma + L) > 0$ and we set $R = \delta - \kappa_d$. Then $v \in B_R(t_{\text{opt}})$ implies $v \in B_\delta(t_*)$. Assuming further $\|x\|_d = R$, the monotonicity and Lipschitz continuity of f then yield

$$\begin{aligned} \langle h_d(x), x \rangle &= \sum_{j=1}^d \langle f(t_{\text{opt}} + v), b_j \rangle x_j = \langle f(t_{\text{opt}} + v) - f(t_{\text{opt}}), v \rangle + \langle f(t_{\text{opt}}) - f(t_*), v \rangle \\ &\geq \gamma \|v\|_{\mathcal{V}_{\text{ext}}^{(d)}}^2 + L \kappa_d \|v\|_{\mathcal{V}_{\text{ext}}^{(d)}} = R(\gamma R + L \kappa_d). \end{aligned}$$

Since $\gamma R - L \kappa_d = \gamma \delta - \kappa_d(\gamma + L) \geq 0$, we conclude $\langle h_d(x), x \rangle = R(\gamma R - L \kappa_d) \geq 0$. By Theorem 30 this yields $h_d(x_*) = 0$ for some x_* with $\|x_*\|_d \leq R$, which is equivalent to $t_d = t_{\text{opt}} + v_*$ solving the projected problem $f|_{\mathcal{V}_{\text{ext}}^{(d)}} = 0$. The uniqueness follows from Theorem 9 applied to $f|_{\mathcal{V}_{\text{ext}}^{(d)}}$.

In the sequel, we assume that $\mathcal{V}_{\text{ext}}^{(d)}$ is a sufficiently good approximation of \mathcal{V}_{ext} as guaranteed by Eq. (32). We note that the Lagrangian (29) is nonsymmetric, consequently we cannot expect the error to be quadratic with respect to the error of the wavefunction. However, we see that the dual variable z enters in (29). Indeed, in the analysis that will follow, the solution z_* of the dual problem enters the error estimates. In the spirit of [37], we start the estimation of $\delta\varepsilon$ with a lemma that concerns the dual solution.

LEMMA 31. *Let f be strongly monotone on $B_\delta(t_*)$, then there exists a unique dual solution $z_* \in \mathcal{V}_{\text{ext}}$ determined by t_* such that (t_*, z_*) is a stationary point of the Lagrangian $\mathcal{L}(\cdot, \cdot)$, i.e., (t_*, z_*) solves (30). Additionally, there exists a corresponding unique $z_d \in \mathcal{V}_{\text{ext}}^{(d)}$ such that (t_d, z_d) solves the discretized equation (31) and approximates the exact dual solution quasi-optimally in the sense that*

$$(33) \quad \|z_d - z_*\|_{\mathcal{V}_{\text{ext}}} \leq c_1 \Theta_d + c_2 \Theta_d^2,$$

with $\Theta_d = \max \{d(t_*, \mathcal{V}_{\text{ext}}^{(d)}), d(z_*, \mathcal{V}_{\text{ext}}^{(d)})\}$.

Proof. By definition t_* solves the second component of (30). Therefore it remains to show the first equation. To that end we use Lax–Milgram [9], for which we need to establish boundedness and coercivity of $f'(t_*)^\dagger$. First, we note that the boundedness of $f'(t_*)$ was shown in Theorem 19. Secondly, we expand f into a Taylor series at t_* , i.e., $f(t_* + w) - f(t_*) = f'(t_*)w + \mathcal{O}(\|w\|_{\mathcal{V}_{\text{ext}}}^2)$ with $w \in B_\delta(t_*)$. The strong monotonicity estimate then yields $\langle f'(t_*)w, w \rangle \geq \gamma \|w\|_{\mathcal{V}_{\text{ext}}}^2 - \mathcal{O}(\|w\|_{\mathcal{V}_{\text{ext}}}^3)$. For an arbitrary u we choose $c \in \mathbb{R}$ sufficiently large such that $w = u/c \in B_\delta(t_*)$. This implies the coercivity of $f'(t_*)$. Thirdly, we remark that boundedness and coercivity of $f'(t_*)$ are transferred straightforwardly to the adjoint operator $f'(t_*)^\dagger$. We set $a(z_*, u) = \langle f'(t_*)^\dagger z_*, u \rangle$ and apply Lax–Milgram to the equation $a(z_*, u) = \mathcal{E}'(t_*)(u)$ for all $u \in \mathcal{V}_{\text{ext}}$. This yields the existence and uniqueness of $z_* \in \mathcal{V}_{\text{ext}}$.

This argumentation holds whenever f is strongly monotone. Hence, the existence and uniqueness of z_d follows by the assumption that $\mathcal{V}_{\text{ext}}^{(d)}$ is a sufficiently good ap-

proximation to \mathcal{V}_{ext} . To show Eq. (33) we decompose $z_d - z_* = z_d - \tilde{z}_d + \tilde{z}_d - z_*$, where $\tilde{z}_d \in \mathcal{V}_{\text{ext}}^{(d)}$ solves

$$(34) \quad (\mathcal{E}'(t_*))(u_d) = \langle f'(t_*)u_d, \tilde{z}_d \rangle, \quad \forall u_d \in \mathcal{V}_{\text{ext}}^{(d)}.$$

Note that this is not the discrete problem since it uses the solution t_* instead of t_d . In the same manner as we previously defined $a(\cdot, \cdot)$ we define a bilinear form from (34). Because $f'(t_*)$ is a bounded and coercive linear map, C ea's lemma [48] implies the quasi optimal approximation by \tilde{z}_d to z_* , i.e., $\|\tilde{z}_d - z_*\|_{\mathcal{V}_{\text{ext}}} \leq C d(z_*, \mathcal{V}_{\text{ext}}^{(d)})$.

To estimate $\|z_d - \tilde{z}_d\|_{\mathcal{V}_{\text{ext}}}$ we use the coercivity of $f'(t_d)$. From Eqs. (34) and (31) we deduce

$$\begin{aligned} \gamma \|z_d - \tilde{z}_d\|_{\mathcal{V}_{\text{ext}}}^2 &\leq \langle f'(t_d)(z_d - \tilde{z}_d), z_d - \tilde{z}_d \rangle \\ &= (\mathcal{E}'(t_d) - \mathcal{E}'(t_*))(z_d - \tilde{z}_d) + \langle (f'(t_*) - f'(t_d))(z_d - \tilde{z}_d), \tilde{z}_d \rangle \\ &\leq L_{\mathcal{E}'} \|t_d - t_*\|_{\mathcal{V}_{\text{ext}}} \|z_d - \tilde{z}_d\|_{\mathcal{V}_{\text{ext}}} + L_{f'} \|t_d - t_*\|_{\mathcal{V}_{\text{ext}}} \|z_d - \tilde{z}_d\|_{\mathcal{V}_{\text{ext}}} \|\tilde{z}_d\|_{\mathcal{V}_{\text{ext}}} \\ &= (L_{\mathcal{E}'} + L_{f'} \|\tilde{z}_d\|_{\mathcal{V}_{\text{ext}}}) \|t_d - t_*\|_{\mathcal{V}_{\text{ext}}} \|z_d - \tilde{z}_d\|_{\mathcal{V}_{\text{ext}}}. \end{aligned}$$

Using the quasi optimality of $\|\tilde{z}_d - z_*\|_{\mathcal{V}_{\text{ext}}}$ we find that $\|\tilde{z}_d\|_{\mathcal{V}_{\text{ext}}}$ is bounded by $\|z_*\|_{\mathcal{V}_{\text{ext}}} + Cd(z_*, \mathcal{V}_{\text{ext}}^{(d)})$ and therefore

$$\begin{aligned} \|z_d - \tilde{z}_d\|_{\mathcal{V}_{\text{ext}}} &\leq \frac{1}{\gamma} \left[L_{\mathcal{E}'} + L_{f'} (\|z_*\|_{\mathcal{V}_{\text{ext}}} + C d(z_*, \mathcal{V}_{\text{ext}}^{(d)})) \right] \|t_d - t_*\|_{\mathcal{V}_{\text{ext}}} \\ &\lesssim c_1 d(t_*, \mathcal{V}_{\text{ext}}^{(d)}) + c_2 d(t_*, \mathcal{V}_{\text{ext}}^{(d)}) d(z_*, \mathcal{V}_{\text{ext}}^{(d)}). \quad \square \end{aligned}$$

In order to estimate the error $\delta\varepsilon = |\mathcal{E}(t_*) - \mathcal{E}(t_d)|$ we define the primal residual $\rho(t_d)(\cdot) : \mathcal{V}_{\text{ext}}^{(d)} \rightarrow \mathbb{R}; u \mapsto -\langle f(t_d), u \rangle$ and the dual residual $\rho^*(t_d, z_d)(\cdot) : \mathcal{V}_{\text{ext}}^{(d)} \rightarrow \mathbb{R}; u \mapsto \mathcal{E}'(t_d)(u) - \langle Df(t_d)(u), z_d \rangle$. The following error characterization is based on the results of Bangerth and Rannacher [3] formulated in a suitable way for this article.

THEOREM 32 (Bangerth–Rannacher, 2003). *For any solution of Eqs. (27) and (28), we have the error representation*

$$(35) \quad 2(\mathcal{E}(t_*) - \mathcal{E}(t_d)) = \mathcal{R}_d^{(3)} + \rho(t_d)(z_* - v_d) + \rho^*(t_d, z_d)(t_* - w_d),$$

with arbitrary $v_d, w_d \in \mathcal{V}_{\text{ext}}^{(d)}$. The remainder term $\mathcal{R}_d^{(3)}$ is cubic in the primal and dual error $e = t_* - t_d$ and $e^* = z_* - z_d$,

$$\begin{aligned} \mathcal{R}_d^{(3)} &= \int_0^1 \left(\mathcal{E}^{(3)}(t_d + se)(e, e, e) - \langle f^{(3)}(t_d + se)(e, e, e), z_d + se^* \rangle \right. \\ &\quad \left. - 3\langle f^{(2)}(t_d + se)(e, e), e^* \rangle \right) s(s-1) ds. \end{aligned}$$

Similarly to the approach in [37] we are able to conclude with the following error estimates for the TCC energy.

THEOREM 33. *Let $\mathcal{V}_{\text{ext}}^{(d)}$ be a sufficiently large subspace of \mathcal{V}_{ext} in the sense that $\Theta_d < c$ (see Lemma 31) for a suitable $c \in (0, 1)$, and denote by (t_*, z_*) and (t_d, z_d) the solutions of Eqs. (30) and (31). If f is strongly monotone at t_* , we have*

$$(36) \quad \delta\varepsilon \leq \|t_d - t_*\|_{\mathcal{V}_{\text{ext}}} (c_1 \|t_d - t_*\|_{\mathcal{V}_{\text{ext}}} + c_2 \|z_d - z_*\|_{\mathcal{V}_{\text{ext}}}),$$

and further

$$(37a) \quad \delta\varepsilon \lesssim \left(d(t_*, \mathcal{V}_{\text{ext}}^{(d)}) + d(z_*, \mathcal{V}_{\text{ext}}^{(d)}) \right)^2,$$

$$(37b) \quad \delta\varepsilon \lesssim \|(e^{T_d} - e^{T_*})\phi_{\text{CAS}}\|_{H^1} (\|(e^{T_d} - e^{T_*})\phi_{\text{CAS}}\|_{H^1} + \|(e^{Z_d} - e^{Z_*})\phi_{\text{CAS}}\|_{H^1}),$$

$$(37c) \quad \delta\varepsilon \lesssim \left(\inf_{\psi \in \mathcal{H}_{\text{ext}}} \|\psi - e^{T_*}\phi_{\text{CAS}}\|_{H^1}^2 + \inf_{\psi \in \mathcal{H}_{\text{ext}}} \|\psi - e^{Z_*}\phi_{\text{CAS}}\|_{H^1}^2 \right)^2.$$

Proof. Using Eq. (30) we can rewrite the dual residual as follows:

$$\begin{aligned} \rho^*(t_d, z_d)(s) &= (\mathcal{E}'(t_d))(s) - \langle f'(t_d)(s), z_d \rangle \\ &= (\mathcal{E}'(t_d) - \mathcal{E}'(t_*))(s) + \langle (f'(t_*) - f'(t_d))(s), z_* \rangle + \langle f'(t_d)(s), z_* - z_d \rangle, \end{aligned}$$

for an arbitrary $s \in \mathcal{V}_{\text{ext}}$. Using Eq. (35) in Theorem 32 we obtain

$$\begin{aligned} 2\delta\varepsilon &\leq |\mathcal{R}_d^{(3)}| + |\langle f(t_d) - f(t_*), z_* - v_d \rangle| + |(\mathcal{E}'(t_d) - \mathcal{E}'(t_*))(t_* - w_d)| \\ &\quad + |\langle (f'(t_*) - f'(t_d))(t_* - w_d), z_* \rangle| + |\langle f'(t_d)(t_* - w_d), z_* - z_d \rangle|. \end{aligned}$$

Exploiting the different Lipschitz continuities further implies

$$(38) \quad \begin{aligned} 2\delta\varepsilon &\leq |\mathcal{R}_d^{(3)}| + L_f \|t_d - t_*\|_{\mathcal{V}_{\text{ext}}} \|z_* - v_d\|_{\mathcal{V}_{\text{ext}}} + L_{\mathcal{E}'} \|t_d - t_*\|_{\mathcal{V}_{\text{ext}}} \|t_* - w_d\|_{\mathcal{V}_{\text{ext}}} \\ &\quad + L_{f'} \|t_d - t_*\|_{\mathcal{V}_{\text{ext}}} \|t_* - w_d\|_{\mathcal{V}_{\text{ext}}} \|z_*\|_{\mathcal{V}_{\text{ext}}} \\ &\quad + C \|t_* - w_d\|_{\mathcal{V}_{\text{ext}}} \|z_* - z_d\|_{\mathcal{V}_{\text{ext}}}. \end{aligned}$$

This yields $2\delta\varepsilon \leq \|t_d - t_*\|_{\mathcal{V}_{\text{ext}}} (c_1 \|t_* - t_d\|_{\mathcal{V}_{\text{ext}}} + c_2 \|z_* - z_d\|_{\mathcal{V}_{\text{ext}}}) + |\mathcal{R}_d^{(3)}|$ for $w_d = t_d$ and $v_d = z_d$. By straightforward computations we estimate

$$|\mathcal{R}_d^{(3)}| \leq L_{\mathcal{E}^{(3)}} \|t_* - t_d\|_{\mathcal{V}_{\text{ext}}}^3 + \zeta L_{f^{(3)}} \|t_* - t_d\|_{\mathcal{V}_{\text{ext}}}^3 + 3L_{f^{(2)}} \|t_* - t_d\|_{\mathcal{V}_{\text{ext}}} \|z_* - z_d\|_{\mathcal{V}_{\text{ext}}},$$

with $\zeta = \max_{s \in [0,1]} \|z_d + se^*\|_{\mathcal{V}_{\text{ext}}}$. Hence, by Lemma 31, $|\mathcal{R}_d^{(3)}| \in \mathcal{O}(\Theta_d^3)$, i.e., we can control the remainder term $\mathcal{R}_d^{(3)}$ by means of Θ_d^3 . Since by assumption $\mathcal{V}_{\text{ext}}^{(d)}$ is a sufficiently large subspace of \mathcal{V}_{ext} in the sense that $\Theta_d < c$, this shows Eq. (36).

The bound in Eq. (37a) follows from inserting the optimal approximations t_{opt} , $z_{\text{opt}} \in \mathcal{V}_{\text{ext}}^{(d)}$ in (38) and applying Theorem 9, Lemma 31 and the fact that $\Theta_d < 1$. Then $\|z_d - z_*\|_{\mathcal{V}_{\text{ext}}} \lesssim \Theta_d$ as the term in $\mathcal{O}(\Theta_d^2)$ becomes negligible. The inequalities (37b) and (37c) follow from Proposition 13. \square

We remark that this error estimate derivation does not require the uniqueness of the solution. In cases with not unique solutions, the *a priori* assumption $t_d \rightarrow t_*$ makes the result meaningful as then the remainder term can be assumed to be small.

We conclude this section by combining previous results to prove Theorem 25, the main result of Subsection 4.2.

Proof of Theorem 25. From Eq. (25), we recall that $\delta E \leq \delta\varepsilon + \delta\varepsilon_{\text{CAS}} + \delta\varepsilon_{\text{CAS}}^*$. Then using Lemma 28, Lemma 29 and Eq. (36) in Theorem 33, the desired result now follows. \square

5. Concluding Remarks and Outlook. In this article, we presented a first analysis of the TCC method, proving locally unique and quasi-optimal solutions in Theorems 19 and 22, and a direct error estimate given by Theorem 25. The conceptual change from the HOMO-LUMO gap to the CAS-ext gap ε_0 is a key aspect

of this article. The definition of ε_0 is tailored for existence, uniqueness and error estimate results that are widely applicable, in particular also to excited state approximations. For merely ground-state studies, better bounds in Section 4.1 are obtainable by changing the considered CAS-ext gap to $\tilde{\varepsilon}_0$. The extended CAS-ext gap $\tilde{\varepsilon}_0$ is larger, and in general increases with the size of the CAS. Since the gap assumption enters directly in the norm estimates, a connection between the constants involved in the norm estimates in Section 4.1 and the size of \mathcal{B}_{CAS} seems likely but remains to be proven. Given the presented analysis, it appears reasonable to assume that the results can be generalized to the continuous formulation of the Schrödinger equation, without reference to a finite-dimensional single-particle basis for external space. This corresponds to $K \rightarrow \infty$, in which case many of the concepts used in Section 4.1 may be generalized. Apparently, the main problem in this generalization is that several properties of the Fock operator do not hold for $K \rightarrow \infty$. In particular, its spectrum is not purely discrete. Even though it may appear that we use the Fock operator and its properties excessively, the presented analysis may also be performed for a different one-particle operator, i.e., not necessarily the Fock operator. Section 4.2 is based on achievements for general variational problems, implying the validity of Theorem 25 for infinite dimensions. The currently most important application of our analysis is the DMRG-TCCSD method [42, 43, 44, 2]. In a recent publication, we investigated its numerical performance in light of the results in this article [10]. Using tensor factorization methods—to obtain a well-chosen basis splitting and an approximation to the FCI solution on \mathcal{H}_{CAS} —simplifies the error estimate in Theorem 25 since the methodological error becomes negligible and δE_{CAS} is quadratically bound. This yields a Galerkin-typical quadratic error estimate for the DMRG-TCC method.

6. Acknowledgements. We would like to thank Rolf Heilemann Myhre, Jiří Pittner, Mihály András Csirik and Christian Schilling for valuable discussions and input to this project.

REFERENCES

- [1] L. ADAMOWICZ, P. PIECUCH, AND K. B. GHOSE, *The state-selective coupled cluster method for quasi-degenerate electronic states*, Mol. Phys., 94 (1998), pp. 225–234, <https://doi.org/10.1080/002689798168510>.
- [2] A. ANTALÍK, L. VEIS, J. BRABEC, Ö. LEGEZA, AND J. PITTNER, *Towards the efficient local tailored coupled cluster approximation and the peculiar case of oxo-mn (salen)*, arXiv preprint arXiv:1905.06833, (2019), <https://arxiv.org/abs/1905.06833>.
- [3] W. BANGERTH AND R. RANNACHER, *Adaptive Finite Element Methods for Differential Equations*, Birkhäuser, 2013, <https://doi.org/10.1007/978-3-0348-7605-6>.
- [4] R. J. BARTLETT AND M. MUSIAL, *Coupled-cluster theory in quantum chemistry*, Rev. Mod. Phys., 79 (2007), pp. 291–352, <https://doi.org/10.1103/RevModPhys.79.291>.
- [5] R. J. BARTLETT, J. WATTS, S. KUCHARSKI, AND J. NOGA, *Non-iterative fifth-order triple and quadruple excitation energy corrections in correlated methods*, Chem. Phys. Lett., 165 (1990), pp. 513–522, [https://doi.org/10.1016/0009-2614\(90\)87031-L](https://doi.org/10.1016/0009-2614(90)87031-L).
- [6] M. BORN AND R. OPPENHEIMER, *Zur quantentheorie der molekeln*, Ann. Phys., 389 (1927), pp. 457–484, <https://doi.org/10.1002/andp.19273892002>.
- [7] G. K.-L. CHAN AND M. HEAD-GORDON, *Highly correlated calculations with a polynomial cost algorithm: A study of the density matrix renormalization group*, J. Chem. Phys., 116 (2002), pp. 4462–4476, <https://doi.org/10.1063/1.1449459>.
- [8] E. EMMRICH, *Gewöhnliche und Operator-Differentialgleichungen: Eine Integrierte Einführung in Randwertprobleme und Evolutionsgleichungen für Studierende*, Springer-Verlag, 2013, <https://doi.org/10.1007/978-3-322-80240-8>.
- [9] L. C. EVANS, *Partial Differential Equations*, American Mathematical Society, 2010, <https://doi.org/10.1090/gsm/019>.
- [10] F. M. FAULSTICH, M. MÁTÉ, A. LAESTADIUS, M. A. CSIRIK, L. VEIS, A. ANTALIK, J. BRABEC,

- R. SCHNEIDER, J. PITTNER, S. KVAAL, AND Ö. LEGEZA, *Numerical and theoretical aspects of the dmrg-tcc method exemplified by the nitrogen dimer*, J. Chem. Theory Comput., (2019), <https://doi.org/10.1021/acs.jctc.8b00960>.
- [11] J. GARCKE AND M. GRIEBEL, *On the computation of the eigenproblems of hydrogen and helium in strong magnetic and electric fields with the sparse grid combination technique*, J. Comput. Phys., 165 (2000), pp. 694–716, <https://doi.org/10.1006/jcph.2000.6627>.
- [12] S. J. GUSTAFSON AND I. M. SIGAL, *Mathematical concepts of quantum mechanics*, Springer Science & Business Media, 2011, <https://doi.org/10.1007/978-3-642-21866-8>.
- [13] T. HELGAKER, P. JORGENSEN, AND J. OLSEN, *Molecular Electronic-Structure Theory*, John Wiley & Sons, 2014, <https://doi.org/10.1002/9781119019572>.
- [14] M. HJORTH-JENSEN, M. P. LOMBARDO, AND U. VAN KOLCK, *An advanced course in computational nuclear physics*, in Lecture Notes in Physics, Berlin Springer Verlag, vol. 936, Springer, 2017, <https://doi.org/10.1007/978-3-319-53336-0>.
- [15] J. HUBBARD, *The description of collective motions in terms of many-body perturbation theory*, Proc. R. Soc. Lond. A, 240 (1957), pp. 539–560, <https://doi.org/10.1098/rspa.1957.0106>.
- [16] N. HUGENHOLTZ, *Perturbation approach to the fermi gas model of heavy nuclei*, Physica, 23 (1957), pp. 533–545, [https://doi.org/10.1016/S0031-8914\(57\)93009-4](https://doi.org/10.1016/S0031-8914(57)93009-4).
- [17] T. KINOSHITA, O. HINO, AND R. J. BARTLETT, *Coupled-cluster method tailored by configuration interaction*, J. Chem. Phys., 123 (2005), p. 074106, <https://doi.org/10.1063/1.2000251>.
- [18] A. KÖHN, M. HANAUER, L. A. MUECK, T.-C. JAGAU, AND J. GAUSS, *State-specific multireference coupled-cluster theory*, Wiley Interdiscip. Rev.: Comput. Mol. Sci., 3 (2013), pp. 176–197, <https://doi.org/10.1002/wcms.1120>.
- [19] K. KOWALSKI, *Properties of coupled-cluster equations originating in excitation sub-algebras*, J. Chem. Phys., 148 (2018), p. 094104, <https://doi.org/10.1063/1.5010693>.
- [20] A. LAESTADIUS AND F. M. FAULSTICH, *The coupled-cluster formalism—a mathematical perspective*, Mol. Phys., (2019), pp. 1–12, <https://doi.org/10.1080/00268976.2018.1564848>.
- [21] A. LAESTADIUS AND S. KVAAL, *Analysis of the extended coupled-cluster method in quantum chemistry*, SIAM J. on Numer. Anal., 56 (2018), pp. 660–683, <https://doi.org/10.1137/17M1116611>.
- [22] J. LANG, A. ANTALÍK, L. VEIS, J. BRABEC, Ö. LEGEZA, AND J. PITTNER, *Towards the linear scaling in dmrg-based tailored coupled clusters: An implementation of dlpmo-tccsd*, arXiv preprint arXiv:1907.13466, (2019), <https://arxiv.org/abs/1907.13466>.
- [23] T. J. LEE AND G. E. SCUSERIA, *Achieving chemical accuracy with coupled cluster methods*, in Quantum Mechanical Electronic Structure Calculations with Chemical Accuracy, S. R. Langhof, ed., Springer, Dordrecht, 1995, ch. 2, pp. 47–108, https://doi.org/10.1007/978-94-011-0193-6_2.
- [24] E. H. LIEB AND B. SIMON, *The hartree-fock theory for coulomb systems*, Commun. Math. Phys., 53 (1977), pp. 185–194, <https://doi.org/10.1007/BF01609845>.
- [25] P.-L. LIONS, *Solutions of hartree-fock equations for coulomb systems*, Commun. Math. Phys., 109 (1987), pp. 33–97, <https://doi.org/10.1007/BF01205672>.
- [26] D. I. LYAKH, M. MUSIAL, V. F. LOTRICH, AND R. J. BARTLETT, *Multireference nature of chemistry: The coupled-cluster view*, Chem. Rev., 112 (2012), pp. 182–243, <https://doi.org/10.1021/cr2001417>.
- [27] R. M. MARTIN, *Electronic structure: basic theory and practical methods*, Cambridge university press, 2004, <https://doi.org/10.1017/CBO9780511805769>.
- [28] H. J. MONKHORST, *Calculation of properties with the coupled-cluster method*, Int. J. Quantum Chem., 12 (1977), pp. 421–432, <https://doi.org/10.1002/qua.560120850>.
- [29] M. A. NIELSEN, *The fermionic canonical commutation relations and the jordan-wigner transform*, School of Physical Sciences The University of Queensland, (2005).
- [30] P. PIECUCH, *Active-space coupled-cluster methods*, Mol. Phys., 108 (2010), pp. 2987–3015, <https://doi.org/10.1080/00268976.2010.522608>.
- [31] P. PIECUCH AND L. ADAMOWICZ, *State-selective multireference coupled-cluster theory employing the single-reference formalism: Implementation and application to the h8 model system*, J. Chem. Phys., 100 (1994), pp. 5792–5809, <https://doi.org/10.1063/1.467143>.
- [32] P. PIECUCH, N. OLIPHANT, AND L. ADAMOWICZ, *A state-selective multireference coupled-cluster theory employing the single-reference formalism*, J. Chem. Phys., 99 (1993), pp. 1875–1900, <https://doi.org/10.1063/1.466179>.
- [33] K. RAGHAVACHARI, G. W. TRUCKS, J. A. POPLE, AND M. HEAD-GORDON, *A fifth-order perturbation comparison of electron correlation theories*, Chem. Phys. Lett., 157 (1989), pp. 479–483, [https://doi.org/10.1016/S0009-2614\(89\)87395-6](https://doi.org/10.1016/S0009-2614(89)87395-6).
- [34] M. RENARDY AND R. C. ROGERS, *An Introduction to Partial Differential Equations*, vol. 13, Springer Science & Business Media, 2006, <https://doi.org/10.1007/b97427>.

- [35] J. RISSLER, R. M. NOACK, AND S. R. WHITE, *Measuring orbital interaction using quantum information theory*, Chem. Phys., 323 (2006), pp. 519–531, <https://doi.org/10.1016/j.chemphys.2005.10.018>.
- [36] T. ROHWEDDER, *The continuous coupled cluster formulation for the electronic schrödinger equation*, ESAIM-Math. Model. Num., 47 (2013), pp. 421–447, <https://doi.org/10.1051/m2an/2012035>.
- [37] T. ROHWEDDER AND R. SCHNEIDER, *Error estimates for the coupled cluster method*, ESAIM-Math. Model. Num., 47 (2013), pp. 1553–1582, <https://doi.org/10.1051/m2an/2013075>.
- [38] R. SCHNEIDER, *Analysis of the projected coupled cluster method in electronic structure calculation*, Numer. Math., 113 (2009), pp. 433–471, <https://doi.org/10.1007/s00211-009-0237-3>.
- [39] I. SHAVITT AND R. J. BARTLETT, *Many-body methods in chemistry and physics: MBPT and Coupled-Cluster Theory*, Cambridge, 2009, <https://doi.org/10.1017/CBO9780511596834>.
- [40] J. C. SLATER AND A. RUSSEK, *Quantum theory of molecules and solids, vol. 1: Electronic structure of molecules*, Am. J. Phys, 32 (1964), pp. 65–66, <https://doi.org/10.1119/1.1970097>.
- [41] S. SZALAY, G. BARCZA, T. SZILVÁSI, L. VEIS, AND Ö. LEGEZA, *The correlation theory of the chemical bond*, Sci. Rep., 7 (2017), <https://doi.org/10.1038/s41598-017-02447-z>.
- [42] L. VEIS, A. ANTALÍK, J. BRABEC, F. NEESE, Ö. LEGEZA, AND J. PITTNER, *Coupled cluster method with single and double excitations tailored by matrix product state wave functions*, J. Phys. Chem. Lett., 7 (2016), pp. 4072–4078, <https://doi.org/10.1021/acs.jpcllett.6b01908>.
- [43] L. VEIS, A. ANTALÍK, J. BRABEC, F. NEESE, Ö. LEGEZA, AND J. PITTNER, *Correction to coupled cluster method with single and double excitations tailored by matrix product state wave functions*, J. Phys. Chem. Lett., 8 (2017), pp. 291–291, <https://doi.org/10.1021/acs.jpcllett.6b02912>.
- [44] L. VEIS, A. ANTALÍK, Ö. LEGEZA, A. ALAVI, AND J. PITTNER, *The intricate case of tetramethyleneethane: A full configuration interaction quantum monte carlo benchmark and multireference coupled cluster studies*, J. Chem. Theory Comput., 14 (2018), pp. 2439–2445, <https://doi.org/10.1021/acs.jctc.8b00022>.
- [45] S. R. WHITE AND R. L. MARTIN, *Ab initio quantum chemistry using the density matrix renormalization group*, J. Chem. Phys., 110 (1999), pp. 4127–4130, <https://doi.org/10.1063/1.478295>.
- [46] J. WLOKA, *Partial differential equations*, Cambridge University, (1987), <https://doi.org/10.1002/zamm.19880680621>.
- [47] H. YSERENTANT, *Regularity and Approximability of Electronic Wave Functions*, Springer, Berlin, Heidelberg, 2010, <https://doi.org/10.1007/978-3-642-12248-4>.
- [48] E. ZEIDLER, *Nonlinear functional analysis and its applications, vol. ii/b*, 1990, <https://doi.org/10.1007/978-1-4612-0981-2>.