

**Quantum Monte Carlo and Molecular Fragmentation Methods for the Treatment of
Electron Correlation in Molecules and Solids**

by

Amanda Dumi

Master of Chemistry, Duquesne University, 2017

Submitted to the Graduate Faculty of
the Dietrich School of Arts and Sciences in partial fulfillment
of the requirements for the degree of

Doctor of Philosophy

University of Pittsburgh

2022

UNIVERSITY OF PITTSBURGH
DEPARTMENT OF CHEMISTRY

This dissertation was presented

by

Amanda Dumi

It was defended on

December 4, 2022

and approved by

Kenneth D. Jordan, Department of Chemistry

Dr. Daniel Lambrecht, Department of Chemistry & Physics, Florida Gulf Coast University

Dr. Jennifer Laaser, Department of Chemistry

Dr. Geoffrey Hutchison, Department of Chemistry

Copyright © by Amanda Dumi
2022

Quantum Monte Carlo and Molecular Fragmentation Methods for the Treatment of Electron Correlation in Molecules and Solids

Amanda Dumi, PhD

University of Pittsburgh, 2022

The quantum mechanical treatment of molecular systems in computational chemistry can offer insight into the nature of chemical bonding, reaction mechanisms, and many experimental observables. An improved treatment of the electronic structure of a molecule is tied to an increase of computational resources which can often become a limiting factor for the size of system that is treatable or the accuracy that can be obtained. To find a balance between these opposing needs, this work explores two approaches which attempt to approximate the electronic structure treatment through using only a subspace of the full system representation.

One approach is through the use of quantum Monte Carlo (QMC) methods, which stochastically solve the Schrödinger equation. Here, diffusion Monte Carlo is used to provide insight into chemical systems that exhibit challenging physics which are unable to be described by the commonly utilized Density Functional Theory. Two systems are explored: a hydrogen atom chemisorbed to the surface of graphene and a model non-valence correlation-bound anion. Additionally, the use of selected CI methods to generate QMC trial wave functions is investigated, as these methods choose the subset of Slater determinants which are most important for a specific system.

The second approach is the development of a fragment selection scheme through unsupervised machine learning approaches. Fragmentation approaches are motivated by the short range nature of correlation effects. The full system is approximated by subsystems which are each treated at a certain level of theory and an estimate of the interactions between them. For this approach to return valuable results, the chemical domains need to capture the most important physics of the desired problem. In this work an unsupervised machine learning based method is developed which will allow the systematic identification of important chemical domains with minimal quantum mechanical data which can improve transferability and automation of fragmentation approaches.

Table of Contents

1.0 Introduction	1
1.1 Project summaries	2
1.1.1 Quantum Monte Carlo Approaches	2
1.1.2 Real Space Partition Methods	4
2.0 Theoretical Background	6
2.1 Mean Field Methods	9
2.1.1 Hartree-Fock	9
2.1.2 Density Functional Theory	11
2.1.2.1 Exchange-Correlation Functionals	12
2.2 Wave Function Methods	14
2.2.1 Configuration Interaction	14
2.2.1.1 Truncation of the FCI Space	15
2.3 Coupled-Cluster	17
2.3.1 Equation-of-Motion	18
2.4 Stochastic methods	20
2.4.1 Considerations for QMC	21
2.4.2 Variational Monte Carlo	21
2.4.3 Projector Monte Carlo	23
2.5 Fragmentation Methods	25
2.5.1 Fragment treatments	26
3.0 The binding of atomic hydrogen on graphene from density functional theory and diffusion Monte Carlo calculations	28
3.1 Summary	28
3.2 Introduction	28
3.3 Methods	30
3.3.1 Density Functional Theory Calculations	30

3.3.2	Quantum Monte Carlo Calculations	31
3.4	Results & Discussion	33
3.4.1	Binding Energy	33
3.4.2	Binding Density	35
3.4.3	Charge Density Differences Between DMC and PBE	37
3.5	Conclusions	37
3.6	Acknowledgements	39
3.7	Data Availability Statement	39
4.0	The Role of High-Order Electron Correlation Effects in a Model System for Non-valence Correlation-bound Anions	40
4.1	Summary	40
4.2	Introduction	41
4.3	Computational Details	43
4.3.1	EOM Coupled Cluster	43
4.3.2	DMC	45
4.3.3	Restricted CI and CIPSI-generated Trial Wave Functions for DMC Calculations	46
4.3.4	AFQMC	48
4.3.5	Radial Orbital Densities	48
4.4	Results and Discussion	49
4.4.1	Results for $R = 4 \text{ \AA}$: the Correlation Bound Region	49
4.4.2	Results for $R = 7 \text{ \AA}$: the Electrostatically Bound Region	55
4.5	Conclusion	56
5.0	Characterizing unsupervised learning methods ability to determine important chemical domains.	59
5.1	Summary	59
5.2	Introduction	60
5.3	Methods	62
5.3.1	Clustering approaches	63
5.3.2	Molecular Representations	66

5.3.3 Validation	68
5.3.4 Systems	69
5.4 Results and Discussion	71
5.4.1 Water Clusters	71
5.4.2 Methylthiophenes	73
5.4.3 Silyl Ketene	74
5.5 Conclusion	76
6.0 Conclusions	78
A.1 Additional works	78
A.1.1 Q-GPU	78
A.1.2 AFMC applications Towards Regional Embedding.	79
A.1.3 Characterization of silyl ketenes	80
A.2 Supplemental Material for Chapter 3	80
A.3 Supplemental Material for Chapter 4	85
A.4 Supplemental Material for Chapter 5	104
A.5 Supplemental Material for Chapter 5	111
Appendix. Bibliography	118

List of Tables

Table 1:	Binding energy (meV) of a hydrogen atom chemisorbed on graphene calculated with various DFT functionals and with DMC.	34
Table 2:	Dependence of the total energies and the EBE of the model $(\text{H}_2\text{O})_4$ cluster at $R = 4 \text{ \AA}$ on the supplemental diffuse basis functions. Results obtained using the EOM-CCSD method.	44
Table 3:	Sensitivity of the EBE of the $(\text{H}_2\text{O})_4$ model to the “core” basis set. Results obtained using the EOM-MP2 method.	45
Table 4:	EBCs of the $(\text{H}_2\text{O})_4$ model calculated using HF, EOM, and AFQMC methods and employing the aug-cc-pVTZ+7s7p basis set.	51
Table 5:	EBCs of the $(\text{H}_2\text{O})_4$ model calculated using the DMC method and various trial wave functions ¹	52
Table 6:	PBE total energies of each system and binding energies (eV) with respect to increasing the Monkhorst-Pack k -point grid. The hydrogen total energy value used to calculate the binding energy was for a 1x1x1 k -grid with a 250 Ry wave function cutoff.	81
Table 7:	PBE total and binding energies (eV) with respect to varying the kinetic energy cutoff values of the wave function on a 4x4x1 k -point grid. The hydrogen total energy value used to calculate the binding energy was for a 1x1x1 k -grid with a 250 Ry wave function cutoff.	81
Table 8:	PBE0 and HSE total energies and binding energies (eV) with respect to increasing the Monkhorst-Pack k -point grid.	82
Table 9:	Non-self-consistent field calculations: system and binding energies used to calculate the differences in Table II of the manuscript, reported without basis set superposition error correction.	82
Table 10:	Total and binding energies (eV) of a hydrogen atom on graphene calculated with QMC.	83

Table 11: Total energies (eV) at each twist angle used for graphene DMC treatment. . .	83
Table 12: Total energies (eV) at each twist angle used for hydrogen atom chemisorbed on graphene DMC treatment.	84
Table 13: AFQMC total energies at used to calculate the electron binding energies. AFQMC calculation details can be found in the manuscript.	101

List of Figures

Figure 1: Schematic of a random walk	22
Figure 2: Model of hydrogen chemisorbed to graphene	32
Figure 3: electron density change due to the adsorption of the H atom to the distorted graphene sheet	36
Figure 4: electron density difference of H on graphene: DMC-PBE	38
Figure 5: The model $(\text{H}_2\text{O})_4$ system	43
Figure 6: Radially integrated charge densities for R=4 in $(\text{H}_2\text{O})_4$	54
Figure 7: Radially integrated charge densities for R=7 in $(\text{H}_2\text{O})_4$	57
Figure 8: Visual representations of the three clustering approach explored in this work .	64
Figure 9: The methylthiophene monomer 2-D structure	70
Figure 10: Structures of silyl ketene systems studied	71
Figure 11: Visualization of preservation of covalent bonds in $(\text{H}_2\text{O})_4$	72
Figure 12: fragment size variance for $((\text{H}_2\text{O})_{21})$	73
Figure 13: Performance of clustering for the methylthiophene test set	74
Figure 14: Assessment of fragmentation schemes on silyl ketene trimers	75
Figure 15: Representative visualization of fragmentation for silyl ketene trimer	75
Figure 16: Electron density difference for Hgraphene: DMC-HSE	85
Figure 17: Water cluster test set treated with affinity propagation	105
Figure 18: Water cluster test set treated with mean shift clustering	106
Figure 19: Detailed results for Methylthiophene test set	107
Figure 20: Methylthiophene test set treated with k -means clustering. The structures were optimized at either the Hartree-Fock level of theory (left) or ω -B97X-D level of theory (right).	107

Figure 21: Methylthiophene test set treated with spectral clustering, which worked well for all descriptors. The structures were optimized at either the Hartree-Fock level of theory (left) or ω -B97X-D level of theory (right). All molecular representations are shown and their markings indicated in the legend. Two molecules resulting in different clustering upon the improved level of theory.	108
Figure 22: Methylthiophene test set treated with affinity propagation clustering. The structures optimized were at either the Hartree-Fock level of theory (left) or ω -B97X-D level of theory (right). All molecular representations are shown and their markings indicated in the legend.	109
Figure 23: Methylthiophene test set treated with mean shift clustering. The structures were optimized at either the Hartree-Fock level of theory (left) or ω -B97X-D level of theory (right). The explored representations were not formulated for mean shift, and the clustering results were not viable.	109
Figure 24: Percent error of the energy and the speedup over the supermolecular calculation for the SK dimer to assess the performance of the fragment approaches and descriptor quality.	111
Figure 25: Water cluster test set treated with affinity propagation	112
Figure 26: Water cluster test set treated with mean shift clustering	113
Figure 27: Detailed results for Methylthiophene test set	114
Figure 28: Methylthiophene test set treated with k -means clustering. The structures were optimized at either the Hartree-Fock level of theory (left) or ω -B97X-D level of theory (right).	114
Figure 29: Methylthiophene test set treated with spectral clustering, which worked well for all descriptors. The structures were optimized at either the Hartree-Fock level of theory (left) or ω -B97X-D level of theory (right). All molecular representations are shown and their markings indicated in the legend. Two molecules resulting in different clustering upon the improved level of theory.	115

Figure 30: Methylthiophene test set treated with affinity propagation clustering. The structures optimized were at either the Hartree-Fock level of theory (left) or ω -B97X-D level of theory (right). All molecular representations are shown and their markings indicated in the legend.	116
Figure 31: Methylthiophene test set treated with mean shift clustering. The structures were optimized at either the Hartree-Fock level of theory (left) or ω -B97X-D level of theory (right). The explored representations were not formulated for mean shift, and the clustering results were not viable.	116
Figure 32: Percent error of the energy and the speedup over the supermolecular calculation for the SK dimer to assess the performance of the fragment approaches and descriptor quality.	117

1.0 Introduction

“The general theory of quantum mechanics is now almost complete, The underlying physical laws necessary for the mathematical theory of a large part of physics and the whole of chemistry are thus completely known, and the difficulty is only that the exact application of these laws leads to equations much too complicated to be soluble. It therefore becomes desirable that approximate practical methods of applying quantum mechanics should be developed, which can lead to an explanation of the main features of complex atomic systems without too much computation.”-Dirac¹

Electronic structure methods aim to accurately describe the electrons within a molecule or solid. This becomes computationally expensive when trying to describe a potential between all electrons which are dynamically responding to the changing potential of the others, generating a competing effort between obtaining accurate results and feasible calculations.

One can simply avoid this complexity by treating the system in a mean-field approach such as Hartree-Fock, however electron correlation is unaccounted for. Neglecting electron correlation is not always an appropriate approximation for property prediction and ways of incorporating these effects are necessary. Full Configuration Interaction, which is formally exact in a complete basis, has prohibitive $N!$ scaling. Whereas coupled-cluster singles doubles with perturbative triples methods is often regarded as the “gold standard” level of theory, but incurs an N^7 scaling. The higher scaling restricts the application of these methods to small systems while many of the chemical processes of interest involve large molecules.

On the other hand, one of the most frequently applied methods is Density Functional Theory (DFT), a mean field approach which attempts to approximate the missing correlation energy through an exchange-correlation functional that has favorable scaling with system size. Although DFT has demonstrated the ability to successfully describe many molecular systems, there are many occasions in which this approximation falls short, e.g. systems which exhibit strong correlation, systems in which the physics used to construct the functional may differ from that being treated, and in situations targeting more than just total energies as exchange-correlation functionals are often benchmarked against energies only.

In this work, we probe various approaches to focusing efforts in the electronic structure calculations in an intuitive way to allow an accurate description including correlation energy while minimizing or reducing the computational demand. We approach this from two different direc-

tions. The first is through the application of is a stochastic sampling to solve the Schrödinger equation through quantum Monte Carlo (QMC) approaches. The second is through an attempt to choose molecular partitioning for fragmentation approaches through the utilization of unsupervised machine learning methods. These two projects are introduced briefly in the Section 1.1.1 and Section 1.1.2. In the chapters to follow, the QMC related projects are discussed in Chapter 3 and Chapter 4 and the fragmentation work can be found in Chapter 5.

1.1 Project summaries

1.1.1 Quantum Monte Carlo Approaches

One way to allow for the treatment of the electronic structure for larger systems is to use a stochastic approach to solve the Schrödinger equation. When sampled with a Monte Carlo-type approach, these methods are collectively known as quantum Monte Carlo (QMC) methods, though there are many formulations. We demonstrate the ability of QMC to accurately describe two challenging chemical problems where DFT is unreliable: hydrogen chemisorbed to a graphene surface and a model $(\text{H}_2\text{O})_4$ system.

The first application is describing a hydrogen atom chemisorbed to graphene. The unique electronic, optical, and transport properties of graphene make it an important system for a wide range of applications, many of which involve or are impacted by the adsorption of atoms or molecules. To bring these applications to fruition, a deeper understanding of the interaction of atoms and molecules with graphene is required, and has been the subject of several experimental and theoretical studies.^{2–14} Most computational studies of adsorption of atoms and molecules on graphene have employed DFT, primarily due to its favorable scaling with system size. However, a reliable theoretical description of interactions at the graphene surface has proven to be challenging for DFT.^{2,7,8,15} The diffusion Monte Carlo (DMC)¹⁶ method, which is a real-space stochastic approach to solving the many-body Schrödinger equation is particularly attractive given its low scaling with the number of electrons and high parallelizability. DMC also has the advantages of being systematically improvable and its energy being much less sensitive to the basis set employed than methods

that work in the space of Slater determinants. DMC has been used to describe the adsorption of various species on graphene including O₂¹⁰, a water molecule^{9,17}, and a platinum atom.¹¹ In the present work, we use the DMC method to calculate the binding energy of H to graphene in the chemisorbed state.

The second challenging chemical system we will explore using QMC methods is a model (H₂O)₄ system, which is an example of a non-valence correlation-bound (NVCB) anion. In recent years, there has been growing interest in NVCB anions in which long-range correlation effects are crucial for the binding of the excess electron.^{18–31} By definition, NVCB anions are unbound in the HF approximation. Usual methods for treating these systems are restricted by a rather large scaling with system size, so it is of interest to know if DMC can provide an accurate description. As NVCBs are relevant to image potential states of larger species such as C₆₀ and graphitic systems, a method which enables an accurate ab initio treatment will allow for furthered understandings of these systems. This particular model system is useful as there is no net dipole and thus electrostatics will not play a prominent role in binding the excess electron.^{18,25} Additionally, the coordinates of the water molecules can be tuned to regimes in which the NVCB character dominates to regimes in which the NVCB character is lessened or non-existent. This model system will allow us to also understand the quality of the trial wave function needed to allow for an accurate DMC treatment. In DMC, many of the approximations are systematically improvable. However, the fixed node error discussed in Section 2.4 is one that needs care in ensuring the effects of the approximation are minimal. The nodes of the wave function are set by the trial wave function. Various sources of trial wave functions are explored, from single determinant methods like HF and DFT, to multi-Slater determinant expansions which result in a nodal surface more representative of the ground state. The multi-Slater determinant trial wave functions are sourced from a hand-selected truncation scheme and through an automated approach of selected CI methods such as perturbatively selected configuration interaction scheme (CIPSI)^{32–36}. From a method development standpoint, there is interest in understanding how multideterminant trial wave functions can impact a QMC treatment.^{37–39}

1.1.2 Real Space Partition Methods

To compensate for high computational scaling, fragmentation approaches estimate the energy and other properties of large molecular systems by partitioning the system into small subsystems, where the final estimate of the energy becomes the accumulation of the parts. This can be very successful as electron correlation is often a short range phenomena. The accuracy of the fragmentation approach hinges on the electronic structure treatment of each fragment, the approach used to describe the interaction between the fragments, and the way in which the molecule is partitioned. In the fragmentation schemes, the best case scenario for scaling becomes $\mathcal{O}(N^p) \rightarrow N_{frag} \mathcal{O}(f^p)$, where N_{frag} is the number of fragments, f is representative of the fragment size, and p is the exponential value dependent on the level of electronic structure theory utilized.^{40,41} This partitioning of a single, very costly calculation into N_{frag} smaller calculations achieves two important objectives: 1) Computational scaling with system size is reduced with reasonable fragment definition and 2) trivial parallelization is possible by treating subsystems separately, with the potential to efficiently utilize high performance computing resources.

As pointed out by Herbert, the choice of fragments for a system is not well-defined, but affects the quality of results obtained.⁴⁰ In some systems, a natural approach towards partitioning arises when there is a stark difference in the types of bonding present in the system. However, the choice of fragments is not always as clear-cut. In such cases, fragmentation requires the comparison of total energies, dipoles, or polarizabilities. Often the fragment definition is based on predefined functional groups or chosen manually.⁴² Some methods of energy estimation are defined based on specific fragmentation schemes such as the systematic molecular fragmentation (SMF) and systematic molecular fragmentation by annihilation (SMFA).^{43–47} In these methods, fragments are built around functional groups with larger fragments made from their groupings. Ultimately, the level of fragmentation is at the discretion of the user to achieve the desired level of accuracy. However, functional group definition may become ambiguous. For example, there is no set number of monomers to include from a polymer backbone to acquire an accurate description capture the chemical behavior. Additionally, a fragment definition based only on functional groups may not consider the interacting chemical environment.

If instead a more general approach could be developed which can consider the molecular sys-

tem in an automatic and unbiased fashion, work would be made possible related to various fields such as quantum embedding approaches, Quantum Mechanics-Molecular Mechanics (QM/MM) approaches, quantum computing efforts, and the description of large systems with ab initio methods. A desirable approach to choosing fragments would have low computational cost and prioritize keeping associated molecular components intact to treat fully with quantum mechanics while the estimation of their interactions should occur only at the most weakly bound points. Here we present a proof of concept for an approach utilizing unsupervised machine learning methods to identify the strongly interacting substructures of the system. Clustering methods are a form of unsupervised machine learning used to identify substructures in data sets, as a result these approaches are fundamental to data-mining procedures. Previous work in chemistry utilized clustering methods to identify structure-property relationships in large databases⁴⁸, to determine the number of residues to treat in quantum refinement methods⁴⁹, and to partition large proteins into peptides using an amino acids representation using graph based methods.⁵⁰ The application of clustering methods to produce logical fragments of individual molecules at an atomic level is an unexplored direction.

2.0 Theoretical Background

Electronic Structure methods are theoretical approaches which attempt to describe the behavior of electrons in a chemical system to gain insight into the properties of that system. The electrons are described by a wave function Ψ , which contains all information about the state. We specifically denote the wave function of all electrons as $\Psi(r_1 \dots r_N)$ where N is the number of electrons. When Ψ is acted on by an operator, observables of the system can be obtained. Specifically, let us consider the non-relativistic Schrödinger equation:

$$\hat{H}\Psi(\vec{r}_1 \dots \vec{r}_N) = E\Psi(\vec{r}_1 \dots \vec{r}_N) \quad (1)$$

where \hat{H} is the Hamiltonian, an operator which is applied to the wave function, Ψ , to return the ground state energy of the system, E , and Ψ . The full molecular Hamiltonian and has the following form:

$$\hat{H} = \sum_A^M -\frac{1}{2m_A} \nabla_A^2 + \sum_{i=1}^N -\frac{\hbar^2}{2m_e} \nabla_i^2 + \frac{e^2}{4\pi\epsilon_0} \left(\sum_{i < j} \frac{1}{|\vec{r}_i - \vec{r}_j|} - \sum_{i=1}^N \sum_{A=1}^M \frac{Z_A}{|\vec{r}_i - \vec{R}_A|} + \sum_{A < B} \frac{Z_A Z_B}{|\vec{R}_A - \vec{R}_B|} \right) \quad (2)$$

where the first summation is over M ions, m_A is the mass of the ion, the second summation is over N electrons, \hbar is the reduced Plank's constant, m_e is the mass of an electron, e is the elementary charge, ϵ_0 is the vacuum permittivity, Z is the atomic number of the nuclei, and R_A is the position of the ions. The first is the kinetic energy operator of the ions, the second term is the electronic kinetic energy operator, while the last three terms are potential energy terms resulting from Coulomb interactions between electron-electron, electron-ion, and ion-ion. To simplify this, Hartree atomic units will be used throughout this work which reduces the electron mass (m_e), the expression $4\pi\epsilon_0$, reduced Planck's constant, and the elementary charge to 1.

In practice, finding an exact solution to Equation 2 is unfeasible except for systems containing a few particles. Describing larger systems requires approximations to be introduced. One of the most commonly utilized approximations is the Born-Oppenheimer approximation, which enforces that the electrons are moving within a fixed field of ions. This is a reasonable approximation as

the larger weight of ions relative to electrons will cause the relative speed of the ions to be much slower. This allows a simplification of the Hamiltonian to:

$$\hat{H} = \sum_i \left[-\frac{1}{2} \nabla_i^2 - \sum_{iA} \frac{Z_A}{|\vec{r}_i - \vec{R}_A|} \right] + \sum_{i < j} \frac{1}{|\vec{r}_i - \vec{r}_j|} \quad (3)$$

Separating this expression into the one-electron and two-electron terms leaves us with:

$$\hat{H} = \sum_i \hat{h}(i) + \sum_{i < j} \hat{v}(i, j) \quad (4)$$

where $\hat{h}(i)$ is the one-electron operator which contains the electron kinetic energy and the electron-ion potential, and $\hat{v}(i, j)$ is the two-electron operator which describes the electron-electron repulsion energy.

In addition to the Born-Oppenheimer approximation, it is also useful to approximate the many electron wave function as a product of single particle functions, i.e. potential between electrons, V_{ee} , is ignored. This is known as the Hartree product wave function. This simple product is physically correct for Bosonic systems, but will fail to describe a Fermionic system as the Hartree product will not enforce that the electrons be indistinguishable, i.e. that our wave function exhibits antisymmetry. To enforce antisymmetry of the wave function, Slater determinants are used. The Slater determinant has the form of:

$$\Psi = \frac{1}{\sqrt{N!}} \begin{vmatrix} \chi_1(\mathbf{x}_1) & \chi_2(\mathbf{x}_1) & \cdots & \chi_N(\mathbf{x}_1) \\ \chi_1(\mathbf{x}_2) & \chi_2(\mathbf{x}_2) & \cdots & \chi_N(\mathbf{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \chi_1(\mathbf{x}_N) & \chi_2(\mathbf{x}_N) & \cdots & \chi_N(\mathbf{x}_N) \end{vmatrix}$$

where χ_i is a single particle orbital and \mathbf{x}_i is the position of the i^{th} electron. Here, the single particle spin orbitals can be understood and x denoting a collective variable to describe position and spin. The single particle orbitals in computational chemistry can have many forms. Plane waves are often used for materials calculation, but one of the most frequently used for molecular calculations are molecular orbitals (MOs) that are represented by a linear combination of atomic orbitals centered on atoms:

$$|\chi_i\rangle = \sum_v c_{vi} \phi_v \quad (5)$$

where χ are atomic atomic orbitals and c is a matrix of coefficients where the columns represent the molecular orbitals and the rows are atomic orbitals. Once working within a single particle basis we can express the one- and two-electron terms as integrals over the orbitals. The one-electron term becomes

$$\langle i | \hat{h} | j \rangle = \int d\mathbf{x}_1 \chi_i^*(\mathbf{x}_1) \hat{h} \chi_j(\mathbf{x}_1) \quad (6)$$

The two-electron term can be expressed as

$$\langle ij | \hat{v} | kl \rangle = \int d\mathbf{x}_1 d\mathbf{x}_2 \chi_i^*(\mathbf{x}_1) \chi_j^*(\mathbf{x}_2) \hat{v} \chi_k(\mathbf{x}_1) \chi_l(\mathbf{x}_2) \quad (7)$$

which is known as the physicists' notation, or one can express this in chemists' notation as:

$$[ij|kl] = \int d\mathbf{x}_1 d\mathbf{x}_2 \chi_i^*(\mathbf{x}_1) \chi_j(\mathbf{x}_1) \hat{v} \chi_k(\mathbf{x}_2)^* \chi_l(\mathbf{x}_2) \quad (8)$$

For a wave function to be exact and lead to an exact energy, an infinite basis is required. This is not feasible in computational chemistry and instead the wave function is approximated by a finite basis. The variational principle can be used to inform the quality of the result obtained by the approximate wave function. The variational principle states that an approximate wave function which satisfies the boundary conditions of the problem will give an energy that is greater than the true ground state energy:

$$E_v = \frac{\int dR \Psi^*(R) \hat{H} \Psi(R)}{\int dR \Psi^*(R) \Psi(R)} \geq E_0 \quad (9)$$

Here, $R = (r_1, r_2, \dots, r_N)$ representing the positions of all electrons, E_v is the energy of the approximate wave function, and E_0 is the true ground state energy. Thus if a method is said to be variational then the improvement of an approximate wave function will lead to a minimization of the variational energy.

Currently, the problem is proposed as placing identical particles in orbitals. This idea lends itself well to another expression of these equations known as second quantization. This expression of the integrals will be exceptionally helpful when discussing wave function methods such as Configuration Interaction. The basics of the formalism are given here, but details of can be found in Reference [51]. Second quantization is best understood if one thinks of a Slater determinant as an occupation vector.

$$|n_\alpha\rangle = |n_1, n_2, n_3, \dots, n_\alpha, \dots\rangle \quad (10)$$

where the values of the occupation will be 1 or 0 indicated whether an orbital is occupied or unoccupied. This means the sum over all occupied orbitals in the vector will result in the total number of electrons $\sum_{\alpha} n_{\alpha} = N$. The occupation of specific orbitals within the vector can be indicated through the use of creation a_i^{\dagger} and annihilation operators a_i . Creation operators will place a particle within the i^{th} orbital, while the annihilation operator will remove a particle from the i^{th} orbital. The molecular Hamiltonian (non-relativistic and within the Born-Oppenheimer approximation) in second quantization form is:

$$\begin{aligned}\hat{H} &= \sum_i \hat{h}(i) + \sum_{i < j} \hat{v}(i, j) \\ \hat{H} &= \sum_{pq} h_{pq} a_p^{\dagger} a_q + \sum_{pq, rs} v_{pqrs} a_p^{\dagger} a_q^{\dagger} a_r a_s + V_{NN}\end{aligned}\quad (11)$$

where h_{pq} are the one electron Hamiltonian terms for orbital pair pq and v_{pqrs} are the two-electron terms for $pqrs$ orbitals, and V_{NN} is the potential resulting from the fixed ions.

2.1 Mean Field Methods

2.1.1 Hartree-Fock

The Hartree-Fock (HF) approach is one of the fundamental approximations for quantum chemistry applications. This method provides a description of an electron in the mean field of the other particles. The formulation of the method can be understood as finding the best set of spin orbitals for a single Slater determinant which give the best estimation of our ground state wave function. Following from the variational principle, the best orbitals will be those which result in the lowest ground state energy. The energy expression for a single Slater determinant in the HF formalism is given by:

$$E_{HF} = \sum_i \langle i | h | i \rangle + \frac{1}{2} \sum_{ij} [ii|jj] - [ij|ji] \quad (12)$$

where i and j are the occupied orbitals. These expressions are the one and two electron integrals described in Equation 6 and Equation 8:

$$\langle i | \hat{h} | i \rangle = \int d\mathbf{x}_1 \chi_i^*(\mathbf{x}_1) \hat{h}(1) \chi_i(\mathbf{x}_1) \quad (13)$$

$$[ii|jj] = \langle ii | \hat{v} | jj \rangle = \int \int d\mathbf{x}_1 d\mathbf{x}_2 \chi_i^*(\mathbf{x}_1) \chi_j(\mathbf{x}_1) \frac{1}{r_{12}} \chi_j^*(\mathbf{x}_2) \chi_j(\mathbf{x}_2) \quad (14)$$

$$[ij|ji] = \langle ii | \hat{v} | jj \rangle = \int \int d\mathbf{x}_1 d\mathbf{x}_2 \chi_i^*(\mathbf{x}_1) \chi_j(\mathbf{x}_1) \frac{1}{r_{12}} \chi_j^*(\mathbf{x}_2) \chi_i(\mathbf{x}_2) \quad (15)$$

The HF equations can then be solved by employing the method of Lagrange multipliers with additional constraints to ensure that the MOs remain orthogonal. A Fock operator is defined which will have eigenvalues that are the optimal MOs to minimize E_{HF} .

$$f(\mathbf{x}_1) = h(\mathbf{x}_1) + \sum_{j \neq i} (\mathcal{J}_j(\mathbf{x}_1) - \mathcal{K}_j(\mathbf{x}_1)) \quad (16)$$

where \mathcal{J} is Coulomb interaction of electron in spatial orbital χ_i with the average charge distribution of all other electrons shown in Equation 18. The second term, \mathcal{K} , is the exchange operator which does not have a classical analog, but follows from the antisymmetry requirement of the wave function (Equation 19).

$$\langle i | h(1) | j \rangle = \int d\mathbf{x}_1 \chi_i^*(\mathbf{x}_1) h(1) \chi_i(x_1) \quad (17)$$

$$\mathcal{J}(x_1) = \int d\mathbf{x}_2 |\chi_i(\mathbf{x}_2)|^2 r_{12}^{-1} \quad (18)$$

$$\mathcal{K}(x_1) \chi_i(x_1) = \left[\int dx_2 \chi_j^*(x_2) r_{12}^{-1} \chi_i(x_2) \right] \chi_j(x_1) \quad (19)$$

The restriction on the summation of $i \neq j$ would require a Fock matrix defined for each orbital, but since cases where $i = j$, $[\mathcal{J} - \mathcal{K}] \chi_i = 0$ the restriction on the sum can be removed and a single Fock matrix can be used for all orbitals. The HF equations can then be defined within a basis set, which result in the Roothaan equations. If we expand our basis as in Equation 5 the matrix form of the equations become:

$$\sum_v F_{\mu v} C_{vi} = \epsilon_i \sum_v S_{\mu v} C_{vi} \quad (20)$$

where μv are indices of the atomic orbitals and i is the index of the molecular orbital. This expression of the energy is in accordance with the variational principle. The HF approach is an example of a self-consistent method where we must iterate through steps of forming the Fock matrix, Diagonalizing the Fock matrix, and checking for convergence. The convergence measures are set to a

convergence in the energy or in the density, though there exist other convergence checks. This will be in a self-consistent fashion, where the values of the Fock matrix depends on it's own solution.

The HF approach is often the theoretical starting point of other theoretical approaches. Due to the fact that HF is a mean-field approach, electrons are not correlated to one another. Electron correlation can be important to describe many chemical phenomena. In fact, the difference in energy of the true system within the nonrelativistic treatment from the energy described with the HF method is the correlation energy, $E_{corr} = E_{True} - E_{HF}$. Methods discussed in the next sections this are attempting to recover E_{corr} through approximate approaches to treating interacting electrons.

2.1.2 Density Functional Theory

Density Functional Theory (DFT) is another mean field approach but attempts to estimate the missing correlation energy. DFT is one of the most common methods used in solving the electronic structure in chemistry applications as it offers an approximation of the effects of electron correlation, while still having low computational scaling ($\sim N^3$). The fundamentals of the method are described here as it can provide reliable results in many cases, but also understanding the reasons of failure are a motivating facet of this work.

Fundamentally, the energy of a system can be understood as the interactions of the electrons and an external potential which is induced by the interaction of electrons with nuclei and additional sources:

$$E = \langle \Psi | \hat{H} | \Psi \rangle = \langle \Psi | \hat{T} + \hat{V}_{ee} | \Psi \rangle + \int d\mathbf{r} v_{ext}(\mathbf{R}) n(\mathbf{R}) \quad (21)$$

where \hat{T} and \hat{V}_{ee} are the kinetic energy and electron-electron potential operators. Since the wave function, Ψ is dependent on v_{ext} , all terms in this equation are also dependent on v_{ext} . In other words, the energy can be understood as a functional of the external potential, which can be represented as $E[v_{ext}]$. Hohenberg and Kohn postulated that the many-particle wave function can be represented by a reduced variable of the electron density^{52,53} It follows that if one can find the exact ground state electron density, the exact ground state energy is recovered. The formulation of a universal functional was defined which will be valid for any V_{ext} :

$$F_{HK}[n] = \langle T \rangle + \langle V_{ee} \rangle \quad (22)$$

and an energy functional defined as:

$$E_{HK}[n, v_{ext}] = F[n] + \int dr v_{ext}(r)n(r) \quad (23)$$

. For quantum chemistry applications, a further approximation of DFT is often used to simplify finding a good energy functional. This is motivated by the fact that from the previous formulation we are able to describe the external potential and the average Coulomb Hartree term, but the rest of the functional is unknown. Introduced by Kohn and Sham, the theorem states that a non-interacting system can be used to represent the electron density of the interacting, where the auxiliary system is chosen to reproduce some but not all properties of a system.⁵⁴ This changes our expression of the energy to

$$E_{KS}[n] = T_{ip}[n] + \int dr v_{ext}(r)n(r) + E_H[n] + E_{xc}[n] \quad (24)$$

where T_{ip} is now the independent particle kinetic energy, E_H is the Hartree potential and E_{xc} is the exchange correlation energy, E_{xc} , defined as:

$$E_{xc}[n] = \langle T \rangle - T_{ip}[n] + \langle V_{ee} \rangle - E_H \quad (25)$$

where $\langle T \rangle$ is the true kinetic energy, T_{ip} is the kinetic energy of the system in absence of any electron-electron interactions, $\langle V_{ee} \rangle$ is the true electron interaction potential, and E_H is the Hartree potential. From Equation 24, all terms can be solved exactly except for E_{xc} , thus much of the following discussion will revolve around how approximate E_{xc} values are determined.⁵⁵

2.1.2.1 Exchange-Correlation Functionals

The E_{xc} value is determined through the use of an exchange-correlation functional. Since the exact form of the functional is not known, there have been many approximations. There exists families of functionals characterized by the information incorporated. The families of functionals are colloquially referred to as rungs making up Jacob's ladder⁵⁶ and are discussed briefly here.⁵⁷ As more ingredients are incorporated, the accuracy of the functional generally increases. However, it must be emphasized that DFT is not systematically improvable.

The first rung is the local density approximation (LDA) functional, which is based on the uniform electron gas (UEG). Since there is no analytical form for the LDA correlation functional,

parameterized data is used from QMC results for the UEG from Ceperly and Alder.⁵⁸ The UEG is a good approximation for homogeneous systems and as a result the LDA functional can give reasonable descriptions in cases which resemble the UEG, but this functional often falls short for many molecules, surfaces, and defect solids. which are inhomogenous. The second rung in the ladder are generalized gradient approximation (GGA) functionals which improve upon LDA functionals by introducing inhomogeneity correction factors through the density gradient. GGA functionals correct for the overestimation of the interaction energy of the LDA functional. The next rung of Jacob's ladder contains meta-generalized gradient approximation (mGGA) functionals, which are an extension of GGA functionals that include further improvements to the describe the density. Two different improvements can be included, the Laplacian of the density or the kinetic energy density. The flexibility in the functional form provided by the inclusion of the kinetic energy density make it a more common functional ingredient. mGGA can improve upon GGA results, but are often more sensitive to the integration grid.

The functionals listed above all exhibit self interaction error, which results from the approximate exchange formulation. Hybrid functionals introduce exact exchange from HF in an attempt to correct for self interaction error.^{59,60} Hybrid functionals will be explored in Chapter 3 as tool to describe hydrogen chemisorbed to the surface of graphene. Further improvements to the functional form to correct for self interaction error results from range separated functionals.⁶¹

Beyond these main rungs of Jacob's ladder, there are other improvements to functional accuracy which can be based on dispersion corrections or inclusion of data from wave function methods.^{62–65}

The choices of functional are vast, but the lack of a systematically improvable framework, this can bring variability into the reliability of DFT results. Additionally, there are many systems which are a formidable challenge for DFT. Systems which exhibit strong correlation are poorly described by DFT. This is due to strong correlation requiring multireference wave functions to be captured accurately. The single reference nature of DFT leads to shortcomings in the description of these systems. The metal-insulator transition of Mott^{66,67} is a well-studied example.

There have been attempts to incorporate strong correlation into various functionals.^{68–70} Another approach is the use of the Hubbard model as a way to capture the effects of strong correlation.^{71,72} These methods have been successful in many cases⁷³. Though there are still questions

about the selection of the U and J parameters as often times these parameters are not transferable between systems and are selected based of semi-empirical data or more recently through machine learning approaches of QMC-based searches.⁶⁴

There have also been effort to go beyond standard DFT to describe strongly correlated systems.⁷⁴ One such approach for materials is the use of model Hamiltonians.⁷⁵ These approaches involve downfolded Hamiltonians, tight-binding models, dynamical mean-field theory (DMFT). DMFT methods map a many-body lattice problem to a local model called an impurity model.^{76,77} This approach has been used to successfully describe the Mott insulator transition mentioned above.⁷⁸

2.2 Wave Function Methods

2.2.1 Configuration Interaction

In the Configuration Interaction (CI) approach, the wave function becomes a linear expansion of a set of Slater determinants.

$$|\Psi_{CI}\rangle = \sum_k c_k^{CI} |D_k^{CI}\rangle \quad (26)$$

$$E_{CI} c_k^{CI} = \sum_l \langle D_k^{CI} | \hat{H} | D_l^{CI} \rangle \quad (27)$$

The level of excitation contained in the CI can be truncated to only include certain levels of excitations.

$$|\Psi_{CI}\rangle = (1 + \sum_{ia} c_i^a a_a^\dagger a_i + \sum_{ijab} c_{ij}^{ab} a_a^\dagger a_b^\dagger a_i a_j + \dots) |\Phi_{HF}\rangle \quad (28)$$

where a and a^\dagger are discussed in 2 and i and j are occupied orbitals and a and b are virtual orbitals. If all possible excitations are chosen, this is known as Full Configuration Interaction (FCI). FCI provides the exact solution within a single particle basis. The approach is size extensive (energy grows linearly with the number of electrons) and size consist (energy of the system is strictly separable.). However, the drawback of this method is the scaling of $N!$ as it considers all N electrons in $2n$ orbitals. This scaling can be slightly reduced by enforcing an S_z value (z component of the angular momenta), but the scaling still remains exponential.

2.2.1.1 Truncation of the FCI Space

In many instances, many of the determinants in the Hilbert space are not necessary for an accurate descriptions of the system.^{79,80} For instance, in many systems the FCI expansion, the HF determinant will have the largest coefficient. If the coefficient on the HF determinant is large enough, a small determinant expansion can often provide accurate results. Once one is considering a truncated CI, the methods are no longer size-consistent or size-extensive. Interacting with a subspace of the full Hilbert space can be done in two ways:

- choosing only certain orbitals and/or allowing only certain excitations
- selected CI approaches: a family of algorithms which look to incorporate determinants in an informed way

When including only certain excitations, the number of determinants is reduced. If only single excitations are considered, this is referred to as Configuration Interaction Singles (CIS), where Configuration Interaction Singles and Doubles (CISD) incorporates singles and doubles, and this pattern continues for including triple CISDT or quadruple CISDTQ excitations. Restrictions on the orbitals also allows for a significant reduction in the CI cost. This approach revolves around determining an active space which defines the orbitals which excitations are allowed to and from. One form of this type of restriction is freezing the core orbitals. While on the other end of things, the high-lying virtual orbitals can be frozen since the excitations into this space offer only small contribution to the wave function. However, this does lead to an ambiguity in the choice of active space for a given system. If one is defining a subset of orbitals and allowing all excitations within that subset, this approach is known as Complete Active Space Configuration Interaction (CASCI). The idea of reducing excitation level and orbital spaces are also a possibility.

In selected CI methods, the determinants are not chosen *a priori* but instead are selected during a process based on how the determinant contributes to the FCI wave function. There are various formulations of selected CI approaches, some of which were first presented in 1960,⁸¹ but selected CI methods are experiencing a renaissance in recent years. Some of the formulations are heat bath-CI^{82,83}, full configure interaction quantum Monte Carlo (FCIQMC)⁸⁴, adaptive CI⁸⁵, and perturbatively selected configuration interaction scheme (CIPSI)^{32,33}. The work presented here takes advantage of CIPSI.

In the CIPSI method, an initial zeroth-order wave function is defined as an initial determinant expansion:

$$|\Psi^{(0)}\rangle = \sum_I c_I |D_I\rangle \quad (29)$$

where I can begin at 1 if using a single determinant or many if a small multi-determinant wave function was generated from a truncated CI, such as CIS. The variational energy, $E^{(0)}$, can be evaluated as:

$$E^{(0)} = \frac{\langle \Psi^{(0)} | \hat{H} | \Psi^{(0)} \rangle}{\langle \Psi^{(0)} | \Psi^{(0)} \rangle} \quad (30)$$

From the $|\Psi^{(0)}\rangle$, determinants external to this reference space are generated, $|\alpha\rangle$. The set of $|\alpha\rangle$ determinants are evaluated using a perturbative criteria, the Epstein-Nesbet perturbative energy (Equation 31) to give an estimate of the correlation energy recovered by the inclusion of $|\alpha\rangle$ into the variational wave function.

$$e_\alpha^{(2)} = \frac{\langle \Psi^{(0)} | \hat{H} | \alpha \rangle^2}{E^{(0)} - \langle \alpha | \hat{H} | \alpha \rangle}. \quad (31)$$

This can provide a second order energy correction to the variational energy

$$E^{(2)} = \sum_\alpha e_\alpha^{(2)} = \sum_\alpha \frac{\langle \alpha | \hat{H} | \Psi^{(0)} \rangle^2}{E^{(0)} - \langle \alpha | \hat{H} | \alpha \rangle} \quad (32)$$

where $E^{(2)}$ gives an approximation of the correlation energy that $\Psi^{(0)}$ is neglecting. A CIPSI iteration consists of the following steps:

1. Define a variational wave function which is an expansion over a set of *internal determinants*, \mathcal{I}

$$|\Psi^{(0)}\rangle = \sum_{I \in \mathcal{I}} c_I |D_I\rangle$$

2. Assess the variational energy according to Equation 29
3. Generate \mathcal{A}_n , the set of all *external determinants* $\alpha, \notin \mathcal{I}_n$ and $\langle \Psi^{(0)} | \hat{H} | \alpha \rangle \neq 0$
4. Compute the perturbative contribution according to Equation 32
5. Take the α determinants with the largest $e^{(2)}$ contribution and incorporate them into \mathcal{I}_{n+1}
6. Iterate until the number of determinants included in the space exceeds a certain threshold or the perturbative contribution to the energy from the external set of determinants contributes less than a certain threshold.

The CIPSI method has its own formulation of calculating excited states, forces, and dipoles.^{36,86} However, it is often used a trial wave function for quantum Monte Carlo techniques.^{37,87}

2.3 Coupled-Cluster

The CC approach uses an exponential Ansatz of the form

$$\Psi_{CC} = e^{\hat{T}} |\Phi_0\rangle \quad (33)$$

where Φ_0 is a reference function. The \hat{T} is a cluster operator which contains degrees of excitations in the form $\hat{T} = \hat{T}_1 + \hat{T}_2 + \hat{T}_3 + \dots$ and $\hat{T}_1 = \sum_{ia} t_{ia} a_i^\dagger a_a$. The amplitudes of the cluster coefficients and the ground state energy, E_0 are found by solving the following set of nonlinear equation:

$$\langle \Phi_0 | e^{-\hat{T}} \hat{H} e^{\hat{T}} | \Phi_0 \rangle = E_0, \langle \Phi^* | e^{-\hat{T}} \hat{H} e^{\hat{T}} | \Phi_0 \rangle = 0 \quad (34)$$

where Φ^* is a specific excitation. This exponential Ansatz guarantees the CC wave function is size consistent and size extensive.

If one considered all excitation operators, the FCI wave function would be recovered, though this is just as computationally prohibitive. In practice, only certain excitation operators are used. In Coupled Cluster Singles and Doubles (CCSD) only single and double excitations are considered. This approach scales as n^6 where n is the number of spatial orbitals. However, failings of CCSD in certain cases motivated further excitations to be considered.[†] If triple excitations are included, the CC scales as n^8 , which becomes prohibitive for many systems. Instead, the effect of the triple excitations can be approximated by a perturbative correction which is a method known as Coupled Cluster Singles, Doubles, and Perturbative Triples (CCSD(T)), which scales as n^7 . This method is often referred to as the gold standard in computational quantum chemistry.

Additionally as one makes trade-offs of treating larger systems sizes with the CC methods which contain a lower number of excitations to achieve a lower scaling, some chemical phenomena will be poorly described. A method which attempts to overcome having to make this trade off is domain-based local pair natural orbital coupled-cluster (DLPNO)[88], which is a linear scaling approximation to the CC methods.

2.3.1 Equation-of-Motion

CC methods can also be used to obtain excited states for other approximate methods such as linear response theory or equation-of-motion (EOM), which can help explore photochemistry and spectroscopy questions.

The typical methods for excited states are ΔE methods are solving separately for the states of interest. equation-of-motion (EOM) methods instead allow us to solve for the excited state from a ground state reference.

$$\Psi_{ex} = \mathcal{R}|\Psi_0\rangle \quad (35)$$

In particular, this discussion focuses on EOM-CC methods. The CC equations can be written to express the impact of the cluster expansion on the reference state in the form of a normal Hamiltonian⁸⁹

$$\hat{H}_N = \hat{H} - \langle \Phi_0 | \hat{H} | \Phi_0 \rangle \quad (36)$$

where the cc equation becomes:

$$\hat{H}_N e^T |\Phi_0\rangle = \Delta E_0 e^T |\Phi_0\rangle \quad (37)$$

where the correlation energy is $\Delta E_0 = E_0 - \langle \Phi_0 | e^T \hat{H} e^{-T} | \Phi_0 \rangle$. These same expressions can be used for the excited state, Ψ_k . The solution to the Schrödinger equation becomes

$$\hat{H}_N |\Psi_k\rangle = \Delta E_k |\Psi_k\rangle \quad (38)$$

$$\hat{H}_N \mathcal{R} |\Psi_0\rangle = \Delta E_k \mathcal{R} |\Psi_0\rangle \quad (39)$$

where $\Delta E_k = E_k - \langle \Phi_0 | \hat{H} | \Phi_0 \rangle$ and E_k is the energy of the targeted state.

Subtracting the EOM Ansatz inserted into the Schrödinger equation ($H\mathcal{R}|\Psi_0\rangle = E_{ex}\mathcal{R}|\Psi\rangle$) from the application of the \mathcal{R} operator onto the ground state ($\mathcal{R}H|\Psi_0\rangle = E_0\mathcal{R}|\Psi\rangle$) results in:

$$[H_N, \mathcal{R}] |\Psi_0\rangle = \Delta \omega_k \nabla |\Psi_0\rangle \quad (40)$$

where $\omega_k = E_k - E_0$ and E_0 is energy of a single reference determinant. The CC exponential ansatz of Equation 33 can be substituted and the exponential cluster operator introduced through left multiplication which results in

$$e^{-\hat{T}} [H, \mathcal{R}] e^{\hat{T}} |\Phi_0\rangle = \Delta \omega_k e^{-\hat{T}} \mathcal{R}_k e^{\hat{T}} |\Phi_0\rangle. \quad (41)$$

If the normal-ordered Hamiltonian is slightly reordered to:

$$\bar{H}_N = \langle \Phi_0 | e^{-\hat{T}} \hat{H} e^{\hat{T}} | \Phi_0 \rangle \quad (42)$$

allows us to express the EOM equations as an eigenvalue problem of \bar{H}_N :

$$\bar{H}_N \mathcal{R} |\Psi_0\rangle = \omega_k \mathcal{R} |\Psi_0\rangle \quad (43)$$

However, \hat{H}_N is not Hermitian and thus, the right and left eigenvectors do not need to agree, but do exhibit bi orthogonality. The left hand solution is a de-excitation operator, \hat{L} .

The EOM solution will often target a specific property such as an excitation energy, ionization potential, electron affinity, and others. The solution results from a diagonalization of the \hat{H}_N operator within the appropriate configurational subspace for the property. For example, the EE-EOM-CCSD subspace would be for $|\Phi_i^a\rangle$ and $|\Phi_{ij}^{ab}\rangle$

The excitation operator has a similar form the the cluster operator in CC methods;

$$\mathcal{R} = r_0 + \sum_i \sum_a r_i^a a_a^\dagger a_i + \sum_{i,j} \sum_{a,b} r_{ij}^{ab} a_a^\dagger a_b^\dagger a_i a_j + \dots \quad (44)$$

The r coefficients can be solved through a coupled cluster approach. Due to the similarity in form, this operator can use similar truncation schemes to those used in CC. The EOM method shares similar scaling to CC when the same level of excitations are included in the operator. This means that although the method provides a clear path to excited states, the application will be restricted by scaling. DLPNO implementations of EOM exist implementations of to overcome the scaling restrictions.

2.4 Stochastic methods

The Hilbert space of quantum systems containing over a few atoms and a few orbitals is so large for describing interacting electrons that many-body integration becomes impossible. This is part of the problem referred to as the *curse of dimensionality*. A possible work around is to sample the phase space through a stochastic Monte Carlo (MC) method.

Stepping away from the electronic structure perspective, Monte Carlo integration methods sample a multidimensional integral at a randomly selected set of points, instead of assessing the multidimensional integral on a grid as is done in standard integration techniques. To generalize this, the naive MC approach will take a function, sample random points within a volume element that one is trying to integrate over, and have values for this function at random points returned. Due to the law of large numbers, the expected value of the integral can be understood as the mean result of the samples. In these methods, the statistical error and statistical variance must be known to assess the reliability.

Given a probability distribution, Π , a random walk can be taken through the possible states contributing to this distribution. This problem can be cast as a Markov chain where the samples become a series of events which are described as a transition matrix $T(s \rightarrow s')$, where a transition to s' only depend on the current state (s) and no others. This characteristic describes the concept of ergodicity if the starting state information is unknown after a certain number of steps. This requires that 1) One can move from one state to any other state in a finite number of steps with non-zero probability, 2) the transition matrix is non-periodic, and 3) average return time to any state is finite.

One can also enforce detailed balance to ensure only the intended distribution is sampled. In other words, the flux from state "A" to state "B" is equal to the flux of going from "B" to "A". In practice, the Metropolis-Hastings algorithm is used to enforce detailed balance. The Metropolis-Hastings algorithm is given in Algorithm 2.4 and a simple visualization shown in Figure 1.

Algorithm 1 Metropolis-Hastings algorithm

1. Describe the transition matrix and probability distribution
 2. Set an initial state
 3. iterate M steps
 - Sample s' from $T(s^{(n)} \rightarrow s')$, where $s^{(n)}$ is state at current iteration.
 - calculate the acceptance ratio $A = \frac{\Pi(s')T(s' \rightarrow s^{(n)})}{\Pi(s^{(n)})T(s^{(n)} \rightarrow s')}$
 - accept or reject by drawing a random number from a uniform distribution between (0,1),
$$u_n. s^{n+1} = \begin{cases} s', & A > u_n \\ s^{(n)}, & A \leq u_n \end{cases}$$
-

2.4.1 Considerations for QMC

In this work and in many QMC methods, a correlated Slater-Jastrow wave function is used. [90–93].

$$\Psi_{SJ}(R) = \det(\varphi_k(x_i)) e^{\sum_{i < j} u(x_i, x_j)} \quad (45)$$

where φ_k are the orbitals, and the exponential term is the Jastrow correlation factor. This can be reformulated over all electrons and expressed as correlations through 1-, 2-, and 3-body terms as

$$U(R) = \sum_i^I u(\mathbf{r}_{iI}) + \sum_{i < j}^N u(\mathbf{r}_{ij}) + \sum_{i < j}^N \sum_l^M u(\mathbf{r}_{il}, \mathbf{r}_{jl}, \mathbf{r}_{ij}) \quad (46)$$

The Jastrow factor can have various forms.[94, 95] The forms used within this work are described in subsequent the chapters.

2.4.2 Variational Monte Carlo

The variational Monte Carlo (VMC) approach is based on the variational theorem. The method is used to optimize the Jastrow parameters of the trial wave function. In VMC a random walk is used to sample the distribution:

$$\Pi(R) = \frac{|\Psi_T(R)|^2}{\int |\Psi_T(R)|^2 dr} \quad (47)$$

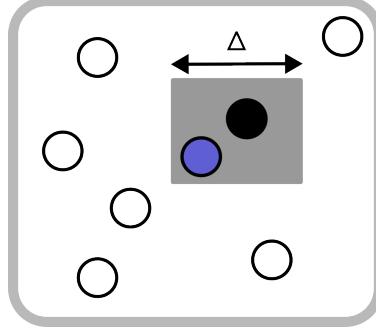


Figure 1: Schematic of a random walk, where the black filled dot is a proposed move within the gray square to the blue position. The acceptance of the move will be determined by the Metropolis-Hastings algorithm. Figure adapted from [53]

where Ψ_T is the trial wave function. Equation 9 can then be rewritten as an average value of the “local energy” over the distribution $\Pi_T(R)$.

$$E_V = \int dR \Pi_T(R) E_L(R) = \langle E_L(R) \rangle_{\Pi} \quad (48)$$

where the local energy is defined as:

$$E_L = Re \frac{\hat{H}\Psi_T(R)}{\Psi_T(R)} \quad (49)$$

and an average variational energy over a number of random walk steps can be defined as:

$$\bar{E}_V = \frac{1}{M} \sum_n E_L(R^{(n)}) \quad (50)$$

As discussed, the error bar information of the random walk is a necessary assessment of the reliability of the results. For VMC, the variance of the energy can be assessed through:

$$\nu_T = \frac{\int dR |\Psi_T|^2 (E_L(R) - E_V)^2}{\int dR |\Psi_T|^2} \quad (51)$$

which as an average variance value over M steps has the form:

$$\bar{\nu}_T = \frac{1}{M} \sum_{n=1}^M (E_L(R^{(n)}) - \bar{E}_V)^2 \quad (52)$$

This method is restricted by the quality of the Ψ_T . In other words, even with sufficient sampling, VMC will produce an exact answer only if $\Psi_T = \Psi_{exact}$.

2.4.3 Projector Monte Carlo

Projector Monte Carlo is an umbrella term for a variety of methods that all filter out the ground state by repeatedly applying a many-body projector with a random walk[53, 96]. This formulation can be understood if one begins with the time-dependent Schrödinger equation:

$$i\frac{d\Psi(r,t)}{dt} = \hat{H}\Psi(r,t) \quad (53)$$

and introduced the transformation of substitution $t \rightarrow i\tau$ so that the wave function is evolves through “imaginary time” as

$$\frac{d\Psi(r,\tau)}{d\tau} = \hat{H}\Psi(r,\tau). \quad (54)$$

Expressing Ψ as a linear combinations of eigenfunctions:

$$\Psi(R) = \sum_j c_j \varphi_j(R) \quad (55)$$

$$\hat{H}\varphi_i = E_i \varphi_i \quad (56)$$

allows one to write a formal solution to the imaginary time Schrödinger equation:

$$|\Psi(\tau_1 + \delta\tau)\rangle = e^{-H\delta\tau} |\Psi(\tau_1)\rangle \quad (57)$$

this can then be expanded into energy eigenstates, as is done in Equation 55

$$|\Psi(\delta\tau)\rangle = \sum_{i=0}^{\infty} c_i e^{-\varepsilon_i \delta\tau} |\varphi_i\rangle. \quad (58)$$

Finally, a shift by E_0 is introduced to ensure a stationary solution as $\tau \rightarrow \infty$

$$|\Psi(\delta\tau)\rangle = \sum_{i=0}^{\infty} c_i e^{-\varepsilon_i - E_0 \delta\tau} |\varphi_i\rangle \quad (59)$$

Thus as imaginary time progresses through a repeated application of the projector, any state that is not orthogonal to the ground state and is higher in energy will be damped out leaving the ground state projected out.

Here, the discussion will focus on diffusion Monte Carlo (DMC) as that is the main projector Monte Carlo approach used in this work. The DMC projector is

$$\hat{G}_D = e^{-\Delta\tau(\hat{H} - E_T)} \quad (60)$$

where E_T is the trial energy and is introduced as a way to control the normalization. $\Delta\tau$ is the time step in imaginary time. The projection is carried out through a random walk through the space of electron coordinates. However, \hat{G}_D is needed to be in a spatial representation. This is accomplished through a short time expansion, which is made possible through Trotter's formula[97, 98]:

$$E^{-\tau(\hat{T}+\hat{V})} = \lim_{n \rightarrow \infty} (e^{-\Delta\tau\hat{T}} e^{-\Delta\tau\hat{V}})^n \quad (61)$$

The common approach to sampling is a branching scheme where one is duplicating and terminating random walks. One can consider the effect of the random walk through the impact of the Hamiltonian components individually. The kinetic energy term in the absence of the potential term leads to a diffusive step and the random walk progresses from R to R' according to the diffusion probability:

$$P_d(R' \leftarrow R) = (2\pi\Delta\tau)^{\frac{3N}{2}} \exp((|R - R'|^2)/(\Delta\tau)) \quad (62)$$

If the potential term is considered in the absence of the kinetic term, branching steps are defined according to $\exp(V(R') - E_T)$. However, since \hat{T} and \hat{V} do not commute, an error is incurred. The error can be controlled through alternating between branching and diffusive steps as well as by utilizing the second order Trotter-Suzuki formula.[99]

Beyond the formulation of the method, there are a few practical considerations that must be made to the algorithm to produce reliable results with efficient sampling.

Importance Sampling: Given a long enough propagation through imaginary time the ground state will be projected out for any initial guess wave function which adheres to the boundary conditions of the problem. However this is often unstable and inefficient. Instead, one can utilize importance sampling, which introduces a guiding wave function, Ψ_G , that approximates the true ground state wave function. The projector becomes

$$G(R \leftarrow R') = \Psi_G(R) \langle R | e^{-\Delta\tau(\hat{H}-E_T)} | R' \rangle \Psi_G(R)^{-1} \quad (63)$$

which results in higher sampling in areas where Ψ_G is large.

Fixed Node Approximation: The Fermion wave function is not a probability distribution as it will have positive and negative regions. The most common way to address the sign problem is to forbid moves that would change the sign of the trial wave function. The fixed-node energy is an upper

bound to the exact ground state energy. The incorporation of this approximation is referred to as FN-DMC. The DMC investigations presented in Chapters 3 and 4 will utilize FN-DMC approach.

2.5 Fragmentation Methods

To overcome the high scaling of computational cost with system size, fragment approaches approximate the energy and other properties of the system by partitioning into subsystems, where the final result becomes the accumulation of the parts. One potential advantage of fragmentation is to reduce the computational scaling from $\mathcal{O}(N^p)$ to $\mathcal{O}(N_{\text{frag}}^q f^p)$, where p is the scaling exponent for the chosen electronic structure method, where N_{frag} is the number of fragments, f is representative of the fragment size, and q is an exponential specific to the fragmentation approach. Judicious choices in the fragment approach can lead to a reduction in overall computational scaling ($q < p$), and in best cases linear scaling in the number of fragments can be obtained, $\mathcal{O}(N_{\text{frag}} f^p) \approx \mathcal{O}(N_{\text{frag}})$, where constant fragment size f is assumed. Another appealing feature of fragment approaches is that parallelization is facile since much of the workload for the individual subsystems can be treated separately, which provides the potential to efficiently utilize high performance computing resources to reduce the overall wall time of calculations. Beyond these computational advantages, fragment approaches can quantify magnitudes and roles of inter-fragment interactions that can deepen chemical understanding similar to analysis techniques such as energy decomposition analysis or symmetry adapted perturbation theory.[100–104] A perhaps surprising application of fragment approaches was proposed for quantum computing, where it may offer a route to overcome the challenges of system-environment interactions which hamper the integrity of results for quantum simulations of large systems. Instead of treating an entire (large) system on a quantum computer, fragment approaches were suggested so that the system could be partitioned into subsystems which lead to circuits that are less susceptible to errors, while incorporating the most important chemical interactions in the quantum domain.[105] In summary, fragment approaches offer useful features ranging from computational efficiency over chemical insights to potential applications in unprecedented fields such as quantum computing.

The accuracy and efficiency of fragment approaches hinges on a number of factors, including

the choice of electronic structure treatment for each fragment, approaches to incorporate inter-fragment interactions, and the specific choices for partitioning of the system. Significant attention has been devoted to the treatment of inter- and intra-fragment interactions. These approaches will be discussed briefly, but are beyond the scope of the current work and interested readers are directed towards a number of helpful reviews.[40, 106] In contrast, this work investigates how a number of strategies for the selection of fragments impact the accuracy and efficiency of the results.

2.5.1 Fragment treatments

Although improving the treatment of inter- and intrafragment interactions is not a focus of the current body of work, a brief overview of the main families of methods are given for context. These methods can be categorized as, many-body expansion (MBE) approaches, and multilayer composite approaches. (“ONIOM” scheme), approximate functional group additivity approaches.

The many body expansion approaches can have many flavors, but they essentially look to take the energy of each fragment and estimate the energy of including many body terms into the energy. The estimation of many-body effects to the energy is done by studying groups of fragments according to Equation 64

$$E = \sum_I E_I + \sum_I \sum_{J>I} \Delta E_{IJ} + \sum_I \sum_{J>I} \sum_{K>J} \Delta E_{IJK} + \dots \quad (64)$$

where the ΔE terms represent additions of the many body terms to the initial fragment energy, i.e., the energy of the isolated fragments subtracted from the energy of combined fragment. For a concrete example, the two body form of this term is

$$\Delta E_{IJ} = E_{IJ} - E_I - E_J \quad (65)$$

The popular quantum chemistry method of fragment molecular orbitals fall under the umbrella of a MBE. The many-body terms will be truncated at a certain level of interaction.

A different family of methods are known as multilayer approaches which break the system into areas of interest to estimate the full system as a combination of subsystems treated at varying levels of theory. This in theory would allow one to employ a high level of theory to a realm

which one knows interesting physics is occurring, but then switch to a medium or low level theory for those areas participating less in the area of interest. Methods that fall under this family are ONIOM[107–109] and molecule-in-molecule (MIM) approaches[110–112]. The third type of fragmentation method attempts to determine the size of the fragments through an iterative scheme until convergence of the property with changing fragment size.

3.0 The binding of atomic hydrogen on graphene from density functional theory and diffusion Monte Carlo calculations

The test and figures in this chapter have been adapted from A. Dumi, S. Upadhyay, L. Bernasconi, H. Shin, A. Benali, and K.D. Jordan The binding of atomic hydrogen on graphene from density functional theory and diffusion Monte Carlo calculations, *J. Chem. Phys.*, **2022**, 156, 144702, DOI:10.1063/5.0085982.with the permission of AIP Publishing.with permission from The author's contribution to the work included performing trial wave function generation calculations, generation of images, setting up the QMC workflow with S. Upadhyay, along with performing analysis and assisting with figures, and editing/revising the manuscript.

3.1 Summary

In this work density functional theory (DFT) and diffusion Monte Carlo (DMC) methods are used to calculate the binding energy of a H atom chemisorbed on the graphene surface. The DMC value of the binding energy is about 16% smaller in magnitude than the Perdew–Burke–Ernzerhof (PBE) result. The inclusion of exact exchange through the use of the Heyd–Scuseria–Ernzerhof (HSE) functional brings the DFT value of the binding energy closer in line with the DMC result. It is also found that there are significant differences in the charge distributions determined using PBE and DMC approaches.

3.2 Introduction

The unique electronic, optical, and transport properties of graphene make it an important system for a wide range of applications, many of which involve or are impacted by the adsorption of atoms or molecules. To bring these applications to fruition, a deeper understanding of the interaction of atoms and molecules with graphene is required, and, not surprisingly, this has been the

subject of several experimental and theoretical studies.^{2–14}

The adsorption of H atoms on graphene has been the subject of multiple studies.^{3–6,8,113} It is known that there is both a weakly absorbed state in which barriers for diffusion are small and a much more strongly bound chemisorbed state^{114,115}, which is the focus of this work. Chemisorbed H atoms open up the band gap and allow for tuning of electronic properties.¹¹⁶ It has been demonstrated that even a single chemisorbed hydrogen atom causes an extended magnetic moment in the graphene sheet.^{117,118} On the other hand, there is evidence that given the ready diffusion of H in the physisorbed state, the H atoms tend to pair up on the surface leading to non-magnetic species.¹¹⁸ Finally, interest in the hydrogen/graphene system has also been motivated by the potential use of graphene and graphitic surfaces for hydrogen storage.¹⁴ In spite of the interest in H chemisorbed on graphene, we are unaware of experimental values of the binding energy.

Most computational studies of adsorption of atoms and molecules on graphene have employed density functional theory (DFT), primarily due to its favorable scaling with system size, allowing for the treatment of larger periodic structures. However, a reliable theoretical description of interactions at the graphene surface has proven to be challenging for DFT.^{2,7,8,15} In recent years considerable progress has been made in extending correlated wave function methods to periodic systems.^{87,119–123} Among these methods, the diffusion Monte Carlo (DMC)¹⁶ method, which is a real-space stochastic approach to solving the many-body Schrödinger equation is particularly attractive given its low scaling with the number of electrons and high parallelizability. DMC also has the advantages of being systematically improvable and its energy being much less sensitive to the basis set employed than methods that work in the space of Slater determinants. In DMC calculations, the atomic basis set is important only to the extent that it impacts the nodal surface. DMC has been used to describe the adsorption of various species on graphene including O₂¹⁰, a water molecule^{9,17}, and a platinum atom.¹¹ In a study of a physisorbed H atom on graphene, Ma et al. found that different DFT functionals gave binding energies ranging from 5 to 97 meV, while DMC calculations gave a value of only 5 ± 5 meV.⁸ Various DFT calculations utilizing the Perdew-Burke-Ernzerhof (PBE)¹²⁴ and Perdew-Wang (PW91)¹²⁵ functionals predict the chemisorbed H atom species to be bound by 480 to 1,440 meV.^{126–134} However, this large spread is primarily a result of some calculations employing small supercells resulting in an unphysical description of the low-coverage situation, too small a *k*-point grid, or small atom-localized basis sets that do not

adequately describe the binding and introduce large basis set superposition error (BSSE). In the present work, we use the DMC method to calculate the binding energy of H to graphene in the chemisorbed state.

3.3 Methods

All calculations reported in this study used a 5x5x1 supercell of graphene, as it was large enough to make inconsequential the interaction between periodic images of the adsorbed hydrogen atom and to assure that there are essentially unperturbed C atoms between the buckled regions in adjacent images in the *x* and *y* directions. The geometries of graphene, both pristine and with a chemisorbed H atom, were provided by Kim et al.,¹³⁵ and were obtained using the PBE+D3 DFT method.^{124,136} For all systems, a vacuum spacing of 16 Å was used.

3.3.1 Density Functional Theory Calculations

The single particle orbitals used in the trial wave functions for variational Monte Carlo (VMC) and DMC calculations were calculated using the PBE functional with the correlation consistent electron core potential (ccECP)^{137,138} pseudopotentials and a plane wave basis with an energy cut-off of 3,400 eV. Monkhorst-Pack *k*-point grid meshes¹³⁹ were employed with a 13.6 meV Marzari-Vanderbilt-DeVita-Payne cold smearing of the occupations.¹⁴⁰ The PBE results were converged at a 6x6x1 *k*-point grid to 1 meV for graphene and graphene with an adsorbed hydrogen atom. The hydrogen atom trial was generated using a 1x1x1 *k*-point grid. Convergence studies can be found in Table S1 and S2 of the Supplementary Material.

In addition to the PBE calculations used to generate the trial wave functions for DMC, DFT calculations were carried out with the PBE0¹⁴¹ and Heyd–Scuseria–Ernzerhof (HSE) functionals¹⁴² to determine if inclusion of exact exchange proves important for the adsorption energy. Due to the inclusion of exact exchange, these calculations would be computationally prohibitive in a plane wave basis, particularly with the high energetic cutoff required by the ccECP pseudopotential. For this reason, they were carried out all-electron with the POB-TZVP Gaussian type orbital

(GTO) basis set.¹⁴³ Due to the use of GTOs, these calculations suffer from basis set superposition error (BSSE), which we corrected using Grimme's geometry-dependent counterpoise correction scheme.^{144,145} This correction resulted in a 113 meV reduction in the magnitude of the binding energy when using the PBE0 functional. For the PBE0 and HSE calculations, a 12x12x1 k -point grid was used to assure binding energies converged to within 2meV. Convergence data are supplied in Table S3 of the Supplementary Material.

The plane wave DFT calculations were carried out with the QUANTUM ESPRESSO version 6.3 code.^{146–148} The Gaussian basis DFT calculations were carried out with CRYSTAL17,^{149,150} save for the HSE calculation of the isolated hydrogen atom which was carried out using NWChem version 6.8¹⁵¹ using the same basis as the calculations in CRYSTAL17.

3.3.2 Quantum Monte Carlo Calculations

DMC is a projector quantum Monte Carlo (QMC) method, solving the Schrödinger equation in imaginary time $\tau = it$; any initial state $|\psi\rangle$, that is not orthogonal to the true ground state $|\phi_0\rangle$, will evolve to the ground state in the long time limit. When dealing with Fermionic particles, the DMC method requires the use of the fixed-node approximation¹⁵² to maintain the antisymmetric property of the wave function. For efficient sampling and to reduce statistical fluctuations, we use a Slater-Jastrow trial wave function fixing the nodes through a Slater determinant comprised of single-particle orbitals, which, in this work, are expanded in a B-spline basis. The Jastrow factor is a function that reduces the variance by explicitly describing dynamic correlation. The Jastrow factor contains terms for one-body (electron-ion), two-body (electron-electron) and three-body (electron-electron-ion) interactions. The one- and two-body terms were described with spline functions¹⁵³, while the three-body terms were represented by polynomials.¹⁵⁴ 10 parameters were used for the one-body terms per atom type, and 10 parameters were employed per spin-channel for the two-body terms. The cutoffs on the one- and two-body terms were fixed to the Wigner-Seitz radius of the simulation cell. The three-body terms were comprised of 26 parameters per term with a cutoff of 10 Bohr. The parameters in the Jastrow factor were separately optimized for each geometry with the linear method¹⁵⁵ using VMC. To reduce the cost of the DMC calculations as well as to reduce the fluctuations near the ionic core regions, ccECP pseudopotentials were used to

replace the core electrons.^{137,138} The ccECP pseudopotentials were designed to be used with high-accuracy many-body methods. The non-local effects due to the pseudopotentials were addressed

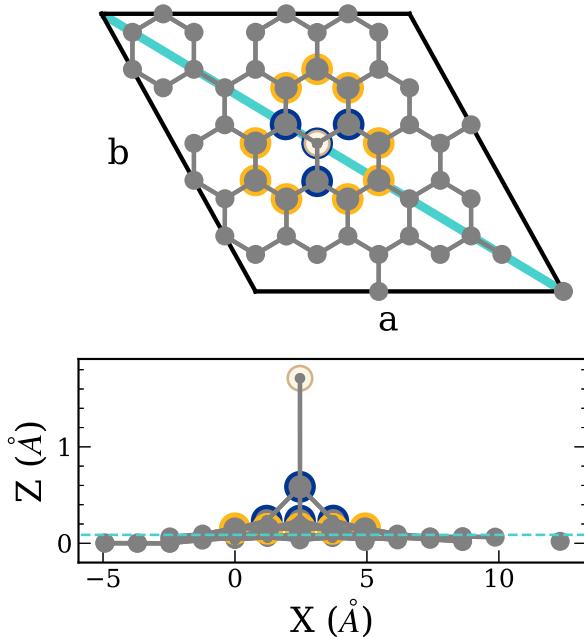


Figure 2: Perpendicular view of the simulation cell (top) and a parallel view obtained by projection onto the xz -plane (bottom). The carbon atoms are colored gray and the hydrogen atom is denoted as white. For the perpendicular view, the cyan line represents the slice of the cell used to visualize electron density differences. For the parallel view, the dotted cyan line represents the mean carbon z position. Blue outlined atoms are greater than one standard deviation away from the mean carbon z position, whereas yellow atoms are between $0.5\text{--}1.0 \sigma$.

using the determinant-localization approximation along with the t-moves method (DLTM).^{156,157} Finite size effects were addressed using twist averaging.¹⁵⁸ The twist angles were chosen to be the symmetry unique points of the $6\times 6\times 1$ k -point grid shifted by half a grid step away from the gamma point in each direction.

The DMC calculations were performed using the branching scheme proposed by Zen et al. (ZSGMA)¹⁵⁹ with a population control target of 8,192 walkers and a time step of 0.005 a.u., which represented a balance between computational cost and finite timestep error in previous work.¹⁰

We define the binding energy as

$$E_b = E_{dgr+H} - (E_{gr} + E_H) \quad (66)$$

where E_{dgr+H} is the energy of the distorted graphene sheet with a chemisorbed atomic hydrogen, E_H is the energy of a hydrogen atom, and E_{gr} is the energy of a pristine graphene sheet. In the chemisorbed state, the hydrogen atom bonds directly over a carbon atom, causing this carbon to be pulled out of the sheet towards the hydrogen.^{160,161} The adjacent carbons are also pulled in the direction of the hydrogen leading to a distorted graphene sheet.

The QMC calculations were carried out using the QMCPACK code, with the workflow between QUANTUM ESPRESSO and QMCPACK managed by Nexus.^{162–164} Figures 2 and 4 were rendered with matplotlib¹⁶⁵ and the density plots were generated using VESTA.¹⁶⁶

3.4 Results & Discussion

3.4.1 Binding Energy

Table 1 contains a summary of the binding energies of a hydrogen atom chemisorbed on graphene from this work and selected values from previous publications using the PW91 and PBE functionals. These literature values range from -790 to -980 meV. This wide spread of binding energies is caused by (1) the use in some studies of small supercells for which there are sizable interactions between the CH groups in adjacent cells, and (2) the use in some studies of small atom-centered basis sets without corrections for BSSE. Our calculations with the PBE functional in conjunction with a plane wave basis set give a binding energy of -821 meV. This should be contrasted with our -691 ± 19 meV DMC result. There are several possible sources for the difference between the PBE and DMC values of the binding energy. These include errors in the DFT calculations due to self interaction and planar graphene having more multiconfigurational character than H/graphene, with this being better described with DMC than with PBE. We note that the inclusion of the D3 dispersion correction with the PBE functional only changes the magnitude of the binding energy by 0.03 eV.

Table 1: Binding energy (meV) of a hydrogen atom chemisorbed on graphene calculated with various DFT functionals and with DMC.

Method	Binding energy
This Work	
PBE ¹	-821
PBE ²	-871
PBE0 ²	-851 (-800)
HSE ²	-794 (-743)
DMC	-691 \pm 19
Previous Work	
PW91	-810 to -830 ¹²⁷ , -870 ¹²⁸
PBE	-790 ¹²⁹ , -840 ¹³⁰ , -980 ¹³¹

¹ Calculation was done in the plane wave basis

² Calculation was done in the Gaussian basis set with corrections for BSSE. Values in parentheses include a correction for the basis set incompleteness as described in the text.

The PBE binding energy is 51 meV lower in magnitude in the plane wave than in the GTO basis set when the same k -point grid is used, and this value is used as a correction for the basis set incompleteness error for the results with other functionals in Table 1. The calculations in the GTO basis set give a slightly smaller in magnitude binding energy with PBE0 than with PBE. However, with HSE, we obtain a binding energy 77 meV smaller in magnitude than the PBE result. Applying the correction for the basis set incompleteness error, we obtain -800 meV for the PBE0 binding energy and -743 meV for the HSE binding energy, with the latter being in reasonable agreement with the DMC result of -691 meV. Although the 130 meV difference between the plane-wave PBE and DMC values of the binding energy may appear to be small, this energy difference, of that magnitude is consistent with an order of magnitude change in the hydrogen evolution current at

room temperature on graphene electrodes.¹³⁵

In order to better understand the origin of the difference in the PBE and HSE H-atom adsorption energies, we also carried out non-self-consistent calculations, using PBE densities to evaluate the HSE energies. These calculations gave a binding energy only 21 meV smaller in magnitude than obtained from the self-consistent HSE calculations. This demonstrates that the functional is more important than the density in establishing the binding energy. Detailed information can be found in Table S4 of the Supplementary Material.

Detailed results of the DMC calculations can be found in Tables S5-S7 in the Supplementary Material.

3.4.2 Binding Density

It is instructive to examine the change in the electron density associated with the binding of the H atom to the distorted graphene as determined from the PBE and DMC calculations. The density change is given by

$$\rho_b = \rho_{dgr+H} - (\rho_{dgr} + \rho_H), \quad (67)$$

where ρ_H is the charge density of the hydrogen atom, and ρ_{dgr+H} and ρ_{dgr} are the charge densities of the distorted graphene sheet with and without hydrogen, respectively. For the QMC density, the density was accumulated during the VMC and DMC calculations, the mixed estimator bias was found to be insignificant, and was thus not corrected.

The ρ_b density differences for both DMC and PBE are shown in Figure 3. The dark blue and gold regions represent a loss and gain of electron density, respectively. As expected, there is a shift in electron density from the carbon atom participating in the carbon-hydrogen bond as well as to the three adjacent carbon atoms. These qualitative changes in the density are consistent with previous theoretical and experimental studies.^{160,161} The rehybridization from sp^2 to sp^3 of the carbon participating in the CH bond and the weakening of the π bonds due to the distortion of the graphene lead to the electron density shift. The change in the charge distribution is similar for PBE and DMC, with the most noticeable difference being a greater increase of density at remote C atoms in the DMC than in the PBE calculations.

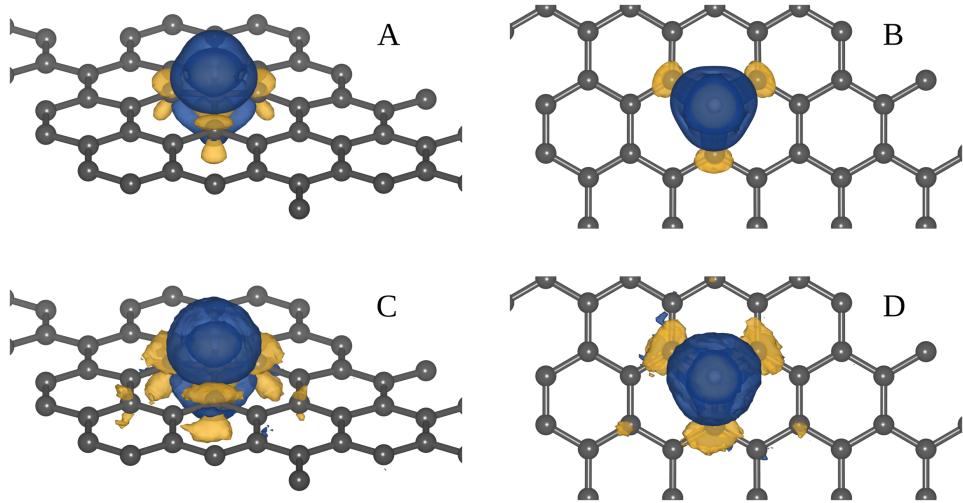


Figure 3: Change of the electron density due to the adsorption of the H atom to the distorted graphene sheet (Eq. 67). ρ_b from PBE calculations is shown from an oblique angle (A) and aligned along the c axis (B). ρ_b from DMC calculations (C) and (D) is shown from the same perspectives. Gold and blue represent a gain and loss of electron density, respectively. Note that there is a region of increased charge density at the C-H bond that is enveloped by a region of loss in the charge density. The binding density was visualized using an isovalue of 2.8×10^{-5} for DMC and 3.9×10^{-5} for PBE, in both cases capturing 95% of the differential charge density.

3.4.3 Charge Density Differences Between DMC and PBE

In this section, the difference between the DMC and PBE charge densities for distorted graphene with the adsorbed hydrogen atom as well as for pristine planar graphene without the adsorbed hydrogen atom are considered. The charge density difference for each system is calculated according to

$$\Delta\rho_{\text{system}} = \rho_{\text{system}}^{\text{DMC}} - \rho_{\text{system}}^{\text{PBE}}, \quad (68)$$

where $\rho_{\text{system}}^{\text{DMC}}$ is the DMC charge density of a given system (either distorted graphene with the adsorbed hydrogen or pristine graphene) and $\rho_{\text{system}}^{\text{PBE}}$ is the corresponding PBE charge density. $\Delta\rho_{gr}$ and $\Delta\rho_{dgr+H}$ are reported in Figure 4 along the 110 slice through the unit cell, which captures the carbon-hydrogen bond. From the top-down perspective in Figure 2, the 110 lattice plane bisects the cell diagonally through the longer of the two diagonals and is indicated by the solid cyan line. In Figure 4, blue represents areas where the PBE density is larger, while gold areas represent areas where the DMC density is larger. The DMC density, in comparison with the PBE density, has greater weight in the bonding region between atoms. We note that the HSE density displays similar differences as the PBE density. Figure S2 of the Supplementary Material includes a visualization of the DMC-HSE density difference. This is the case for both the planar graphene without hydrogen and the system with hydrogen chemisorbed to graphene. Even though there are significant differences between the PBE and DMC densities for both systems, the difference is similar in the two systems, consistent with it not introducing a large error in the PBE value of the binding energy.

3.5 Conclusions

Calculations of the binding energy of a hydrogen atom on a graphene sheet were carried out using various DFT methods and with DMC. The DMC calculations provide a benchmark value of the binding energy. Our best estimate of the binding energy from DMC calculations is -691 ± 19 meV. The PBE result obtained with a plane-wave basis set gives a binding energy about 20% larger in magnitude than the DMC result. The global hybrid functional, PBE0, gives a binding energy close to that of PBE. In comparison, HSE, a range-separated hybrid functional, gives a smaller

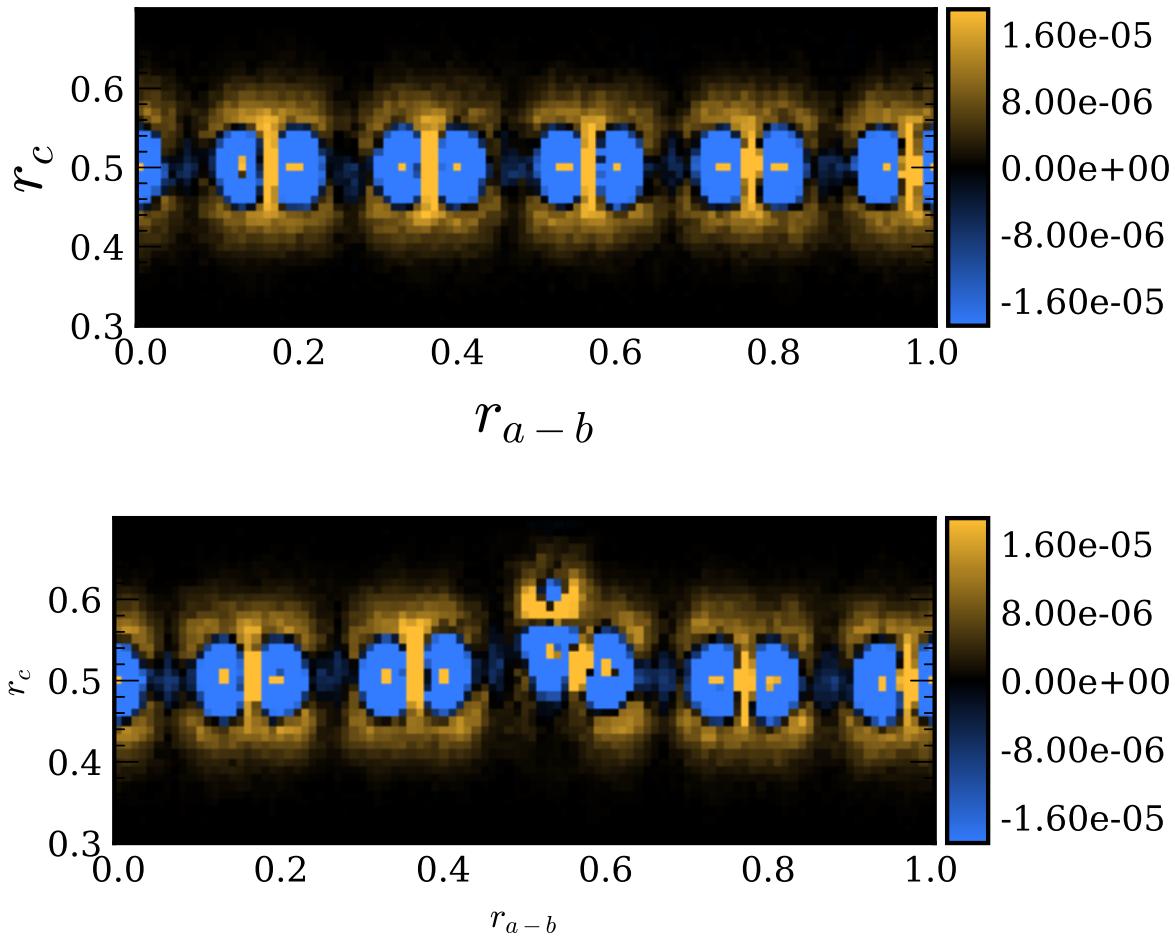


Figure 4: Visualization of the difference of PBE and DMC densities sliced along the 110 lattice plane of the unit cell for the graphene sheet, $\Delta\rho_{gr}$, (top) and H adsorbed onto graphene, $\Delta\rho_{dgr+H}$, (bottom). The abscissa represents traversing the 110 plane in fractional coordinates, while the ordinate represents traversing the c axis in fractional coordinates. Blue regions represent places where the PBE density is larger, while the gold color represents regions where the DMC density is larger.

binding energy of -743 meV, after a correction applied for the basis set incompleteness error, and is much closer to the value from DMC calculations. Interestingly, there are significant differences in the DMC and PBE charge densities of both graphene and H/graphene.

3.6 Acknowledgements

We thank Dr. Dan Sorescu for helpful discussion and for sharing the coordinates of his calculations. A.B. and H.S were supported by the U.S. Department of Energy, Office of Science, Basic Energy Sciences, Materials Sciences and Engineering Division, as part of the Computational Materials Sciences Program and Center for Predictive Simulation of Functional Materials. An award of computer time was provided by the Innovative and Novel Computational Impact on Theory and Experiment (INCITE) program. The DMC and plane wave DFT calculations used resources of the Argonne Leadership Computing Facility, which is a DOE Office of Science User Facility supported under contract DE-AC02-06CH11357. The DFT calculation using Gaussian orbitals were carried out on computing resources in the University of Pittsburgh’s Center for Research Computing. K.D.J. acknowledges NSF (CBET-2028826) for partial support of this work. S.U. was supported in part by the Pittsburgh Quantum Institute (PQI) Graduate Quantum Leader Award.

3.7 Data Availability Statement

The data that support the findings of this study are openly available on the Materials Database Facility at https://acdc.alcf.anl.gov/mdf/detail/dumi_dmc_hgraphene_v1.3, with the following DOI: 10.18126/s1wc-tya.

4.0 The Role of High-Order Electron Correlation Effects in a Model System for Non-valence Correlation-bound Anions

The test and figures in this chapter have been adapted from The Role of High-Order Electron Correlation Effects in a Model System for Non-valence Correlation-bound Anions, *J. Chem. Phys.*, **2020**, 153, 224118, DOI: 10.1063/5.0030942 with the permission of AIP. The author's contribution to the work included performing CIPSI calculations and the subsequent DMC calculations, along with performing analysis and assisting with figures, and editing/revising the manuscript.

4.1 Summary

Non-valence correlation bound anions (NVCB) are molecules in which the primary binding interaction comes from electron correlation, thus methods which capture electron correlation are necessary for an accurate description. The diffusion Monte Carlo (DMC), auxiliary field quantum Monte Carlo (AFQMC), and equation-of-motion coupled cluster (EOM-CC) methods are used to calculate the electron binding energy (EBE) of the non-valence anion state of a model $(\text{H}_2\text{O})_4$ cluster. The viability of the CIPSI procedure to select determinants to describe the binding of the electron for the NVCB system is explored. Two geometries are considered, one at which the anion is unbound and the other at which it is bound in the Hartree-Fock (HF) approximation. It is demonstrated that DMC calculations can recover from the use of a HF trial wave function that has collapsed onto a discretized continuum solution, although larger electron binding energies are obtained when using a trial wave function for the anion that provides a more realistic description of the charge distribution, and, hence, of the nodal surface. For the geometry at which the cluster has a non-valence correlation-bound anion, both the inclusion of triples in the EOM-CC method and the inclusion of supplemental diffuse d functions in the basis set are important. DMC calculations with suitable trial wave functions give EBE values in good agreement with our best estimate EOM-CC result. AFQMC using a trial wave function for the anion with a realistic electron density gives a value of the EBE nearly identical to the EOM-CC result when using the same basis set. For the

geometry at which the anion is bound in the HF approximation, the inclusion of triple excitations in the EOM-CC calculations is much less important. The best estimate EOM-CC EBE value is in good agreement with the results of DMC calculations with appropriate trial wave functions.

4.2 Introduction

In recent years, there has been growing interest in a class of anions known as non-valence correlation-bound (NVCB) anions in which long-range correlation effects are crucial for the binding of the excess electron.^{18–31} By definition, NVCB anions are unbound in the Hartree-Fock (HF) approximation. Due to their highly spatially extended charge distributions, large, flexible basis sets are required for the theoretical characterization of NVCB anions. However, with such basis sets, the wave function from Hartree-Fock (HF) calculations on the excess electron system collapses onto the neutral plus an electron in an orbital that can be viewed as a discretized representation of a continuum solution.¹⁸ Methods that start from the HF wave function including second-order Møller-Plesset perturbation theory (MP2)¹⁶⁷ or coupled-cluster singles and doubles with perturbative triples (CCSD(T))¹⁶⁸ do not recover from this collapse onto the continuum, while methods such as orbital-optimized MP2 (OOMP2)¹⁶⁹ or Bruckner coupled-cluster¹⁷⁰ can overcome this problem.¹⁸ The majority of calculations of NVCB anions have employed the equation-of-motion coupled-cluster singles and doubles (EOM-CCSD) method.¹⁷¹ Among the NVCB anions studied computationally to date are C₆₀, C₆F₆, TCNE, (NaCl)₂, Xe_n clusters, large polyaromatic hydrocarbons, and certain (H₂O)_n clusters.^{18–25,30,31}

The EOM-CCSD method displays an $\mathcal{O}(N^6)$ scaling with system size, and higher order EOM-CC methods are even more computationally demanding. As a result, most of the calculations of NVCB anions carried out to date have not been fully converged with respect to basis set or the level of excitations treated in the EOM procedure. We note, however, that by using domain-based local pair natural orbitals (DLPNO), electron affinity EOM-CCSD calculations have recently been carried out on systems described by up to 4,500 basis functions.⁸⁸

In the present work, we apply two quantum Monte Carlo (QMC) methods to the problem of calculating the electron binding energy (EBE) of the non-valence anion of a model (H₂O)₄ cluster.

The first approach considered is fixed-node diffusion Monte Carlo (DMC),^{172–175} using various single Slater determinant (SD) and multideterminant (MD) trial wave functions. DMC is a real-space method, with the major sources of error resulting from the use of finite time steps and the fixed-node approximation. The finite time step error can be largely eliminated by running calculations at different time steps and then extrapolating to the zero time step limit. The fixed-node error results from imposition of a nodal surface via a trial wave function, which is necessary to ensure Fermionic behavior, and can be addressed by a variety of means including expanding the number of Slater determinants in the trial wave function or by applying the backflow transformation.¹⁷⁶ It is important to note that, by virtue of working in real space, fixed-node DMC energies are much less sensitive to the choice of the atomic basis set than methods such as EOM-CCSD that operate in a space of Slater determinants.

The second QMC approach considered is the auxiliary field QMC (AFQMC) method.^{177–183} AFQMC calculations sample an over-complete space of nonorthogonal Slater determinants. The finite time step error can be mitigated as in DMC. The error that arises from constraining the phase of the wave function to zero can be systematically reduced by improving the trial wave function. Phaseless AFQMC is additionally subject to the limitations of the atomic basis set employed. DMC scales as $\sim \mathcal{O}(N^3)$ with system size, while AFQMC displays an $\sim \mathcal{O}(N^4)$ scaling in most implementations. One of the goals of these calculations is to determine whether DMC calculations can recover from the use of a trial wave function that has collapsed onto a discretized continuum orbital in the case of the excess electron. Additionally, we explore whether correlation effects that are missing in EOM-CCSD are important for electron binding.

In our calculations, we employ a model $(\text{H}_2\text{O})_4$ cluster that has been investigated in earlier studies by our group.^{18,25} In this model, depicted in Figure 5, the monomers are arranged so that the net dipole moment is zero. If the distance R is varied, with all other geometrical parameters held fixed, the system can be tuned from a regime (large R) that the excess electron weakly binds in the HF approximation to one (small R) at which it is not bound in the HF approximation. i.e., at which it is NVCB in nature.

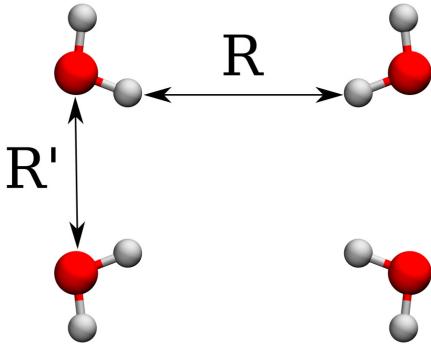


Figure 5: The model $(\text{H}_2\text{O})_4$ system considered in this study. R' held fixed at 3.46105 Å, and R is either 4 Å or 7 Å. Image generated using VMD.¹⁸⁴

4.3 Computational Details

4.3.1 EOM Coupled Cluster

The EOM methods considered in this study are EOM-MP2,¹⁸⁵ EOM-CCSD,¹⁷¹ EOM-CCSD(T)(a)*,¹⁸⁶ and EOM-CCSDT,^{187,188} listed in order of increasing sophistication in terms of treatment of correlation effects. In the EOM-MP2 and EOM-CCSD methods, the neutral molecule is treated at the MP2 and CCSD levels, respectively, and the amplitudes from these calculations are used to perform unitary transformation of the Hamiltonian. This "dressed" Hamiltonian is then used to carry out a 1-particle plus 2-particle-1-hole CI calculation on the anion. In the EOM-CCSDT method, the neutral species is first treated at the CCSDT level, and the transformed Hamiltonian is used to do CI calculation on the anion that includes up to 3-particle-2-hole configurations. The EOM-CCSD(T)(a)* method includes in an approximate manner both triple excitations in the ground state coupled cluster calculations and 3-particle-2-hole excitations in the treatment of the anion.¹⁸⁶

The main basis set used for the EOM calculations reported in this study is aug-cc-pVTZ+7s7p, formed by supplementing the aug-cc-pVTZ Gaussian-type orbital (GTO) basis set^{189,190} with a 7s7p set of diffuse functions centered at the middle of the cluster and similar to the set from Ref.¹⁸. The exponents of the supplemental functions start at 0.023622, with each successive exponent

Table 2: Dependence of the total energies and the EBE of the model $(\text{H}_2\text{O})_4$ cluster at $R = 4 \text{ \AA}$ on the supplemental diffuse basis functions. Results obtained using the EOM-CCSD method.

basis set	neutral (Ha)	anion (Ha)	EBE (meV)
aug-cc-pVTZ	-305.327947	-305.331344	92.4
aug-cc-pVTZ+1s	-305.327953	-305.332359	119.9
aug-cc-pVTZ+2s	-305.327957	-305.334226	170.6
aug-cc-pVTZ+3s	-305.327958	-305.334460	176.9
aug-cc-pVTZ+7s	-305.327958	-305.334462	177.0
aug-cc-pVTZ+7s1p	-305.327979	-305.334604	180.3
aug-cc-pVTZ+7s7p	-305.327987	-305.334622	180.6
aug-cc-pVTZ+3s1p	-305.327979	-305.334602	180.2

being smaller by a factor of 3.2. However, as seen from Table 2, the supplemental 7s7p set of diffuse functions can be truncated to 3s1p without significantly impacting the EBE as calculated at the EOM-CCSD level. Moreover, as shown in Table 3, expanding the main basis set (i.e., the non-supplemented portion) from aug-cc-pVTZ to aug-cc-pVQZ^{189,190} makes only a small impact on the EBE (4% at $R = 4 \text{ \AA}$) . In contrast, reducing the main basis set to aug-cc-pVDZ^{189,190} leads to a 14% reduction in the EBE. (These results were obtained using the EOM-MP2 method, but as seen from comparison of the results in Tables 2 and 3, using the aug-cc-pVTZ+3s1p basis set in both cases, the EBEs from the calculations with the EOM-CCSD and EOM-MP2 methods agree to within 0.5 meV.) The smaller aug-cc-pVDZ+3s1p basis will be used in the EOM-CCSDT calculations, which would have been computationally prohibitive with aug-cc-pVTZ+7s7p or aug-cc-pVTZ+3s1p basis sets. Finally, EOM-CCSD(T)(a)* calculations were carried out with aug-cc-pVTZ+3s1p3d basis sets, where the exponents of the d functions match those of the s and p functions, to assess the importance of supplemental d functions on the EBEs. The EOM calculations utilized the frozen core approximation and were carried out using the Coupled-Cluster techniques for Computational Chemistry (CFOUR) program package.^{191,192}

⁰SD/X indicates that the trial wave function employed a single Slater determinant with X (either HF or B3LYP)

Table 3: Sensitivity of the EBE of the $(\text{H}_2\text{O})_4$ model to the “core” basis set. Results obtained using the EOM-MP2 method.

	Neutral (Ha)	Anion (Ha)	EBE (meV)
$R = 4.0 \text{ \AA}$			
aug-cc-pVDZ+3s1p	-305.0371957	-305.0428558	154.0
aug-cc-pVTZ+3s1p	-305.3092869	-305.3159306	180.8
aug-cc-pVQZ+3s1p	-305.4008845	-305.4078074	188.4
$R = 7.0 \text{ \AA}$			
aug-cc-pVDZ+3s1p	-305.0383747	-305.0432259	132.0
aug-cc-pVTZ+3s1p	-305.3104923	-305.3157472	143.0
aug-cc-pVQZ+3s1p	-305.4021640	-305.4075716	147.1

4.3.2 DMC

The DMC calculations were carried out using trial wave functions represented as products of one or more Slater determinants with a Jastrow factor with one-, two-, and three-body terms.^{162,193,194} The parameters in the Jastrow factors were optimized using variational Monte Carlo (VMC), and the resulting trial wave functions were then employed in subsequent DMC calculations. Three types of SD trial wave functions were employed. These used HF orbitals, Becke-Lee-Yang-Parr (B3LYP) DFT orbitals,^{195–198} and natural orbitals (NOs) from small restricted single plus double excitation configuration interaction (SDCI) calculations designed to bind the excess electron when it is not bound in the HF approximation. In addition, DMC calculation were carried out using MD trial wave functions, with the determinants being determined either from the restricted SDCI procedure or from configuration interaction using a perturbative selection made iteratively (CIPSI) calculations.³² Details on these calculations are provided below.

To reduce the computational cost of the DMC calculations, the ccECP pseudopotentials^{137,138} were employed together with GTO basis sets that we designate as cc-pVDZ / ccECP, aug-cc-pVDZ

orbitals. When different types of trial wave functions are used for the neutral (N) and anion (A) this is indicated by the double slash.

/ ccECP, aug-cc-pVDZ / ccECP+3s1p, and aug-cc-pVDZ / ccECP+7s7p. The "core" cc-pVDZ / ccECP^{137,138} basis set was designed for use with the ccECP pseudopotentials; the "aug" indicates that the diffuse aug functions from the aug-cc-pVDZ basis sets of Dunning and co-workers are included; and the 7s7p set of diffuse functions are those described above in the Section 4.3.1.¹⁹⁰ The T-moves scheme was used to control the localization error for nonlocal pseudopotentials.¹⁹⁹

The double-zeta rather than the larger triple-zeta basis set was used as the core basis set due to the relative insensitivity of DMC calculations to the choice of the atomic basis set. For most of the DMC calculations a fixed population of 16,000 walkers and time steps of 0.001, 0.003, and 0.005 a.u. were employed, with the reported results obtained by linear extrapolation to zero time step. However, this population is much larger and the time steps much smaller than what is actually required to achieve well converged energies with minimized finite time step and fixed population errors. Indeed, DMC calculations using Hartree-Fock trial wave functions, larger time steps (specifically 0.05, 0.1, and 0.2 a.u.) and a smaller population of only 1,000 walkers produce an electron binding energy within error bars of that obtained using the smaller time steps and larger populations. Additionally, a DMC calculation with a B3LYP trial wave function with a time step of 0.05 is in agreement with the values obtained with the smaller time steps and larger populations suggesting that these parameters do not depend strongly on the choice of starting orbitals. In light of this, the 0.05 a.u. time step and smaller walker population were employed in the DMC calculations using CIPSI trial wave functions to mitigate the additional cost associated with the MD space. The VMC and DMC calculations were carried out using the QMCPACK code.^{162,163} The orbitals for the SD-based trial wave functions and the restricted SDCI MD wave function were both generated using the General Atomic and Molecular Electronic Structure System (GAMESS) program package,^{200–202} whereas the CIPSI wave functions were generated using the Quantum Package 2.0 code.³³

4.3.3 Restricted CI and CIPSI-generated Trial Wave Functions for DMC Calculations

The restricted SDCI procedure employed the HF wave function for the neutral molecule and a specially tailored SDCI wave function for the anion, which included all symmetry-allowed single and double excitations, with the latter restricted so that one of the electrons excited is from the

orbital occupied by the excess electron in the HF wave function. This approach, when used with a flexible basis, gives a bound anion. NOs were generated from the SDCI wave function of the anion and were used in a SD trial wave function for subsequent DMC calculations. In addition, the SDCI wave function itself (expanded in terms of HF orbitals) was used in MD DMC calculations on the anion for $R = 4 \text{ \AA}$. In this case, a threshold of 0.001 on the magnitude of coefficients in the CI expansion was used in choosing the retained determinants. This resulted in a wave function with 1,392 Slater determinants.

By design, the restricted SDCI wave function does not allow for change of the correlation energy of the valence electrons due to the presence of the excess electron. This possibility is allowed for in the CIPSI MD trial wave functions. The CIPSI calculations were carried out using B3LYP orbitals rather than Hartree-Fock orbitals because the former avoids the problem of collapse onto a discretized continuum solution at $R = 4 \text{ \AA}$.^{195–197} Since the CIPSI calculations have not approached the full configuration interaction limit as indicated by the second-order perturbative correction to the energy, a judicious choice of starting orbitals is required to construct a physically meaningful trial wave function. In order to generate compact wave functions for both the anion and the neutral, NOs were iteratively refined through successive CIPSI calculations, each beginning from a single determinant reference of natural orbitals from the previous iteration. For each NO-generating CIPSI calculation, approximately 100,000 determinants were retained and used to generate NOs for the next iteration, for a total of six NO generation cycles. With the determinant of resulting NOs as a reference, a final CIPSI calculation was carried out, stopping when at least 150,000 determinants were included in the variational space for the anion and at least 100,000 determinants for the neutral. The resulting determinant spaces were used as the DMC trial wave functions.

Both the restricted SDCI and the CIPSI calculations used to generate the trial wave functions for subsequent DMC calculations were carried out using the ccECP pseudopotentials. The aug-cc-pVDZ/ccECP+7s7p and aug-cc-pVDZ/ccECP+3s1p basis sets were used for the CIPSI and CIPSI calculations, respectively.

4.3.4 AFQMC

AFQMC^{177–183} utilizes the Hubbard-Stratonovich transformation²⁰³ to represent the imaginary-time propagator as a multi-dimensional integral over auxiliary-fields. Ground-state properties are sampled from a random walk in the space of non-orthogonal Slater determinants subject to the phaseless constraint¹⁸⁰ introducing a bias which can be systematically reduced based on the quality of the nodal surface of the trial wave function employed. While sophisticated trial wave functions generated from regularized orbital-optimized MP2 (κ -OOMP2)²⁰⁴ or complete active space self-consistent field (CASSCF)^{205–207} are required to obtain quantitative predictions for some bi-radicaloids and transition metals, high accuracy has been obtained, even for systems exhibiting non-trivial electron correlation such as dipole-bound anions,²⁰⁸ with single-determinant trial wave functions consisting of HF or Kohn-Sham orbitals.^{208,209}

In this work, we perform calculations with a graphics processing unit (GPU) implementation of AFQMC,²¹⁰ utilizing single-precision floating-point arithmetic and two-electron integrals decomposed via a modified Cholesky decomposition (10^{-5} cutoff).²¹¹ These calculations made use of the aug-cc-pVTZ+7s7p basis set, a small imaginary-time step of 0.005 a.u., and correlated all electrons. For the neutral species and electrostatically bound anion ($R = 7 \text{ \AA}$), the Hartree-Fock wave function was used as the trial wave function. For the NVCB anionic species ($R = 4 \text{ \AA}$), a SD trial wave function comprised of natural orbitals from the restricted SDCI calculation as detailed in Section 4.3.3 (but now carried out without pseudopotentials) was used.

4.3.5 Radial Orbital Densities

To compare the description of the charge distribution of the excess electron as calculated using different theoretical methods, we generate radial electron density plots. This choice is motivated by the fact that the excess electron occupies an orbital belonging to the totally symmetric representation. The radial electron densities are generated by numerically integrating over the angular components of the singly occupied molecular or natural orbital. First, Molden files are created from the output data from the various generating programs using cclib when supported.²¹² With the Molden files as input, PySCF is used to generate the electron density on a uniform radial grid and 5810 point Lebedev-Laikov angular grid as tabulated in quadpy.^{213–216} Finally, a numerical

integration is performed over the angular components. An example of this workflow is presented in detail in the Supplementary Information.

4.4 Results and Discussion

The EBEs obtained from the EOM and AFQMC calculations are summarized in Table 4, and the results from the various DMC calculations are summarized in Table 5. We consider first the results obtained for $R = 4 \text{ \AA}$, for which HF calculations do not bind the excess electron.

4.4.1 Results for $R = 4 \text{ \AA}$: the Correlation Bound Region

From Table 4, it is seen that the EOM-CCSD/aug-cc-pVTZ+7s7p calculations give a value of the EBE of 181 meV for the $(\text{H}_2\text{O})_4$ cluster model at $R = 4 \text{ \AA}$. This increases to 196 meV with the EOM-CCSD(T)(a)^{*} method. The AFQMC calculations using the same basis set and for the anion a single determinant of NOs from the restricted SDCI calculation for the trial wave function produce an EBE value of 194 ± 10 meV, comparable to the EOM-CCSD(T)(a)^{*} result. The EOM-CCSD(T)(a)^{*} and EOM-CCSD(T) EBE values calculated with this basis set are nearly identical, demonstrating that the approximate treatment of triples in the former procedure introduces a negligible error in the EBE. The contribution of supplemental diffuse functions was checked using the EOM-CCSD(T)(a)^{*} method and the aug-cc-pVTZ+3s1p2d basis set. These calculations reveal that the inclusion of the supplemental diffuse d functions leads to a ~ 10 meV increase in the EBE. With the inclusion of this correction, we obtain an estimated EOM-CCSDT EBE of 212 meV. It is expected that the inclusion of the supplemental d functions in the basis set used for the AFQMC calculations would lead to a similar increase in the EBE obtained using that method.

The restricted SDCI procedure, by itself, is not expected to give an accurate value of the EBE and is designed to generate appropriate trial wave functions for DMC or AFQMC calculations on the anion. In fact, the EBE resulting from the HF treatment of the neutral and the restricted SDCI treatment of the anion using the aug-cc-pVTZ+7s7p basis set is 345 meV, appreciably larger than the EOM and AFQMC values. This over-binding is due in part to the fact that the restricted SDCI

wave function, like the HF wave function, overestimates the magnitude of the dipole moment of the water molecules, resulting in a too favorable electrostatic interaction. We also constructed a single determinant trial wave function for the anion using the natural orbitals of the restricted SDCI expansion. We note also that the single determinant of NOs generated from the restricted SDCI wave function and using the aug-cc-pVTZ+7s7p basis set places the anion 160 meV above the neutral when the latter is treated in the HF approximation. This is not surprising since this calculation neglects correlation effects other than those incorporated in the determination of the orbitals. What is important is that the approaches based on the restricted SDCI procedure provide a realistic description of the orbital occupied by the excess electron and avoid the collapse onto the discretized continuum as was observed with the HOMO in the HF calculations.

In light of the close agreement between the EOM-CCSD(T)(a)^{*} and AFQMC values of the EBE of the (H₂O)₄ model at R = 4 Å, when using a comparable basis sets in the two approaches it is relevant to determine whether DMC calculations with sufficiently flexible trial wave functions give an EBE close to the AFQMC and EOM values consistent with these results. DMC calculations using HF trial wave functions together with the aug-cc-pVDZ/ccECP+7s7p basis set give an EBE of 183 ± 10 meV, appreciably smaller than the EOM-CCSD(T)(a)^{*} and AFQMC values. Interestingly, essentially the same EBE is obtained from the DMC calculations using a Slater determinant of HF orbitals expanded in the aug-cc-pVDZ/ccECP basis set without the 7s7p supplemental set of diffuse functions. However, if the aug diffuse functions are also removed, the DMC calculations fail to bind the excess electron. We believe that this is a consequence of the fact that with the cc-pVDZ basis set there is a near zero probability of sampling regions of space at large distances from the molecule, which are important for describing the charge distribution of the excess electron.

A significantly larger value of the EBE is obtained from SD DMC calculations using B3LYP orbitals in place of HF orbitals. The resulting EBE of 212 ± 11 meV, within statistical error, agrees with the EOM-CCSD(T)(a)^{*} and AFQMC values. A similar value of the EBE is obtained from DMC calculations using a single determinant of HF orbitals for the neutral cluster and a single

¹SD/X indicates that the trial wave function employed a single Slater determinant with X (either HF or B3LYP) orbitals. When different types of trial wave functions are used for the neutral (N) and anion (A) this is indicated by the double slash.

²SD/X indicates that the trial wave function employed a single Slater determinant with X (either HF or B3LYP) orbitals. When different types of trial wave functions are used for the neutral (N) and anion (A) this is indicated by the double slash.

Table 4: EBEs of the $(\text{H}_2\text{O})_4$ model calculated using HF, EOM, and AFQMC methods and employing the aug-cc-pVTZ+7s7p basis set.

Method	EBE (meV)
$R = 4.0 \text{ \AA}$	
HF	-0.4
EOM-CCSD	180.6
EOM-CCSD(T)(a)*	195.8
EOM-CCSDT	197.5 ¹ (212.0) ²
AFQMC SD/HF(N)//SD/NO SDCI(A)	194 \pm 10
$R = 7.0 \text{ \AA}$	
HF	41.3
EOM-CCSD	140.2
EOM-CCSD(T)(a)*	141.7
EOM-CCSDT	143.3 ¹ (154.2) ²
AFQMC SD/HF	181 \pm 5

¹ This EOM-CCSDT/aug-cc-pVTZ+7s7p value was estimated by adding the difference of EBEs from the EOM-CCSD(T)(a)* and EOM-CCSDT calculations with the aug-cc-pVDZ+3s1p basis set to the value from EOM-CCSD(T)(a)*/aug-cc-pVTZ+7s7p.

² The EOM-CCSDT/aug-cc-pVTZ+7s7p3d value was estimated by adding the difference between the EBEs calculated with the EOM-CCSD(T)(a)* with the aug-cc-pVTZ+3s1p and aug-cc-pVTZ+3s1p3d basis sets to the EOM-CCSDT/aug-cc-pVTZ+7s7p estimated value in footnote [1] to assess the effect of incorporating diffuse d functions into the basis.

determinant of natural orbitals from the restricted SDCI procedure described in Section 4.3.3 for the anion. DMC calculations using a SD of HF orbitals for trial wave function of the neutral and

Table 5: EBEs of the $(\text{H}_2\text{O})_4$ model calculated using the DMC method and various trial wave functions²

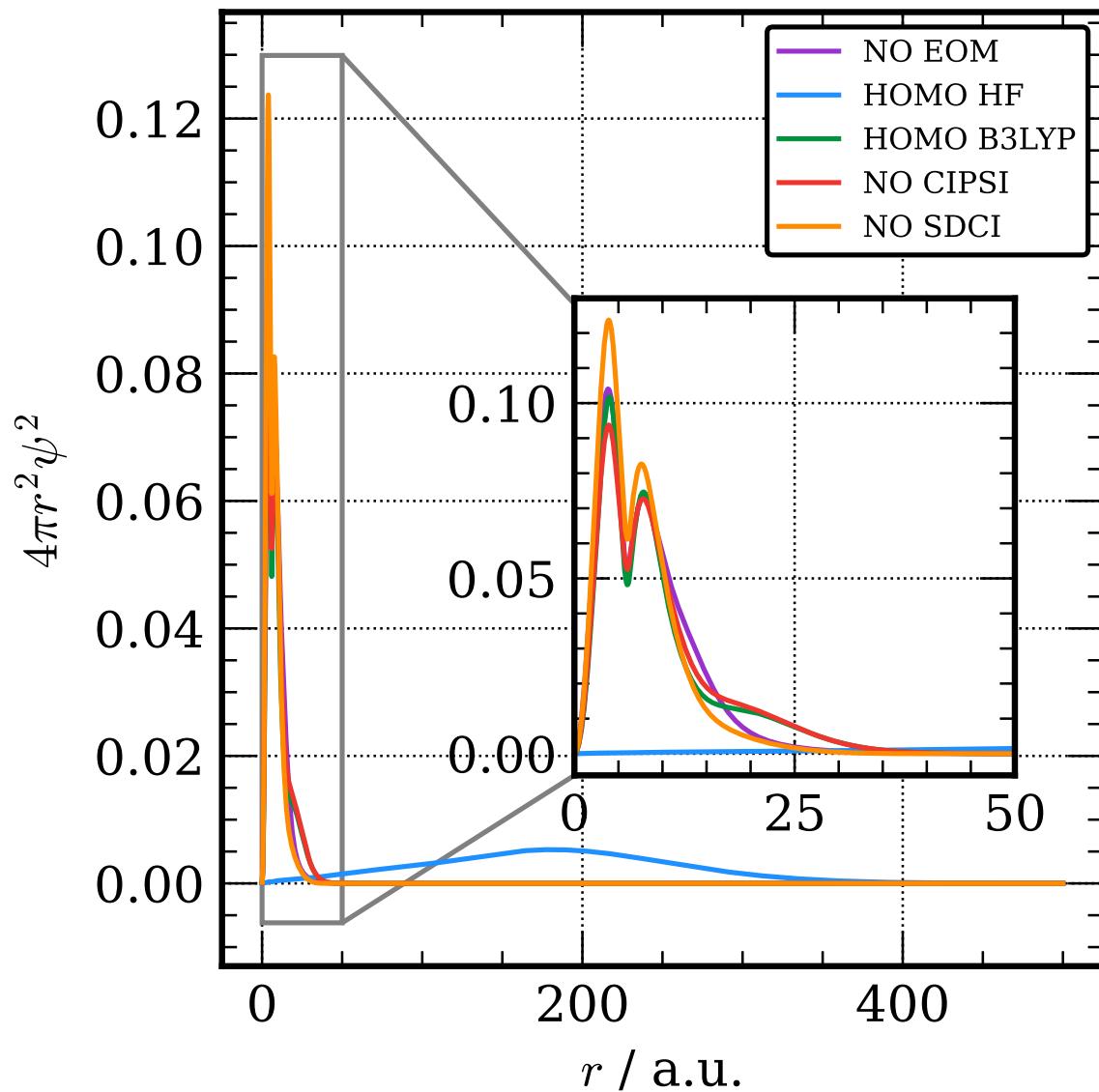
wave function	basis set	EBE (meV)
$R = 4.0 \text{ \AA}$		
SD/HF	aug-cc-pVDZ+7s7p	183 ± 10
SD/HF	aug-cc-pVDZ	176 ± 12
SD/HF	cc-pVDZ	-528 ± 25
SD/B3LYP	aug-cc-pVDZ+7s7p	212 ± 11
SD/HF(N)//SD/NO SDCI(A)	aug-cc-pVDZ+7s7p	205 ± 10
SD/HF(N)//MD/NO SDCI(A)	aug-cc-pVDZ+7s7p	202 ± 12
MD/CIPSI NO	aug-cc-pVDZ+3s1p	190 ± 9
$R = 7.0 \text{ \AA}$		
SD/HF	aug-cc-pVDZ+7s7p	141 ± 14
SD/B3LYP	aug-cc-pVDZ+7s7p	164 ± 9
SD/HF(N)//SD/NO SDCI(A)	aug-cc-pVDZ+7s7p	160 ± 9
MD/CIPSI NO	aug-cc-pVDZ+3s1p	159 ± 8

a trial wave function for the anion retaining 1,392 of the most important determinants from the restricted SDCI calculation gives an EBE of 202 ± 12 meV, close to the values obtained using the single determinants B3LYP orbitals or of NOs from the SDCI calculation (for the anion). The DMC value of the EBE resulting from the anionic trial wave function using a SD of NOs from the restricted SDCI MD calculation results is 205 ± 10 meV, similar to that from DMC calculations using as trial wave functions the MD restricted SDCI wave function for the anion and the HF wave function for the neutral.

Figure 6 compares the radial charge distributions of the singly occupied orbital from the HF and B3LYP calculations on the excess electron system as well as of the NOs associated with the excess electron from EOM-CCSD, restricted SDCI and CIPSI calculations. The collapse of the singly occupied orbital from the HF calculations onto a discretized continuum orbital is readily apparent. In contrast, the NOs from the EOM-CCSD and restricted SDCI calculations and the singly occupied orbital from the B3LYP calculation on the anion are more localized and are qualitatively similar to one another. These results are consistent with the nodal surface for the anion being significantly improved when using a SD trial wave function that has a physically reasonable charge distribution for the orbital occupied by the excess electron. Thus, although DMC calculations do recover from the collapse of the HF trial wave function onto a discretized continuum solution in the case of the anion, starting with such a trial function leads to a greater nodal surface error for the anion than for the neutral cluster. However, we also note that the radial distribution function of the singly occupied orbital from the B3LYP calculation on the anion has a spurious peak near 25 atomic units from the center of the cluster. This is likely a consequence of the self-interaction error in the B3LYP functional. The relevant NO extracted from the CIPSI calculations, which were carried using B3LYP orbitals, exhibits a similar shoulder.

Our final set of DMC calculations at $R = 4 \text{ \AA}$ used MD trial wave functions determined from CIPSI calculations for the neutral and anionic clusters. The strategy used in performing the CIPSI calculations was presented in Section 4.3.3, where it was noted that these calculations, unlike those with the restricted SDCI wave functions, allow for the correlation between the valence electrons change due to the presence of the excess electron. The DMC calculations using the CIPSI trial wave function resulted in an EBE of 190 ± 9 meV for $R = 4 \text{ \AA}$, slightly under-binding compared to the single determinant DMC value of the EBE obtained using B3LYP orbitals though in close

Figure 6: Radially integrated charge densities of the singly occupied orbitals from HF and B3LYP calculations and the singly occupied natural orbital from EOM-CCSD, SDCI, and CIPSI calculations of the model ($(\text{H}_2\text{O})_4$) cluster anion at $R = 4 \text{ \AA}$. All plots generated using Matplotlib.¹⁶⁵



agreement with the results of DMC calculations carried out with the restricted SDCI trial wave function.

4.4.2 Results for $R = 7 \text{ \AA}$: the Electrostatically Bound Region

We now consider the results obtained for the $(\text{H}_2\text{O})_4$ cluster model at $R = 7 \text{ \AA}$, for which HF calculations with the aug-cc-pVTZ+7s7p basis set bind the excess electron by 41 meV. In this case, the EOM-CCSD and EOM-CCSD(T)(a)* calculations give EBEs of 140 meV and 142 meV, respectively. Thus unlike the situation for $R = 4 \text{ \AA}$, the inclusion of triples in the EOM-CC procedure is relatively unimportant at $R = 7 \text{ \AA}$. The DMC calculations using SD HF trial wave functions give an EBE of 141 ± 14 meV, while the DMC calculations using as trial wave functions single determinants of B3LYP orbitals, single determinants generated using the restricted SDCI procedure, or MD trial wavefunctions generated using the CIPSI procedure give similar EBEs values ranging from 159 ± 8 to 164 ± 9 meV.

Since the anion is bound in the HF approximation at $R = 7 \text{ \AA}$, we also were able to calculate EBEs using separate, frozen-core coupled-cluster calculations for the neutral and anion with the following coupled-cluster methods: coupled-cluster singles, doubles, and a perturbative treatment of triples $\Delta\text{CCSD}(\text{T})^{168}$, coupled-cluster singles, doubles, and triples ($\Delta\text{CCSDT})^{217-220}$, and CCSDT with the perturbative treatment of quadruple excitations ($\Delta\text{CCSDT}(\text{Q})^{221}$) methods. The Δ indicates that the EBE is derived from the energy difference between the separate calculations on the neutral and anion. The ΔCCSDT and $\Delta\text{CCSDT}(\text{Q})$ calculations were carried out with only the aug-cc-pVDZ+3s1p basis set. These calculations indicate that full treatment of the triples, and even approximate treatment of the quadruple excitation contributions, has less than a 1 meV effect on the EBE of the $(\text{H}_2\text{O})_4$ cluster model at $R = 7.0 \text{ \AA}$. On the other hand, the inclusion of diffuse d function in the supplemental set of functions leads to a 12 meV increase in the EBE. With this correction we obtain an estimated EOM-CCSDT EBE of 154 meV, which is in good agreement with the DMC results using suitable trial wave functions.

The AFQMC calculations give an EBE of 181 ± 5 meV, significantly larger than the EOM-CC results or DMC values. This most likely reflects an inadequacy of the HF wave function used for the anion in the AFQMC calculations. Support for this interpretation is provided by examination

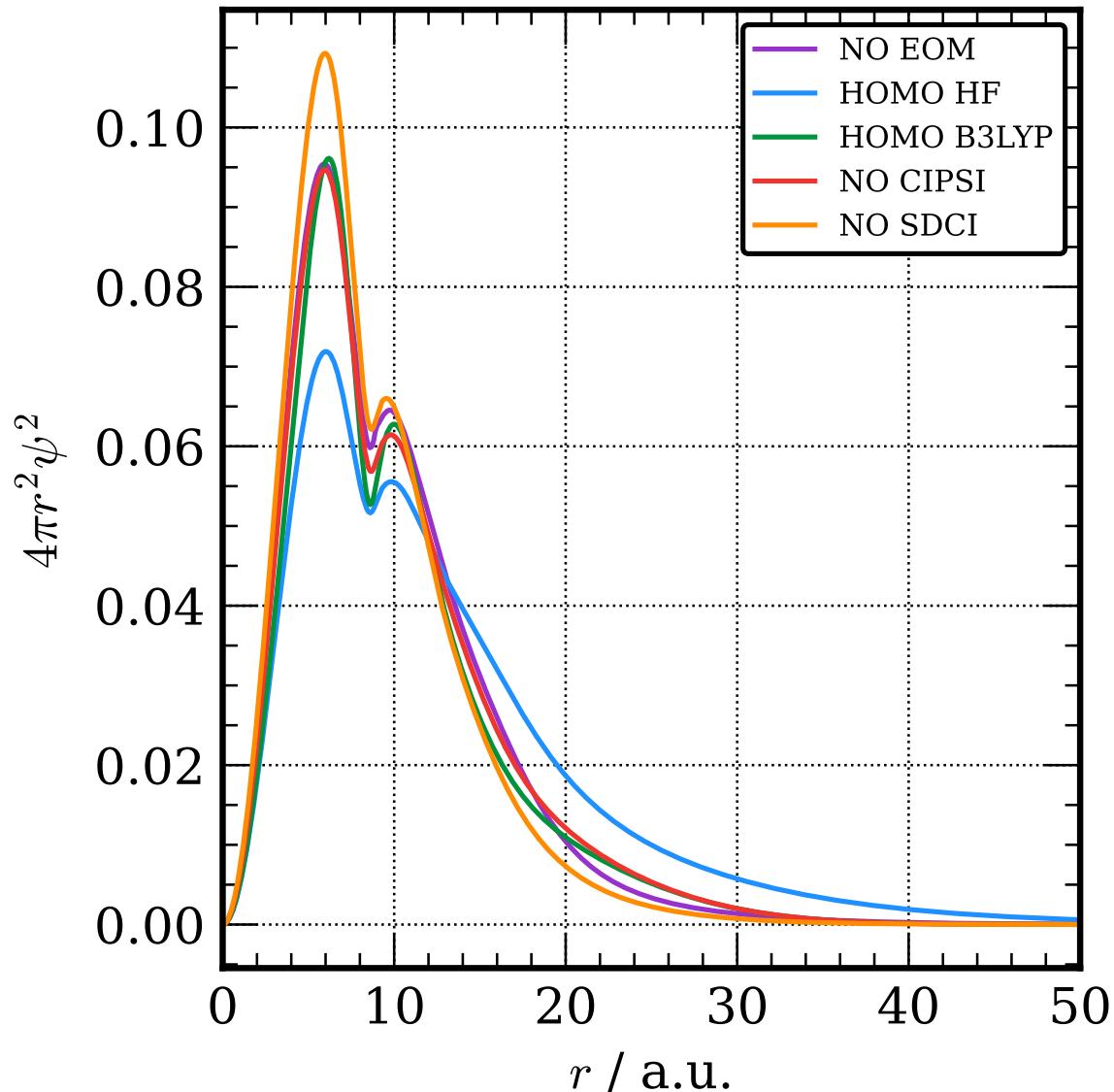
of Figure 7, which shows the radial charge distribution of the excess electron for the $(\text{H}_2\text{O})_4$ model at $R = 7 \text{ \AA}$. From this figure it is seen that although that the HF wave function has not collapsed onto the continuum as it did in the $R = 4 \text{ \AA}$ cluster, it is still much more diffuse than that from calculations that include correlation effects. It is also seen from comparisons of Figures 6 and 7 that the charge distribution associated with the NO occupied by the excess electron in the EOM-CCSD calculations for the cluster with $R = 7 \text{ \AA}$, is more radially extended than that at $R = 4 \text{ \AA}$. Another noticeable difference between the charge density plots for $R = 7 \text{ \AA}$ and 4 \AA is the reduction of the long-range shoulder in the radial charge distribution of the HOMO from the B3LYP calculations on the anion and in the relevant NO from the CIPSI calculations on the anion carried out using B3LYP orbitals, suggesting that self-interaction errors are less problematical at $R = 7 \text{ \AA}$.

4.5 Conclusion

In this study we have applied various EOM-CC methods and two different quantum Monte Carlo methods to calculate the EBE of a model $(\text{H}_2\text{O})_4$ cluster at two geometries, one at which the anion is bound in the HF approximation and the other at which it is not. Diffusion Monte Carlo calculations using single determinant trial functions based on Hartree-Fock orbitals are shown to bind the excess electron even when the initial wave function for the anion has collapsed onto the neutral plus discretized continuum orbital. However, such calculations significantly underestimate the EBE, whereas SD DMC calculations using trial wave functions for the anion with a more realistic charge distribution for the excess electron give larger EBE values that are in close agreement with our best estimate EOM-CCSDT values for both geometries considered.

For $R = 4 \text{ \AA}$, at which the anion is correlation bound, use of such trial wave functions accurately reflecting the physical charge density resulted in AFQMC-predicted EBE values in agreement with the EOM-CCSD(T)(a)* result (when using comparable basis sets). However, at $R = 7 \text{ \AA}$, AFQMC calculations with HF trial wave functions significantly overestimate the EBE compared to EOM-CC and DMC values, suggesting the need for an improved trial wave functions in this case. For the $(\text{H}_2\text{O})_4$ model system, the restricted SDCI represents an economical way to create trial wave functions for QMC calculations on non-valence anions that are not bound in the Hartree-Fock

Figure 7: Radially integrated charge densities of the singly occupied orbitals from HF and B3LYP calculations and the singly occupied natural orbital from EOM-CCSD, restricted SDCI, and CIPSI calculations of the model ($(\text{H}_2\text{O})_4$) cluster anion at $R = 7 \text{ \AA}$.



approximation. However, it remains to be seen if this strategy will be as effective for systems in which the neutral species is more strongly correlated than the model $(\text{H}_2\text{O})_4$ cluster.

Finally, we note that at $R = 4 \text{ \AA}$, for which the anion is NVCB in nature, the most frequently used method to characterize such anions, EOM-CCSD, underestimates the EBE by about 10% compared to the result of EOM-CCS DT calculations. Both DMC and AFQMC are viable alternatives to high order EOM methods, and while more computationally demanding for the $(\text{H}_2\text{O})_4$ cluster, they demonstrate lower scaling with system size than EOM methods, making them attractive for the characterization of non-valence anions of much larger systems.

5.0 Characterizing unsupervised learning methods ability to determine important chemical domains.

The author's contribution to the work included development of descriptors, writing of code, running calculations, performing analysis, and writing/editing the manuscript.

5.1 Summary

In order to reduce the computational cost of large calculations, we look to fragmenting large molecules into smaller subsystems which can best represent the whole. This work proposes a scheme for automatic molecular fragmentation through unsupervised learning approaches in which each fragment is chosen to best retain the important features of the bonding environment. This chapter highlights the efforts at benchmarking the performance of our proposed method on a set of test systems. To this end, a set of clustering algorithms (spectral, agglomerative, k -means, and affinity propagation) were studied in combination with various molecular representations, including those incorporating bonding information derived from quantum mechanics. The performance for the clustering/descriptor combinations was assessed for test systems spanning a range from easily distinguishable fragments such as non-covalently bound water cluster to oligomers in which lowest-loss fragmentation is ambiguous. Overall, it is found that spectral clustering works well in all systems tested, showing very little sensitivity to the representation employed. Spectral, agglomerative, and k -means clustering produce reasonable fragments for systems with clear fragmentation patterns. Though in the oligomer system, spectral clustering achieves the best performance as assessed by offering a trade-off between lowest error and highest speed-up and is thus recommended as the most robust clustering approach for molecular systems. The approach has the potential to improve reproducibility and transferability by replacing manual fragmentation with quantitative partitioning criteria.

5.2 Introduction

Common quantum chemistry methods that provide an accurate description of molecules are often restricted to small systems in terms of number of atoms or basis set size due to high scaling of the computational cost with system size, N . For example, coupled-cluster singles doubles with perturbative triples methods is often regarded as the “gold standard” level of theory, but incurs an N^7 scaling. Full configuration interaction (FCI), which is formally exact in a complete basis, has prohibitive $N!$ scaling. The higher scaling restricts the application of these methods to small systems while many of the chemical processes of interest involve large molecules. To compensate for this high scaling, fragmentation approaches estimate the energy and other properties of large molecular systems by partitioning the system into small subsystems, where the final estimate of the energy becomes the accumulation of the parts. The accuracy of this approach hinges on the electronic structure treatment of each fragment, the approach used to describe the interaction between the fragments, and the way in which the molecule is partitioned. The possible inter- and intra-fragment treatment approaches are vast, but beyond the scope of the current article, though interested readers are directed towards a number of helpful reviews.^{40,106} In the fragmentation schemes, the best case scenario for scaling becomes $\mathcal{O}(N^p) \rightarrow N_{frag} \mathcal{O}(f^p)$, where N_{frag} is the number of fragments, f is representative of the fragment size, and p is the exponential value dependent on the level of electronic structure theory utilized.^{40,41} This partitioning of a single, very costly calculation into N_{frag} smaller calculations achieves two important objectives: 1) Computational scaling with system size is reduced with reasonable fragment definition and 2) trivial parallelization is possible by treating subsystems separately, with the potential to efficiently utilize high performance computing resources. In addition to enabling the treatment of larger systems, fragmentation methods can provide detailed insight into interfragment interactions when combined with analysis techniques such as energy decomposition analysis or symmetry adapted perturbation theory.^{100–104} Fragmentation approaches can also assist recent efforts of enabling quantum chemistry via quantum computers. Current quantum computation is restricted as the integrity of the results can be sacrificed due to interactions of the hardware with the environment. Fragmentation methods have been suggested as ways to treat larger molecular systems on quantum computers. By breaking molecules into smaller domains before treatment on a quantum computer, the most important chemical interactions can

be efficiently described without succumbing to errors resulting from the quantum computing hardware.²²²

As pointed out by Herbert, the choice of fragments for a system is not well-defined, but affects the quality of results obtained.⁴⁰ In some systems, a natural approach towards partitioning arises when there is a stark difference in the types of bonding present in the system, such as in non-covalent molecular clusters. For covalent systems, however, the choice of fragments is not always as clear-cut. In such cases, fragmentation requires the comparison of total energies, dipoles, or polarizabilities. Often the fragment definition is based on predefined functional groups or chosen manually.⁴² Some methods of energy estimation are defined based on specific fragmentation schemes such as the systematic molecular fragmentation (SMF) and systematic molecular fragmentation by annihilation (SMFA).^{43–47} In these methods, fragments are built around functional groups or larger fragments made from their groupings. Ultimately, the level of fragmentation is at the discretion of the user to achieve the desired level of accuracy. However, functional group definition may become ambiguous. For example, there is no set number of monomers to include from a polymer backbone to acquire an accurate description capture the chemical behavior. Additionally, a fragment definition based only on functional groups may not consider the interacting chemical environment.

A desirable approach to choosing fragments would have low computational cost and prioritize keeping associated molecular components intact to treat fully with quantum mechanics while the estimation of their interactions should occur only at the most weakly bound points. *To this end, we propose an approach utilizing clustering methods to identify the strongly interacting substructures of the system, which we term Automatic Fragmentation of Molecules using Clustering (AFMC) approach.* Clustering methods are a form of unsupervised machine learning used to identify substructures in data sets, as a result these approaches are fundamental to data-mining procedures. Previous work in chemistry utilized clustering methods to identify structure-property relationships in large databases⁴⁸, to determine the number of residues to treat in quantum refinement methods⁴⁹, and to partition large proteins into peptides using an amino acids representation using graph based methods.⁵⁰ The application of clustering methods to produce logical fragments of individual molecules at an atomic level is an unexplored direction. The motivation is that these unsupervised machine learning algorithms (UML), given a certain level of molecular information, will be able

to group the atoms interacting most strongly with each other, ensuring that segmentation occurs between atoms which are the most weakly connected. This approach is able to operate independent from functional group definition which will become useful for capturing non-covalent interactions and non-local interactions in materials or biomolecules. Additionally, this approach is expected to overcome shortcomings in other fragmentation approaches such as severing of double bonds or ring structures, since the molecular representation should be designed to avoid this.

The article is organized as follows, Section 5.3.1 will discuss the clustering approaches used, Section 5.3.2 will describe the representation of chemical data used as input for the clustering approach, with validation methods described in Section 5.3.3, results and a discussion on clustering performance follow in Section 5.4.

5.3 Methods

In UML approaches, the clustering depends fundamentally on two factors: the features used to describe each data point and the algorithm used to identify domains within the data. This section introduces UML (clustering) approaches along with molecular representations (features) upon which fragment selection is based. Additionally, metrics are discussed to assess the quality of a chosen partitioning and molecular test systems are presented. Four representations are explored in this work and are described in Section 5.3.2. Several clustering approaches selected from a range of different families of algorithm were considered. The main article focuses on agglomerative, k -means, and spectral clustering; additional algorithms tested can be found in the supporting information. These clustering methods are described in Section 5.3.1. The code used to explore automatic molecular fragmentation using clustering can be found at <https://github.com/amandadumi/molfrag>. Clustering approaches are utilizing the implementations contained within the scikit-learn Python package.²²³

5.3.1 Clustering approaches

This subsection describes the clustering methods explored in this work. This work explores the application of three clustering algorithms, representing different approaches for the selection of fragments: agglomerative, k -means, and spectral clustering.^{224–228} At a minimal level, all chosen algorithms require only one user input: the number of subsets to identify within the data, i.e. the desired number of fragments. Starting from this user input, the clustering approaches automatically determine fragments based on selected molecular representations, as outlined below.

Agglomerative clustering²²⁹ is performed in a bottom-up fashion where in each iteration the most similar clusters are merged. In the initial iteration, all of the atoms are regarded as individual clusters, which are then merged into larger clusters. Merging occurs according to a linkage criterion which describes the similarity (or distance) between clusters. In this work, the Ward linkage criterion is used, which chooses which clusters so that the variance of Euclidean distances within each cluster is minimized. Here, the variance of a cluster is calculated as the residual sum of squares (RSS) of all variables in a cluster C ,

$$d_C = \sum_{i,j \in C} ||x_i - x_j||^2, \quad (69)$$

where i, j are observations (atoms) within a cluster with associated data points x . The Ward criterion²²⁹ results in a more regular distribution of cluster sizes compared to other choices, which is advantageous for the speed in a fragment calculation, as even distribution of fragment sizes is associated with equal computational cost distribution. This process is repeated until the user-specified number of fragments has been obtained. A visualization of this algorithm is shown in Figure 8, where the progression of this algorithm can be understood as a dendrogram. The dots are grouped through iterations until the desired number of clusters are found. The time complexity of the agglomerative clustering method is $\mathcal{O}(n^3)$, but can be reduced to $\mathcal{O}(n^2)$ with various optimizations, where n is the number of data points.²²⁴ A shortcoming of this method is that it is a greedy algorithm, or in other words once an atom is assigned to a fragment only considering the local environment and this assignment will not be reassessed.

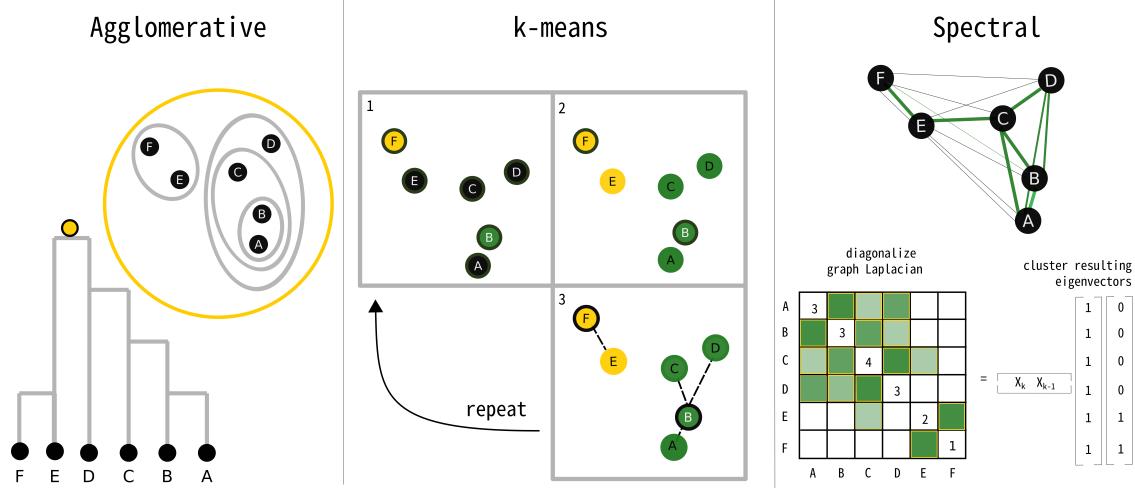


Figure 8: Visual representations of the three clustering approach explored in this work. The dots are a minimal example of a set of data. Agglomerative clustering (left), k -means clustering (center), and spectral clustering (right). Detailed descriptions of each method are found in the text.

The k -means clustering approach which iteratively minimizes the distance between the data points x_i and the cluster centers \bar{x}_c according to:

$$\text{Minimize } \left\{ J \equiv \sum_{c=1}^k \sum_{i=1}^{n_c} \|x_i^{(c)} - \bar{x}_c\|^2 \right\} \quad (70)$$

where c denotes a particular cluster, $x_i^{(c)}$ a data point within this cluster, and \bar{x}_c is the centroid associated with cluster c . k is the number of fragments requested by the user. During the first iteration, the centroids are chosen randomly among the data points. Data points are then assigned to different centroids (clusters) based on shortest distances; centroids are updated by averaging over the data points associated with a given cluster,

$$\bar{x}_c = \sum_{j \in c} x_j^{(c)}. \quad (71)$$

This procedure of assigning points based on proximity to the centroids and update of centroids is repeated iteratively until changes between centroids falls below a threshold. An illustration of this approach is given in Figure 8.

It is interesting to note some differences and similarities between agglomerative and k -means clustering. First, it is noted that at first glance the sum of squares used as objective function within k -means clustering is related to Ward's linkage criteria used within agglomerative clustering. However, an important difference is that the k -means approach is non-greedy, meaning that a reassignment of points between different clusters is possible between different iterations. This property implies that k -means has the potential to generate a global solution, whereas agglomerative clustering aims at a locally optimal solution. Another interesting difference is that k -means clustering is an iterative procedure, meaning that termination and thus runtime depend on the specific data, whereas agglomerative clustering is guaranteed to terminate after a certain number of steps when the user-requested number of clusters has been identified. The convergence of k -means clustering implementations, specifically of Lloyd's algorithm²²⁶, varies significantly between average-case and worst-case scenarios. For typical applications, k -means is observed to converge in few iterations, leading to an average-case observed complexity that is linear, $O(n)$, where we use O to denote an approximate, observed scaling. In practice, k -means is therefore found to be more readily applicable to larger data sets compared to agglomerative clustering. However, we do note that under worst-case scenarios the number of iterations required for k -means convergence scales in a superpolynomial fashion, leading to lower bounds of $2^{\Omega(\sqrt{n})}$ for the worst-case computational scaling, where Ω denotes a lower bound.²³⁰ Lastly, we note that the k -means procedure is typically repeated for different (random) choices of initial centroids to ensure convergence to a global minimum, since initial centroids may impact the convergence to a specific solution. For this reason, this work reports k -means results over ten runs with different initial centroids for each calculation. In contrast, agglomerative clustering always produces the same results for a given set of input data except in the case of identical cluster distances; in such a case the final clusters would depend on the sorting of the input data.²³¹ For molecular representations using floating point numbers to map molecular structure, it is unlikely to encounter exact degeneracies. However, some of the molecular representations rounding these representations to integers, it would be possible to encounter degeneracies and thus a dependency of agglomerative clustering results on sorting of the input data.

The third clustering algorithm explored here is spectral clustering.^{10.1093/imaiai/iay008, 232,233 A} A visualization of this method can be seen in Figure 8. Spectral clustering begins by defining an

affinity matrix A_{ij} which describes the similarity of each pair of data points, ij . In the present context, the similarity between two data points can be understood as the presence of a bond or the strength of the bond. From the affinity matrix, a degree matrix, D , is built, which sums the rows of A onto the diagonal of D . The off-diagonal elements of A combined with D are combined to form the graph Laplacian of the data as $L = D - A$. The eigenvectors and eigenvalues resulting from the diagonalization of L represent the data in a lower dimension space that leads to clearer separation for linear cuts. Following the spectral decomposition, the eigenvalues are then clustered by another method such as k -means, a discretized approach, or others.^{232–234}. In this work, the k -means approach is utilized. The time complexity of the spectral clustering method is $\mathcal{O}(n^3)$, where n is the number of data points.

Additionally, the application of affinity propagation and mean shift clustering were explored, but were unsuccessful in producing useful molecular fragments. These methods are of interest since the number of clusters are chosen automatically through an a variety of approaches to analyze the density of a given set of data. However, these methods did not produce consistent results with the molecular representations explored in this work and, in many cases, no viable subsystems resulted from the fragmentation for the representations explored in this work. Although affinity propagation had some success, it was not consistent across test sets. Results for these two additional methods are included in the SI. These finding do not rule out the possible use of these clustering methods for chemical fragmentation as a tailored the molecular representation for this method may be needed.

5.3.2 Molecular Representations

This subsection describes the molecular representation used to describe the chemical system. The success of clustering depends on the representation, i.e. the features used to describe the relationship between the data points (here atoms). This work explores the application of four different representations. Two descriptors are derived from structure information alone: Cartesian-based and a covalent radii based bond matrix-based descriptors. The structure-only derived representation provide a low-cost descriptor as no quantum mechanical information is incorporated. Alternatively, incorporating bonding information from a quantum mechanical treatment should provide a

more detailed descriptor, though at a higher computational cost. Two descriptors which incorporate quantum mechanic information are explored: the Mayer bond matrix descriptor and the rounded Mayer bond matrix descriptor.

The descriptors are presented in a way that they represent an affinity or similarity matrix between objects. The clustering methods utilize this information in different ways. Spectral clustering uses the affinity matrix to perform the subspace search, agglomerative clustering will invert the affinity matrix to indicate that those data points more strongly interacting are closer in space, and the k -mean approach will use each row of the affinity matrix as a description of the dimensions in which the vector norm is measured. The representation is the reciprocal of each element which is handled automatically within the `molfrag` code.

The Cartesian representation describes the position of each atom as x , y , and z components. Distances between any two atoms, A and B , are calculated via the conventional Euclidean distance, $R_{AB} = \sqrt{(x_A - x_B)^2 + (y_A - y_B)^2 + (z_A - z_B)^2}$. If used as a precomputed similarity matrix to describe the strength of the interactions between molecules, the representation is given as $G_{AB}^{xyz} = \frac{1}{R_{AB}}$. In k -means and agglomerative clustering, the Cartesian coordinates are fed directly to the clustering algorithm. In the spectral clustering algorithms, the Cartesian data must be represented by an affinity matrix. The affinity matrix, A , defined as describes the strength of interactions between data pairs.

When a descriptor exhibits a more block diagonal structure in the descriptor, many clustering methods including those explored here are able to distinguish easily between these sections. To explore if this idea can assist in molecular fragmentation definition, a representation is constructed around a bond matrix based on the covalent radii of the atoms (cr). The covalent radii has been used in supervised machine learning as a way to numerically capture a chemical environment, thus exploring this as a feature for unsupervised machine learning is sensible.^{235,236} This is a Boolean matrix in which one indicates the presence of a bond as determined by

$$G_{AB}^{cr} = \begin{cases} 1 & \text{if } R_{AB} \leq 1.1(A_{cr} + B_{cr}) \\ 0 & \text{otherwise} \end{cases}$$

where the G_{AB}^{cr} is the descriptor entry, R_{AB} is the distance between atom A and atom B , and A_{cr} (B_{cr}) is the covalent radii of A (B).^{237,238}

The remaining two representations incorporate information of the bonding environment via the Mayer bond order as a surrogate for density matrix and thereby quantum mechanical bonding information. The Mayer bond order is defined in terms of spin orbitals as:

$$G_{AB}^{Mbm} = 2 \sum_{\mu \in A} \sum_{\nu \in B} (\mathbf{PS})_{\mu\nu} (\mathbf{PS})_{\nu\mu}, \quad (72)$$

where P is the density matrix and S is overlap matrix in an atomic orbital basis μ and ν .²³⁹ The Mayer bond order matrix (G^{Mbm}) representation is a slight modification of the form of the values into an affinity matrix, where the magnitude of the Mayer bond matrix element represents the similarity between two atoms. As previously mentioned, some clustering methods benefit from a more block diagonal structure of the descriptor, we also look to coarse grain the descriptor through rounding the values of the matrix. This results in the rounded Mayer bond order matrix (G^{rMbm}) representation which rounds G_{AB}^{Mbm} according to:

$$G_{AB}^{rMbm} = \begin{cases} \lceil B_{AB} \rceil, & \text{if } \{B_{AB}\} \geq 0.5 \\ \lfloor B_{AB} \rfloor, & \text{otherwise.} \end{cases} \quad (73)$$

The intention of this representation is to dampen out insignificant pairs, allowing only the most strongly interacting pairs to be considered in the descriptor and making the cuts between clusters more obvious. In general, it is expected that incorporating the bond order into the descriptor will enable the fragmentation approach to preserve bonds between the most strongly interacting parts of the molecules.

5.3.3 Validation

In this section, the means of defining successful partitioning of a molecular system is outlined. In this work, two different approaches are used: one which ensures the expected fragmentation is produced and another that quantifies the recovered energy of the full system. The set of water clusters and methylthiophene tetramers were chosen that should produce very clear clustering for a requested number of clusters. The performance of the clustering methods on these test cases

can be assessed through external validation, which compares the resulting cluster labels to a expected/correct cluster labels.^{240,241} The external validation statistic used in this work is the Adjusted Rand Index (ARI). The Rand index, 74, measures the frequency of occurrence of agreement over the total pairs.

$$\mathcal{R}(\mathcal{C}, \mathcal{C}') = \frac{N_{11} + N_{00}}{n(n-1)/2} \quad (74)$$

Here \mathcal{C} is the resulting clustering and \mathcal{C}' is the expected clustering consisting of n total points. N_{00} and N_{11} are the number of data point pairs in the same clusters and different clusters respectively for C and C' . However, since there is a small probability these data points could end up in the same cluster by chance, the ARI is used. The ARI corrects for chance by using a baseline of expected similarity from a random model. A value close to unity represents total agreement between the expected and the actual fragmentation. The ARI approaching zero reflects an increase in the difference between the expected and returned fragmentation. In the oligomer systems, the correct fragmentation is ambiguous and thus correct clustering is not known ahead of time. for these systems the percent error of the fragmented system from supermolecular energy is used as the metric of success.

5.3.4 Systems

To measure the success of the AFMC approach, families of molecules are considered to explore the performance in cases where the partitioning introducing the lowest error is apparent for noncovalently and covalently bound molecules and a case of more ambiguous fragmentation in oligomers.

Water Clusters: The performance of the AFMC with the various descriptors on noncovalently bound molecules, a set of water clusters were explored. Water clusters, $(\text{H}_2\text{O})_n$ for $n=2$ to $n=21$ optimized with the TIP4P water model were used from the Wales cluster database.^{242,243}

Methylthiophenes: To benchmark performance of the fragmentation methods on covalently bound systems with a clear desired fragmentation, a set of methylthiophene tetramers were explored. When partitioning these systems, 4 fragments were requested. Coordinates for the tetramers were generated with Open Babel by providing a SMILES representation with defined linkage atoms, indicated by a box in Figure 9.²⁴⁴ The geometry of each molecule was determined at two lev-

els of theory to study the sensitivity of the fragmentation to small perturbations in the structure. The levels of theory used optimize the structures were Hartree-Fock/6-311G** and ω B97X-D/6-311G**.^{245,246} Introduction of broken bonds within the molecule were treated by hydrogen capping. The hydrogen cap contributions were then treated by subtracting the energy of all hydrogen atoms from that of the fragment calculation.

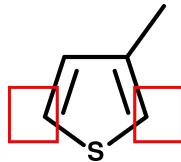


Figure 9: The methylthiophene monomer 2-D structure with the linkage atoms highlighted in red.

Silyl ketene oligomers: Silyl ketene oligomers (SKs) provide a set of systems where ideal fragmentation is ill-defined due to the backbone and side group interactions which may be essential to describe their chemistry and as a result, not clear fragmentation pattern. SKs have the general form of $(RC=C=O)_n$, where R is a SiR_3 group (Figure 10(left)). This class of molecule is a candidate for chain polymerization and can avoid the undesired 2+2 cycloadditions observed in aryl and alkyl ketenes. Previous work in our group has aimed to predict stable structures and polymerization mechanisms.²⁴⁷ As large polymer units are considered, the computational cost grows and fragmentation becomes an attractive and possibly necessary option. In this work, oligomers of the SK monomers act as a test system to explore the performance of the clustering approaches in terms of the clustering ability to reduce computational time while minimize the difference in error when compared to the supermolecular calculation. The systems explored consist of a dimer and trimer of the tert-Butyldiphenylsilyl monomer, with a methylonate nucleophile to begin the polymerization displayed in Figure 10 (center, right). The SK structures were generated with Avogadro2 and optimized at the Becke-3 Parameter-Lee-Yang-Parr (B3LYP) including the Becke-Johnson dispersion correction (-D3(BJ)) level of theory with the pc-1 basis set.¹⁹⁵⁻¹⁹⁸ The work by Mardirossian et al. suggests this optimization level represents a balance of computational cost and accuracy.⁵⁷ Introduction of broken bonds within the molecule were treated by hydrogen capping. The hydrogen cap contributions were then treated by subtracting the energy of all hydrogen atoms from that of the fragment calculation.

This test case also aimed to determine the sensitivity of the clustering methods to the level of theory used for the Mayer bond order calculations. Using the optimized structure above, the Mayer bond matrix was also calculated at varying levels of theory: Hartree-Fock, B3LYP^{195–198}, and ω B97M-V²⁴⁸ each level of theory is paired with four different basis sets: STO-3G^{249,250}, 3-21G^{251,252}, cc-pVDZ, cc-pVTZ.^{253,254}

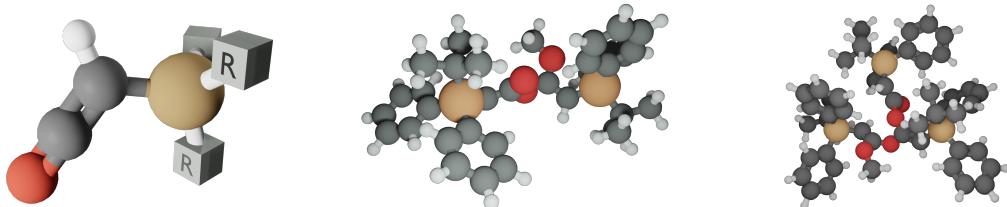


Figure 10: Structure of a generic silyl ketene unit (left), the dimer (center) and trimer (right) used for the clustering benchmark.

5.4 Results and Discussion

5.4.1 Water Clusters

The agglomerative, k -means, and spectral clustering return the expected clustering as a comparison of the expected and true clustering results with an average ARI of 1 for all descriptors. This indicates the non-covalently bound water structures are correctly clustered when the number of clusters is predefined. This is expected as non-covalently bound systems have apparent non-interacting or weakly interacting points are clearly represented by the representation and resolved by the clustering algorithm. These results suggests these clustering methods could reliably find reasonable fragments for non-covalently bonded systems.

In order to recover the total energy of a water cluster, often more than one water monomer should be included in a fragment, placing importance on clustering methods which allow for the selection of cluster numbers are valuable. Ideally, the clustering algorithms should resolve to keep individual water molecules intact within each cluster as opposed to separating a molecule between

two fragments.⁴⁰ Figure 11 demonstrates for spectral clustering with the Cartesian descriptor cases where the number of clusters is lower than the number of water monomers. Though in addition to ensuring only non-covalent bonds are broken, another important metric is that the method return balanced clusters. To explore this, the $(\text{H}_2\text{O})_{21}$ molecule was partitioned $N_{frag} = 2 \dots 21$. The results shown in Figure 12 demonstrate that spectral clustering can maintain fragments with a balanced number of water molecules in each fragment over the range tested, while agglomerative and k -means have trouble for a smaller number of requested fragments. If one considers the variance in the water molecules per cluster for In this case the returned fragments are defined as including one or more water monomers without any segmentation of O-H bonds. This successful behavior is observed for k -means, agglomerative, and spectral clustering algorithms with all descriptors.

The timing of the clustering itself is another important metric to consider, and the AFMC has negligible cost for generating fragments. Once the features of the descriptor are calculated, the routine to set up the molecule, generate fragments, and create output files took approximately 2 seconds for all clustering methods on an 4 core laptop with a i5-5200 CPU. The AFMC approach thus offers a computationally efficient way to generate fragments which can then be used in combination with any interfragment treatment approaches to estimate the energy and properties of large molecules.

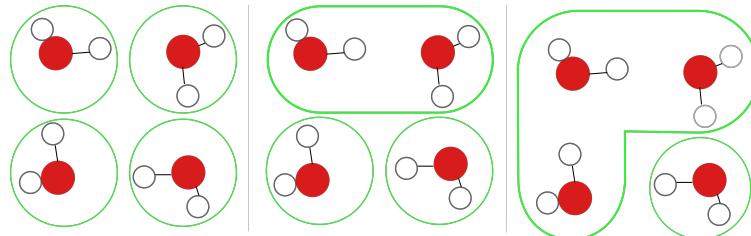


Figure 11: A demonstrative result for the preservation of covalent bonds in cases when the N_c is less than the N_m . The N_c requested decreases from 4 (requested on the left) to 2 (requesting on the right), clusters are designated by the green outlines. This example is performed with spectral clustering on the Cartesian representation for $(\text{H}_2\text{O})_4$. Covalent bonds are preserved for all descriptor/clustering algorithm combinations.

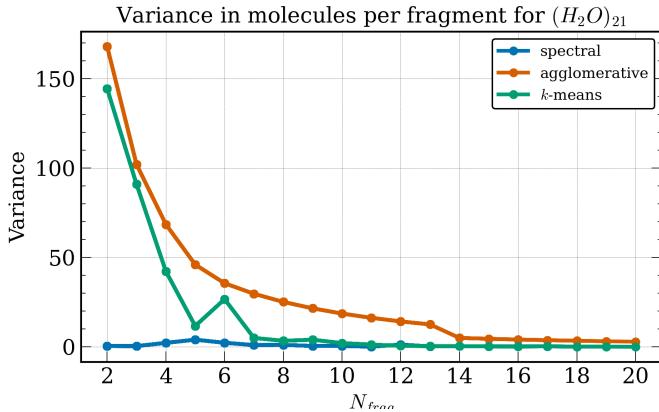


Figure 12: The variance in cluster size for cases when the number of requested clusters is less than the number of monomers where the covalent bonds intact. For for $(H_2O)_{21}$, the N_{frag} requested increases from 2 to 20. Shown here for the G_{Mbm} descriptor.

5.4.2 Methylthiophenes

The clustering method performance for the covalently bound methylthiophene tetramers into four fragments are shown in Figure 13; The effects of geometry optimization on the fragmentation results are also presented. Optimization are performed at the Hartree-Fock/6-311G** level of theory subfigure a and ω B97X-D/6-311G** level of theory was used in subfigure b. Spectral clustering performs well for all molecular representations with no dependence on the descriptor, level of theory, or basis set used. Other clustering methods have a strong dependence on the descriptor used and variation as the level of theory used in the geometry optimization changes. These structures are challenging due to the descriptors maintaining less of a block diagonal structure, i.e more non-local interactions. A less block diagonal structure in the representation means cluster boundaries are much less clear leading to problems for certain unsupervised learning algorithms. However, as spectral clustering first embeds the representation into a lower dimensional space before clustering, it is able to resolve the primary interactions. Notably, agglomerative, k -means and spectral clustering perform well with the Cartesian descriptor.

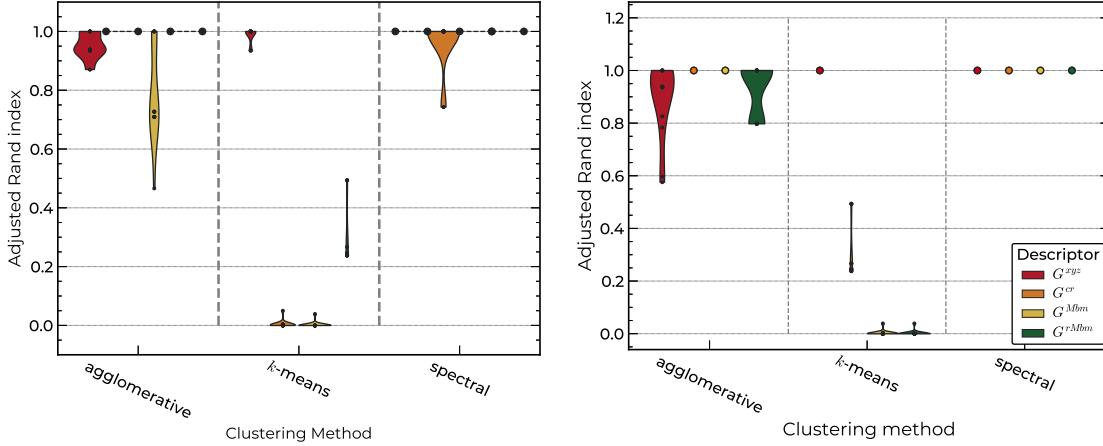


Figure 13: Performance of clustering for the methylthiophene test set, measured as the average ARI across the set of methylthiophenes. The colored dots at 1 indicate successful clustering of the full test set, as the average ARI is 1 for the descriptor/clustering algorithm combination.

5.4.3 Silyl Ketene

Silyl ketenes, realistic systems with more ambiguity in the choice of partitioning, are challenging applications for these clustering approaches. The dimer results are found in Figure A.5 and trimer results can be seen in Figure 5.4.3. The findings across the dimers and trimers are comparable and suggest spectral clustering is the most robust clustering algorithm, achieving the lowest errors with the largest speedups across all descriptor/method pairs. Interestingly, a fragmentation pattern with the lowest error is found when the covalent radii bond matrix is used for the dimer or the Cartesian descriptor is used for the trimer. The favorable fragmentation using these purely structure-based descriptors suggests it may be possible to select fragments that result in high accuracies in molecular properties without relying on the incorporation quantum mechanical information in the descriptor for some systems.

Agglomerative clustering produces acceptable fragments, though in general, a larger deviation from the supermolecular result is observed. On the other hand, only the Cartesian descriptor with the k -means clustering algorithm yields reasonable result, which is likely due to the conserved spatial information allowing for an accurate choice of centroids opposed to representations which

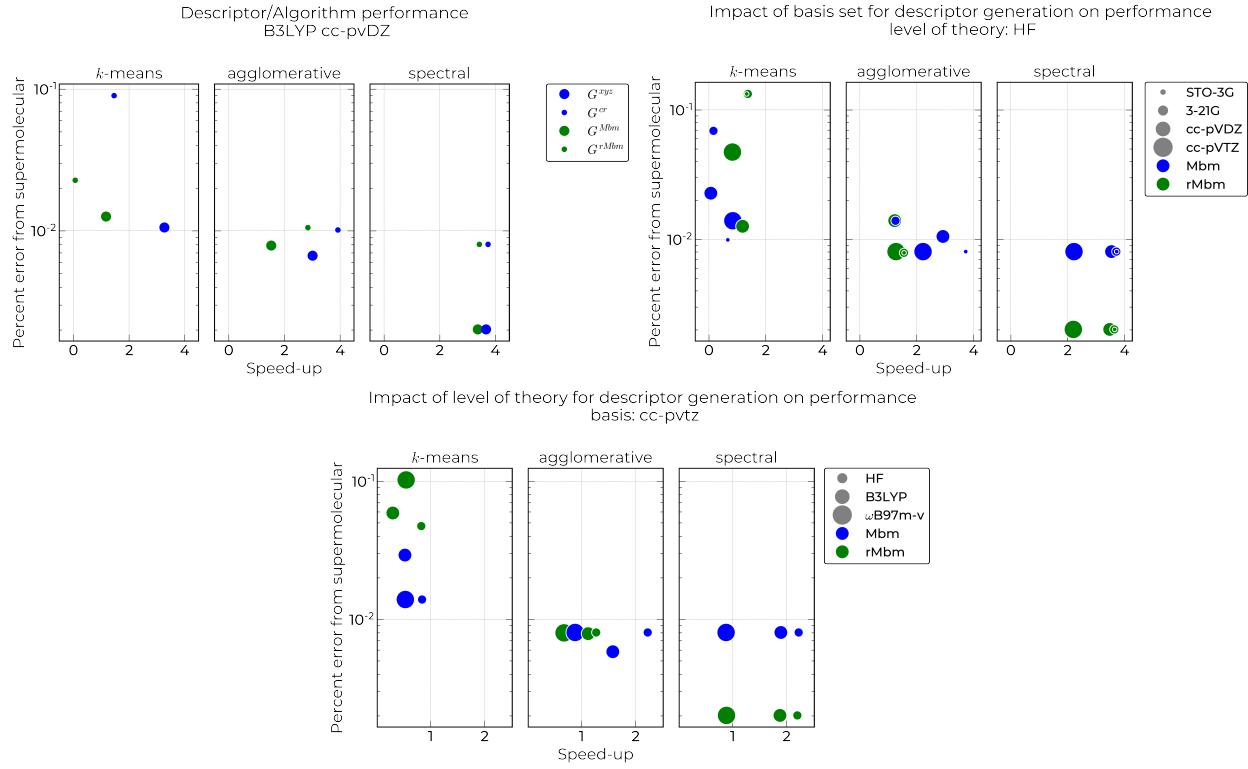


Figure 14: Assessment of fragmentation schemes on silyl ketene trimers: Presented is the percent error of the energy and the speedup over the supermolecular calculation for the SK trimer to assess the performance of the fragment approaches and descriptor quality.



Figure 15: Representative visualization of resulting fragmentation for silyl ketenes. Results from agglomerative (left) and spectral (right) clustering on the G^{Mbm} descriptor. Colors represent fragment identity.

employ bonding environment as the features. Overall, the clustering approaches investigated are relatively insensitive to the level of theory and basis set used to generate the Mayer bond matrix, which is encouraging as it may allow for future computational savings in future applications on larger more complex systems where a purely structure base descriptor may not incorporate the necessary interaction important to the structure. If cases arise in larger systems where the quantum mechanics based descriptors become necessary, the low level approximation to the bond order will suffice as long as the bond matrix is meaningful.

5.5 Conclusion

In this work we explore the Automatic Fragmentation of Molecules through Clustering (AFMC) as an automated partitioning scheme for molecular systems. AFMC approach utilizes UML methods to determine molecular domains in computationally efficient and transferable ways. Three classes of systems were studied to assess the performance of the AFMC approach: water clusters to investigate non-covalently bonded molecules, methylthiophene tetramers to probe the behavior of AFMC with ring structures, and two silyl ketene oligomers to explore a realistic chemical case with ambiguous fragmentation choices. Several clustering approaches and molecular representations were tested. We find that this approach with various clustering methods can accurately identify meaningful molecular domains for non-covalently bound molecular systems (water clusters) and in the case of covalently bonded systems with aromatic units. Overall, a spectral clustering approach was able to produce balanced and sensible molecular fragments as can be seen by the low error and high speed up for the silyl ketene structures. The clustering is performed on molecular representations derived from either the molecular structure alone or a low-level quantum mechanics prediction of the Mayer bond matrix. Both classes of representations performed well for the systems studied, though it is not yet clear whether more intricate, correlation-dependent, bonding schemes will benefit from the quantum mechanical informed descriptors and we aim to explore this in future work. The combination of Spectral clustering with the Cartesian descriptor provides the most reliable clustering with minimal preparation cost, thus we recommend this combination for molecules alike to those studied here. The combination of these descriptors and UML techniques

provide a low cost way to determine acceptable fragments for computational chemistry for further, more accurate, quantum mechanical calculations.

6.0 Conclusions

In this work, we probe various approaches to focusing efforts in the electronic structure calculations in an intuitive way to allow an accurate description including correlation energy while minimizing or reducing the computational demand. We approach this from two different directions. The first is through the application of is a stochastic sampling to solve the Schrödinger equation through quantum Monte Carlo (QMC) approaches. The second is through an attempt to choose molecular partitioning for fragmentation approaches through the utilization of unsupervised machine learning methods.

The first is through the application of is a stochastic sampling to solve the Schrödinger equation through quantum Monte Carlo (QMC) approaches.

The second is through an attempt to choose molecular partitioning for fragmentation approaches through the utilization of unsupervised machine learning methods.

A.1 Additional works

As some of the research directions are not tied directly to the dissertation theme or are still being developed, they are mentioned here. These works can briefly be described as contribution to improving quantum computing emulators, an extension of the AFMC approach to regional embedding fragment definition, and a density functional theory analysis of silyl ketene systems to inform experimental synthesis.

A.1.1 Q-GPU

One such project was an opportunity to collaborate with Dr. Tang and others on a Quantum Computing emulation project which aimed to improve the emulation of quantum computing quantum chemistry simulation by increasing the number of qubits one can emulate through the utilization of graphical processing units.[255] The improvements through Q-GPU are Q-GPU

features i) proactive state amplitude transfer, ii) zero state amplitude pruning, iii) delayed qubit involvement, and iv) lossless nonzero state amplitude compression. Experimental results across nine representative quantum circuits indicate that Q-GPU significantly reduces the execution time of the state-of-the-art GPU-based QCS by 71.89% ($3.55\times$ speedup). Q-GPU also outperforms the state-of-the-art OpenMP CPU implementation, the Google Qsim-Cirq simulator, and the Microsoft QDK simulator by $1.49\times$, $2.02\times$, and $10.82\times$, respectively. The author's contribution were in collaborating on the hydrogen chain circuit, which was one of 9 quantum circuit used to test Q-GPU, as well as collecting timing information with the improvements for the circuits.

A.1.2 AFMC applications Towards Regional Embedding.

Collaborators: Shiv Upadhyay and Daniel S. Lambrecht

In the previous chapter, we looked at partitioning an entire molecule into important domains in systems where there is no focal interaction. However, there are many cases in which a focal area of the calculation is known such as an adsorbate to a surface, a polymer chain growth site, or an enzyme in a binding pocket. One way in which these calculation are made more affordable is through localization schemes. A recent scheme was proposed for a regional embedding approach in which localization of the occupied and virtual molecular orbitals (MOs) occurs separately for a fragment area of interest. This fragment region was not defined in the original work requiring a set of convergence test for each system studied, taking away from the speed-ups granted by the method. In this chapter, the utility of the AFMC approach to determining this fragment area is presented. Specifically, the spectral clustering approaches as proposed in our previous work with a newly-developed descriptor that is tailored to target an area of interest within a chemical system. This approach allows for a reliable definition of the fragment area over which localization of the molecular orbitals occurs, which removes the necessity of convergence testing for fragment size.

Though these preliminary results are brief, they do suggest a path forward in using fragmentation approaches with regional embedding schemes, or other localization schemes. To move this work forward, the characterization of the current descriptors must be applied to a wider variety of systems to understand how system agnostic the descriptors can be. Additionally, other forms of descriptors should be investigated.

A.1.3 Characterization of silyl ketenes

Collaborators Krista Schoonover, Ian Baster, Sarah Mitchell, Emily Pentzer and Daniel S. Lambrecht.

A joint experimental and computational study to explore the preference for nucleophilic addition or deprotonation of combinations of different anions and different silyl ketene units. Specifically, the silyl ketene with TBDPS functional groups paired with various nucleophiles were explored with DFT. The addition reaction was thermodynamically favored, while deprotonation was primarily kinetically favored. The author's contribution to this work were in the calculation of the molecular electrostatic potential and small contributions in discussion to the machine learning study which was led and performed by Daniel. S. Lambrecht. This manuscript is currently in preparation with collaborators Krista Schoonover, Ian Baster, Sarah Mitchell, Emily Pentzer and Daniel S. Lambrecht. This manuscript is currently in preparation.

A.2 Supplemental Material for Chapter 3

The Supplementary Material document includes the total energies and error bars for the quantum Monte Carlo calculations, the total energies for the DFT calculations, and details of the convergence of the DFT total energies with respect to the k -point grid and kinetic energy cutoff of the plane wave basis, and a comparison of the density difference of DMC-PBE and DMC-HSE.

Convergence of PBE total and binding energies with respect to relevant parameters

In this section, the sensitivity of PBE total and binding energies with respect to k -point grid and cutoff of the kinetic energy of the wave function for the plane wave basis.

Table 6: PBE total energies of each system and binding energies (eV) with respect to increasing the Monkhorst-Pack k -point grid. The hydrogen total energy value used to calculate the binding energy was for a $1 \times 1 \times 1$ k -grid with a 250 Ry wave function cutoff.

k -point grid	hydrogen on graphene	graphene	H	binding energy
2x2x1	-7759.736	-7745.293	-13.604	-0.839
3x3x1	-7759.751	-7745.339	-13.604	-0.808
4x4x1	-7759.749	-7745.323	-13.604	-0.823
5x5x1	-7759.749	-7745.322	-13.604	-0.824
6x6x1	-7759.749	-7745.325	-13.604	-0.820

Table 7: PBE total and binding energies (eV) with respect to varying the kinetic energy cutoff values of the wave function on a $4 \times 4 \times 1$ k -point grid. The hydrogen total energy value used to calculate the binding energy was for a $1 \times 1 \times 1$ k -grid with a 250 Ry wave function cutoff.

kinetic energy cutoff	hydrogen on graphene	graphene	H	binding energy
150	-7758.435	-7744.019	-13.603	-0.813
250	-7759.749	-7745.323	-13.604	-0.823
300	-7759.790	-7745.362	-13.604	-0.824
350	-7759.799	-7745.372	-13.604	-0.824

Hybrid functional total and binding energies (eV) with respect to k -point grid

Table 8: PBE0 and HSE total energies and binding energies (eV) with respect to increasing the Monkhorst-Pack k -point grid.

functional	k -point grid	hydrogen on graphene	graphene	H	binding energy
PBE0	6x6x1	-51820.644	-51806.227	-13.537	-0.880
PBE0	10x10x1	-51820.644	-51806.255	-13.537	-0.851
PBE0	12x12x1	-51820.644	-51806.253	-13.537	-0.854
HSE	6x6x1	-51820.264	-51805.892	-13.565	-0.807
HSE	10x10x1	-51820.264	-51805.907	-13.565	-0.792
HSE	12x12x1	-51820.264	-51805.905	-13.565	-0.794

Non-self-consistent field calculations details

Table 9: Non-self-consistent field calculations: system and binding energies used to calculate the differences in Table II of the manuscript, reported without basis set superposition error correction.

density functional	energy functional	graphene	H	Hgraphene	binding energies
pbe	pbe	-1903.884	-0.496	-1904.415	-.962
hse	pbe	-1903.883	-0.496	-1904.414	-.947
pbe	hse	-1903.989	-0.498	-1904.520	-.879
hse	hse	-1903.989	-0.498	-1904.521	-.905

DMC Calculation Details

Total Energies for DMC calculations

Table 10: Total and binding energies (eV) of a hydrogen atom on graphene calculated with QMC.

Trial wave function	hydrogen on graphene	graphene	hydrogen	binding energy
PBE	-7748.200 ± 0.012	-7733.893 ± 0.014	-13.616 ± 0.002	-0.691 ± 0.019

Total Energies per twist angle

Table 11: Total energies (eV) at each twist angle used for graphene DMC treatment.

Twist number	Total Energy
1	-7733.872 ± 0.046
2	-7734.059 ± 0.057
3	-7734.074 ± 0.058
4	-7733.851 ± 0.041
5	-7733.492 ± 0.052
6	-7733.467 ± 0.048
7	-7733.080 ± 0.057
8	-7734.137 ± 0.047
9	-7734.080 ± 0.041
10	-7733.853 ± 0.047
11	-7734.444 ± 0.050
12	-7733.983 ± 0.048

Table 12: Total energies (eV) at each twist angle used for hydrogen atom chemisorbed on graphene DMC treatment.

Twist number	Total Energy
1	-7748.443 ± 0.057
2	-7748.474 ± 0.045
3	-7748.283 ± 0.046
4	-7748.294 ± 0.047
5	-7748.006 ± 0.050
6	-7747.736 ± 0.056
7	-7748.748 ± 0.045
8	-7748.092 ± 0.063
9	-7747.601 ± 0.042
10	-7748.488 ± 0.047
11	-7747.941 ± 0.041
12	-7747.779 ± 0.050
13	-7747.761 ± 0.051
14	-7748.275 ± 0.045
15	-7748.655 ± 0.054
16	-7748.380 ± 0.053
17	-7748.037 ± 0.058
18	-7748.615 ± 0.047

DMC-DFT density difference

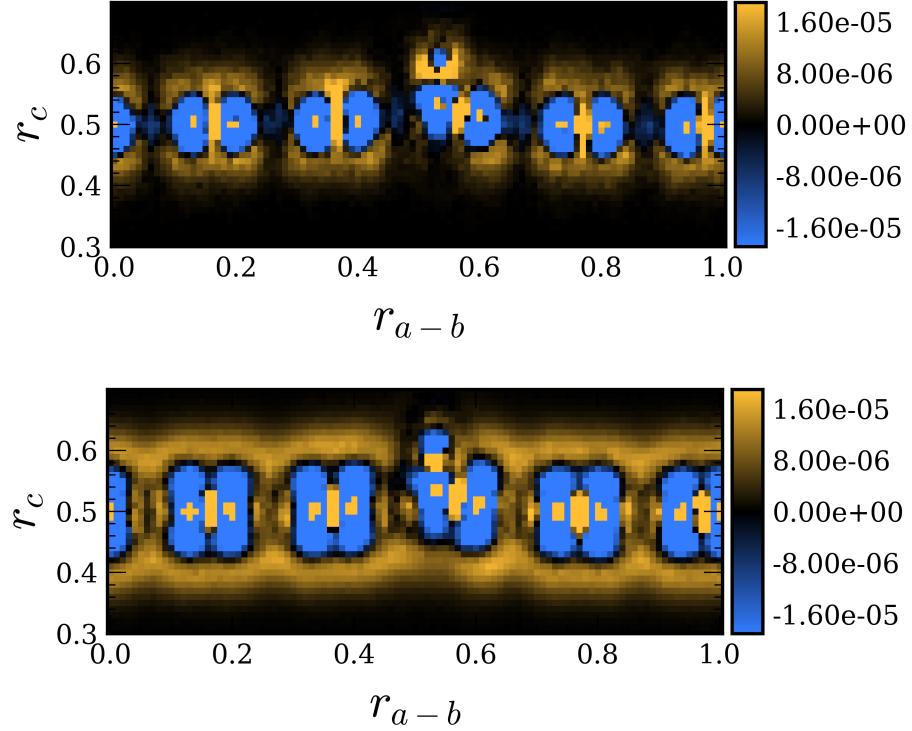


Figure 16: Similar to Figure 3 in the main text. DMC-PBE density (top) is reproduced from the main text and compared with the DMC-HSE density (bottom). We attribute the differences to primarily the difference in basis sets, plane waves vs Gaussian orbitals, and the use of pseudopotentials in the plane wave calculation.

A.3 Supplemental Material for Chapter 4

Additional simulation details such as input files are available at: https://github.com/shivupa/Water4_JCP_Special_Issue_Supplemental_Material.

Geometries

The geometries are given in xyz format for the (H₂O)₄ structures in Angstroms.

Geometry at R = 4 Å

13

WATER 4 SYSTEM 4.0 Angstroms

X	0.000000	0.000000	0.000000
O	0.000000	1.730527	2.893437
H	0.000000	1.387569	2.000001
H	0.000000	2.681413	2.785434
O	0.000000	-1.730527	2.893437
H	0.000000	-1.387569	2.000001
H	0.000000	-2.681413	2.785434
O	0.000000	1.730527	-2.893437
H	0.000000	1.387569	-2.000001
H	0.000000	2.681413	-2.785434
O	0.000000	-1.730527	-2.893437
H	0.000000	-1.387569	-2.000001
H	0.000000	-2.681413	-2.785434

Geometry at R = 7 Å

13

WATER 4 SYSTEM 7.0 Angstroms

X	0.000000	0.000000	0.000000
O	0.000000	1.730527	4.393437
H	0.000000	1.387569	3.500001
H	0.000000	2.681413	4.285434
O	0.000000	-1.730527	4.393437
H	0.000000	-1.387569	3.500001
H	0.000000	-2.681413	4.285434
O	0.000000	1.730527	-4.393437
H	0.000000	1.387569	-3.500001
H	0.000000	2.681413	-4.285434

O	0.000000	-1.730527	-4.393437
H	0.000000	-1.387569	-3.500001
H	0.000000	-2.681413	-4.285434

Supplemental Basis Functions

The exponents for the 7s7p and 3s1p diffuse Gaussian type orbitals are given below in the GAMESS format.

7s7p

```
S    1
  1   0.02362232   1.0
S    1
  1   0.00738198   1.0
S    1
  1   0.00230687   1.0
S    1
  1   0.00072090   1.0
S    1
  1   0.00022528   1.0
S    1
  1   0.00007040   1.0
S    1
  1   0.00002200   1.0
P    1
  1   0.02362232   1.0
P    1
  1   0.00738198   1.0
P    1
  1   0.00230687   1.0
P    1
  1   0.00072090   1.0
P    1
  1   0.00022528   1.0
P    1
  1   0.00007040   1.0
```

P 1
1 0.00002200 1.0

3s1p

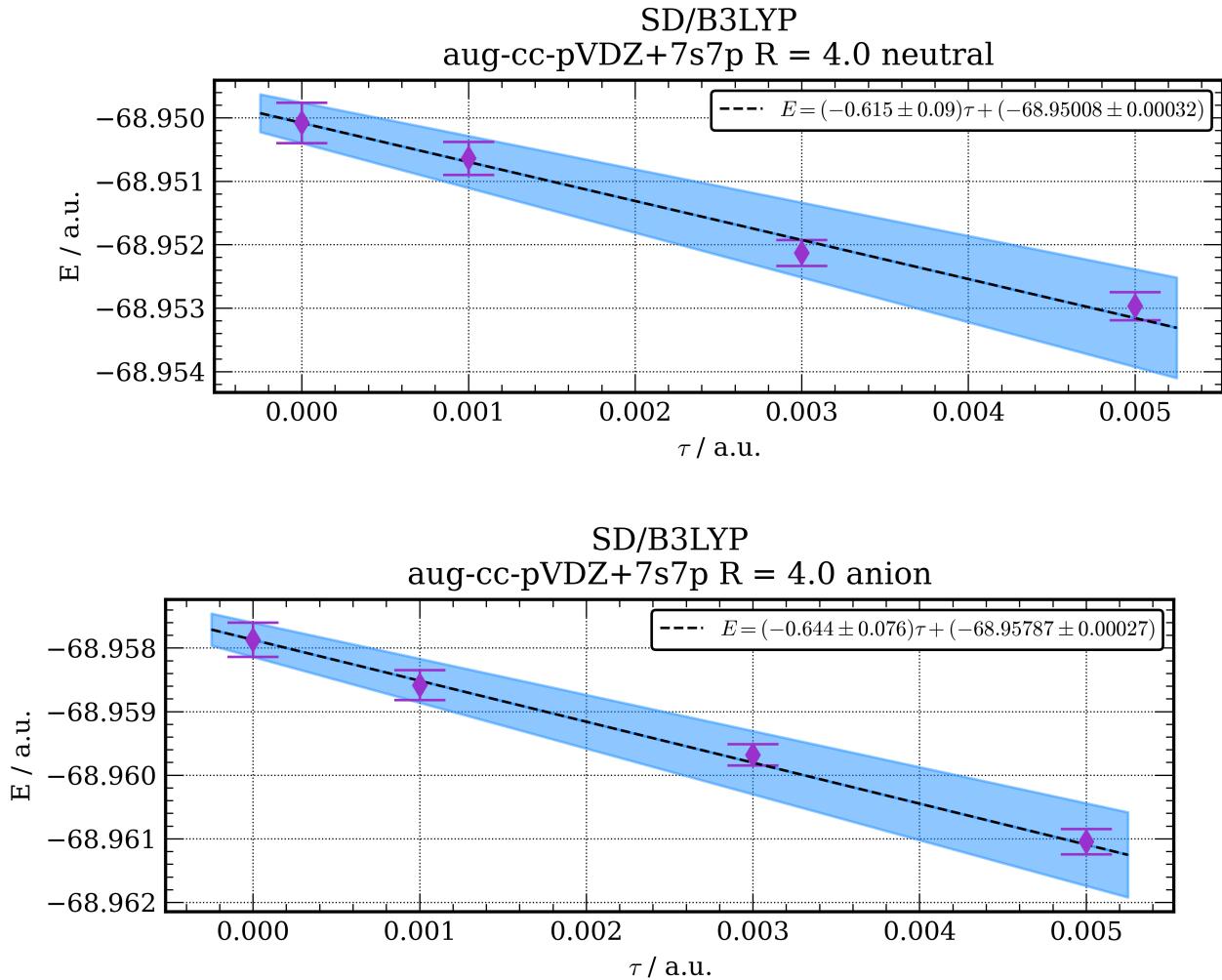
S 1
1 0.02362232 1.0
S 1
1 0.00738198 1.0
S 1
1 0.00230687 1.0
P 1
1 0.02362232 1.0

3s1p3d

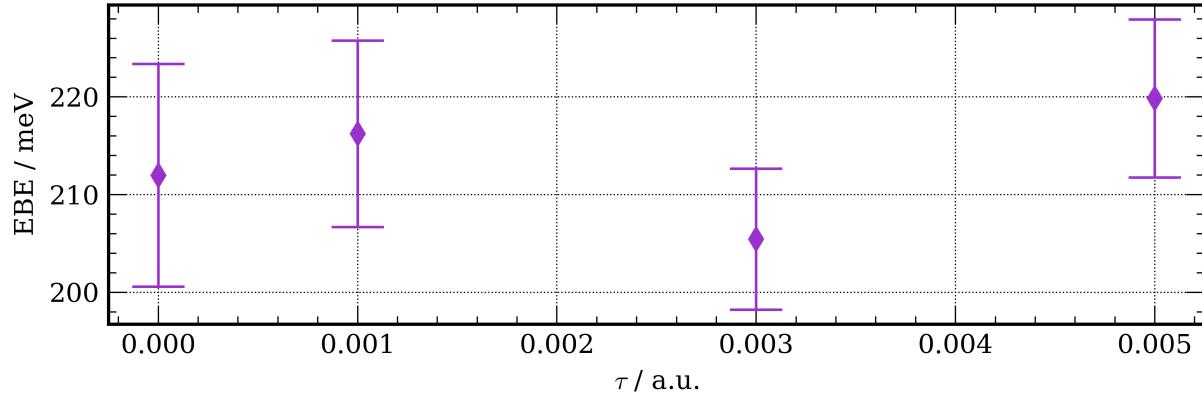
S 1
1 0.02362232 1.0
S 1
1 0.00738198 1.0
S 1
1 0.00230687 1.0
P 1
1 0.02362232 1.0
D 1
1 0.02362232 1.0
D 1
1 0.00738198 1.0
D 1
1 0.00230687 1.0

DMC Extrapolation

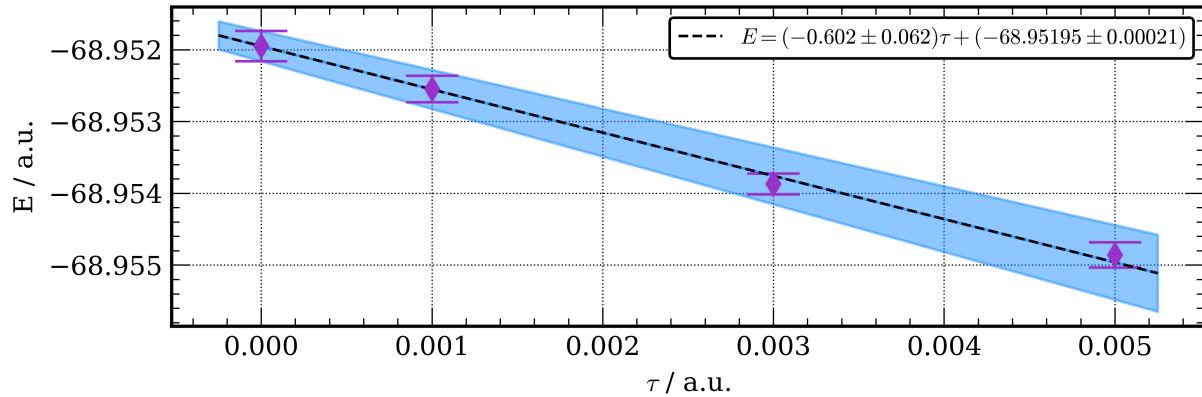
Summaries of the zero time step linear extrapolation are plotted below for the w(H₂O)₄ system, the R parameter is indicated in each plot title. The blue shaded region corresponds to the error in the fit of the DMC energies at the three timesteps (0.001, 0.003, 0.005).



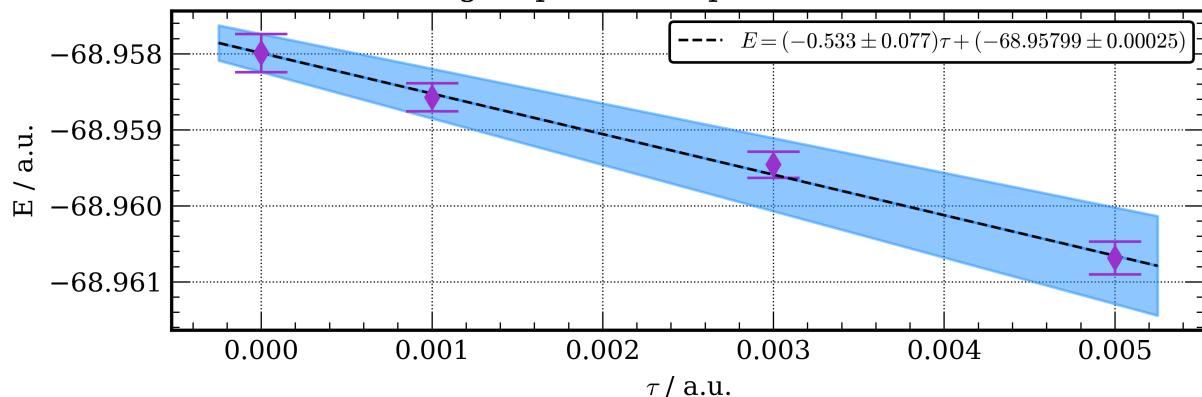
SD/B3LYP
aug-cc-pVDZ+7s7p R = 4.0 EBE



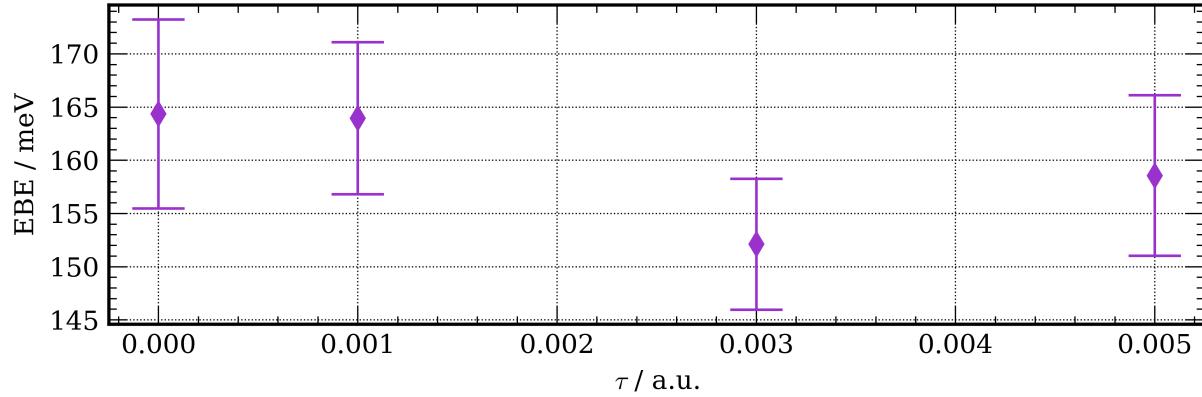
SD/B3LYP
aug-cc-pVDZ+7s7p R = 7.0 neutral



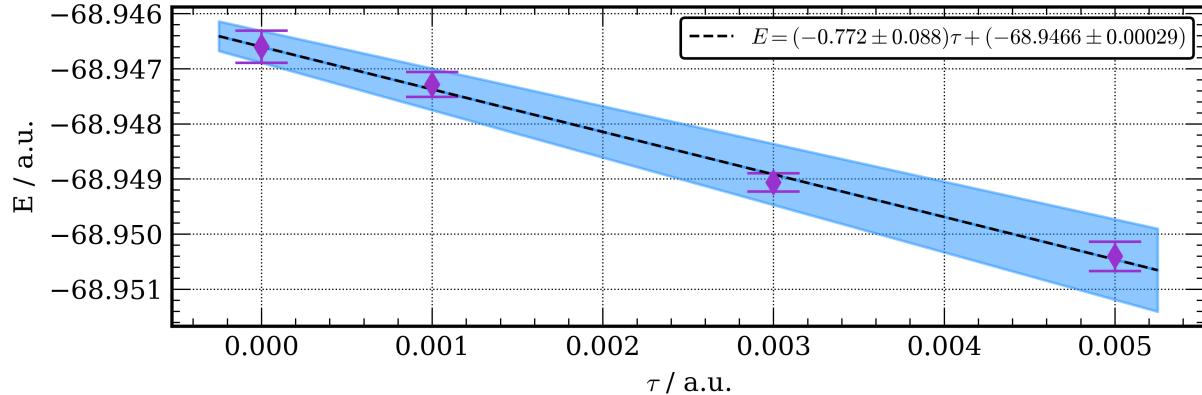
SD/B3LYP
aug-cc-pVDZ+7s7p R = 7.0 anion



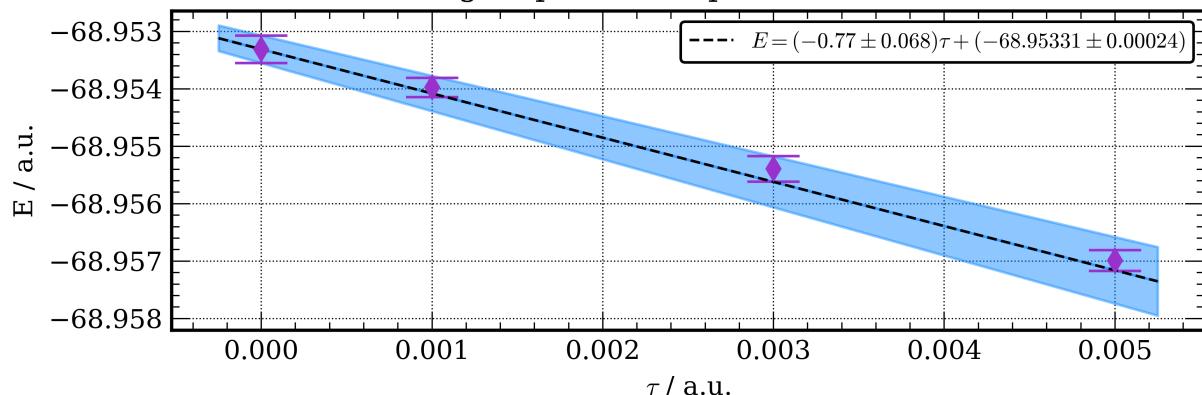
SD/B3LYP
aug-cc-pVDZ+7s7p R = 7.0 EBE



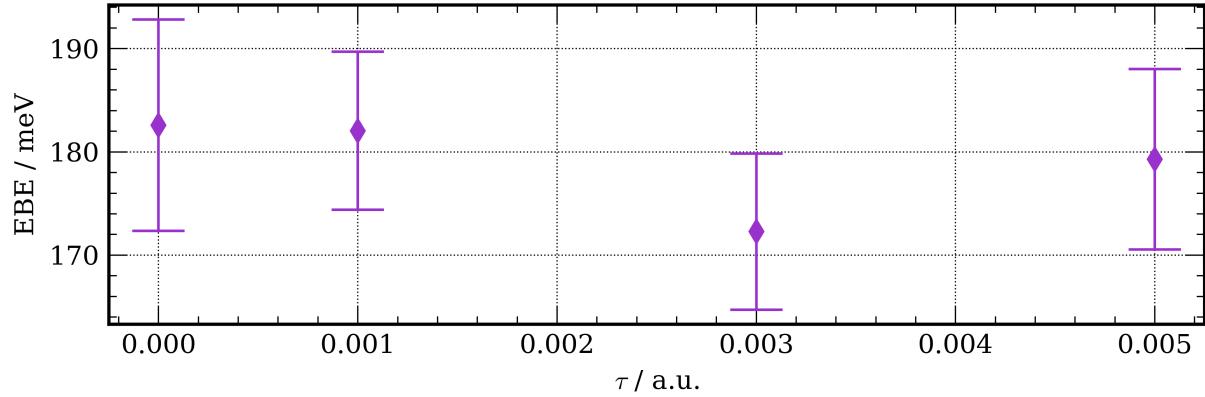
SD/HF
aug-cc-pVDZ+7s7p R = 4.0 neutral



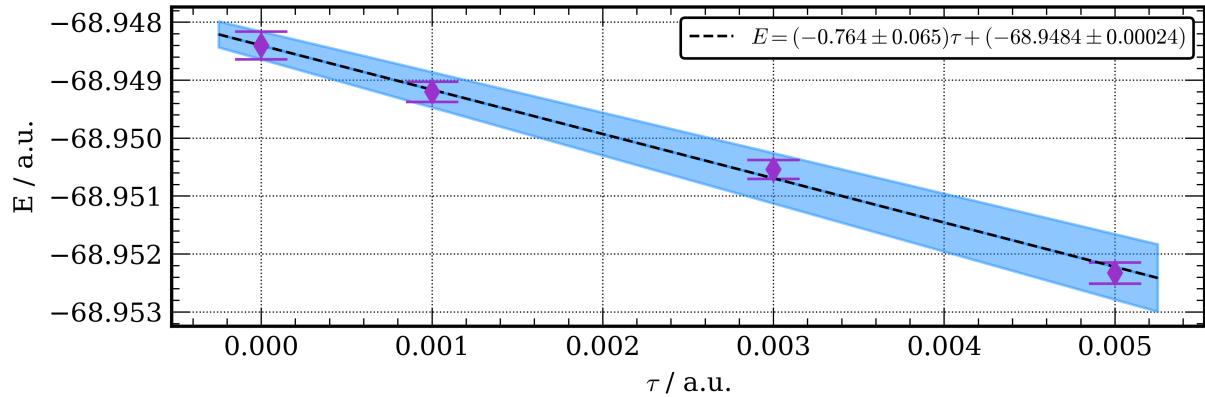
SD/HF
aug-cc-pVDZ+7s7p R = 4.0 anion



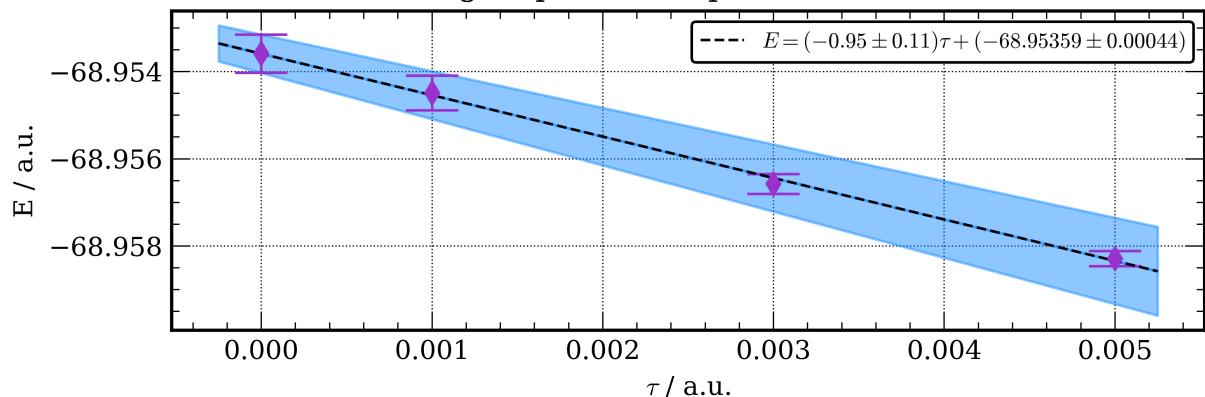
SD/HF
aug-cc-pVDZ+7s7p R = 4.0 EBE



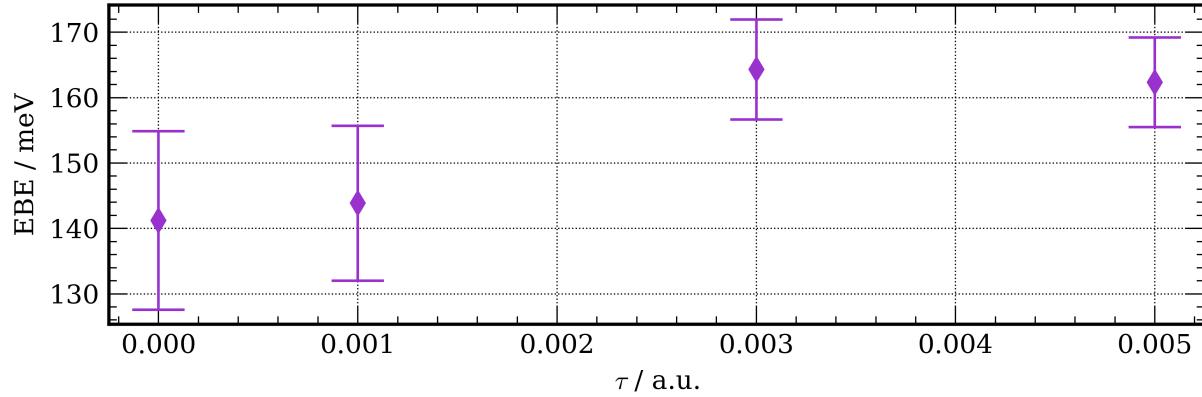
SD/HF
aug-cc-pVDZ+7s7p R = 7.0 neutral



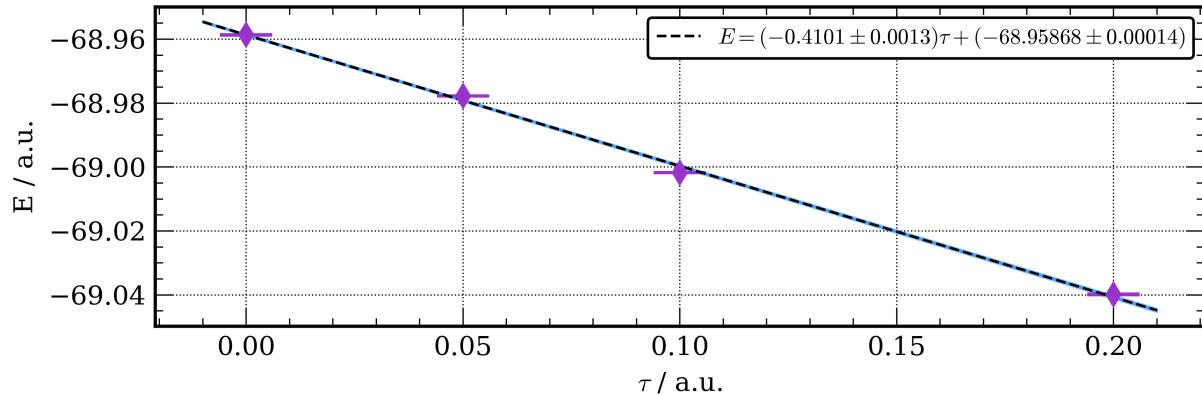
SD/HF
aug-cc-pVDZ+7s7p R = 7.0 anion



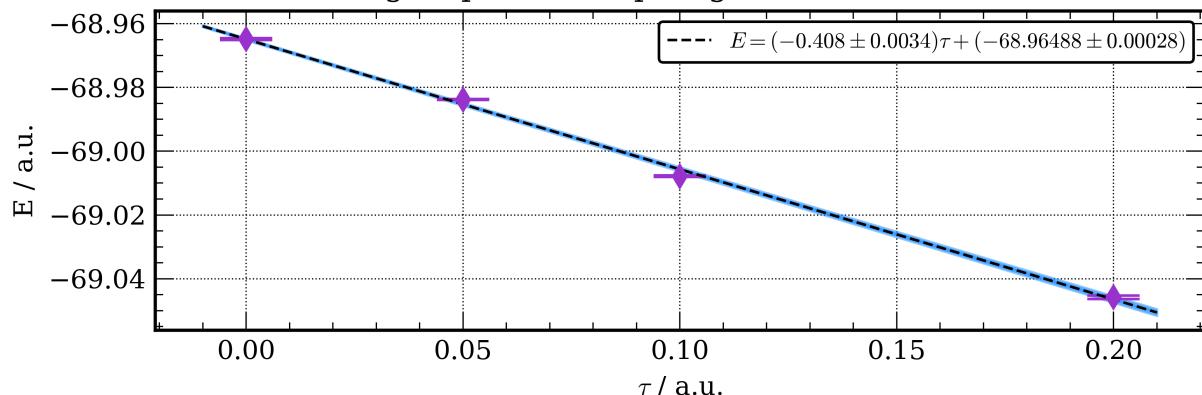
SD/HF
aug-cc-pVDZ+7s7p R = 7.0 EBE



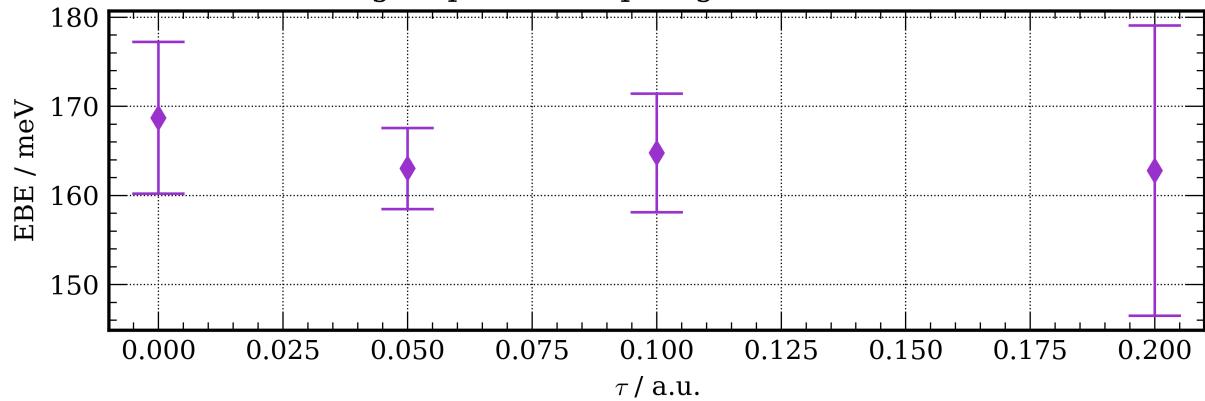
SD/HF
aug-cc-pVDZ+7s7p large ts R = 4.0 neutral



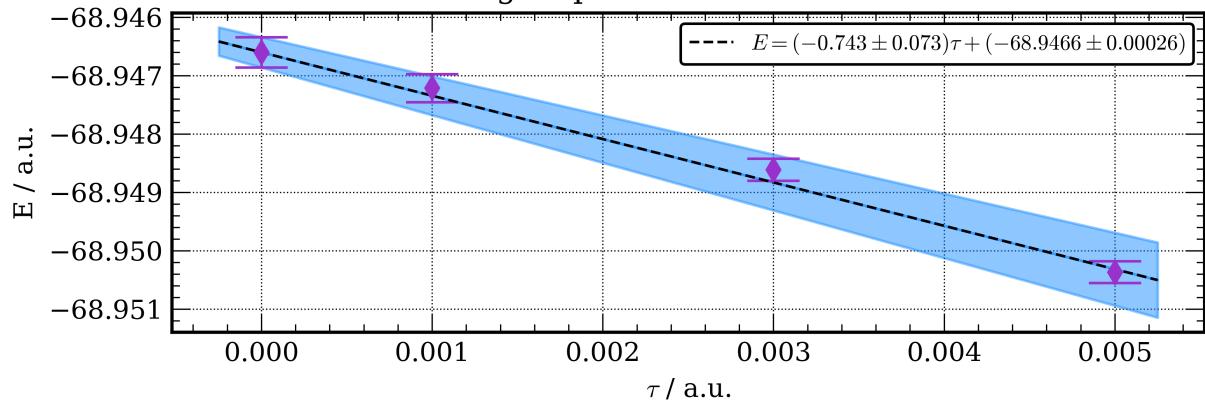
SD/HF
aug-cc-pVDZ+7s7p large ts R = 4.0 anion



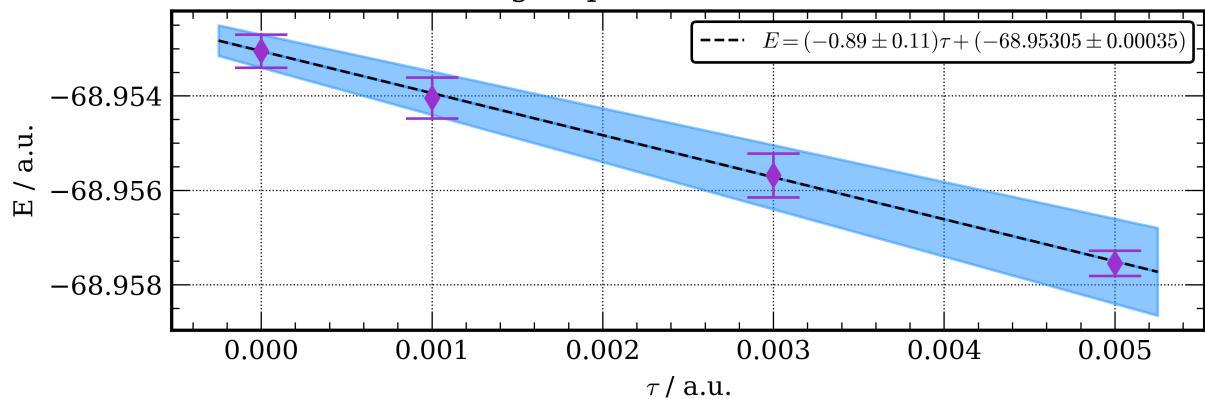
SD/HF
aug-cc-pVDZ+7s7p large ts R = 4.0 EBE



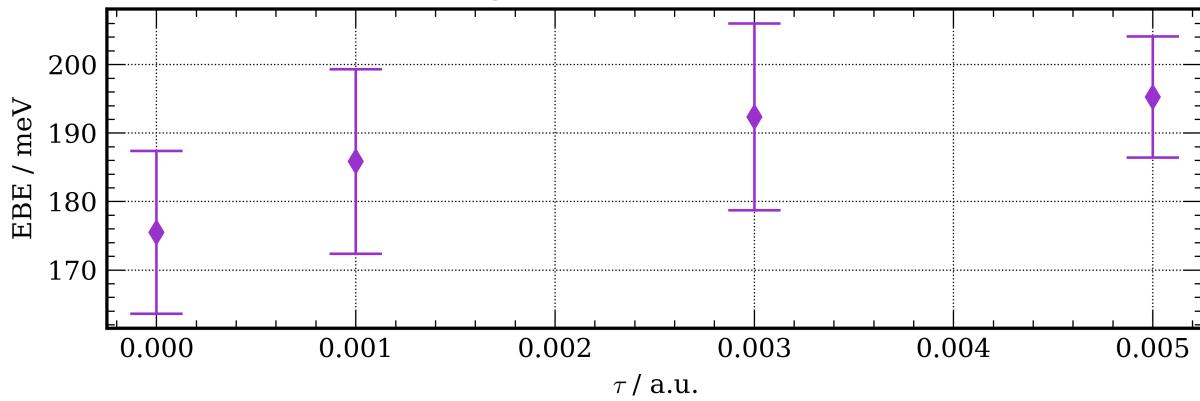
SD/HF
aug-cc-pVDZ R = 4.0 neutral



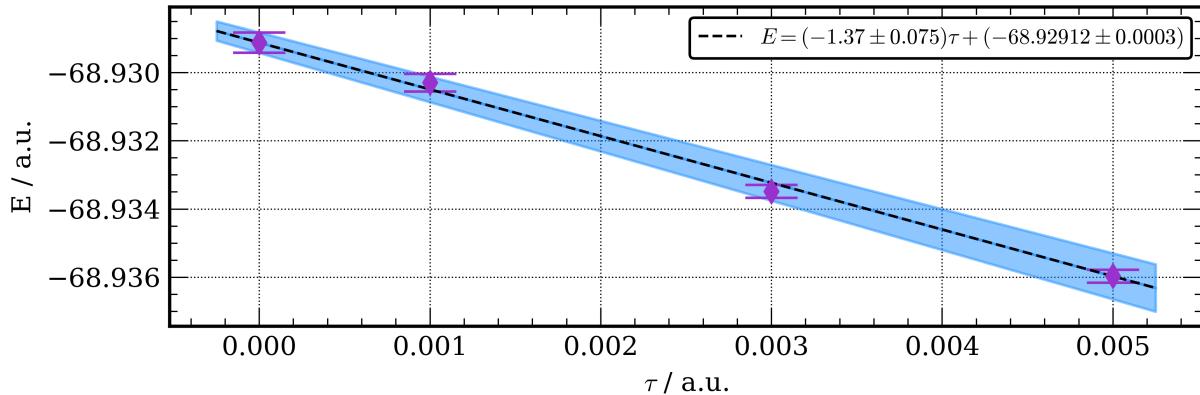
SD/HF
aug-cc-pVDZ R = 4.0 anion



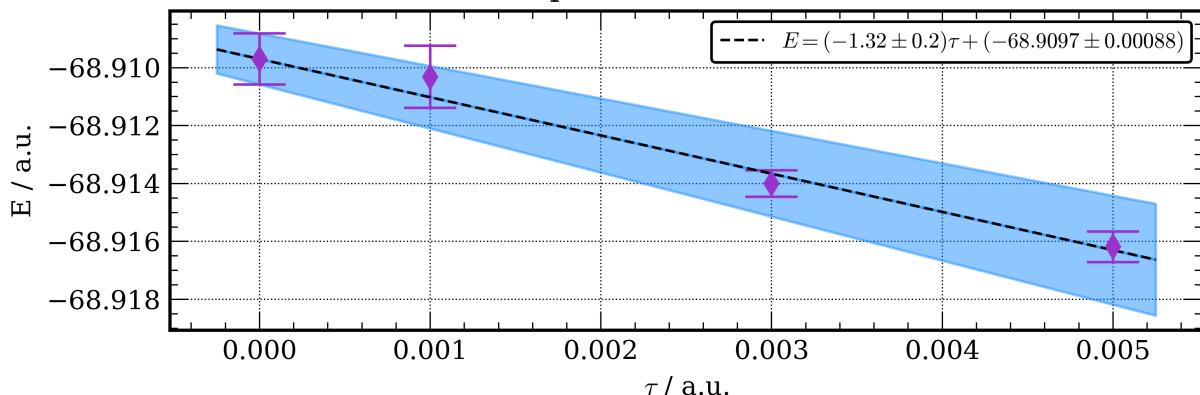
SD/HF
aug-cc-pVDZ R = 4.0 EBE

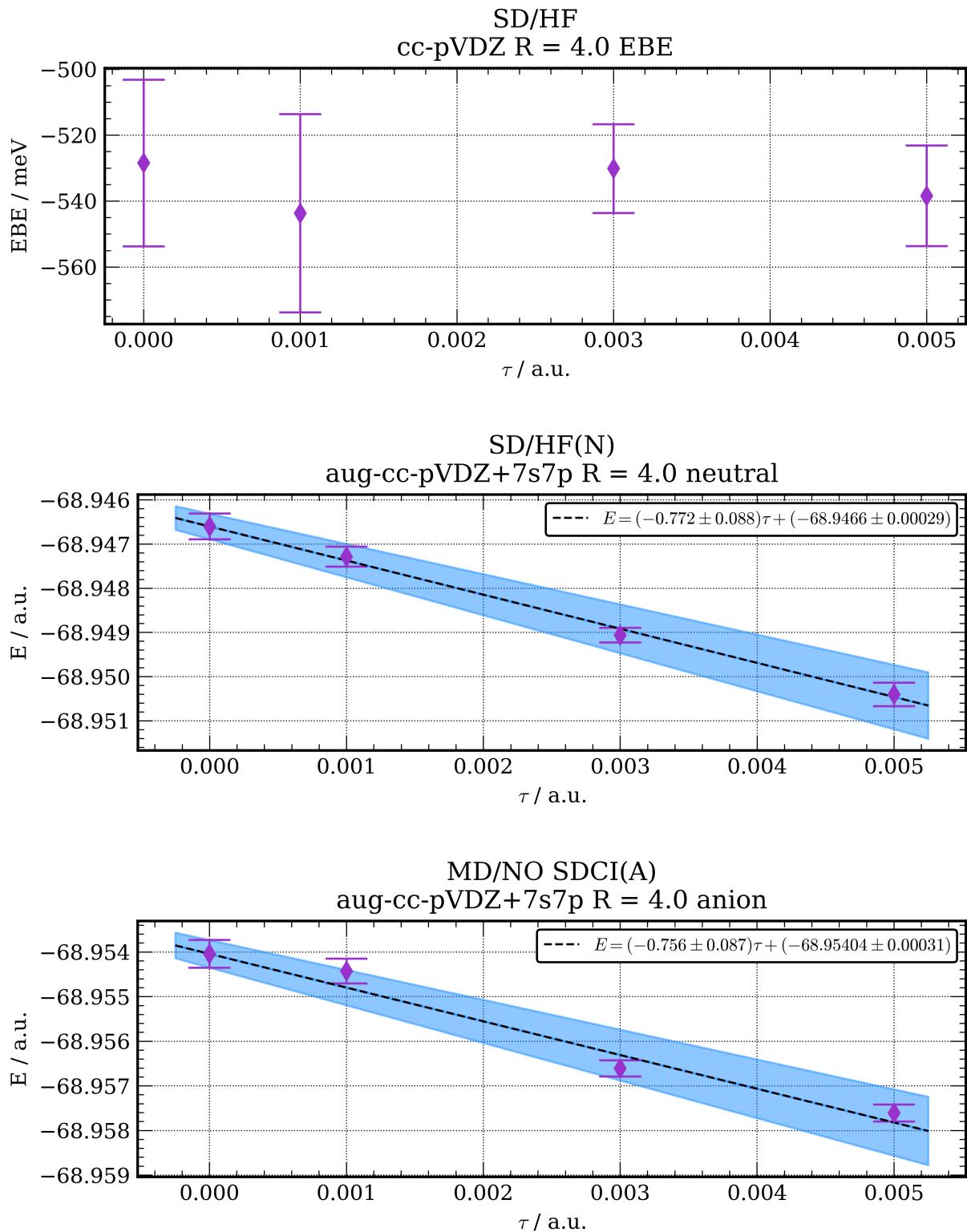


SD/HF
cc-pVDZ R = 4.0 neutral

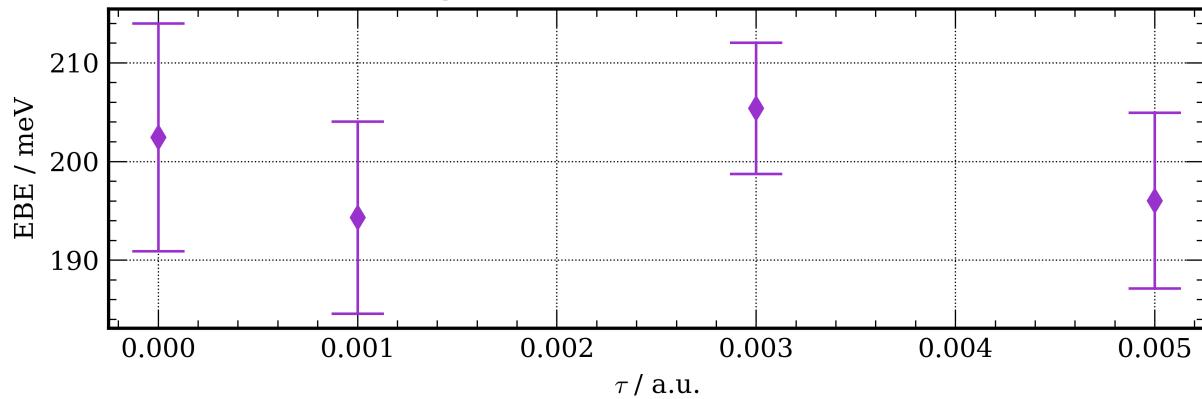


SD/HF
cc-pVDZ R = 4.0 anion

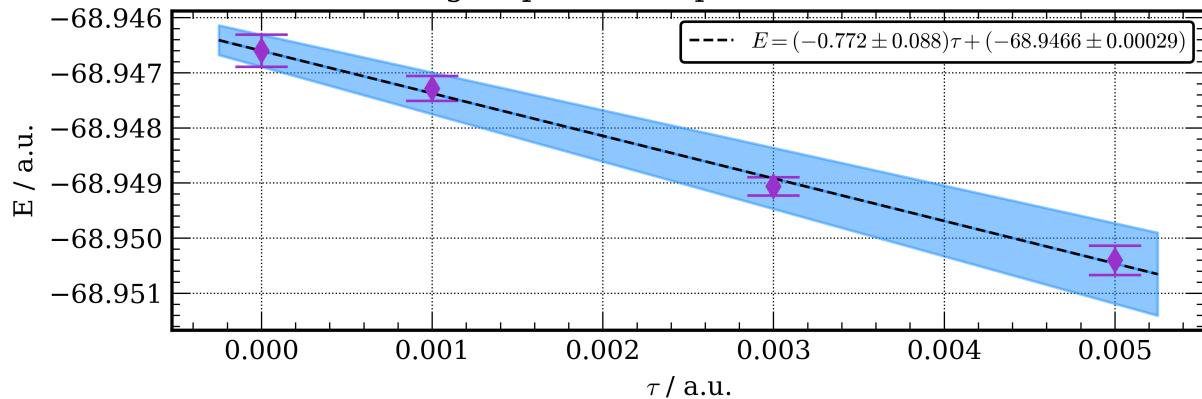




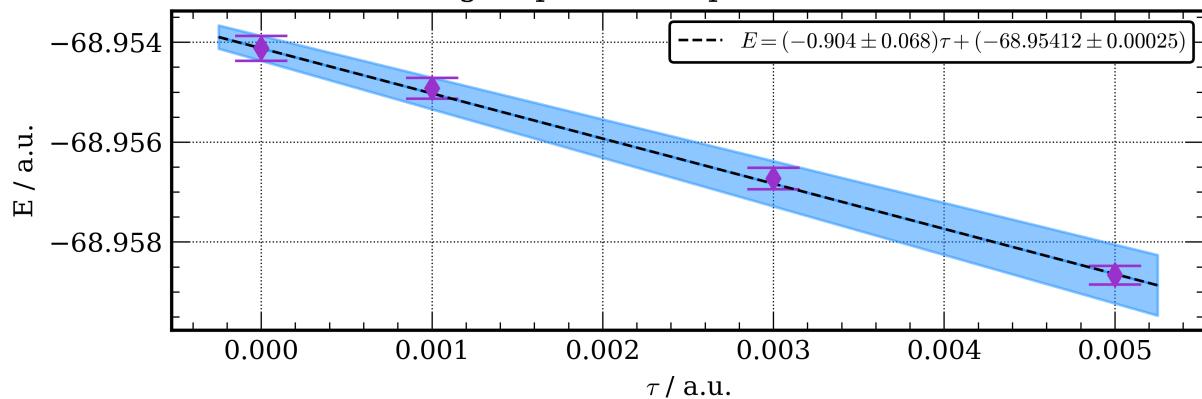
SD/HF(N)//MD/NO SDCI(A)
aug-cc-pVDZ+7s7p R = 4.0 EBE



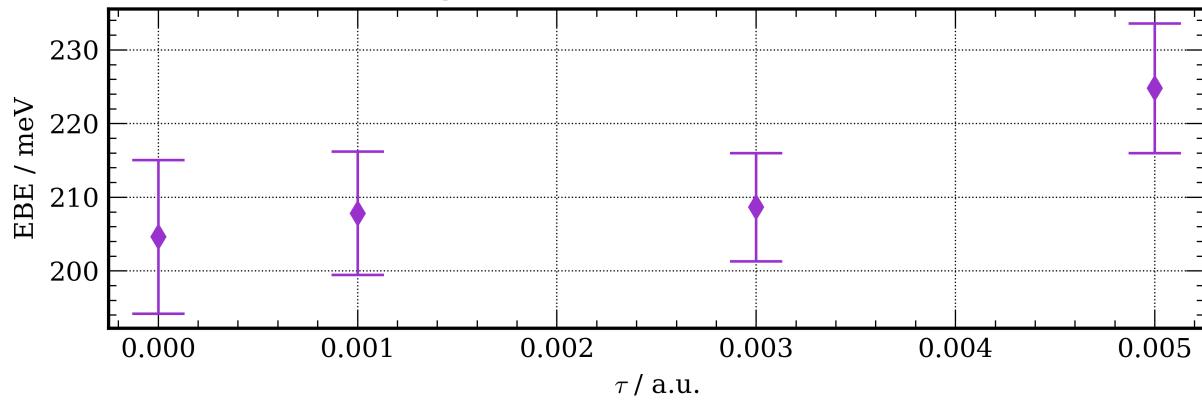
SD/HF(N)
aug-cc-pVDZ+7s7p R = 4.0 neutral



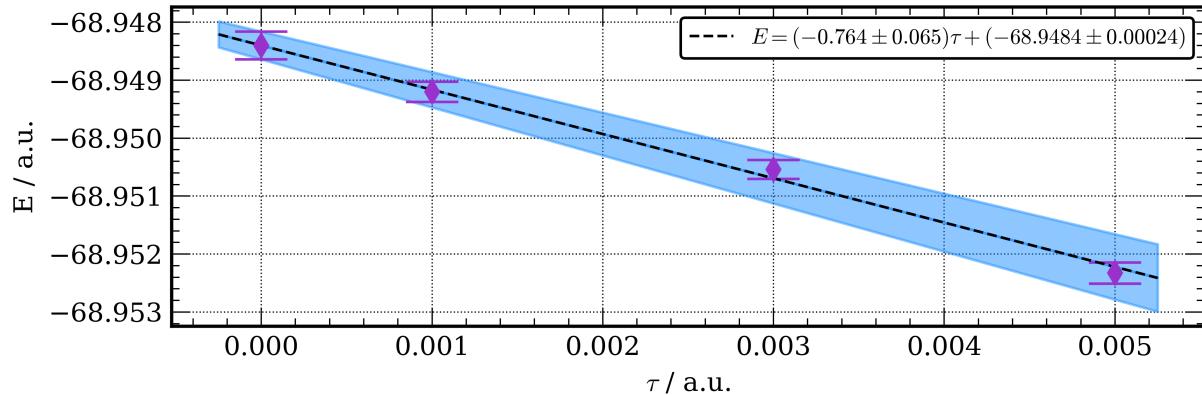
SD/NO SDCI(A)
aug-cc-pVDZ+7s7p R = 4.0 anion



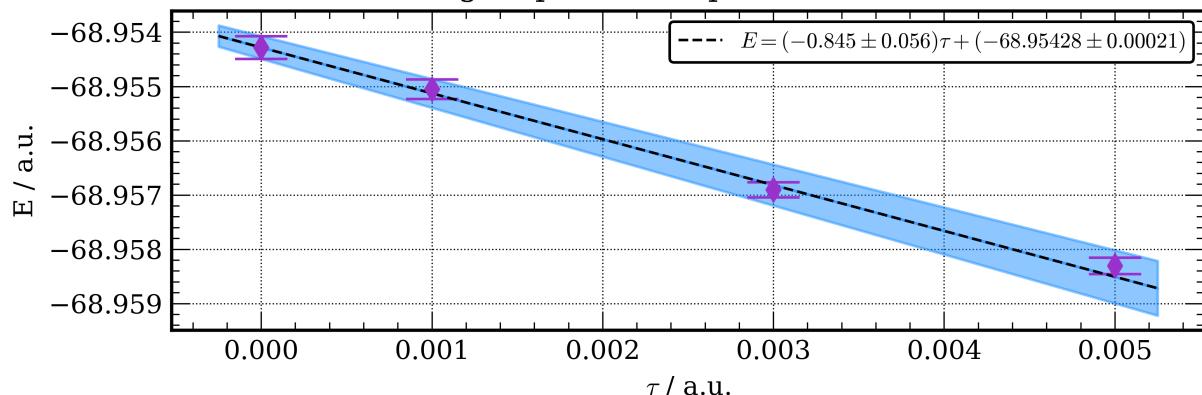
SD/HF(N)//SD/NO SDCI(A)
aug-cc-pVDZ+7s7p R = 4.0 EBE



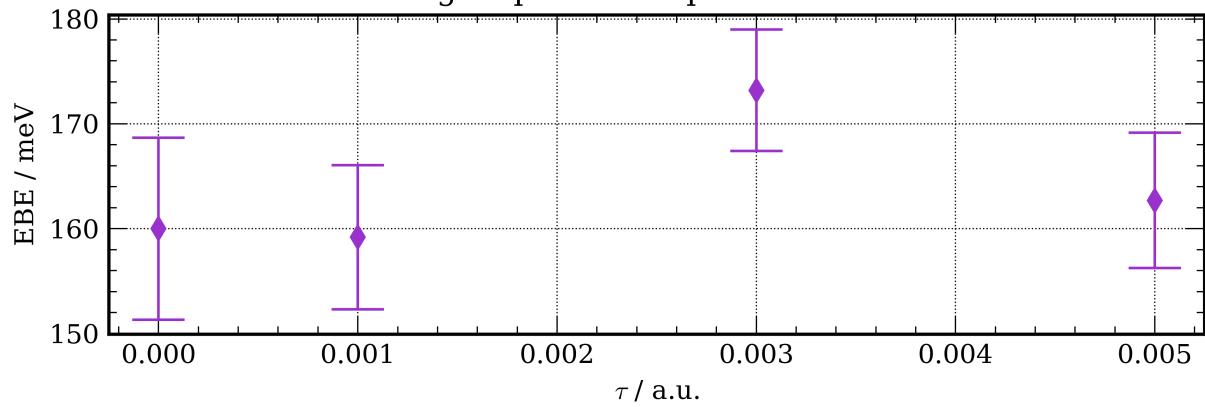
SD/HF(N)
aug-cc-pVDZ+7s7p R = 7.0 neutral



SD/NO SDCI(A)
aug-cc-pVDZ+7s7p R = 7.0 anion



SD/HF(N)//SD/NO SDCI(A)
aug-cc-pVDZ+7s7p R = 7.0 EBE



AFQMC energies

Table 13: AFQMC total energies at used to calculate the electron binding energies. AFQMC calculation details can be found in the manuscript.

R (Å)	neutral energy (Ha)	anion energy (Ha)
4	305.464859 ± 0.000274	-305.472004 ± 0.000266
7	-305.432678 ± 0.000121	-305.439315 ± 0.000121

Radial orbital density plots

The radial orbital density, $\psi^2(r)$, plots are created by integrating over angular portion of the norm of the single particle wave function.

$$\psi^2(r) = \int_0^{2\pi} \int_0^\pi \psi^2(r, \theta, \phi) r^2 \sin(\theta) dr d\theta d\phi \quad (75)$$

Discretizing this expression using a uniform radial grid and a Lebedev-Laikov quadrature for the angular components, yields a form that can be readily evaluated.

$$\psi^2(r_i) = 4\pi r_i^2 \sum_j^{N^{ang}} w_j^{ang} \psi^2(r_i, \theta_j, \phi_j) \quad (76)$$

The function $\psi^2(r_i)$ from Equation 76 can be plotted with the points r_i serving as the abscissa. Since the singly occupied orbitals are normalized, the proximity of the sum of the radial quadrature to unity is used as a check.

$$\sum_i^{N^{rad}} \psi^2(r_i) w_i^{rad} = \sum_i^{N^{rad}} \psi^2(r_i) \Delta r \approx 1 \quad (77)$$

Required software versions

Required software	version
numpy	1.18.4
quadpy	0.16.2
pyscf	1.7.0
cclib	1.6.3

Step 1: Generating a Molden file Molden files were generated using `cclib`, with the exception of the natural orbital from the CIPSI calculations. Since QuantumPackage is not supported by `cclib`, Molden files were created using the native utility in QuantumPackage 2.0. For the Molden files generated with `cclib`, the `-g/-ghost` flag indicates the presence of a ghost atom. By default the only molecular orbitals can be written to a Molden file, therefore the `-n/-naturalorbitals` flag was created to allow natural orbitals to be written in place of molecular orbitals. This flag is not yet available in the official distribution, but a request to incorporate it in the official distribution has been opened (<https://github.com/cclib/cclib/pull/948>).

```
$ ccwrite molden -g "X" -n QUANTUM_CHEMISTRY_OUTPUT_FILE
```

Step 2: Integrating over the angular components of the singly occupied orbital

quadpy was used to generate the Lebedev-Laikov integration weights and points. The singly occupied molecular/natural orbital was evaluated at these points using PySCF.

```
import numpy
import quadpy
import pyscf
import pyscf.tools

filename = "FILENAME.molden"
r_max = 100
num_radial_pts = 1000
mo_idx = 20
# use pyscf to load the molden
mol, mo_energy, mo_coeff, mo_occ, irrep_labels, spins = pyscf.tools.molden.load(
    filename
)
# extract the singly occupied orbital coefficients
singly_occ_orb = mo_coeff[:, mo_idx]
# generate the angular points and weights using quadpy
lebedev_laikov = quadpy.u3.schemes["lebedev_131"]()
angular_pts = lebedev_laikov.theta_phi
angular_weights = lebedev_laikov.weights
num_angular_pts = len(angular_pts[0])
# generate the radial points and weights using numpy
radial_pts = numpy.linspace(r_max, 0, num_radial_pts, endpoint=False)[::-1]
radial_weights = numpy.ones_like(radial_pts) * (radial_pts[1] - radial_pts[0])
# a helper function to convert radial and angular points to cartesian
def sph2cart(r, theta_phi):
```

```

theta = theta_phi[0]
phi = theta_phi[1]

x = r * numpy.cos(theta) * numpy.sin(phi)
y = r * numpy.sin(theta) * numpy.sin(phi)
z = r * numpy.cos(phi)

return numpy.vstack((x, y, z)).T

# integrate over the angular points for each radial point

values = []

for r in radial_pts:

    r_pts = r * numpy.ones(num_angular_pts)
    coords = sph2cart(r_pts, angular_pts)
    ao = mol.eval_gto("GT0val_cart", coords)
    value = angular_weights @ ao @ singly_occ_orb
    values.append(4 * numpy.pi * r * r * value ** 2)

# output the values

values = numpy.array(values)

numpy.savetxt("{}_values.txt".format(filename), values)
numpy.savetxt("{}_r.txt".format(filename), radial_pts)

# check the norm of the orbital

print(radial_weights @ values)

```

A.4 Supplemental Material for Chapter 5

This appendix includes both a more detailed view of the averaged data presented in the manuscript and presents the data for methods that did not produce viable fragments when paired with the current molecular representations. For the test systems below, the ARI for each molecule in the test set is given. Each plot corresponds to a clustering algorithm, where each molecular representation is included. The data for the water clusters is presented in Section A.5 followed by additional information for the methylthiophenes in Section A.5. For the raw data, structures, and

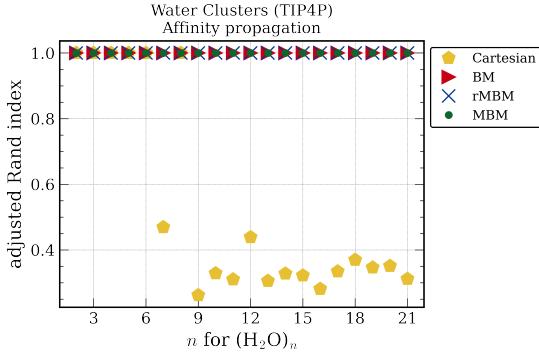


Figure 17: Water cluster test set treated with affinity propagation clustering. All molecular representations are shown and their markings indicated in the legend.

production scripts please see at GitHub hosted data at https://github.com/amandadumi/a1_supporting_info.

Water Clusters For the test set of water clusters, the ARI across the test set are reported for each structure. A more detailed view of agglomerative, spectral, and k -means clustering are provided. Additionally, data is provided for the clustering methods that did not provide chemically-relevant fragments, either inconsistently or not at all.

The additional methods explored which did not produce viable fragments with the current system, either inconsistently or not at all, are also included. Affinity propagation, shown in Figure 25. Mean shift, shown in Figure 26. For affinity propagation, the Cartesian descriptor resulted in a cluster number greater than the number of water monomers for systems that included more than 6 monomers. This resulted in segmented covalent bonds. This likely results from the fact that the difference in the representation between bonded and non-bonded pairs was not sensitive enough. Affinity propagation can be sensitive to the preference of each data-point, which basically describes how likely a given point is to be a cluster center. In this work, we also looked to tuning this preference value by trying a few different options. The typical affinity matrix construction was attempted by placing the sum of distances for a given molecule along the diagonal. A second approach was explored which placed the sum of the covalent-radii determined bond matrix valued for a given atom along the diagonal. We then looked to instead incorporating the Mayer

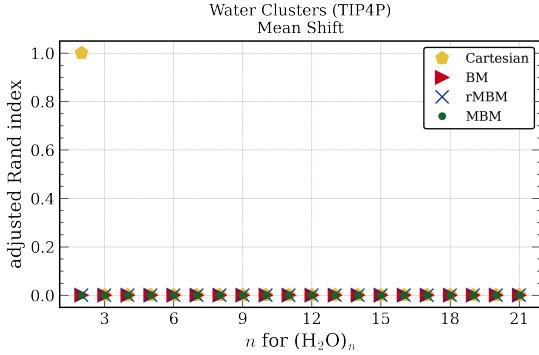


Figure 18: Water cluster test set treated with mean shift clustering. All molecular representations are shown and their markings indicated in the legend. This clustering algorithm with our current descriptor formulation did not produce meaningful fragments.

bond matrix of values for a given molecule along the diagonal. The later did slightly improve the results, though not quite enough to give us reliable results. The quantum mechanics-based (QM) descriptors, rMbm and Mbm, when used with with affinity propagation performed well, though this performance was not maintained in the methylthiophene molecular test sets as seen in the next section. The mean shift clustering results are shown in Figure 26. This algorithm did not perform well with the current descriptors as can be seen by the ARI of 0 for all molecules. Again, more clusters than the number of water monomers were chosen.

Methylthiophenes

The result below are for the methylthiophene test set of molecules at both the HF and ω -B97X-D level of theory. Methods which did not produce molecular fragments in accordance with chemical intuition when used to cluster the current molecular representations are also included.

Agglomerative clustering results are shown in Figure 27, k -means clustering results are shown in Figure 28, spectral clustering shown in Figure 29. The effect of the level of theory used to optimize the structure had on the clustering/molecular representation can be seen by comparing the right and left image from each figure. Generally, agglomerative clustering was able to perform well for the molecular representation explored, except the rounded Mayer bond matrix which had lost too much detail in the descriptor, making any meaningful groupings difficult for this algorithm.

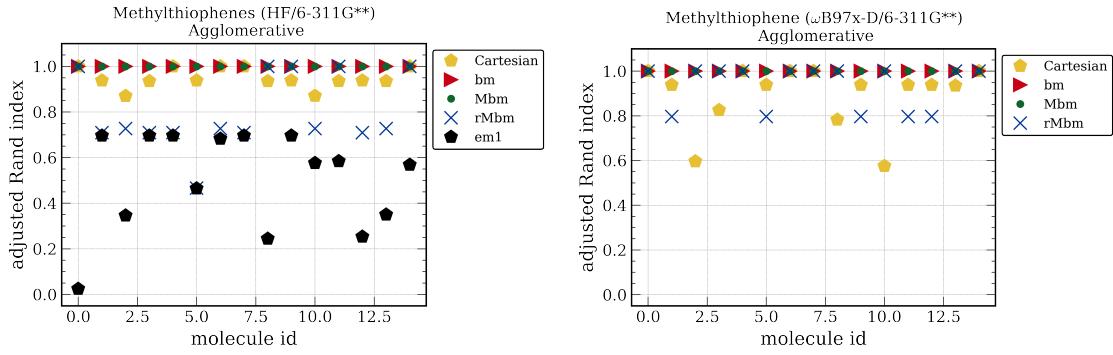


Figure 19: Methylthiophene test set treated with agglomerative clustering. The structures were optimized at either the Hartree-Fock level of theory (left) or ω -B97X-D level of theory (right).

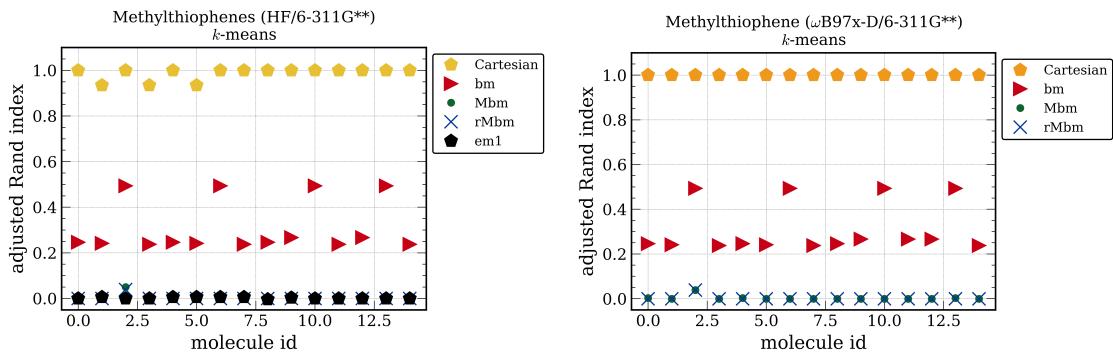


Figure 20: Methylthiophene test set treated with k -means clustering. The structures were optimized at either the Hartree-Fock level of theory (left) or ω -B97X-D level of theory (right).

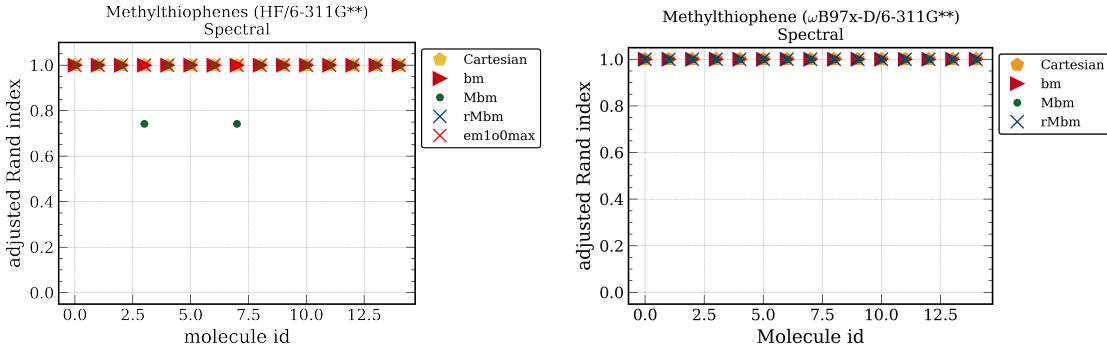


Figure 21: Methylthiophene test set treated with spectral clustering, which worked well for all descriptors. The structures were optimized at either the Hartree-Fock level of theory (left) or ω -B97X-D level of theory (right). All molecular representations are shown and their markings indicated in the legend. Two molecules resulting in different clustering upon the improved level of theory.

Worth noting is that agglomerative clustering was improved when a better level of theory was used to generate the Cartesian descriptor, while more variation was introduced for the rounded Mayer bond Matrix. The k -means clustering with the Cartesian descriptor show a few structures form fragments with slight variations of the correct fragments, while other descriptors show very little or no variation. Generally performance is poor for all but the Cartesian descriptor. Spectral clustering was the most robust, showing a high ARI for all representations. The effect of the level of theory is small, but two molecules in the test set did depend on the higher level of theory to produce the correct clusters.

Results for affinity propagation clustering results are included in Figure 30. Oddly, for affinity propagation, the Cartesian descriptor which performed poorly in the water cluster test set, is the only descriptor that produces reliable fragments for the methylthiophene test set. Although we do see good performance of affinity propagation in both test sets, the reliance on molecular representation warrants further exploration and potentially a more widely applicable descriptor generated for this descriptor. Since a single representation was not viable for both test sets, this method was not pursued further. Mean shift clustering results are shown in ???. Mean shift clustering on the ex-

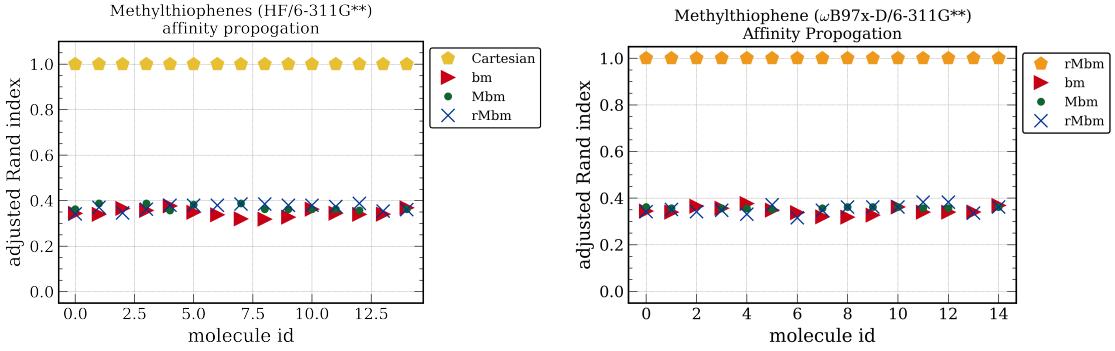


Figure 22: Methylthiophene test set treated with affinity propagation clustering. The structures optimized were at either the Hartree-Fock level of theory (left) or ω -B97X-D level of theory (right). All molecular representations are shown and their markings indicated in the legend.

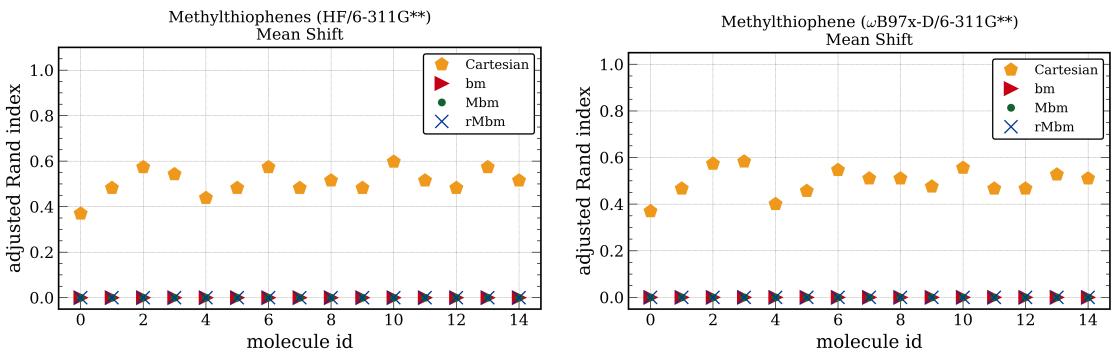


Figure 23: Methylthiophene test set treated with mean shift clustering. The structures were optimized at either the Hartree-Fock level of theory (left) or ω -B97X-D level of theory (right). The explored representations were not formulated for mean shift, and the clustering results were not viable.

plored descriptors results in a low ARI value for all representation indicating the current molecular descriptors are not constructed in a way with compliments this clustering method.

Silyl Ketenes

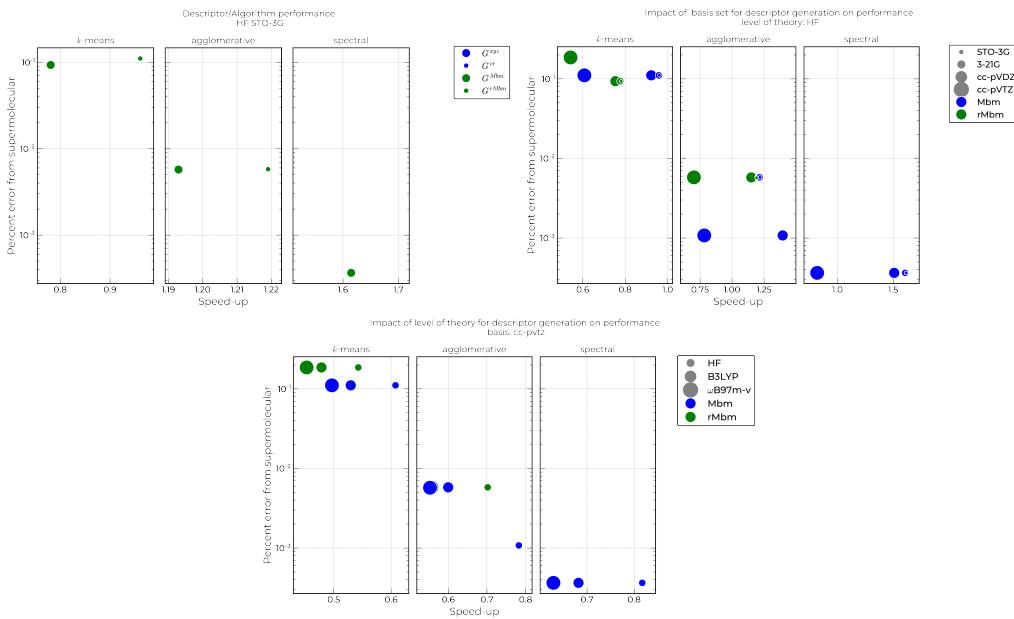


Figure 24: Percent error of the energy and the speedup over the supermolecular calculation for the SK dimer to assess the performance of the fragment approaches and descriptor quality.

A.5 Supplemental Material for Chapter 5

This appendix includes both a more detailed view of the averaged data presented in the manuscript and presents the data for methods that did not produce viable fragments when paired with the current molecular representations. For the test systems below, the ARI for each molecule in the test set is given. Each plot corresponds to a clustering algorithm, where each molecular representation is included. The data for the water clusters is presented in Section A.5 followed by additional information for the methylthiophenes in Section A.5. For the raw data, structures, and production scripts please see at GitHub hosted data at https://github.com/amandadumi/a1_supporting_info.

Water Clusters For the test set of water clusters, the ARI across the test set are reported for each structure. A more detailed view of agglomerative, spectral, and *k*-means clustering are provided.

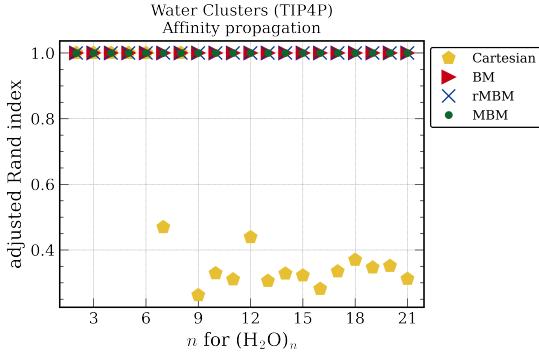


Figure 25: Water cluster test set treated with affinity propagation clustering. All molecular representations are shown and their markings indicated in the legend.

Additionally, data is provided for the clustering methods that did not provide chemically-relevant fragments, either inconsistently or not at all.

The additional methods explored which did not produce viable fragments with the current system, either inconsistently or not at all, are also included. Affinity propagation, shown in Figure 25. Mean shift, shown in Figure 26. For affinity propagation, the Cartesian descriptor resulted in a cluster number greater than the number of water monomers for systems that included more than 6 monomers. This resulted in segmented covalent bonds. This likely results from the fact that the difference in the representation between bonded and non-bonded pairs was not sensitive enough. Affinity propagation can be sensitive to the preference of each data-point, which basically describes how likely a given point is to be a cluster center. In this work, we also looked to tuning this preference value by trying a few different options. The typical affinity matrix construction was attempted by placing the sum of distances for a given molecule along the diagonal. A second approach was explored which placed the sum of the covalent-radii determined bond matrix valued for a given atom along the diagonal. We then looked to instead incorporating the Mayer bond matrix of values for a given molecule along the diagonal. The later did slightly improve the results, though not quite enough to give us reliable results. The quantum mechanics-based (QM) descriptors, rMBM and MBM, when used with affinity propagation performed well, though this performance was not maintained in the methylthiophene molecular test sets as seen in the next

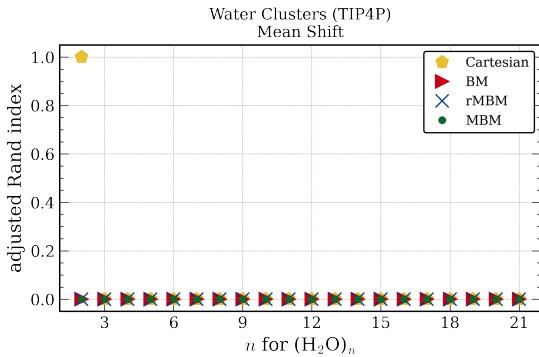


Figure 26: Water cluster test set treated with mean shift clustering. All molecular representations are shown and their markings indicated in the legend. This clustering algorithm with our current descriptor formulation did not produce meaningful fragments.

section. The mean shift clustering results are shown in Figure 26. This algorithm did not perform well with the current descriptors as can be seen by the ARI of 0 for all molecules. Again, more clusters than the number of water monomers were chosen.

Methylthiophenes

The result below are for the methylthiophene test set of molecules at both the HF and ω -B97X-D level of theory. Methods which did not produce molecular fragments in accordance with chemical intuition when used to cluster the current molecular representations are also included.

Agglomerative clustering results are shown in Figure 27, k -means clustering results are shown in Figure 28, spectral clustering shown in Figure 29. The effect of the level of theory used to optimize the structure had on the clustering/molecular representation can be seen by comparing the right and left image from each figure. Generally, agglomerative clustering was able to perform well for the molecular representation explored, except the rounded Mayer bond matrix which had lost too much detail in the descriptor, making any meaningful groupings difficult for this algorithm. Worth noting is that agglomerative clustering was improved when a better level of theory was used to generate the Cartesian descriptor, while more variation was introduced for the rounded Mayer bond Matrix. The k -means clustering with the Cartesian descriptor show a few structures form fragments with slight variations of the correct fragments, while other descriptors show very little or

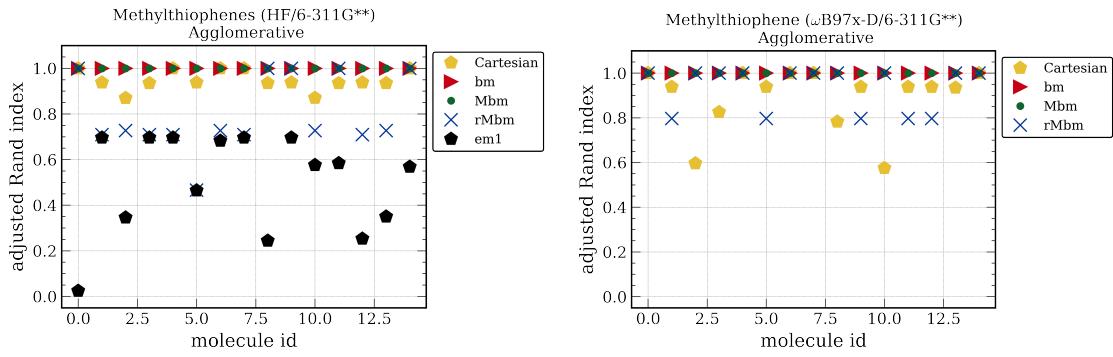


Figure 27: Methylthiophene test set treated with agglomerative clustering. The structures were optimized at either the Hartree-Fock level of theory (left) or ω -B97X-D level of theory (right).

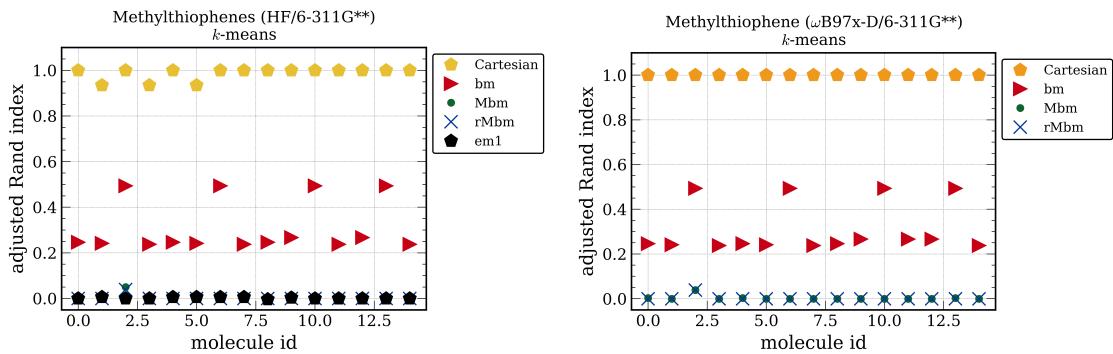


Figure 28: Methylthiophene test set treated with k -means clustering. The structures were optimized at either the Hartree-Fock level of theory (left) or ω -B97X-D level of theory (right).

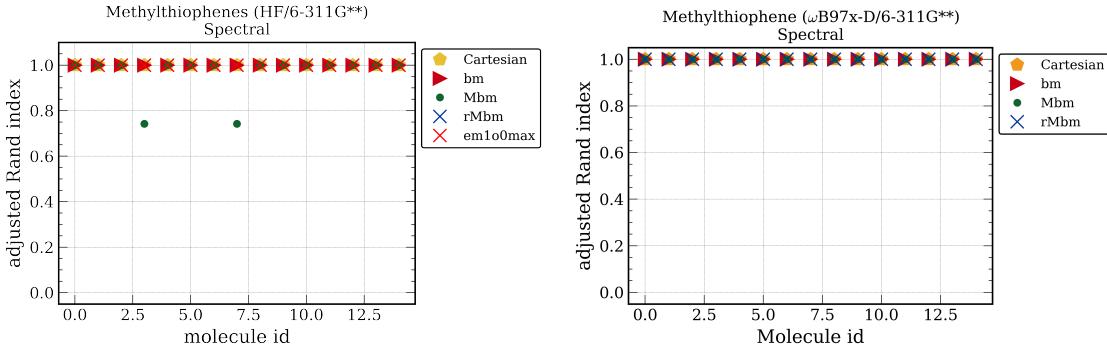


Figure 29: Methylthiophene test set treated with spectral clustering, which worked well for all descriptors. The structures were optimized at either the Hartree-Fock level of theory (left) or ω -B97X-D level of theory (right). All molecular representations are shown and their markings indicated in the legend. Two molecules resulting in different clustering upon the improved level of theory.

no variation. Generally performance is poor for all but the Cartesian descriptor. Spectral clustering was the most robust, showing a high ARI for all representations. The effect of the level of theory is small, but two molecules in the test set did depend on the higher level of theory to produce the correct clusters.

Results for affinity propagation clustering results are included in Figure 30. Oddly, for affinity propagation, the Cartesian descriptor which performed poorly in the water cluster test set, is the only descriptor that produces reliable fragments for the methylthiophene test set. Although we do see good performance of affinity propagation in both test sets, the reliance on molecular representation warrants further exploration and potentially a more widely applicable descriptor generated for this descriptor. Since a single representation was not viable for both test sets, this method was not pursued further. Mean shift clustering results are shown in ???. Mean shift clustering on the explored descriptors results in a low ARI value for all representation indicating the current molecular descriptors are not constructed in a way with compliments this clustering method.

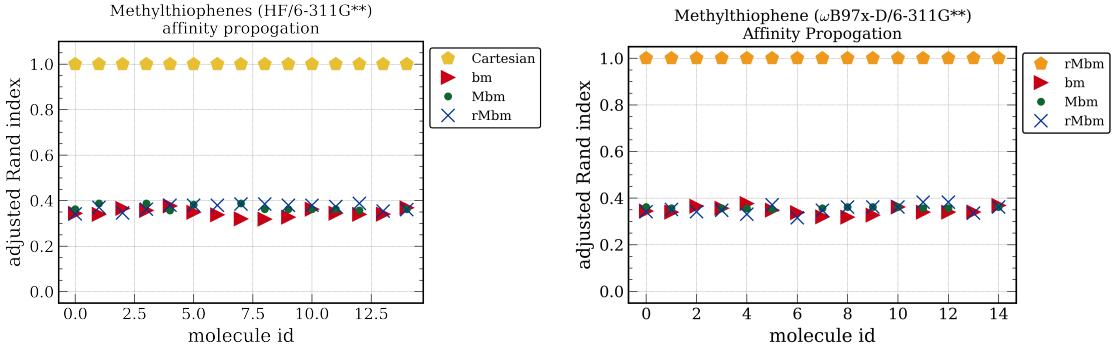


Figure 30: Methylthiophene test set treated with affinity propagation clustering. The structures optimized were at either the Hartree-Fock level of theory (left) or ω -B97X-D level of theory (right). All molecular representations are shown and their markings indicated in the legend.

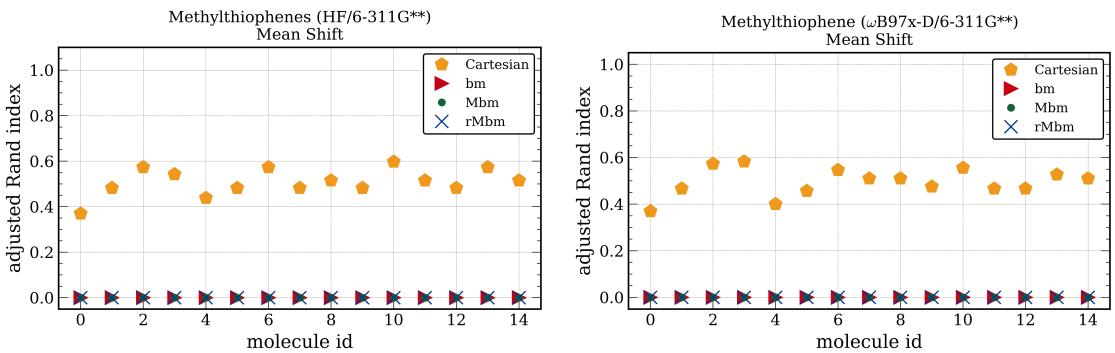


Figure 31: Methylthiophene test set treated with mean shift clustering. The structures were optimized at either the Hartree-Fock level of theory (left) or ω -B97X-D level of theory (right). The explored representations were not formulated for mean shift, and the clustering results were not viable.

Silyl Ketenes

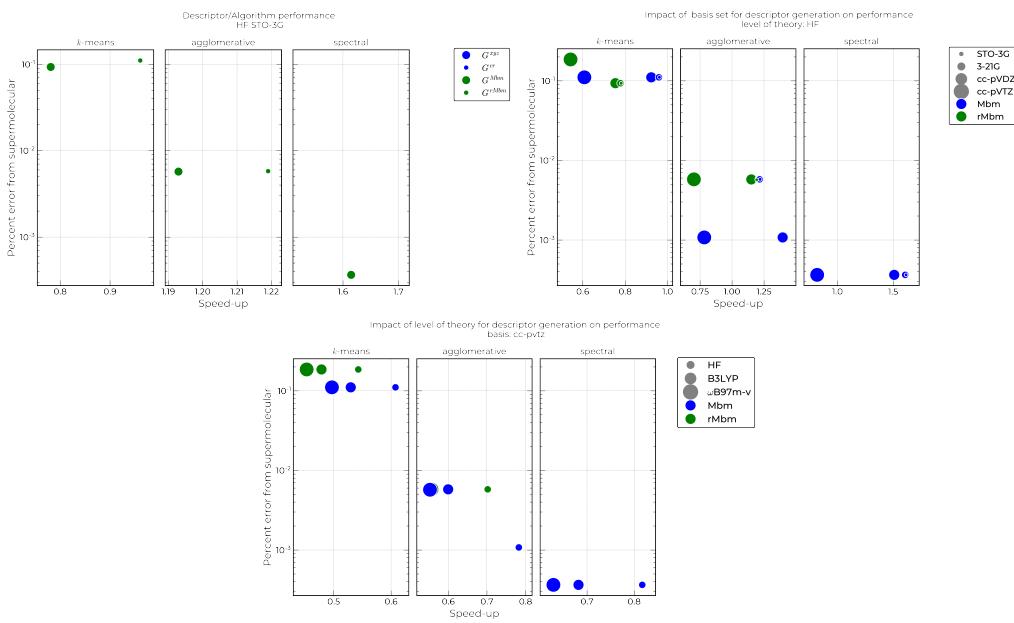


Figure 32: Percent error of the energy and the speedup over the supermolecular calculation for the SK dimer to assess the performance of the fragment approaches and descriptor quality.

Appendix Bibliography

- ¹P. A. M. Dirac and R. H. Fowler, “Quantum mechanics of many-electron systems”, Proceedings of the Royal Society of London. Series A, Containing Papers of a Mathematical and Physical Character **123**, 714–733 (1929).
- ²T. J. Hughes, R. A. Shaw, and S. P. Russo, “Computational investigations of dispersion interactions between small molecules and graphene-like flakes”, *J. Phys. Chem. A* **124**, 9552–9561 (2020).
- ³S. Niaz, T. Manzoor, and A. H. Pandith, “Hydrogen storage: Materials, methods and perspectives”, *Renew. Sustain. Energy Rev* **50**, 457–469 (2015).
- ⁴A. Fomkin, A. Pribylov, I. Men’shchikov, A. Shkolin, O. Aksyutin, A. Ishkov, K. Romanov, and E. Khozina, “Adsorption-based hydrogen storage in activated carbons and model carbon structures”, *Reactions* **2**, 209–226 (2021).
- ⁵C. Ataca, E. Aktürk, S. Ciraci, and H. Ustunel, “High-capacity hydrogen storage by metallized graphene”, *Appl. Phys. Lett.* **93**, 043123 (2008).
- ⁶G. K. Dimitrakakis, E. Tylianakis, and G. E. Froudakis, “Pillared graphene: A new 3-D network nanostructure for enhanced hydrogen storage”, *Nano Lett.* **8**, 3166–3170 (2008).
- ⁷M. Bartolomei, M. I. Hernández, J. Campos-Martínez, R. Hernández-Lamoneda, and G. Giorgi, “Permeation of chemisorbed hydrogen through graphene: A flipping mechanism elucidated”, *Carbon* **178**, 718–727 (2021).
- ⁸J. Ma, A. Michaelides, and D. Alfè, “Binding of hydrogen on benzene, coronene, and graphene from quantum Monte Carlo calculations”, *J. Chem. Phys.* **134**, 134701 (2011).
- ⁹J. Ma, A. Michaelides, D. Alfè, L. Schimka, G. Kresse, and E. Wang, “Adsorption and diffusion of water on graphene from first principles”, *Phys. Rev. B* **84**, 033402 (2011).
- ¹⁰H. Shin, Y. Luo, A. Benali, and Y. Kwon, “Diffusion Monte Carlo study of O₂ adsorption on single layer graphene”, *Phys. Rev. B* **100**, 075430 (2019).
- ¹¹J. Ahn, I. Hong, G. Lee, H. Shin, A. Benali, and Y. Kwon, “Adsorption of a single Pt atom on graphene: spin crossing between physisorbed triplet and chemisorbed singlet states”, *Phys. Chem. Chem. Phys.* **23**, 22147–22154 (2021).

- ¹²A. K. Geim, “Graphene: Status and prospects”, *Science* **324**, 1530–1534 (2009).
- ¹³K. S. Novoselov, A. K. Geim, S. V. Morozov, D. Jiang, M. I. Katsnelson, I. V. Grigorieva, S. V. Dubonos, and A. A. Firsov, “Two-dimensional gas of massless Dirac fermions in graphene”, *Nature* **438**, 197–200 (2005).
- ¹⁴O. K. Alekseeva, I. V. Pushkareva, A. S. Pushkarev, and V. N. Fateev, “Graphene and graphene-like materials for hydrogen energy”, **15**, 273–300 (2020).
- ¹⁵F. Hummel, T. Tsatsoulis, and A. Grüneis, “Low rank factorization of the Coulomb integrals for periodic coupled cluster theory”, *J. Chem. Phys.* **146**, 124105 (2017).
- ¹⁶W. M. C. Foulkes, L. Mitas, R. J. Needs, and G. Rajagopal, “Quantum monte carlo simulations of solids”, *Rev. Mod. Phys.* **73**, 33–83 (2001).
- ¹⁷J. G. Brandenburg, A. Zen, M. Fitzner, B. Ramberger, G. Kresse, T. Tsatsoulis, A. Grüneis, A. Michaelides, and D. Alfè, “Physisorption of water on graphene: Subchemical accuracy from many-body electronic structure methods”, *J. Phys. Chem. Lett.* **10**, 358–368 (2019).
- ¹⁸V. K. Voora, A. Kairalapova, T. Sommerfeld, and K. D. Jordan, “Theoretical approaches for treating non-valence correlation-bound anions”, *J. Chem. Phys.* **147**, 214114 (2017).
- ¹⁹V. K. Voora and K. D. Jordan, “Nonvalence Correlation-Bound Anion State of C₆F₆: Doorway to Low-Energy Electron Capture”, *J. Phys. Chem. A* **118**, 7201–7205 (2014-09, 2014).
- ²⁰T. H. Choi and K. D. Jordan, “Model potential study of non-valence correlation-bound anions of (C₆₀)_n clusters: the role of electric field-induced charge transfer”, *Faraday Discuss.* **217**, 547–560 (2019-07, 2019).
- ²¹V. K. Voora and K. D. Jordan, “Nonvalence Correlation-Bound Anion States of Polycyclic Aromatic Hydrocarbons”, *J. Phys. Chem. Lett.* **6**, 3994–3997 (2015-10, 2015).
- ²²T. Sommerfeld, B. Bhattacharai, V. P. Vysotskiy, and L. S. Cederbaum, “Correlation-bound anions of NaCl clusters”, *J. Chem. Phys.* **133**, 114301 (2010-09, 2010).
- ²³V. G. Bezchastnov, V. P. Vysotskiy, and L. S. Cederbaum, “Anions of xenon clusters bound by long-range electron correlations”, *Phys. Rev. Lett.* **107**, 133401 (2011-09, 2011).
- ²⁴A. Kairalapova, K. D. Jordan, M. F. Falcetta, D. K. Steiner, B. L. Sutter, and J. S. Gowen, “Prediction of a non-valence temporary anion state of (NaCl)2”, *J. Phys. Chem. B* **123**, 9198–9205 (2019).

- ²⁵A. Kairalapova, K. D. Jordan, D. N. Maienshein, M. C. Fair, and M. F. Falcetta, “Prediction of a nonvalence temporary anion shape resonance for a model (H₂O)₄ system”, *J. Phys. Chem. A* **123**, 2719–2726 (2019).
- ²⁶J. P. Rogers, C. S. Anstöter, and J. R. R. Verlet, “Ultrafast dynamics of low-energy electron attachment via a non-valence correlation-bound state”, *Nat. Chem.* **10**, 341–346 (2018-03, 2018).
- ²⁷J. P. Rogers, C. S. Anstöter, and J. R. R. Verlet, “Evidence of Electron Capture of an Outgoing Photoelectron Wave by a Nonvalence State in (C₆F₆)_n⁻”, *J. Phys. Chem. Lett.* **9**, 2504–2509 (2018-05, 2018).
- ²⁸J. N. Bull, C. S. Anstöter, and J. R. R. Verlet, “Ultrafast valence to non-valence excited state dynamics in a common anionic chromophore”, *Nat. Commun.* **10**, 5820 (2019-12, 2019).
- ²⁹J. R. R. Verlet, C. S. Anstöter, J. N. Bull, and J. P. Rogers, “Role of Nonvalence States in the Ultrafast Dynamics of Isolated Anions”, *J. Phys. Chem. A* **124**, 3507–3519 (2020-05, 2020).
- ³⁰V. K. Voora, L. S. Cederbaum, and K. D. Jordan, “Existence of a Correlation Bound s-Type Anion State of C₆₀”, *J. Phys. Chem. Lett.* **4**, 849–853 (2013-03, 2013).
- ³¹V. K. Voora and K. D. Jordan, “Nonvalence Correlation-Bound Anion States of Spherical Fullerenes”, *Nano Lett.* **14**, 4602–4606 (2014-08, 2014).
- ³²B. Huron, J. P. Malrieu, and P. Rancurel, “Iterative perturbation calculations of ground and excited state energies from multiconfigurational zeroth-order wavefunctions”, *J. Chem. Phys.* **58**, 5745–5759 (1973).
- ³³Y. Garniron, T. Applencourt, K. Gasperich, A. Benali, A. Ferté, J. Paquier, B. Pradines, R. Assaraf, P. Reinhardt, J. Toulouse, P. Barbaresco, N. Renon, G. David, J. P. Malrieu, M. Vérit, M. Caffarel, P. F. Loos, E. Giner, and A. Scemama, “Quantum Package 2.0: An Open-Source Determinant-Driven Suite of Programs”, *J. Chem. Theory Comput.* **15**, 3591–3609 (2019).
- ³⁴M. Dash, S. Moroni, A. Scemama, and C. Filippi, “Perturbatively selected configuration-interaction wave functions for efficient geometry optimization in quantum monte carlo”, *Journal of Chemical Theory and Computation* **14**, 4176–4182 (2018).
- ³⁵F. Kossoski, Y. Damour, and P.-F. Loos, “Hierarchy configuration interaction: Combining seniority number and excitation degree”, *The Journal of Physical Chemistry Letters* **13**, 4342–4349 (2022).

- ³⁶M. Dash, S. Moroni, C. Filippi, and A. Scemama, “Tailoring CIPSI expansions for QMC calculations of electronic excitations: The case study of thiophene”, *Journal of Chemical Theory and Computation* **17**, 3426–3434 (2021).
- ³⁷F. D. Malone, A. Benali, M. A. Morales, M. Caffarel, P. R. C. Kent, and L. Shulenburger, “Systematic comparison and cross-validation of fixed-node diffusion Monte Carlo and phaseless auxiliary-field quantum Monte Carlo in solids”, *Phys. Rev. B* **102**, 161104 (2020).
- ³⁸S. Pathak, B. Busemeyer, J. N. B. Rodrigues, and L. K. Wagner, “Excited states in variational Monte Carlo using a penalty method”, *The Journal of Chemical Physics* **154**, 034101 (2021).
- ³⁹S. D. Pineda Flores and E. Neuscamman, “Excited state specific multi-slater jastrow wave functions”, *The Journal of Physical Chemistry A* **123**, 1487–1497 (2019).
- ⁴⁰J. M. Herbert, “Fantasy versus reality in fragment-based quantum chemistry”, *Journal of Chemical Physics* **151**, 1–38 (2019).
- ⁴¹J. Y. Hasegawa, “Fragment-based configuration interaction wave function to calculate environmental effect on excited states in proteins and solutions”, *Chemical Physics Letters* **571**, 77–81 (2013).
- ⁴²S. Müller, “Flexible heuristic algorithm for automatic molecule fragmentation: application to the UNIFAC group contribution model”, *Journal of Cheminformatics* **11**, 57 (2019).
- ⁴³V. Deev and M. A. Collins, “Approximate ab initio energies by systematic molecular fragmentation”, *The Journal of Chemical Physics* **122**, 154102 (2005).
- ⁴⁴M. A. Collins and V. A. Deev, “Accuracy and efficiency of electronic energies from systematic molecular fragmentation”, *Journal of Chemical Physics* **125**, 10.1063/1.2347710 (2006).
- ⁴⁵M. A. Addicoat and M. A. Collins, “Accurate treatment of nonbonded interactions within systematic molecular fragmentation”, *The Journal of Chemical Physics* **131**, 104103 (2009).
- ⁴⁶M. A. Collins, “Systematic fragmentation of large molecules by annihilation”, *Phys. Chem. Chem. Phys.* **14**, 7744–7751 (2012).
- ⁴⁷R. Kobayashi, M. A. Addicoat, A. T. Gilbert, R. D. Amos, and M. A. Collins, “The SMFA program for quantum chemistry calculations on large molecules”, *WIREs Computational Molecular Science* **9**, e1413 (2019).
- ⁴⁸T. D. Huan, A. Mannodi-Kanakkithodi, C. Kim, V. Sharma, G. Pilania, and R. Ramprasad, “A polymer dataset for accelerated property prediction and design”, *Scientific Data* **3**, 1–10 (2016).

- ⁴⁹M. Zheng, N. W. Moriarty, Y. Xu, J. R. Reimers, P. V. Afonine, and M. P. Waller, “Solving the scalability issue in quantum-based refinement: Q|R#1”, *Acta Crystallographica Section D Structural Biology* **73**, 1020–1028 (2017).
- ⁵⁰M. Wolter, M. von Looz, H. Meyerhenke, and C. R. Jacob, “Systematic Partitioning of Proteins for Quantum-Chemical Fragmentation Methods Using Graph Algorithms”, *Journal of Chemical Theory and Computation* **17**, 1355–1367 (2021).
- ⁵¹T. Helgaker, P. Jørgensen, and J. Olsen, “Second quantization”, in *Molecular electronic-structure theory* (John Wiley & Sons, Ltd, 2000), pp. 1–33.
- ⁵²P. Hohenberg and W. Kohn, “Inhomogeneous electron gas”, *Phys. Rev.* **136**, B864–B871 (1964).
- ⁵³R. M. Martin, L. Reining, and D. M. Ceperley, *Interacting electrons: Theory and computational approaches* (Cambridge University Press, Cambridge, 2016).
- ⁵⁴W. Kohn and L. J. Sham, “Self-consistent equations including exchange and correlation effects”, *Phys. Rev.* **140**, A1133–A1138 (1965).
- ⁵⁵R. Jones, “Density functional theory: Its origins, rise to prominence, and future”, *Reviews of Modern Physics* **87**, 897–923 (2015).
- ⁵⁶J. P. Perdew, A. Ruzsinszky, J. Tao, V. N. Staroverov, G. E. Scuseria, and G. I. Csonka, “Prescription for the design and selection of density functional approximations: more constraint satisfaction with fewer fits”, *J Chem Phys* **123**, 62201 (2005-08-08, 2005).
- ⁵⁷N. Mardirossian and M. Head-Gordon, “Thirty years of density functional theory in computational chemistry: an overview and extensive assessment of 200 density functionals”, *Molecular Physics* **115**, 2315–2372 (2017).
- ⁵⁸D. M. Ceperley and B. J. Alder, “Ground state of the electron gas by a stochastic method”, *Phys. Rev. Lett.* **45**, 566–569 (1980).
- ⁵⁹T. Schmidt and S. Küümel, “One- and many-electron self-interaction error in local and global hybrid functionals”, *Phys. Rev. B* **93**, 165120 (2016-04, 2016).
- ⁶⁰S. Śmiga and L. A. Constantin, “Unveiling the physics behind hybrid functionals”, *J. Phys. Chem. A* **124**, 5606–5614 (2020-07-09, 2020).
- ⁶¹M. Brütting, H. Bahmann, and S. Kümmel, “Hybrid functionals with local range separation: Accurate atomization energies and reaction barrier heights”, *J. Chem. Phys.* **156**, 104109 (2022).

- ⁶²S. Grimme, S. Ehrlich, and L. Goerigk, “Effect of the damping function in dispersion corrected density functional theory”, *Journal of Computational Chemistry* **32**, 1456–1465 (2011).
- ⁶³S. Grimme, A. Hansen, J. G. Brandenburg, and C. Bannwarth, “Dispersion-corrected mean-field electronic structure methods”, *Chemical Reviews* **116**, 5105–5154 (2016).
- ⁶⁴P. Tavadze, R. Boucher, G. Avendaño-Franco, K. X. Kocan, S. Singh, V. Dovale-Farelo, W. Ibarra-Hernández, M. B. Johnson, D. S. Mebane, and A. H. Romero, “Exploring DFT+U parameter space with a bayesian calibration assisted by markov chain monte carlo sampling”, *npj Computational Materials* **7**, 182 (2021-11-11, 2021).
- ⁶⁵M. Yu, S. Yang, C. Wu, and N. Marom, “Machine learning the hubbard U parameter in DFT+U using bayesian optimization”, *npj Computational Materials* **6**, 180 (2020-11-27, 2020).
- ⁶⁶N. F. Mott, “Metal-insulator transition”, *Rev. Mod. Phys.* **40**, 677–683 (1968-10, 1968).
- ⁶⁷M. Imada, A. Fujimori, and Y. Tokura, “Metal-insulator transitions”, *Rev. Mod. Phys.* **70**, 1039–1263 (1998).
- ⁶⁸F. Malet and P. Gori-Giorgi, “Strong Correlation in Kohn-Sham Density Functional Theory”, *Physical Review Letters* **109**, 246402 (2012).
- ⁶⁹A. D. Becke, “Density functionals for static, dynamical, and strong correlation”, *The Journal of Chemical Physics* **138**, 074109 (2013).
- ⁷⁰A. D. Becke, “Communication: Two-determinant mixing with a strong-correlation density functional”, *The Journal of Chemical Physics* **139**, 021104 (2013).
- ⁷¹B. Himmetoglu, A. Floris, S. de Gironcoli, and M. Cococcioni, “Hubbard-corrected DFT energy functionals: The LDA+U description of correlated systems”, *International Journal of Quantum Chemistry* **114**, 14–49 (2014).
- ⁷²S. A. Tolba, K. M. Gameel, B. A. Ali, H. A. Almossalami, and N. K. Allam, “The DFT+U: Approaches, accuracy, and applications”, in *Density functional calculations*, edited by G. Yang (IntechOpen, Rijeka, 2018).
- ⁷³R. Tesch and P. M. Kowalski, “Hubbard U parameters for transition metals from first principles”, *Phys. Rev. B* **105**, 195153 (2022).
- ⁷⁴E. Pavarini, “Solving the strong-correlation problem in materials”, *La Rivista del Nuovo Cimento* **44**, 597–640 (2021).

- ⁷⁵R. Requist and E. K. U. Gross, “Model Hamiltonian for strongly correlated systems: Systematic, self-consistent, and unique construction”, Phys. Rev. B **99**, 125114 (2019-03, 2019).
- ⁷⁶J. Tilly, P. V. Sriluckshmy, A. Patel, E. Fontana, I. Rungger, E. Grant, R. Anderson, J. Tennyson, and G. H. Booth, “Reduced density matrix sampling: Self-consistent embedding and multiscale electronic structure on current generation quantum computers”, Phys. Rev. Research **3**, 033230 (2021).
- ⁷⁷C. J. C. Scott and G. H. Booth, “Extending density matrix embedding: A static two-particle theory”, Phys. Rev. B **104**, 245114 (2021-12, 2021).
- ⁷⁸G. Kotliar, S. Y. Savrasov, K. Haule, V. S. Oudovenko, O. Parcollet, and C. A. Marianetti, “Electronic structure calculations with dynamical mean-field theory”, Rev. Mod. Phys. **78**, 865–951 (2006).
- ⁷⁹L. Bytautas and K. Ruedenberg, “A priori identification of configurational deadwood”, Chemical Physics, Moving Frontiers in Quantum Chemistry: **356**, 64–75 (2009).
- ⁸⁰J. S. M. Anderson, F. Heidar-Zadeh, and P. W. Ayers, “Breaking the curse of dimension for the electronic Schrödinger equation with functional analysis”, Computational and Theoretical Chemistry **1142**, 66–77 (2018).
- ⁸¹C. F. Bender and E. R. Davidson, “Studies in Configuration Interaction: The First-Row Diatomic Hydrides”, Physical Review **183**, 23–30 (1969).
- ⁸²A. A. Holmes, N. M. Tubman, and C. J. Umrigar, “Heat-bath configuration interaction: An efficient selected configuration interaction algorithm inspired by heat-bath sampling”, Journal of Chemical Theory and Computation **12**, 3674–3680 (2016).
- ⁸³A. A. Holmes, C. J. Umrigar, and S. Sharma, “Excited states using semistochastic heat-bath configuration interaction”, Journal of Chemical Physics **147**, 1–11 (2017).
- ⁸⁴G. H. Booth, A. J. W. Thom, and A. Alavi, “Fermion Monte Carlo without fixed nodes: A game of life, death, and annihilation in Slater determinant space”, The Journal of Chemical Physics **131**, 054106 (2009).
- ⁸⁵J. B. Schriber and F. A. Evangelista, “Communication: An adaptive configuration interaction approach for strongly correlated electrons with tunable accuracy”, The Journal of Chemical Physics **144**, 161106 (2016).

- ⁸⁶M. Dash, J. Feldt, S. Moroni, A. Scemama, and C. Filippi, “Excited States with Selected Configuration Interaction-Quantum Monte Carlo: Chemically Accurate Excitation Energies and Geometries”, *Journal of Chemical Theory and Computation* **15**, 4896–4906 (2019).
- ⁸⁷A. Benali, K. Gasperich, K. D. Jordan, T. Applencourt, Y. Luo, M. C. Bennett, J. T. Krogl, L. Shulenburger, P. R. C. Kent, P.-F. Loos, A. Scemama, and M. Caffarel, “Toward a systematic improvement of the fixed-node approximation in diffusion Monte Carlo for solids—A case study in diamond”, *J. Chem. Phys.* **153**, 184111 (2020).
- ⁸⁸A. K. Dutta, M. Saitow, B. Demoulin, F. Neese, and R. Izsák, “A domain-based local pair natural orbital implementation of the equation of motion coupled cluster method for electron attached states”, *J. Chem. Phys.* **150**, 164123 (2019-04, 2019).
- ⁸⁹M. Musiał, “Equation-of-motion coupled-cluster models”, in *Quantum chemistry and dynamics of excited states* (John Wiley & Sons, Ltd, 2020), pp. 77–108.
- ⁹⁰E. A. Hylleraas, “Neue Berechnung der Energie des Heliums im Grundzustande, sowie des tiefsten Terms von Ortho-Helium”, *Z. Physik* **54**, 347–366 (1929).
- ⁹¹A. Bijl, “The lowest wave function of the symmetrical many particles system”, *Physica* **7**, 869–886 (1940).
- ⁹²R. Dingle, “LI. The zero-point energy of a system of particles”, *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science* **40**, 573–578 (1949).
- ⁹³R. Jastrow, “Many-body problem with strong forces”, *Phys. Rev.* **98**, 1479–1484 (1955-06-01, 1955).
- ⁹⁴K. Esler, J. Kim, D. Ceperley, and L. Shulenburger, “Accelerating quantum monte carlo simulations of real materials on GPU clusters”, *Computing in Science and Engineering* **14**, 40–51 (2012).
- ⁹⁵V. Natoli and D. M. Ceperley, “An optimized method for treating long-range potentials”, *Journal of Computational Physics* **117**, 171–178 (1995-03-01, 1995).
- ⁹⁶E. Pavarini, E. Koch, and S. Zhang, eds., *Many-Body Methods for Real Materials*, Vol. 9, Schriften Des Forschungszentrums Jülich. Modeling and Simulation (Forschungszentrum Jülich GmbH Zentralbibliothek, Verlag, Jülich, Sept. 2019).
- ⁹⁷H. F. Trotter, “On the Product of Semi-Groups of Operators”, *Proceedings of the American Mathematical Society* **10**, 545–551 (1959).

- ⁹⁸E. Nelson, “Feynman integrals and the schrödinger equation”, *Journal of Mathematical Physics* **5**, 332–343 (1964).
- ⁹⁹M. Suzuki, “Improved Trotter-like formula”, *Physics Letters A* **180**, 232–234 (1993).
- ¹⁰⁰B. Jeziorski, R. Moszynski, and K. Szalewicz, “Perturbation Theory Approach to Intermolecular Potential Energy Surfaces of van der Waals Complexes”, *Chemical Reviews* **94**, 1887–1930 (1994).
- ¹⁰¹K. Szalewicz, “Symmetry-adapted perturbation theory of intermolecular forces”, *WIREs Computational Molecular Science* **2**, 254–272 (2012).
- ¹⁰²D. S. Lambrecht, “Generalizing energy decomposition analysis to response properties to inform expedited predictive models”, *Computational and Theoretical Chemistry* **1149**, 24–30 (2019).
- ¹⁰³E. J. Berquist and D. S. Lambrecht, “A first principles approach for partitioning linear response properties into additive and cooperative contributions”, *ChemRXiv* **10.26434/c**, 1–25 (2018).
- ¹⁰⁴R. Z. Khaliullin, E. A. Cobar, R. C. Lochan, A. T. Bell, and M. Head-Gordon, “Unravelling the origin of intermolecular interactions using absolutely localized molecular orbitals”, *Journal of Physical Chemistry A* **111**, 8753–8765 (2007).
- ¹⁰⁵Y. Kawashima, E. Lloyd, M. P. Coons, Y. Nam, S. Matsuura, A. J. Garza, S. Johri, L. Huntington, V. Senicourt, A. O. Maksymov, J. H. V. Nguyen, J. Kim, N. Alidoust, A. Zaribafyan, and T. Yamazaki, “Optimizing electronic structure simulations on a trapped-ion quantum computer using problem decomposition”, *Communications Physics* **4**, 245 (2021).
- ¹⁰⁶M. S. Gordon, D. G. Fedorov, S. R. Pruitt, and L. V. Slipchenko, “Fragmentation Methods: A Route to Accurate Calculations on Large Systems”, *Chemical Reviews* **112**, 632–672 (2012).
- ¹⁰⁷L. W. Chung, W. M. C. Sameera, R. Ramozzi, A. J. Page, M. Hatanaka, G. P. Petrova, T. V. Harris, X. Li, Z. Ke, F. Liu, H.-B. Li, L. Ding, and K. Morokuma, “The ONIOM Method and Its Applications”, *Chemical Reviews* **115**, 5678–5796 (2015).
- ¹⁰⁸T. Vreven, K. S. Byun, I. Komáromi, S. Dapprich, J. A. J. Montgomery, K. Morokuma, and M. J. Frisch, “Combining Quantum Mechanics Methods with Molecular Mechanics Methods in ONIOM”, *Journal of Chemical Theory and Computation* **2**, 815–826 (2006).
- ¹⁰⁹L. D. Jacobson, R. M. Richard, K. U. Lao, and J. M. Herbert, “Chapter Two - Efficient Monomer-Based Quantum Chemistry Methods for Molecular and Ionic Clusters”, in , Vol. 9, edited by R. A. Wheeler, *Annual Reports in Computational Chemistry* (Elsevier, 2013), pp. 25–58.

- ¹¹⁰N. J. Mayhall and K. Raghavachari, “Molecules-in-Molecules: An Extrapolated Fragment-Based Approach for Accurate Calculations on Large Molecules and Materials”, *Journal of Chemical Theory and Computation* **7**, 1336–1343 (2011).
- ¹¹¹K. V. J. Jose and K. Raghavachari, “Evaluation of Energy Gradients and Infrared Vibrational Spectra through Molecules-in-Molecules Fragment-Based Approach”, *Journal of Chemical Theory and Computation* **11**, 950–961 (2015).
- ¹¹²A. Saha and K. Raghavachari, “Analysis of Different Fragmentation Strategies on a Variety of Large Peptides: Implementation of a Low Level of Theory in Fragment-Based Methods Can Be a Crucial Factor”, *Journal of Chemical Theory and Computation* **11**, 2012–2023 (2015).
- ¹¹³R. M. Ferullo, N. F. Domancich, and N. J. Castellani, “On the performance of van der Waals corrected-density functional theory in describing the atomic hydrogen physisorption on graphite”, *J. Phys. Chem. A* **500**, 283–286 (2010).
- ¹¹⁴X. Sha and B. Jackson, “First-principles study of the structural and energetic properties of H atoms on a graphite (0001) surface”, *Surf. Sci.* **496**, 318–330 (2002).
- ¹¹⁵L. Jeloaica and V. Sidis, “DFT investigation of the adsorption of atomic hydrogen on a cluster-model graphite surface”, *Chem. Phys. Lett.* **300**, 157–162 (1999).
- ¹¹⁶S. Sahu and G. C. Rout, “Band gap opening in graphene: a short theoretical study”, *Int. Nano Lett.* **7**, 81–89 (2017).
- ¹¹⁷H. González-Herrero, J. M. Gómez-Rodríguez, P. Mallet, M. Moaied, J. J. Palacios, C. Salgado, M. M. Ugeda, J.-Y. Veuillet, F. Yndurain, and I. Brihuega, “Atomic-scale control of graphene magnetism by using hydrogen atoms”, *Science* **352**, 437–441 (2016).
- ¹¹⁸H. González-Herrero, E. C.-d. Río, P. Mallet, J.-Y. Veuillet, J. J. Palacios, J. M. Gómez-Rodríguez, I. Brihuega, and F. Ynduráin, “Hydrogen physisorption channel on graphene: a highway for atomic H diffusion”, *2d Mater.* **6**, 021004 (2019).
- ¹¹⁹T. N. Mihm, A. R. McIsaac, and J. J. Shepherd, “An optimized twist angle to find the twist-averaged correlation energy applied to the uniform electron gas”, *J. Chem. Phys.* **150**, 191101 (2019).
- ¹²⁰J. M. Callahan, M. F. Lange, and T. C. Berkelbach, “Dynamical correlation energy of metals in large basis sets from downfolding and composite approaches”, *J. Chem. Phys.* **154**, 211105 (2021).

- ¹²¹G. H. Booth, A. Grüneis, G. Kresse, and A. Alavi, “Towards an exact description of electronic wavefunctions in real solids”, *Nature* **493**, 365–370 (2013).
- ¹²²T. Schäfer, F. Libisch, G. Kresse, and A. Grüneis, “Local embedding of coupled cluster theory into the random phase approximation using plane waves”, *J. Chem. Phys.* **154**, 011101 (2021).
- ¹²³T. Schäfer, B. Ramberger, and G. Kresse, “Quartic scaling MP2 for solids: A highly parallelized algorithm in the plane wave basis”, *J. Chem. Phys.* **146**, 104101 (2017).
- ¹²⁴J. P. Perdew, K. Burke, and M. Ernzerhof, “Generalized gradient approximation made simple”, *Phys. Rev. Lett.* **77**, 3865–3868 (1996).
- ¹²⁵J. P. Perdew, “Unified theory of exchange and correlation beyond the local density approximation”, in *Electronic structure of solids '91*, Vol. 17, edited by P. Ziesche and H. Eschrig, *Physical Research* (1991), pp. 11–20.
- ¹²⁶T. Roman, W. A. Diño, H. Nakanishi, H. Kasai, T. Sugimoto, and K. Tange, “Hydrogen pairing on graphene”, *Carbon* **45**, 218–220 (2007).
- ¹²⁷. u. \u{S}ljivan\u{c}anin, E. Rauls, L. Hornekær, W. Xu, F. Besenbacher, and B. Hammer, “Extended atomic hydrogen dimer configurations on the graphite(0001) surface”, *J. Chem. Phys.* **131**, 084706 (2009).
- ¹²⁸P. O. Lehtinen, A. S. Foster, Y. Ma, A. V. Krasheninnikov, and R. M. Nieminen, “Irradiation-induced magnetism in graphite: A density functional study”, *Phys. Rev. Lett.* **93**, 187202 (2004).
- ¹²⁹Y. Lin, F. Ding, and B. I. Yakobson, “Hydrogen storage by spillover on graphene as a phase nucleation process”, *Phys. Rev. B* **78**, 041402 (2008).
- ¹³⁰S. Casolo, O. M. Løvvik, R. Martinazzo, and G. F. Tantardini, “Understanding adsorption of hydrogen atoms on graphene”, *J. Chem. Phys.* **130**, 054704 (2009).
- ¹³¹R. H. Miwa, T. B. Martins, and A. Fazzio, “Hydrogen adsorption on boron doped graphene: an ab initio study”, *Nanotechnology* **19**, 155708 (2008).
- ¹³²A. Ishii, M. Yamamoto, H. Asano, and K. Fujiwara, “DFT calculation for adatom adsorption on graphene sheet as a prototype of carbon nanotube functionalization”, *J. Phys. Conf. Ser.* **100**, 052087 (2008).
- ¹³³W. Li, M. Zhao, T. He, C. Song, X. Lin, X. Liu, Y. Xia, and L. Mei, “Concentration dependent magnetism induced by hydrogen adsorption on graphene and single walled carbon nanotubes”, *J. Magn. Magn. Mater.* **322**, 838–843 (2010).

- ¹³⁴D. W. Boukhvalov, M. I. Katsnelson, and A. I. Lichtenstein, “Hydrogen on graphene: Electronic structure, total energy, structural distortions and magnetism from first-principles calculations”, *Phys. Rev. B* **77**, 035427 (2008).
- ¹³⁵M. A. Kim, D. C. Sorescu, S. Amemiya, K. D. Jordan, and H. Liu, “Real-time modulation of hydrogen evolution activity of graphene electrodes using mechanical strain”, *ACS Applied Materials & Interfaces* **14**, 10691–10700 (2022).
- ¹³⁶S. Grimme, J. Antony, S. Ehrlich, and H. Krieg, “A consistent and accurate ab initio parametrization of density functional dispersion correction (DFT-D) for the 94 elements H-Pu”, *J. Chem. Phys.* **132**, 154104 (2010).
- ¹³⁷A. Annaberdiyev, G. Wang, C. A. Melton, M. C. Bennett, L. Shulenburger, and L. Mitas, “A new generation of effective core potentials from correlated calculations: 3d transition metal series”, *J. Chem. Phys.* **149**, 134108 (2018-10, 2018).
- ¹³⁸M. C. Bennett, C. A. Melton, A. Annaberdiyev, G. Wang, L. Shulenburger, and L. Mitas, “A new generation of effective core potentials for correlated calculations”, *J. Chem. Phys.* **147**, 224106 (2017-12, 2017).
- ¹³⁹H. J. Monkhorst and J. D. Pack, “Special points for Brillouin-zone integrations”, *Phys. Rev. B* **13**, 5188–5192 (1976).
- ¹⁴⁰N. Marzari, D. Vanderbilt, A. De Vita, and M. C. Payne, “Thermal contraction and disordering of the al(110) surface”, *Phys. Rev. Lett.* **82**, 3296–3299 (1999).
- ¹⁴¹C. Adamo and V. Barone, “Toward reliable density functional methods without adjustable parameters: The PBE0 model”, *J. Chem. Phys.* **110**, 6158–6170 (1999).
- ¹⁴²A. V. Krukau, O. A. Vydrov, A. F. Izmaylov, and G. E. Scuseria, “Influence of the exchange screening parameter on the performance of screened hybrid functionals”, *J. Chem. Phys.* **125**, 224106 (2006).
- ¹⁴³M. F. Peintinger, D. V. Oliveira, and T. Bredow, “Consistent Gaussian basis sets of triple-zeta valence with polarization quality for solid-state calculations”, *J. Comput. Chem.* **34**, 451–459 (2012).
- ¹⁴⁴H. Kruse and S. Grimme, “A geometrical correction for the inter- and intra-molecular basis set superposition error in Hartree-Fock and density functional theory calculations for large systems”, *J. Chem. Phys.* **136**, 154101 (2012).

- ¹⁴⁵J. G. Brandenburg, M. Alessio, B. Civalleri, M. F. Peintinger, T. Bredow, and S. Grimme, “Geometrical correction for the inter- and intramolecular basis set superposition error in periodic density functional theory calculations”, *J. Phys. Chem. A* **117**, 9282–9292 (2013).
- ¹⁴⁶P. Giannozzi, O. Baseggio, P. Bonfà, D. Brunato, R. Car, I. Carnimeo, C. Cavazzoni, S. de Gironcoli, P. Delugas, F. Ferrari Ruffino, A. Ferretti, N. Marzari, I. Timrov, A. Urru, and S. Baroni, “Quantum ESPRESSO toward the exascale”, *J. Chem. Phys.* **152**, 154105 (2020).
- ¹⁴⁷P. Giannozzi, S. Baroni, N. Bonini, M. Calandra, R. Car, C. Cavazzoni, D. Ceresoli, G. L. Chiarotti, M. Cococcioni, I. Dabo, A. Dal Corso, S. de Gironcoli, S. Fabris, G. Fratesi, R. Gebauer, U. Gerstmann, C. Gougoussis, A. Kokalj, M. Lazzeri, L. Martin-Samos, N. Marzari, F. Mauri, R. Mazzarello, S. Paolini, A. Pasquarello, L. Paulatto, C. Sbraccia, S. Scandolo, G. Sclauzero, A. P. Seitsonen, A. Smogunov, P. Umari, and R. M. Wentzcovitch, “QUANTUM ESPRESSO: a modular and open-source software project for quantum simulations of materials”, *J. Phys.: Condens. Matter* **21**, 395502 (19pp) (2009).
- ¹⁴⁸P. Giannozzi, O. Andreussi, T. Brumme, O. Bunau, M. B. Nardelli, M. Calandra, R. Car, C. Cavazzoni, D. Ceresoli, M. Cococcioni, N. Colonna, I. Carnimeo, A. D. Corso, S. de Gironcoli, P. Delugas, R. A. D. Jr, A. Ferretti, A. Floris, G. Fratesi, G. Fugallo, R. Gebauer, U. Gerstmann, F. Giustino, T. Gorni, J. Jia, M. Kawamura, H.-Y. Ko, A. Kokalj, E. Küçükbenli, M. Lazzeri, M. Marsili, N. Marzari, F. Mauri, N. L. Nguyen, H.-V. Nguyen, A. Otero-de-la-Roza, L. Paulatto, S. Poncé, D. Rocca, R. Sabatini, B. Santra, M. Schlipf, A. P. Seitsonen, A. Smogunov, I. Timrov, T. Thonhauser, P. Umari, N. Vast, X. Wu, and S. Baroni, “Advanced capabilities for materials modelling with QUANTUM ESPRESSO”, *J. Phys.: Condens. Matter* **29**, 465901 (2017).
- ¹⁴⁹R. Dovesi, A. Erba, R. Orlando, C. M. Zicovich-Wilson, B. Civalleri, L. Maschio, M. Rérat, S. Casassa, J. Baima, S. Salustro, and B. Kirtman, “Quantum-mechanical condensed matter simulations with CRYSTAL”, *Wiley Interdiscip. Rev. Comput. Mol. Sci* **8**, e1360 (2018).
- ¹⁵⁰R. Dovesi, F. Pascale, B. Civalleri, K. Doll, N. M. Harrison, I. Bush, P. D’Arco, Y. Noël, M. Rérat, P. Carbonnière, M. Causà, S. Salustro, V. Lacivita, B. Kirtman, A. M. Ferrari, F. S. Gentile, J. Baima, M. Ferrero, R. Demichelis, and M. De La Pierre, “The CRYSTAL code, 1976–2020 and beyond, a long story”, *J. Chem. Phys.* **152**, 204111 (2020).
- ¹⁵¹E. Aprà, E. J. Bylaska, W. A. de Jong, N. Govind, K. Kowalski, T. P. Straatsma, M. Valiev, H. J. J. van Dam, Y. Alexeev, J. Anchell, V. Anisimov, F. W. Aquino, R. Atta-Fynn, J.

Autschbach, N. P. Bauman, J. C. Becca, D. E. Bernholdt, K. Bhaskaran-Nair, S. Bogatko, P. Borowski, J. Boschen, J. Brabec, A. Bruner, E. Cauët, Y. Chen, G. N. Chuev, C. J. Cramer, J. Daily, M. J. O. Deegan, T. H. Dunning, M. Dupuis, K. G. Dyall, G. I. Fann, S. A. Fischer, A. Fonari, H. Fr uchtl, L. Gagliardi, J. Garza, N. Gawande, S. Ghosh, K. Glaesemann, A. W. Götz, J. Hammond, V. Helms, E. D. Hermes, K. Hirao, S. Hirata, M. Jacquelin, L. Jensen, B. G. Johnson, H. Jónsson, R. A. Kendall, M. Klemm, R. Kobayashi, V. Konkov, S. Krishnamoorthy, M. Krishnan, Z. Lin, R. D. Lins, R. J. Littlefield, A. J. Logsdail, K. Lopata, W. Ma, A. V. Marenich, J. Martin del Campo, D. Mejia-Rodriguez, J. E. Moore, J. M. Mullin, T. Nakajima, D. R. Nascimento, J. A. Nichols, P. J. Nichols, J. Nieplocha, A. Otero-de-la-Roza, B. Palmer, A. Panyala, T. Pirojsirikul, B. Peng, R. Peverati, J. Pittner, L. Pollack, R. M. Richard, P. Sadayappan, G. C. Schatz, W. A. Shelton, D. W. Silverstein, D. M. A. Smith, T. A. Soares, D. Song, M. Swart, H. L. Taylor, G. S. Thomas, V. Tipparaju, D. G. Truhlar, K. Tsemekhman, T. Van Voorhis, Á. Vázquez-Mayagoitia, P. Verma, O. Villa, A. Vishnu, K. D. Vogiatzis, D. Wang, J. H. Weare, M. J. Williamson, T. L. Windus, K. Woliński, A. T. Wong, Q. Wu, C. Yang, Q. Yu, M. Zacharias, Z. Zhang, Y. Zhao, and R. J. Harrison, “NWChem: Past, present, and future”, *J. Chem. Phys.* **152**, 184102 (2020).

¹⁵²J. B. Anderson, “Quantum chemistry by random walk: Higher accuracy”, *J. Chem. Phys.* **73**, 3897–3899 (1980).

¹⁵³K. Esler, J. Kim, D. Ceperley, and L. Shulenburger, “Accelerating quantum monte carlo simulations of real materials on GPU clusters”, *Computing in Science Engineering* **14**, 40–51 (2012).

¹⁵⁴N. D. Drummond, M. D. Towler, and R. J. Needs, “Jastrow correlation factor for atoms, molecules, and solids”, *Phys. Rev. B* **70**, 235119 (2004).

¹⁵⁵C. J. Umrigar, J. Toulouse, C. Filippi, S. Sorella, and R. G. Hennig, “Alleviation of the fermion-sign problem by optimization of many-body wave functions”, *Phys. Rev. Lett.* **98**, 110201 (2007).

¹⁵⁶A. Zen, J. G. Brandenburg, A. Michaelides, and D. Alfè, “A new scheme for fixed node diffusion quantum Monte Carlo with pseudopotentials: Improving reproducibility and reducing the trial-wave-function bias”, *J. Chem. Phys.* **151**, 134105 (2019).

- ¹⁵⁷M. Casula, S. Moroni, S. Sorella, and C. Filippi, “Size-consistent variational approaches to nonlocal pseudopotentials: Standard and lattice regularized diffusion Monte Carlo methods revisited”, *J. Chem. Phys.* **132**, 154113 (2010).
- ¹⁵⁸C. Lin, F. H. Zong, and D. M. Ceperley, “Twist-averaged boundary conditions in continuum quantum Monte Carlo algorithms”, *Phys. Rev. E* **64**, 016702 (2001).
- ¹⁵⁹A. Zen, S. Sorella, M. J. Gillan, A. Michaelides, and D. Alfè, “Boosting the accuracy and speed of quantum Monte Carlo: Size consistency and time step”, *Phys. Rev. B* **93**, 241118 (2016).
- ¹⁶⁰D. Cortés-Arriagada, S. Gutiérrez-Oliva, B. Herrera, K. Soto, and A. Toro-Labbé, “The mechanism of chemisorption of hydrogen atom on graphene: Insights from the reaction force and reaction electronic flux”, *J. Chem. Phys.* **141**, 134701 (2014).
- ¹⁶¹H. Jiang, M. Kammler, F. Ding, Y. Dorenkamp, F. R. Manby, A. M. Wodtke, T. F. Miller, A. Kandratsenka, and O. Bünermann, “Imaging covalent bond formation by H atom scattering from graphene”, *Science* **364**, 379–382 (2019).
- ¹⁶²J. Kim, A. D. Baczewski, T. D. Beaudet, A. Benali, M. C. Bennett, M. A. Berrill, N. S. Blunt, E. J. L. Borda, M. Casula, D. M. Ceperley, S. Chiesa, B. K. Clark, R. C. Clay, K. T. Delaney, M. Dewing, K. P. Esler, H. Hao, O. Heinonen, P. R. C. Kent, J. T. Krogel, I. Kylänpää, Y. W. Li, M. G. Lopez, Y. Luo, F. D. Malone, R. M. Martin, A. Mathuriya, J. McMinis, C. A. Melton, L. Mitas, M. A. Morales, E. Neuscamman, W. D. Parker, S. D. P. Flores, N. A. Romero, B. M. Rubenstein, J. A. R. Shea, H. Shin, L. Shulenburger, A. F. Tillack, J. P. Townsend, N. M. Tushman, B. V. D. Goetz, J. E. Vincent, D. C. Yang, Y. Yang, S. Zhang, and L. Zhao, “QMCPACK: an open source ab initio quantum Monte Carlo package for the electronic structure of atoms, molecules and solids”, *J. Phys.: Condens. Matter* **30**, 195901 (2018-04, 2018).
- ¹⁶³P. R. C. Kent, A. Annaberdiyev, A. Benali, M. C. Bennett, E. J. Landinez Borda, P. Doak, H. Hao, K. D. Jordan, J. T. Krogel, I. Kylänpää, J. Lee, Y. Luo, F. D. Malone, C. A. Melton, L. Mitas, M. A. Morales, E. Neuscamman, F. A. Reboreda, B. Rubenstein, K. Saritas, S. Upadhyay, G. Wang, S. Zhang, and L. Zhao, “QMCPACK: Advances in the development, efficiency, and application of auxiliary field and real-space variational and diffusion quantum Monte Carlo”, *J. Chem. Phys.* **152**, 174105 (2020).
- ¹⁶⁴J. T. Krogel, “Nexus: A modular workflow management system for quantum simulation codes”, *Comput. Phys. Commun.* **198**, 154–168 (2016).

- ¹⁶⁵J. D. Hunter, “Matplotlib: A 2D graphics environment”, Computing in Science & Engineering **9**, 90–95 (2007).
- ¹⁶⁶K. Momma and F. Izumi, “VESTA3 for three-dimensional visualization of crystal, volumetric and morphology data”, J. Appl. Crystallogr. **44**, 1272–1276 (2011).
- ¹⁶⁷C. Moller and M. S. Plesset, “Note on an approximation treatment for many-electron systems”, Phys. Rev. **46**, 618–622 (1934-10, 1934).
- ¹⁶⁸K. Raghavachari, G. W. Trucks, J. A. Pople, and M. Head-Gordon, “A fifth-order perturbation comparison of electron correlation theories”, Chem. Phys. Lett. **157**, 479–483 (1989).
- ¹⁶⁹R. C. Lochan and M. Head-Gordon, “Orbital-optimized opposite-spin scaled second-order correlation: An economical method to improve the description of open-shell molecules”, J. Chem. Phys. **126**, 164101 (2007).
- ¹⁷⁰J. Cízek, J. Paldus, and L. Sroubková, “Cluster expansion analysis for delocalized systems”, Int. J. Quantum Chem. **3**, 149–167 (1969).
- ¹⁷¹D. C. Comeau and R. J. Bartlett, “The equation-of-motion coupled-cluster method. Applications to open- and closed-shell reference states”, Chem. Phys. Lett. **207**, 414–423 (1993).
- ¹⁷²R. C. Grimm and R. G. Storer, “Monte-Carlo solution of Schrödinger’s equation”, J. Comput. Phys. **7**, 134–156 (1971).
- ¹⁷³J. B. Anderson, “A random-walk simulation of the Schrödinger equation: H₃⁺”, J. Chem. Phys. **63**, 1499–1503 (1975-08, 1975).
- ¹⁷⁴J. B. Anderson, “Quantum chemistry by random walk. H²P, H₃+D_{3h}¹A[']₁, H₂³Sigma⁺_u, H₄¹Sigma⁺_g, Be¹S”, J. Chem. Phys. **65**, 4121–4127 (1976).
- ¹⁷⁵W. M. C. Foulkes, L. Mitas, R. J. Needs, and G. Rajagopal, “Quantum Monte Carlo simulations of solids”, Rev. Mod. Phys. **73**, 1–51 (2001).
- ¹⁷⁶P. López Ríos, A. Ma, N. D. Drummond, M. D. Towler, and R. J. Needs, “Inhomogeneous backflow transformations in quantum Monte Carlo calculations”, Phys. Rev. E **74**, 066701 (2006-12, 2006).
- ¹⁷⁷R. Blankenbecler, D. J. Scalapino, and R. L. Sugar, “Monte Carlo calculations of coupled boson-fermion systems. I”, Phys. Rev. D **24**, 2278–2286 (1981).
- ¹⁷⁸G. Sugiyama and S. Koonin, “Auxiliary field Monte-Carlo for quantum many-body ground states”, Annals of Physics **168**, 1–26 (1986).

- ¹⁷⁹D. J. Scalapino and R. L. Sugar, “Monte Carlo calculations of coupled boson-fermion systems. II”, Phys. Rev. B **24**, 4295–4308 (1981).
- ¹⁸⁰S. Zhang and H. Krakauer, “Quantum Monte Carlo Method using Phase-Free Random Walks with Slater Determinants”, Phys. Rev. Lett. **90**, 136401 (2003-04, 2003).
- ¹⁸¹S. Zhang, J. Carlson, and J. E. Gubernatis, “Constrained path Monte Carlo method for fermion ground states”, Phys. Rev. B **55**, 7464–7477 (1997-03, 1997).
- ¹⁸²M. Motta and S. Zhang, “Ab initio computations of molecular systems by the auxiliary-field quantum Monte Carlo method”, Wiley Interdiscip. Rev.: Comput. Mol. Sci. **8**, e1364 (2018).
- ¹⁸³S. Zhang, “Ab initio electronic structure calculations by auxiliary-field quantum monte carlo”, in *Handbook of materials modeling : Methods: Theory and modeling* (Springer International Publishing, 2018), pp. 1–27.
- ¹⁸⁴W. Humphrey, A. Dalke, and K. Schulten, “”VMD: Visual molecular dynamics””, J. Mol. Graphics **14**, 33–38 (1996).
- ¹⁸⁵J. F. Stanton and J. Gauss, “Perturbative treatment of the similarity transformed Hamiltonian in equation-of-motion coupled-cluster approximations”, J. Chem. Phys. **103**, 1064–1076 (1995).
- ¹⁸⁶D. A. Matthews and J. F. Stanton, “A new approach to approximate equation-of-motion coupled cluster with triple excitations”, J. Chem. Phys. **145**, 124102 (2016).
- ¹⁸⁷S. A. Kucharski, M. Włoch, M. Musiał, and R. J. Bartlett, “Coupled-cluster theory for excited electronic states: The full equation-of-motion coupled-cluster single, double, and triple excitation method”, J. Chem. Phys. **115**, 8263–8266 (2001).
- ¹⁸⁸K. Kowalski and P. Piecuch, “The active-space equation-of-motion coupled-cluster methods for excited electronic states: Full EOMCCSDt”, J. Chem. Phys. **115**, 643–651 (2001).
- ¹⁸⁹T. H. Dunning, “Gaussian basis sets for use in correlated molecular calculations. I. The atoms boron through neon and hydrogen”, J. Chem. Phys. **90**, 10.1063/1.456153 (1989).
- ¹⁹⁰R. A. Kendall, T. H. Dunning, and R. J. Harrison, “Electron affinities of the first-row atoms revisited. Systematic basis sets and wave functions”, J. Chem. Phys. **96**, 10.1063/1.462569 (1992).
- ¹⁹¹J. F. Stanton, J. Gauss, L. Cheng, M. E. Harding, D. A. Matthews, and P. G. Szalay, *CFOUR, coupled-cluster techniques for computational chemistry, a quantum-chemical program package*.

- ¹⁹²D. A. Matthews, L. Cheng, M. E. Harding, F. Lipparini, S. Stopkowicz, T.-C. Jagau, P. G. Szalay, J. Gauss, and J. F. Stanton, “Coupled-cluster techniques for computational chemistry: The CFOUR program package”, *J. Chem. Phys.* **152**, 214108 (2020).
- ¹⁹³R. Jastrow, “Many-body problem with strong forces”, *Phys. Rev.* **98**, 1479–1484 (1955-06, 1955).
- ¹⁹⁴N. D. Drummond, M. D. Towler, and R. J. Needs, “Jastrow correlation factor for atoms, molecules, and solids”, *Phys. Rev. B* **70**, 235119 (2004-12, 2004).
- ¹⁹⁵A. D. Becke, “Density-functional exchange-energy approximation with correct asymptotic behavior”, *Phys. Rev. A* **38**, 3098–3100 (1988-09, 1988).
- ¹⁹⁶C. Lee, W. Yang, and R. G. Parr, “Development of the Colle-Salvetti correlation-energy formula into a functional of the electron density”, *Phys. Rev. B* **37**, 785–789 (1988-01, 1988).
- ¹⁹⁷S. H. Vosko, L. Wilk, and M. Nusair, “Accurate spin-dependent electron liquid correlation energies for local spin density calculations: a critical analysis”, *Can. J. Phys.* **58**, 1200–1211 (1980).
- ¹⁹⁸P. J. Stephens, F. J. Devlin, C. F. Chabalowski, and M. J. Frisch, “Ab initio calculation of vibrational absorption and circular dichroism spectra using density functional force fields”, *J. Phys. Chem.* **98**, 11623–11627 (1994).
- ¹⁹⁹M. Casula, “Beyond the locality approximation in the standard diffusion Monte Carlo method”, *Physical Review B* **74**, 161102 (2006).
- ²⁰⁰M. W. Schmidt, K. K. Baldridge, J. A. Boatz, S. T. Elbert, M. S. Gordon, J. H. Jensen, S. Koseki, N. Matsunaga, K. A. Nguyen, S. Su, T. L. Windus, M. Dupuis, and J. A. Montgomery, “General atomic and molecular electronic structure system”, *J. Comput. Chem.* **14**, 1347–1363 (1993).
- ²⁰¹M. S. Gordon and M. W. Schmidt, “Chapter 41 - Advances in electronic structure theory: GAMESS a decade later”, in *Theory and Applications of Computational Chemistry*, edited by C. E. Dykstra, G. Frenking, K. S. Kim, and G. E. Scuseria (Elsevier, 2005-01, 2005), pp. 1167–1189.
- ²⁰²G. M. J. Barca, C. Bertoni, L. Carrington, D. Datta, N. De Silva, J. E. Deustua, D. G. Fedorov, J. R. Gour, A. O. Gunina, E. Guidez, T. Harville, S. Irle, J. Ivanic, K. Kowalski, S. S. Leang, H. Li, W. Li, J. J. Lutz, I. Magoulas, J. Mato, V. Mironov, H. Nakata, B. Q. Pham, P. Piecuch, D. Poole, S. R. Pruitt, A. P. Rendell, L. B. Roskop, K. Ruedenberg, T. Sattasathuchana, M. W. Schmidt, J. Shen, L. Slipchenko, M. Sosonkina, V. Sundriyal, A. Tiwari, J. L. Galvez Vallejo,

- B. Westheimer, M. Włoch, P. Xu, F. Zahariev, and M. S. Gordon, “Recent developments in the general atomic and molecular electronic structure system”, *J. Chem. Phys.* **152**, 154102 (2020-04, 2020).
- ²⁰³J. Hubbard, “Calculation of partition functions”, *Phys. Rev. Lett.* **3**, 77 (1959).
- ²⁰⁴J. Lee, F. D. Malone, and M. A. Morales, “Utilizing essential symmetry breaking in auxiliary-field quantum Monte Carlo: Application to the spin gaps of the C₃₆ fullerene and an iron porphyrin model complex”, *J. Chem. Theory Comput.* **16**, 3019–3027 (2020).
- ²⁰⁵B. Rudshteyn, D. Coskun, J. L. Weber, E. J. Arthur, S. Zhang, D. R. Reichman, R. A. Friesner, and J. Shee, “Predicting ligand-dissociation energies of 3d coordination complexes with auxiliary-field quantum monte carlo”, *J. Chem. Theory Comput.* **16**, 3041–3054 (2020).
- ²⁰⁶J. Shee, B. Rudshteyn, E. J. Arthur, S. Zhang, D. R. Reichman, and R. A. Friesner, “On achieving high accuracy in quantum chemical calculations of 3d transition metal-containing systems: A comparison of auxiliary-field quantum monte carlo with coupled cluster, density functional theory, and experiment for diatomic molecules”, *J. Chem. Theory Comput.* **15**, 2346–2358 (2019).
- ²⁰⁷M. Kumar, J. Shee, B. Rudshteyn, D. R. Reichman, R. A. Friesner, C. E. Miller, and J. S. Francisco, “Multiple stable isoprene-ozone complexes reveal complex entrance channel dynamics in the Isoprene+ ozone reaction”, *J. Am. Chem. Soc.*, 10.1021/jacs.0c02360 (2020).
- ²⁰⁸H. Hao, J. Shee, S. Upadhyay, C. Ataca, K. D. Jordan, and B. M. Rubenstein, “Accurate predictions of electron binding energies of dipole-bound anions via quantum Monte Carlo methods”, *J. Phys. Chem. Lett.* **9**, 6185–6190 (2018).
- ²⁰⁹J. Shee, E. J. Arthur, S. Zhang, D. R. Reichman, and R. A. Friesner, “Singlet–Triplet energy gaps of organic biradicals and polyacenes with auxiliary-field quantum monte carlo”, *J. Chem. Theory Comput.* **15**, 4924–4932 (2019).
- ²¹⁰J. Shee, E. J. Arthur, S. Zhang, D. R. Reichman, and R. A. Friesner, “Phaseless auxiliary-field quantum Monte Carlo on graphical processing units”, *J. Chem. Theory Comput.* **14**, 4109–4121 (2018).
- ²¹¹W. Purwanto, H. Krakauer, Y. Virgus, and S. Zhang, “Assessing weak hydrogen binding on Ca+ centers: An accurate many-body study with large basis sets”, *J. Chem. Phys.* **135**, 164105 (2011).

- ²¹²N. M. O'boyle, A. L. Tenderholt, and K. M. Langner, “Cclib: A library for package-independent computational chemistry algorithms”, *J. Comput. Chem.* **29**, 839–845 (2008).
- ²¹³Q. Sun, T. C. Berkelbach, N. S. Blunt, G. H. Booth, S. Guo, Z. Li, J. Liu, J. D. McClain, E. R. Sayfutyarova, S. Sharma, S. Wouters, and G. K.-L. Chan, “PySCF: the Python-based simulations of chemistry framework”, *Wiley Interdiscip. Rev.: Comput. Mol. Sci.* **8**, e1340 (2018).
- ²¹⁴Q. Sun, X. Zhang, S. Banerjee, P. Bao, M. Barbry, N. S. Blunt, N. A. Bogdanov, G. H. Booth, J. Chen, Z.-H. Cui, J. J. Eriksen, Y. Gao, S. Guo, J. Hermann, M. R. Hermes, K. Koh, P. Koval, S. Lehtola, Z. Li, J. Liu, N. Mardirossian, J. D. McClain, M. Motta, B. Mussard, H. Q. Pham, A. Pulkin, W. Purwanto, P. J. Robinson, E. Ronca, E. R. Sayfutyarova, M. Scheurer, H. F. Schurkus, J. E. T. Smith, C. Sun, S.-N. Sun, S. Upadhyay, L. K. Wagner, X. Wang, A. White, J. D. Whitfield, M. J. Williamson, S. Wouters, J. Yang, J. M. Yu, T. Zhu, T. C. Berkelbach, S. Sharma, A. Y. Sokolov, and G. K.-L. Chan, “Recent developments in the PySCF program package”, *J. Chem. Phys.* **153**, 024109 (2020).
- ²¹⁵V. I. Lebedev and D. Laikov, “A quadrature formula for the sphere of the 131st algebraic order of accuracy”, *Doklady Mathematics* **59**, 477–481 (1999).
- ²¹⁶N. Schlömer, N. R. Papior, D. Arnold, M. Ancellin, and R. Zetter, *Nschloe/quadpy v0.15.2*, version v0.15.2, Zenodo, Aug. 2020.
- ²¹⁷J. Lee, D. W. Small, and M. Head-Gordon, “Excited states via coupled cluster theory without equation-of-motion methods: Seeking higher roots with application to doubly excited states and double core hole states”, *J. Chem. Phys.* **151**, 214103 (2019).
- ²¹⁸J. Noga and R. J. Bartlett, “The full CCSDT model for molecular electronic structure”, *J. Chem. Phys.* **86**, 7041–7050 (1987).
- ²¹⁹G. E. Scuseria and H. F. Schaefer, “A new implementation of the full CCSDT model for molecular electronic structure”, *Chem. Phys. Lett.* **152**, 382–386 (1988).
- ²²⁰J. D. Watts and R. J. Bartlett, “The coupled-cluster single, double, and triple excitation model for open-shell single reference functions”, *J. Chem. Phys.* **93**, 6104–6105 (1990).
- ²²¹Y. J. Bomble, J. F. Stanton, M. Kállay, and J. Gauss, “Coupled-cluster methods including non-iterative corrections for quadruple excitations”, *J. Chem. Phys.* **123**, 054101 (2005).

- ²²²T. Yamazaki, S. Matsuura, A. Narimani, A. Saidmuradov, and A. Zaribafyan, *Towards the practical application of near-term quantum computers in quantum chemistry simulations: A problem decomposition approach*, 2018.
- ²²³F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, “Scikit-learn: Machine learning in python”, *Journal of Machine Learning Research* **12**, 2825–2830 (2011).
- ²²⁴R. Sibson, “SLINK: An optimally efficient algorithm for the single-link cluster method”, *The Computer Journal* **16**, 30–34 (1973).
- ²²⁵D. Defays, “An efficient algorithm for a complete link method”, *The Computer Journal* **20**, 364–366 (1977).
- ²²⁶S. Lloyd, “Least squares quantization in PCM”, *IEEE Transactions on Information Theory* **28**, 129–137 (1982).
- ²²⁷H. Steinhaus, “Sur la division des corps matériels en parties”, *Bull. Acad. Pol. Sci., Cl. III* **4**, 801–804 (1957).
- ²²⁸A. Y. Ng, M. I. Jordan, and Y. Weiss, “On spectral clustering: analysis and an algorithm”, *Advances in Neural Information Processing Systems* (2002).
- ²²⁹J. H. Ward, “Hierarchical grouping to optimize an objective function”, *Journal of the American Statistical Association* **58**, 236–244 (1963).
- ²³⁰D. Arthur and S. Vassilvitskii, “How slow is the K-means method?”, in *Proceedings of the twenty-second annual symposium on computational geometry, SCG '06* (2006), pp. 144–153.
- ²³¹X. Yin, I. Neamtiu, S. Patil, and S. T. Andrews, “Implementation-induced Inconsistency and Nondeterminism in Deterministic Clustering Algorithms”, in *2020 IEEE 13th International Conference on Software Testing, Validation and Verification (ICST)* (2020), pp. 231–242.
- ²³²Yu and Shi, “Multiclass spectral clustering”, in *Proceedings ninth IEEE international conference on computer vision* (2003), 313–319 vol.1.
- ²³³J. Shi and J. Malik, “Normalized cuts and image segmentation”, *IEEE Transactions on Pattern Analysis and Machine Intelligence* **22**, 888–905 (2000).

- ²³⁴A. Y. Ng, M. I. Jordan, and Y. Weiss, “On spectral clustering: Analysis and an algorithm”, in Proceedings of the 14th international conference on neural information processing systems: Natural and synthetic, NIPS’01 (2001), pp. 849–856.
- ²³⁵A. B. Georgescu, P. Ren, A. R. Toland, S. Zhang, K. D. Miller, D. W. Apley, E. A. Olivetti, N. Wagner, and J. M. Rondinelli, “Database, Features, and Machine Learning Model to Identify Thermally Driven Metal–Insulator Transition Compounds”, *Chem. Mater.* **33**, 5591–5605 (2021).
- ²³⁶Y. Zhao, Y. Cui, Z. Xiong, J. Jin, Z. Liu, R. Dong, and J. Hu, “Machine Learning-Based Prediction of Crystal Systems and Space Groups from Inorganic Materials Compositions”, *ACS Omega* **5**, 3596–3606 (2020).
- ²³⁷B. Cordero, V. Gómez, A. E. Platero-Prats, M. Revés, J. Echeverría, E. Cremades, F. Barragán, and S. Alvarez, “Covalent radii revisited”, *Dalton Transactions*, 2832 (2008).
- ²³⁸D. G. A. Smith, D. Altarawy, L. A. Burns, M. Welborn, L. N. Naden, L. Ward, S. Ellis, B. P. Pritchard, and T. D. Crawford, “The MolSSI QCArchive project: An open-source platform to compute, organize, and share quantum chemistry data”, *WIREs Computational Molecular Science* **11**, e1491 (2021).
- ²³⁹I. Mayer, “Bond order and valence indices: A personal account”, *Journal of Computational Chemistry* **28**, 18 (2006).
- ²⁴⁰L. Hubert and P. Arabie, “Comparing partitions”, *Journal of Classification* **2**, 193–218 (1985).
- ²⁴¹M. Meilă, “Comparing clusterings—an information based distance”, *Journal of Multivariate Analysis* **98**, 873–895 (2007).
- ²⁴²D. J. Wales, J. P. K. Doye, A. Dullweber, M. P. Hodges, F. Y. Naumkin, F. Calvo, J. Hernández-Rojas, and T. F. Middleton, *The Cambridge Cluster Database*.
- ²⁴³W. L. Jorgensen, J. Chandrasekhar, J. D. Madura, R. W. Impey, and M. L. Klein, “Comparison of simple potential functions for simulating liquid water”, *The Journal of Chemical Physics* **79**, 926–935 (1983).
- ²⁴⁴N. M. O’Boyle, M. Banck, C. A. James, C. Morley, T. Vandermeersch, and G. R. Hutchison, “Open babel: An open chemical toolbox”, *Journal of Cheminformatics* **3**, 1–14 (2011).
- ²⁴⁵J. D. Chai and M. Head-Gordon, “Long-range corrected double-hybrid density functionals”, *Journal of Chemical Physics* **131**, 1–13 (2009).

- ²⁴⁶J. D. Chai and M. Head-Gordon, “Long-range corrected hybrid density functionals with damped atom-atom dispersion corrections”, *Physical Chemistry Chemical Physics* **10**, 6615–6620 (2008).
- ²⁴⁷Y. Xiang, D. J. Burrill, K. K. Bullard, B. J. Albrecht, L. E. Tragesser, J. McCaffrey, D. S. Lambrecht, and E. Pentzer, “Polymerization of silyl ketenes using alkoxide initiators: A combined computational and experimental study”, *Polymer Chemistry* **8**, 5381–5387 (2017).
- ²⁴⁸N. Mardirossian and M. Head-Gordon, “ ω B97M-V: A combinatorially optimized, range-separated hybrid, meta-GGA density functional with VV10 nonlocal correlation”, *The Journal of Chemical Physics* **144**, 214110 (2016).
- ²⁴⁹W. J. Hehre, R. F. Stewart, and J. A. Pople, “Self-consistent molecular-orbital methods. I. Use of gaussian expansions of slater-type atomic orbitals”, *Journal of Chemical Physics* **51**, 2657–2664 (1969).
- ²⁵⁰W. J. Hehre, R. Ditchfield, R. F. Stewart, and J. A. Pople, “Self-consistent molecular orbital methods. IV. Use of gaussian expansions of slater-type orbitals. Extension to second-row molecules”, *Journal of Chemical Physics* **52**, 2769–2773 (1970).
- ²⁵¹E. D. Glendening and D. Feller, “Cation-water interactions: The M+(H₂O)_n clusters for alkali metals, M = li, na, K, rb, and cs”, *Journal of Physical Chemistry* **99**, 3060–3067 (1995).
- ²⁵²J. Andzelm, S. Huzinaga, M. Klobukowski, E. Radzio-Andzelm, Y. Sakai, and T. Hiroshi, “Gaussian basis sets for molecular calculations”, *Physical Sciences Data* **16**, edited by e. S. Huzinaga, 27–426 (1984).
- ²⁵³T. H. Dunning, “Gaussian basis sets for use in correlated molecular calculations. I. The atoms boron through neon and hydrogen”, *The Journal of Chemical Physics* **90**, 1007–1023 (1989).
- ²⁵⁴D. E. Woon and T. H. Dunning, “Gaussian basis sets for use in correlated molecular calculations. III. The atoms aluminum through argon”, *The Journal of Chemical Physics* **98**, 1358–1371 (1993).
- ²⁵⁵Y. Zhao, Y. Guo, Y. Yao, A. Dumi, D. M. Mulvey, S. Upadhyay, Y. Zhang, K. D. Jordan, J. Yang, and X. Tang, “Q-GPU: A recipe of optimizations for quantum circuit simulation using GPUs”, in 2022 IEEE international symposium on high-performance computer architecture (HPCA) (2022), pp. 726–740.