

**ELECTRONIC STRUCTURE METHODS FOR STUDYING NON-COVALENT  
INTERACTIONS IN COMPLEX CHEMICAL ENVIRONMENTS**

A Dissertation  
Presented to  
The Academic Faculty

By

Dominic A. Sirianni

In Partial Fulfillment  
of the Requirements for the Degree  
Doctor of Philosophy in the  
School of Chemistry & Biochemistry

Georgia Institute of Technology

May 2020

Copyright © Dominic A. Sirianni 2020

**ELECTRONIC STRUCTURE METHODS FOR STUDYING NON-COVALENT  
INTERACTIONS IN COMPLEX CHEMICAL ENVIRONMENTS**

Approved by:

Dr. C. David Sherrill, Advisor  
School of Chemistry & Biochemistry  
*Georgia Institute of Technology*

Dr. Jean-Luc Brédas  
School of Chemistry & Biochemistry  
*Georgia Institute of Technology*

Dr. Jesse G. McDaniel  
School of Chemistry & Biochemistry  
*Georgia Institute of Technology*

Dr. Carlos Silva  
School of Chemistry & Biochemistry  
*Georgia Institute of Technology*

Dr. J. C. Gumbart  
School of Physics  
*Georgia Institute of Technology*

Date Approved: March 21, 2020

Nature's evidence in one place seems to contradict that in another, and so far it has not been possible to draw an even halfway coherent picture of the relationship involved. . .

everyone is still groping around in a thick mist,  
and it will probably be a few years before it lifts.

—*Wolfgang Pauli, ca. 1920*

for my wife Theresa,

who always shone a light through the thick mists as we groped along together

## ACKNOWLEDGEMENTS

Before I attempt to acknowledge and express my gratitude to the many individuals who have guided me on my journey, I first would like to recognize that I am, just as we all are, a product both of my experiences and the people with whom I have shared those experiences. In this latter regard I am particularly fortunate: throughout my life, I have been constantly surrounded by exceptionally caring, supportive, and enthusiastic colleagues, collaborators, classmates, friends, and family. For each and every one of you, I am deeply thankful.

\* \* \*

First and foremost I would like to thank my thesis advisor, Prof. C. David Sherrill, for his guidance, instruction, and assistance in my impassioned pursuit of scientific discovery. He has always been patient, understanding, and supportive of me throughout my time at Georgia Tech – first as an NSF REU fellow in his group during the summer of 2014 and then again as a PhD student for the past five years. David’s unbridled enthusiasm for science and fervent interest in learning how the world works has been a constant inspiration, even when the science seems reluctant to be discovered. Not only does David’s zeal for science make working together a real treat, but the wealth of experience in every aspect of academic life that he possesses — and is happy to share with an aspiring academic — has been an invaluable resource for me as I have navigated the various stages of my career while finding my own path.

In addition to David’s guidance and help in my pursuit of scientific progress, I would also like to acknowledge his patience and support in helping me develop as a communicator of science. I love to tell stories, but I often lose my way in details which are not important or by being too long-winded in my approach. David has helped me to realize the beauty of editing, and several strategies which make writing more productive and more enjoyable. I do have a bone to pick, however — when you learn at the feet of the master, and your approach to writing and organizing topics is so profoundly influenced by that tutelage, it

becomes very difficult to organize the introduction of your thesis in a *different* way than the book chapter he has already written on the subject!

I would also like to thank all current (and former!) members of my thesis committee: Prof. Jean-Luc Brédas, Prof. Kenneth Brown, Prof. J.C. Gumbart, Prof. Joseph Perry, Prof. Jesse McDaniel, and Prof. Carlos Silva — your advice, guidance, and enthusiasm have profoundly altered my approach to science for the better. Thank you so very much for all of the help you have given me during my time at Georgia Tech; I knew that I could always pop in to talk about science, or life, or the status of the Georgia Tech football team — which, this past year, was quite pitiful — and it made all the difference to my feeling that I have been on the right path.

In addition to Prof. Sherrill and my thesis committee, I would like to thank a vast array of coworkers and collaborators for always offering a discerning ear and insightful perspective as we worked together to bring projects we were passionate about to fruition. These individuals include Dr. Lori A. Burns, Dr. Jerome Gonthier, Dr. Ryan Richard, Dr. Daniel G. A. Smith, Dr. Jeffrey Shriber, Dr. Daniel Nascimento, Dr. Robert Parrish, Dr. Michael Marshall, Dr. Matthew Kennedy, Mr. Trent Parker, Dr. Brandon Bakr, Mr. Matthew Scheiber, Mr. Asim Alenaizan, Ms. Constance Warden, Mr. Yi Xie, Mr. Derek Metcalf, Mr. Joseph O'Brien, and Mr. Zach Glick from the Sherrill Group; Dr. Daniel Cheney, Dr. Doree Sitkoff, Dr. Sean Zhu, and numerous others at Bristol Myers-Squibb; and Mr. Andrew James, Dr. Justin M. Turney, Dr. Andrew C. Simmonnet, Prof. Wallace Derri-cote, Prof. A. Eugene DePrince III, Prof. Brandon Magers, Prof. Tricia Shepherd, Prof. Ryan Fortenberry, Prof. Ashley Ringer-McDonald, Prof. Stefan Vogt, Prof. Francesco Evangelista, Prof. Konrad Patkowski, and Prof. T. Daniel Crawford in the extended PSI4 and PSI4EDUCATION families. One reason in particular that SETCA and PsiCON are my absolute favorite conferences of the year is because they basically amount to a giant family reunion with all of you, where we have the opportunity to laugh and discuss science together. It has truly been a pleasure to be a part of this family, and I look forward to

introducing even more individuals to our particular brand of collaboration if and when I successfully land a faculty position at a primarily undergraduate institution.

From my heart, I would like to thank all of my friends and colleagues in the School of Chemistry & Biochemistry at Georgia Tech for their part in creating an amazingly supportive and wonderful environment that helped me keep on keeping on. I have tried to pay forward all of the kindness you have shown me throughout my time at Georgia Tech, which I believe is truly an embarrassment of riches. I could fill this entire thesis with examples of your profound impact on my life, but I know that I can be long-winded and I already feel like I'm losing the room. In particular, however, I would like to thank Ms. Courtney Moore, Dr. Abraham Jordan, Dr. Osiris Martinez-Guzman, and Mr. Eric Drew for the honor of their friendship; Profs. M.G. Finn, Joseph Sadighi, Ken Brown, Jean-Luc Brédas, Jesse McDaniel, Josh Kretchmer, Will Gutekunst, Amanda Stockton, Stefan France, Loren Williams, Pete LaPierre, and Jake Soper for their leadership, kindness, and spirit; and Dr. Kenyetta Johnson, Dr. Christine Conwell, Ms. Ashley Edwards, Dr. Cam Tyson, Mr. Robert McCloud, Ms. Ruth Pierre, and Ms. Keisha Harville for all of their help and support, as well as all of the laughs we've shared over common struggle.

I could not have made it to this point in my academic career without the inspiration and guidance given to me through the years by so many teachers, mentors, and tutors, each of whom I have to profoundly thank for their influence on my philosophy as an educator and my development as a lifelong learner. This began from my earliest years, where thanks to the fact that both of my parents and my aunt were teachers in my home town, I grew up as an adopted son of the Kane Area School District. My fascination with chemistry began when, in 6th grade, Mr. Danielson required that we memorize the symbols for the approximately fifty most common elements on the periodic table, and was subsequently fostered by Mr. Bill Ryding and Mr. Gary Johnson in seventh and eighth grade science. Though mathematics and I were early enemies, I first began to love math thanks to Mrs. Stacey Hastings, who taught me Plane Geometry as a freshman in high school, thanks

to her insistence that we *understand* how to construct proofs of, e.g., the side-angle-side relationship for triangles, and was further shaped by Mr. Phil Imbrogno (P.I.) and Mrs. Allie Saquin, whose love for mathematics was absolutely infectious. P.I. was also an amazing example to follow, thanks to the tenacity with which he attacked every aspect of his life, both in the classroom and on the football field. For my early obsession with organic and biochemistry, I blame Mrs. Joanne Perry, whose Honors, A.P., and Organic & Biochemistry courses at Kane Area High School — together with my love of working with my classmates to learn these challenging concepts — inspired me to attend Edinboro University of Pennsylvania with my major declared as Secondary Education in Chemistry in order to enter the family business of teaching.

My time as an education major at Edinboro was not long, however, as I discovered that my love of learning chemistry and mathematics was stronger than my desire to teach high school chemistry. Luckily, I had chosen to attend Edinboro without knowing that both the Department of Chemistry and the Department of Mathematics & Computer Science are incredible, filled with wonderfully kind, dynamic, and encouraging faculty who constantly strive to build an environment for learning where the students' whole person is celebrated, rather than just their enrollment counted. For their support, mentorship, and instruction, I will forever be grateful to Profs. Joanne Smith, Qun Gu, Naod Kebede, Gerry Hoffman, Lisa Unico, Janet Rogers, and Theresa Thewes in the Department of Chemistry and Profs. Corinne Schaeffer, Doug Puharic, Korey Kilburn, Frank Marzano, Rick White, Anne Quinn, Marc Sylvester, Dan Bennet, and John Hoggard in the Department of Mathematics & Computer Science. Additionally, I would like to thank Profs. Smith and Schaeffer for their guidance as my dual academic advisors while at Edinboro, as well as Profs. Naod Kebede and Gerry Hoffman for their mentorship in the research laboratory, which truly lit a fire in my heart that still burns brightly.

Penultimately, I want to thank all of my family and friends from Kane, Moon, Pittsburgh, Edinboro, Georgia Tech, and beyond, without whose presence in my life I would



have been utterly lost. In particular, I want to thank my parents, Jim and Cathy Sirianni, for their constant encouragement and support, and for teaching me how to be a good person; my grandfathers, Joseph Sirianni and Larry Johnson, for being shining examples of taking pride in one's work, and the patience and determination necessary to do so; my aunt and uncle Paula and Bob Mosier and my surrogate siblings Tina, Jeanette, and Robert, for sharing their lives with me and for our constant laughter; Rae-El and Dany Whitman, for their deep and steadfast friendship; Ed Culbertson, Angie Pelehac, Brittany Sanden, and JT Mannozi, for their silliness and all of the times we spent getting lost in the weird part of YouTube at three in the morning when we should have been doing homework.

Finally, and most importantly, I want to thank my wife, Theresa, for everything that she is and aspires to be, and for her help and encouragement in my own journey to be the best version of myself. More clearly than anyone else in my life, you see my whole self and still wish for nothing less than my happiness. I know that I cannot thank you adequately for everything that you are and all that you inspire me to be — I can only hope that these pages are a start, that they mark another milestone in the journey we are sharing together.

## TABLE OF CONTENTS

<b>Acknowledgments</b> . . . . .	v
<b>List of Tables</b> . . . . .	xvi
<b>List of Figures</b> . . . . .	xviii
<b>Chapter 1: Introduction</b> . . . . .	1
1.1 Historical Perspective on the Development of Quantum Mechanics . . . . .	1
1.1.1 Revolutionary Chemistry and Steady-State Physics . . . . .	2
1.1.2 Remaining Challenges in Physics . . . . .	4
1.1.3 The Quantization Revolution . . . . .	7
1.1.4 Atomic Theory Redux . . . . .	8
1.1.5 Towards a General Formulation of Quantum Mechanics . . . . .	11
1.2 Non-Covalent Interactions in Chemistry, Biology, and Physics . . . . .	13
1.3 Prospectus . . . . .	15
<b>I Theoretical Background to the Thesis Project</b>	<b>16</b>
<b>Chapter 2: Introduction to Electronic Structure Theory</b> . . . . .	17
2.1 Basic Formulation . . . . .	18
2.2 Approximate Solution Methods . . . . .	19

2.2.1	The Method of Linear Variations . . . . .	19
2.2.2	Rayleigh-Shrödinger Perturbation Theory . . . . .	20
2.3	The Born–Oppenheimer Approximation . . . . .	22
2.4	Independent Particle Models . . . . .	23
2.5	The Antisymmetry Principle and Slater Determinants . . . . .	25
2.6	Hartree–Fock Molecular Orbital Theory . . . . .	28
2.6.1	The Hartree–Fock Equations . . . . .	28
2.6.2	The Introduction of a Basis Set: The Roothaan Equations . . . . .	31
2.7	Capturing Dynamical Electron Correlation with post-Hartree–Fock Methodologies . . . . .	38
2.7.1	Configuration Interaction & The Exact Electronic Wavefunction . . . . .	40
2.7.2	Møller–Plesset Perturbation Theory . . . . .	46
2.7.3	Coupled-Cluster Theory . . . . .	52
2.8	Practical Considerations for Correlated Computations . . . . .	55
2.8.1	Convergent Methodologies for Electron Correlation . . . . .	56
2.8.2	Accelerating Convergence for Correlated Approaches . . . . .	59
2.9	Density Functional Theory . . . . .	62
2.9.1	The Hohenberg-Kohn Theorems . . . . .	63
2.9.2	The Universal Functional . . . . .	63
2.9.3	Kohn-Sham DFT . . . . .	64
2.9.4	The Exchange-Correlation Functional . . . . .	66
2.9.5	Non-Local Corrections to DFT . . . . .	67
<b>Chapter 3: Theoretical Approaches for Non-Covalent Interactions . . . . .</b>		<b>70</b>

3.1	Defining the “Interaction Energy” . . . . .	70
3.2	The Many-Body Expansion . . . . .	72
3.3	Types of Non-Covalent Interactions . . . . .	73
3.4	Supramolecular Approaches for Non-Covalent Interactions . . . . .	75
3.4.1	Correlation Effects and Size Consistency . . . . .	75
3.4.2	Basis Set Incompleteness and Basis Set Superposition Errors . . . . .	76
3.4.3	Medal Winners for Non-Covalent Interactions . . . . .	78
3.5	Symmetry-Adapted Perturbation Theory . . . . .	80
3.5.1	Basic Formulation of SAPT . . . . .	82
3.5.2	Levels of SAPT . . . . .	83
3.5.3	Additional SAPT Partitions . . . . .	84
<b>II Benchmarking Non-Covalent Interactions: Towards the “Right Answer for the Right Reasons”</b>		<b>88</b>
<b>Chapter 4: Comparison of Explicitly Correlated Methods for Computing High-Accuracy Benchmark Energies for Noncovalent Interactions . . . . .</b>		<b>89</b>
4.1	Abstract . . . . .	89
4.2	Introduction . . . . .	90
4.3	Theoretical & Computational Methods . . . . .	93
4.3.1	Overview of Approximate CCSD(T)-F12 Methods. . . . .	93
4.3.2	Computational Details . . . . .	97
4.4	Results and Discussion . . . . .	99
4.4.1	Basis Set Convergence for A24 Systems . . . . .	100
4.4.2	aXZ vs. XZ-F12: Accuracy Comparison over A24 and Subsets . . . . .	107

4.4.3	aXZ vs. XZ-F12: Extension to the S22 Test Set . . . . .	111
4.4.4	Benchmark Procedures for NCI: Combining Accuracy and Computational Cost . . . . .	114
4.5	Summary and Conclusions . . . . .	118
<b>Chapter 5: Assessment of Density Functional Methods for Geometry Optimization of Bimolecular van der Waals Complexes . . . . .</b>		<b>121</b>
5.1	Abstract . . . . .	121
5.2	Introduction . . . . .	121
5.3	Computational Methods . . . . .	124
5.3.1	Optimized Geometries for A24 systems . . . . .	125
5.3.2	Radial Potential Surface Scans of HBC6 and NBC10x systems . . .	128
5.4	Results and Discussion . . . . .	129
5.4.1	Optimization of A21 Systems . . . . .	129
5.4.2	Prediction of Optimal Intermolecular Separation in NBC7x and HBC6 Interaction Energy Scans . . . . .	134
5.5	Summary & Conclusions . . . . .	138
<b>III Developing Approximate Perturbative Methods for Non-Covalent Interactions</b>		<b>140</b>
<b>Chapter 6: Optimized Damping Parameters for Empirical Dispersion Corrections to Symmetry-Adapted Perturbation Theory . . . . .</b>		<b>141</b>
6.1	Abstract . . . . .	141
6.2	Introduction . . . . .	142
6.3	Theoretical & Computational Methods . . . . .	146
6.3.1	Formulation of SAPT0-D . . . . .	146

6.3.2	Refitting Damping Parameters for SAPT0–D . . . . .	148
6.3.3	Preparation of Test Systems to Evaluate Scaling of Computational Cost . . . . .	150
6.3.4	Preparation of Structures for Application to GPCR Binding . . . . .	151
6.4	Results & Discussion . . . . .	154
6.4.1	Accuracy of SAPT0–D Variants . . . . .	154
6.4.2	Computational Scaling of SAPT0–D . . . . .	157
6.4.3	Differential Binding of Salbutamol to Active vs. Inactive $\beta_1$ AR . . . . .	159
6.5	Conclusions . . . . .	163

## **IV Application to Interesting Chemical Systems 166**

### **Chapter 7: The Influence of Solvation on Non-Covalent Interactions in Biomolecular Complexes: An Intramolecular Symmetry-Adapted Perturbation Study . . . . . 167**

7.1	Abstract . . . . .	167
7.2	Introduction . . . . .	168
7.3	Computational Methods . . . . .	173
7.3.1	Preparation of Geometries for Functionalized Complexes . . . . .	173
7.3.2	Hydration of Functionalized Complexes . . . . .	174
7.3.3	Quantifying Electronic Effect of Solvent on $\pi$ - $\pi$ Interactions and Energy Components via F-/ISAPT . . . . .	176
7.3.4	MP2 Results for Solute-Solute, Solute-Solvent, and Three-Body Interactions . . . . .	178
7.4	Results and Discussion . . . . .	180
7.4.1	Gas-Phase Interactions . . . . .	180
7.4.2	Quantifying ArX–Bz Interactions in Solution via F-/ISAPT . . . . .	181

7.4.3	Many-Body Analysis of Solvated Interactions . . . . .	185
7.4.4	Effect of Multiple Hydration Shells . . . . .	189
7.5	Summary and Conclusions . . . . .	192
7.6	Acknowledgements . . . . .	192
<b>V</b>	<b>Conclusions</b>	<b>193</b>
<b>Chapter 8:</b>	<b>Conclusions and Outlook . . . . .</b>	<b>194</b>
8.1	Conclusions . . . . .	194
8.2	Outlook . . . . .	195
8.2.1	Towards a Multi-Level Embedded SAPT . . . . .	195
8.2.2	The Influence of Long-Range Contacts in Drug–Protein Binding Specificity . . . . .	196
<b>References</b>	<b>. . . . .</b>	<b>210</b>
<b>Vita</b>	<b>. . . . .</b>	<b>211</b>

## LIST OF TABLES

4.1	Interaction energy (kcal/mol) error statistics vs CCSD(T)/CBS for F12n/aXZ and F12n/XZ-F12 methods, gathered for all A24 systems, as well as for the hydrogen bonding (HB), mixed interaction (MX) and dispersion dominated (DD) subsets. . . . .	104
4.2	Interaction energy (kcal/mol) error statistics for F12n/aDZ and F12n/DZ-F12 methods, applied to the S22B test set, as well as for the hydrogen bonding (HB), mixed interaction (MX) and dispersion dominated (DD) subsets. Included for reference are error statistics computed with the DW-CCSD(T**) -F12/aDZ method. . . . .	112
4.3	Interaction energy (kcal/mol) error statistics, error distributions, and timing summaries for benchmark procedure candidates, along with their double- $\zeta$ counterparts, computed over the A24 test set and each subset versus A24B reference energies. . . . .	116
6.1	Datasets utilized in the training and validation sets. All benchmark datasets are of MP2/CBS + $\Delta$ CCSD(T)/aDZ quality or better. For further details of reference levels of theory for each dataset, please refer to Table S8 in the Supplementary Materials of Ref. 39. . . . .	145
6.2	Summary statistics for signed and unsigned interaction energy error distributions computed using SAPT0/jaDZ and variants of SAPT0-D/jaDZ over each of the training and validation sets, as well as over the full set of all systems examined. Values provided for first and third quartiles correspond to the box borders in a box-and-whisker plot, and the second quartile corresponds to the median value for each SE distribution. . . . .	152



6.3	$\Delta\Delta E_{\text{int}}$ values (kcal mol <sup>-1</sup> ) computed between active and inactive forms of the $\beta_1$ AR–salbutamol complex by F-SAPT0–D3M(BJ) and F-SAPT0–D3M(0) in the jun-cc-pVDZ basis set, decomposed into functional group contacts between the full binding pocket of $\beta_1$ AR and fragments of salbutamol. Fragment labels are consistent with those shown in Fig. 6.6, with the row labeled “All” corresponding to the total interaction energy of the $\beta_1$ AR–salbutamol complex. . . . .	160
6.4	Order-1 and order-2 F-SAPT–D analysis quantifying the contributions of contacts predicted by Warne <i>et al.</i> to be important for explaining difference in binding affinity for salbutamol to active vs. inactive states of $\beta_1$ AR ( $\Delta\Delta E_{\text{int}}$ , computed by F-SAPT–D3M(BJ) and F-SAPT0–D3M(0) in the jun-cc-pVDZ basis set. Fragment labels for salbutamol are consistent with those shown in Fig. 6.6. (top) Contributions of residue–functional group pairs hypothesized by Warne <i>et al.</i> to make polar or hydrogen-bonding contacts, as identified in Figure 3 of Ref. 197; (bottom) Contributions of total contact strength between all amino acid sidechains identified in Figure 3 of Ref. 197 and the full salbutamol ligand (labeled “All”). . . . .	161
6.5	Relative contributions of order-1 [i.e., $\beta_1$ AR(fragment)–salbutamol] contacts to the total $\Delta\Delta E_{\text{int}}$ for the full $\beta_1$ AR–salbutamol complex, computed with F-SAPT0–D3M(BJ) and F-SAPT0–D3M(0) in the jun-cc-pVDZ basis set. Also provided for reference are whether or not the particular residue sidechain or peptide bond was hypothesized to be important by Warne <i>et al.</i> in Ref. 197. Fragment labels for $\beta_1$ AR are consistent with the fragmentation procedure described in the Supplementary Information. . . . .	163

## LIST OF FIGURES

2.1	Convergence of correlation energy contribution to non-covalent interaction energies at various truncation orders for the CC series. Values taken from Table S-2 of Ref. 12. . . . .	57
3.1	Comparison of computational classes for NCI. (a) The recommended SAPT model chemistries from Ref. 58 are compared to (b) common DFT approaches and (c) common or recommended wavefunction techniques from Ref. 55 according to both efficiency (purple; time required for adenine-thymine) and accuracy (grey; MAE averaged over S22, HBC6, HSG, and NBC10 databases) metrics. Subset MAE values are shown as inset bars for hydrogen-bonding (red), mixed-influence (green), and dispersion-dominated (blue) NCI motifs. Figure reproduced from Ref. 55. . . . .	79
4.1	Bimolecular complexes included in (a) the A24 and (b) S22 test sets of Hobza and coworkers <sup>14,79</sup> with revised A24B and S22B <sup>30</sup> reference energies, in kcal mol <sup>-1</sup> . Coloring is based on SAPT2+(3)/aTZ results reported previously by Burns <i>et al.</i> <sup>57,77</sup> and indicates interaction type: red for electrostatically dominated interactions (typically hydrogen bonding), blue for dispersion dominated interactions, and yellow-green for interactions of mixed character. . . . .	94
4.2	Convergence of CCSD(T <sup>**</sup> )-F12n/aXZ (n = a, b, c; X = D, T, Q, 5, 6) IEs for (a) the ammonia-water complex, (b) the formaldehyde-ethylene complex, (c) the methane-Ar complex, and (d) the ethylene dimer in forced $\pi$ -stacking geometry. Also plotted are canonical CCSD(T)/aXZ IEs and our revised A24B reference energies (dotted line) obtained at the CCSD(T)/CBS(aQZ, a5Z) [(a) & (c)] or CCSD(T)/CBS(a5Z,a6Z) [(b) & (d)] levels of theory (see text). . . . .	101

- 4.3 Convergence of CCSD(T<sup>\*\*</sup>)-F12n/XZ-F12 (n = a, b, c; X = D, T, Q, 5) IEs for the (a) ammonia-water complex, (b) formaldehyde-ethylene complex, (c) methane-Ar complex, (d) ethylene dimer in forced  $\pi$ -stacking geometry, and (e) methane dimer. Also plotted are canonical CCSD(T)/aXZ IEs and our revised A24B reference energies (dotted line) obtained at the CCSD(T)/CBS(aQZ, a5Z) [a & c] or CCSD(T)/CBS(a5Z, a6Z) [b, d, & e] levels of theory (see text). . . . . 102
- 4.4 Error in IE computed for all bimolecular complexes in the A24 test set relative to A24B reference energies for the (a) CCSD(T<sup>\*\*</sup>)-F12a, (b) CCSD(T<sup>\*\*</sup>)-F12b, and (c) CCSD(T<sup>\*\*</sup>)-F12c methods for both aXZ and XZ-F12 (X = D, T, Q, 5) basis sets. Vertical lines represent individual members of A24, color-coded by interaction type (red = hydrogen bonding, blue = dispersion dominated, yellow/green = mixed interaction). For each level of theory, MAE (black rectangles, given on the left) and MA%E (black ovals, given on the right) are presented. Three shaded error regions are shown: the lightest encompasses  $\pm 0.1$  kcal mol<sup>-1</sup> &  $\sim 4\%$ , next lightest region  $\pm 0.05$  kcal mol<sup>-1</sup> &  $\sim 2\%$ , and darkest region  $\pm 0.01$  kcal mol<sup>-1</sup> &  $\sim 1\%$ . For comparison, errors computed using the current silver-standard DW-CCSD(T<sup>\*\*</sup>)-F12 method, paired with aDZ and aTZ basis sets, are also presented. . . . . 108
- 4.5 Error in IE computed for all bimolecular complexes in the A24 test set relative to Helgaker-extrapolated<sup>20</sup> MP2/CBS(aXZ, a(X+1)Z) reference IEs (analogous  $\zeta$ -levels to A24B), computed using (a) the MP2-F12(3C)/aXZ (X = D, T, Q, 5) and (b) MP2-F12(3C)/XZ-F12 (X = D, T, Q) model chemistries. For each level of theory, MAE (black rectangles, given on the left) and MA%E (black ovals, given on the right) are presented. Shaded error regions are given with identical ranges as in Fig.4.4. . . . . 110
- 4.6 Error in IE computed for all bimolecular complexes in the S22 test set,<sup>79</sup> relative to revised S22B reference energies,<sup>30</sup> computed using (a) the CCSD(T<sup>\*\*</sup>)-F12n/aDZ (n = a, b, c) and (b) CCSD(T<sup>\*\*</sup>)-F12n/DZ-F12 (n = a, b, c) model chemistries. Included for reference are IEs computed using the silver-standard DW-CCSD(T<sup>\*\*</sup>)-F12/aDZ method. For each model chemistry, MAE (black rectangles, given on the left) and MA%E (black ovals, given on the right) are presented. The outer, lightly shaded region encompasses  $\pm 0.5$  kcal mol<sup>-1</sup> &  $\sim 23\%$ , and the inner, darkly shaded region is given to indicate the location of zero error. . . . . 113

- 5.1 Test sets of bimolecular complexes examined here. (a) A21: 21 bound complexes contained in the A24 test set of Hobza and co-workers,<sup>177</sup> (b) NBC7x: seven (recently extended<sup>39</sup>) radial potential scans from the NBC10 test set<sup>140</sup> and (c) HBC6:<sup>30,145</sup> radial potential scans for six doubly hydrogen bonded complexes. Indicated by box coloring [(a)–(c)] or by dot color [(d)] are the noncovalent interaction type for each complex, reported previously:<sup>57,77</sup> red for electrostatic interactions, blue for dispersion interactions, and yellow/green for mixed electrostatic and dispersion interactions. The ternary diagram (d) further indicates the relative magnitude of the interaction energy components for these complexes,<sup>145,178</sup> by placing a colored dot according to the ratios of attractive dispersion/induction and attractive/repulsive electrostatic contributions to the total interaction energy. Proximity to each labeled vertex indicates an increasing fraction of the attraction (repulsion) arising from that component. . . . . 127
- 5.2 Box-and-whisker plots representing both  $\Delta\text{COM}$  signed errors (boxes shaded pink) and LRMSD values (boxes shaded blue) for systems in the A21 test set, optimized using the (a) B3LYP-D3, (b) B97-D3, and (c) M05-2X density functionals together with the DZ, TZ, aDZ, and aTZ basis sets. Boxes encompass the first (Q1) through third (Q3) quartiles of each data set, with values corresponding to the median (Q2) and mean LRMSD and  $\Delta\text{COM}$  signed error indicated as a solid black bar and black square, respectively. Whiskers encompass the full range of LRMSD values and  $\Delta\text{COM}$  signed errors; maximum values are indicated when whiskers surpass the area shown. For reference, a dotted line indicates 0.0 Å, and three levels of shading are provided: light grey encompassing  $\pm 0.1$  Å, medium grey encompassing  $\pm 0.05$  Å, and medium-dark grey encompassing  $\pm 0.01$  Å. . . 130
- 5.3 Values corresponding to individual A21 complexes and box-and-whisker plots detailing test-set-wide distributions for (a) least root mean square displacements (LRMSD) and (b) signed errors in center-of-mass distance ( $\Delta\text{COM}$ ), computed with each density functional using the aug-cc-pVDZ basis set. Values for individual A21 systems, shown with empty circle markers, are grouped and colored according to interaction motif:<sup>57,77</sup> red for electrostatically bound complexes (HB subset), blue for dispersion bound complexes (DD subset), and green for mixed interaction complexes (MX subset). For box-and-whisker plots of each model chemistry, boxes encompass the first (Q1) through third (Q3) quartiles of each data set, with values corresponding to the median (Q2) and mean LRMSD and  $\Delta\text{COM}$  signed error indicated as a solid green bar and black square, respectively. Whiskers encompass the full range of LRMSD values and  $\Delta\text{COM}$  signed errors; maximum values are indicated when whiskers surpass the area shown. For reference, a dotted line indicates 0.0 Å, and three levels of shading are provided: light grey encompassing  $\pm 0.1$  Å, medium grey encompassing  $\pm 0.05$  Å, and medium-dark grey encompassing  $\pm 0.01$  Å. . . . . 133

- 5.4 Scans of the non-counterpoise–corrected interaction energy (unCP IE) along the radial separation coordinate  $R$  in the formamidine dimer (HBC6-3; inset shown) computed with B3LYP-D3 (red), B97-D3 (blue) and M05-2X (green) using the cc-pVDZ basis set. The interpolated optimal intermolecular separation for each curve is indicated with a vertical dotted line in the same colors. For reference, a curve constructed from the CCSD(T)/CBS benchmark IEs at each value of  $R$  is presented in black. . . . . 135
- 5.5 Box-and-whisker plots depicting the distribution of signed error in interpolated optimal center-of-mass displacement ( $\Delta\text{COM}$ ) for radial interaction energy curves in the NBC7x (a & b) and HBC6 (c & d) test sets. For both test sets, box-and-whisker plots representing curves constructed from both counterpoise-corrected (CP; left panels a & c) and uncorrected (unCP right panels, b & d) IEs are given. Whiskers encompass the full range of  $\Delta\text{COM}$  signed errors for the indicated test set, correction scheme, and model chemistry, and boxes illustrate the first (Q1), second (median, black bar), and third quartiles (Q3) of these data; additionally, the mean signed error for each data set is indicated with a black square. For reference, a dotted line indicates 0.0 Å, and three levels of shading are provided: light grey encompassing  $\pm 0.1$  Å, medium grey encompassing  $\pm 0.05$  Å, and medium-dark grey encompassing  $\pm 0.01$  Å. . . . . 136
- 6.1 Ternary diagrams visualizing relative contributions of attractive (-) and/or repulsive (+) electrostatics, induction, and dispersion interaction energy components based on the SAPT0/jun-cc-pVDZ description of IEs for all systems comprising the (a) training and (b) validation sets for SAPT-D parameter training. Each system is represented by a single dot, colored according to the most dominant contribution to the overall SAPT0 IE: red indicating an electrostatically dominated interaction, blue indicating a dispersion dominated interaction, and yellow-green indicating an interaction for which neither electrostatics nor dispersion components dominate. . . . . 146
- 6.2 Violin plots visualizing the distribution of signed errors (SE) of IEs computed for complexes in the validation set with SAPT0/jun-cc-pVDZ (green), SAPT0–D3M(BJ)/jun-cc-pVDZ (orange), and SAPT0–D3M(0)/jun-cc-pVDZ (purple) as compared to CCSD(T)/CBS reference IEs. Violin widths at a given SE correspond to the relative frequency of complexes exhibiting that value of SE. Also provided for convenience in horizontal dotted lines are the first (Q1), second (Q2), and third (Q3) quartiles for each distribution of SE values. . . . . 153

6.3	Violin plots visualizing the distribution of signed errors (SE) of IEs computed for complexes in the full data set (training and validation sets) with SAPT0/jun-cc-pVDZ (green), SAPT0–D3M(BJ)/jun-cc-pVDZ (orange), and SAPT0–D3M(0)/jun-cc-pVDZ (purple) as compared to CCSD(T)/CBS reference IEs. Violin widths at a given SE correspond to the relative frequency of complexes exhibiting that value of SE. Also provided for convenience in horizontal dotted lines are the first (Q1), second (Q2), and third (Q3) quartiles for each distribution of SE values. . . . .	155
6.4	Metrics describing the computational expense of SAPT0–D (blue) as compared to exact SAPT0 (red) for (a) the total wall time for each computation, (b) the total disk space utilized by each computation, and (c) the total memory utilized by each computation for scaling tests using progressively larger subsystems of the 3ACX co-crystal structure. <sup>196</sup> SAPT0 computations on more than 311 atoms failed; however, computations successfully completed for SAPT0–D on 3ACX subsystems with up to 445 atoms. <b>Insets:</b> Structures for selected 3ACX subsystems with 83, 177, 311, and 445 atoms. . . . .	156
6.5	(a) Total interaction energies and components (kcal mol <sup>-1</sup> ) for 3ACX subsystems computed with SAPT0, SAPT0–D3M(0), and SAPT0–D3M(BJ) in the jun-cc-pVDZ basis set. (b) Differences between interaction energies of subsequent 3ACX subsystems (kcal mol <sup>-1</sup> ), illustrating the convergence of SAPT0 and SAPT0–D components and total interaction energies. . . . .	158
6.6	Fragmentation scheme for salbutamol monomer in F-SAPT analysis of interaction energies of $\beta_1$ AR–salbutamol complexes examined here. . . . .	159
7.1	Bimolecular complexes from which the HYD8 test set is constructed. Structures prepared via functionalization of the T-shaped pyridine–benzene complex from the S66 test set, <sup>76,104</sup> before re-optimizing the structures at the B3LYP-D3M(BJ)/aug-cc-pVDZ level of theory within enforced $C_s$ symmetry. Box coloring is based on SAPT0/jun-cc-pVDZ results computed in the gas phase, and indicates the interaction type: blue for dispersion-dominated interactions and yellow-green for mixed electrostatics and dispersion contributions. . . . .	173
7.2	Environment binning schemes employed in this work, illustrated for the HYD8-1 (solvated aniline–benzene) complex. A–B interactions in “EnvC” scheme are computed directly using ISAPT, while A–B interactions in “EnvA” and “EnvB” are computed via F-SAPT post analysis, <sup>222</sup> via accumulation of functional group interactions. . . . .	176

7.3	Total interaction energies and SAPT components for (a) HYD8-3m-w50 (benzene dimer) and (b) HYD8-7m-w50 (pyridinium–benzene) complexes solvated by 50 explicit solvent molecules, computed at the F-/ISAPT0/jun-cc-pVDZ level of theory and averaged over all ten relaxed solvent configurations. “EnvX” labels indicate that explicit solvent molecules are contained within monomer “X” during the SAPT computation (see text). Error bars encompassing the full range of values across all snapshots are also provided for SAPT terms and total IEs. Furthermore, we have provided a set of bars corresponding to the conventional two-body F-SAPT computation in the gas phase, i.e., in the absence of explicit solvent molecules. See Section II C for additional details regarding our nomenclature and details of the F-/ISAPT computations. . . . .	183
7.4	Box-and-whisker plots representing the non-additive three-body correction to total “trimer” energy ( $\Delta E_{ABC}$ ; kcal mol <sup>-1</sup> ), for both relaxed (R) and unrelaxed (U) solvent configurations of each doubly-solvated HYD8 complex, computed at the HF/jun-cc-pVDZ level of theory (see text for details of solvent configuration selection and preparation). Boxes encompass the first (Q1) through third (Q3) quartiles of $\Delta E_{ABC}$ , with values corresponding to the median (Q2) and mean $\Delta E_{ABC}$ indicated as a solid green bar and green triangle, respectively. Additionally, whiskers encompass the full range of $\Delta E_{ABC}$ values for all solvent configurations. . . . .	186
7.5	Box-and-whisker plots representing the total “trimer” interaction energy ( $\Delta E_{ABC}^{\text{IE}}$ ; kcal mol <sup>-1</sup> ) for both relaxed (R) and unrelaxed (U) solvent configurations of benzene dimer (HYD8-3) and pyridinium–benzene (HYD8-7) complexes hydrated by 50 explicit water molecules, computed at the HF/jun-cc-pVDZ (orange boxes) and MP2/jun-cc-pVDZ (blue boxes) levels of theory. Boxes encompass the first (Q1) through third (Q3) quartiles of $\Delta E_{ABC}^{\text{IE}}$ , with values corresponding to the median (Q2) and mean $\Delta E_{ABC}^{\text{IE}}$ indicated as a solid green bar and green triangle, respectively. Additionally, whiskers encompass the full range of $\Delta E_{ABC}^{\text{IE}}$ values for all solvent configurations. . . . .	188
7.6	Mean change in total interaction energies and SAPT components upon first-shell solvation ( $\Delta \text{IE}^{1-\text{G}}$ ; striped bars) and second-shell solvation ( $\Delta \text{IE}^{2-1}$ ; solid bars) for the (a) HYD8-3mX (benzene dimer) and (b) HYD8-7mX (pyridinium–benzene) complexes, averaged over values computed at the F-/ISAPT0/jun-cc-pVDZ level of theory for all ten relaxed solvent configurations. “EnvX” labels indicate that explicit solvent molecules are contained within monomer “X” during the SAPT computation (see text). See Section II C for additional details regarding our nomenclature and details of the F-/ISAPT computations. . . . .	190

## ABSTRACT

Non-covalent interactions (NCI) encompass the quantum mechanical forces felt between atoms and molecules which are not directly bonded to one another. Responsible for governing diverse chemical and physical phenomena, NCI are of fundamental interest in fields including materials design and drug discovery, among others. In order to study NCI accurately, quantum chemical methods must be employed whose computational expense often limits the systems which can be studied to at most 100 atoms. Often, this is challenge is addressed by examining NCI in small, representative subsystems, however this approach neglects the influence of chemical environment on these interactions. Furthermore, the best manner in which to study such environmental effects is still an open question in the field. Meeting these challenges will be the focus of this dissertation: through the development of novel quantum chemical methods, as well as the extension of existing methods, this work will seek to describe the effect of diverse chemical environments on non-covalent interactions. In this way, a more complete understanding of these phenomena will be provided, which can then be exploited to advance various chemical applications.



## SUMMARY

Non-covalent interactions, encompassing the through-space physical forces of attraction and repulsion felt between atoms and molecules separated by finite distances, control numerous chemical and physical phenomena on length and time scales ranging nearly fifteen orders of magnitude: from the association of molecular aggregates and host–guest complexes on the Angström/femtosecond scales and macromolecular secondary structure and dynamics on the micrometer/nanosecond scales to transitions between physical phases of matter on the meter/second scales and even global weather patterns due in part to the increased density of humid air on the thousands of kilometers/millenia scales, understanding NCI and their role in these phenomena is of critical importance to understanding these phenomena themselves. Despite their omnipresence and foundational importance, however, directly submitting NCI to experimental investigation is a significant challenge, due precisely to their nature as being cooperative across these length and time scales.

Computational investigation, however, offers the unique advantage of decoupling this cooperativity, whereby particular interactions within a given system of interest may be studied. The development and application of computational approaches and methodologies to probe NCI has been an ongoing effort in the field of the theoretical and computational molecular sciences for the better part of a century, and thanks to the increased availability of ever-more powerful computing hardware, we are now poised to explore the influence of NCI in systems and on phenomena previously undreamed of. Towards this end, I take in this Thesis a holistic approach to (i) study the fundamental nature of non-covalent interactions (NCI) in small model systems, (ii) develop novel computational tools by which NCI may be accurately investigated in previously inaccessible, extended chemical systems, and finally (iii) apply these novel and existing approaches to examine NCI in complex chemical environments.

In Part I of this Thesis, best practices for the benchmarking of non-covalent interac-

tion energies (IEs) and obtaining accurate structures for bimolecular complexes are established, in an effort to “get the right answers for the right reasons.” Towards this end, Chapter 4 compares interaction energies (IEs) of several non-bonded complexes of various binding motif and interaction strength computed with approximate formulations of explicitly-correlated coupled cluster theory (CCSD(T<sup>\*\*</sup>)-F12n, n = a, b, c; abbreviated F12n) against gold-standard reference IEs to evaluate their performance for generating benchmark-quality descriptions of diverse non-covalent interactions. It was found that, contrary to trends observed for total molecular correlation energies, IEs computed with F12n methods paired with basis sets designed for use with explicitly correlated methods (cc-pVXZ-F12, X = D, T, Q, 5; abbreviated XZ-F12) were less accurate for a given  $\zeta$ -level than when leveraging conventional, correlation consistent basis sets augmented with diffuse functions (aug-cc-pVXZ, X = D, T, Q, 5; abbreviated aXZ). Furthermore, F12n/aXZ model chemistries converged more rapidly towards the complete basis set (CBS) limit than did their F12n/XZ-F12 counterparts, with F12b/aXZ converging most rapidly out of all model chemistries examined. Specifically, F12b/aTZ achieves mean absolute errors of 0.01 kcal mol<sup>-1</sup> for IEs of bimolecular complexes in the A24 test set, rivalling the accuracy of conventional CCSD(T) in the a5Z or even a6Z basis set (where available). Inspired by the performance of these F12n methods, we also established timings for IE computations leveraging F12n/aXZ to assess their computational expense in comparison to canonical benchmark approaches. It was found that while F12b/aXZ was significantly faster than composite approaches based on fully canonical MP2 and CCSD(T), composite approaches leveraging density fitted, frozen natural orbital formulations of CCSD(T) [DF-FNO-CCSD(T)] were even faster than F12b/aXZ while retaining the same level of accuracy. Based on these findings, we have recommended procedures leveraging either DF-FNO-CCSD(T)/[aTQZ;  $\delta$ :aTZ] or F12b/aTZ for obtaining benchmark-quality non-covalent interaction energies.

Furthermore, Chapter 5 explores the suitability of three dispersion-aware density functionals popular for application to non-covalent interactions (B97-D3, B3LYP-D3, and M05-

2X) for generating optimized geometries of non-bonded bimolecular complexes. By comparing against reference geometries for small van der Waals complexes of diverse binding motif generated at the CCSD(T)/CBS level, we establish a reliable protocol for obtaining equilibrium geometries for these complexes using DFT. Each density functional is able to reproduce reference geometries to within  $\pm 0.1$  Å for each of the average monomer center-of-mass displacement and average least root-mean-squared displacement, when paired with the aug-cc-pVDZ (abbreviated aDZ) basis set. Differences in computed equilibrium geometry of this magnitude correspond to differences of interaction energy of only a few tenths of one kcal mol<sup>-1</sup>, equivalent in accuracy to the performance of DFT-D for interaction energies themselves. Additionally, we showed that for both dispersion bound and doubly-hydrogen bonded bimolecular complexes, the optimal intermolecular contact distance interpolated from radial dissociation curves constructed using these DFT-D/aDZ model chemistries were able to reproduce those interpolated from reference curves constructed at the CCSD(T)/CBS level to within  $\pm 0.1$  Å, matching the performance of these DFT-D methods for equilibrium geometries. Due to their favorable performance compared to CCSD(T)/CBS references for estimating both total equilibrium geometries and optimal intermolecular contact distance, we concluded that these DFT-D/aDZ model chemistries were indeed suitable for application to the generation of equilibrium geometries of generic non-bonded complexes.

Part II of this Thesis is concerned with the development of affordable approaches for investigating NCI in extended chemical systems, and particularly the development of approximate, semi-empirical variants to the popular symmetry-adapted perturbation theory (SAPT) approach. SAPT has become a valuable computational tool offering physical insight into the fundamental nature of non-covalent interactions in diverse chemical systems by directly computing the electrostatics, exchange (steric) repulsion, induction (polarization), and London dispersion contributions to the interaction energy using quantum mechanics. Further application of SAPT to novel chemical problems is limited primarily by

its computational expense, where even for its most affordable variant, SAPT0, computing the London dispersion contribution to the interaction energy (IE) scales as the fifth power of system size,  $\mathcal{O}(N^5)$ . In Chapter 6, we optimize damping parameters for the semiempirical  $-D3$  dispersion correction of Grimme and co-workers, so that they are suitable for use as replacements of the computationally expensive dispersion term in SAPT0. Parameters are obtained by fitting to a large set of 2295 interaction energies computed at the CCSD(T)/CBS level of theory. This reduces the algorithmic scaling of SAPT0 from  $\mathcal{O}(N^5) \rightarrow \mathcal{O}(N^4)$  while retaining the physically meaningful interpretation of IE components characteristic of all SAPT methods. This scaling reduction translates into a nearly  $2.5\times$  speedup over conventional SAPT0 for systems with  $\sim 300$  atoms. Furthermore, this allows for SAPT-D computations to be performed on systems with over 450 atoms, while offering nearly equivalent accuracy to SAPT0 when compared against reference IEs for a diverse set of approximately 8,100 bimolecular complexes. We have further extended our formulation of SAPT-D to be consistent with the functional group partition (F-SAPT-D) and applied this method to conclude that the difference in binding affinity for partial agonist salbutamol to the G-protein coupled  $\beta_1$ -adrenergic receptor between active and inactive forms is due to the cooperative effects of both peptide bonds and residues outside the immediate binding pocket, indicating that a local contact model for protein-ligand binding is insufficient to discriminate between binding conformations which possess similar activities.

Finally, in Part III of this Thesis, Chapter 7 examines the extent to which chemical environment “tunes” NCI by leveraging both intramolecular SAPT (ISAPT) and functional-group partitioned SAPT (F-SAPT) to study solvated  $\pi - \pi$  interactions. In doing so, we investigate (i) possible approaches by which to compute non-covalent interactions embedded in a chemical environment, and (ii) quantify the tuning of these interactions due to the environment relative to the interactions in the gas phase. We have applied our approach to quantify the extent to which explicit water solvent modulates  $\pi - \pi$  interactions in several functionalized, T-shaped arene-benzene complexes, hydrated by a statistically

diverse set of solvent configurations. We have found that, for systems wherein no significant non-additive three-body interaction between the monomers and the collective solvent environment are present, the solvent environment does not significantly tune  $\pi - \pi$  interactions, either due to the choice of system partitioning or solvent configuration. For systems where the nonadditive three-body interaction is significant, however — i.e., where it either is greater than  $\sim 2 \text{ kcal mol}^{-1}$  or where it deviates between solvent configurations by greater than  $\sim 1 \text{ kcal mol}^{-1}$  — the solvent environment does tune the interaction, sometimes by up to several  $\text{kcal mol}^{-1}$  for both total interaction energies and F-/ISAPT components. Finally, we have shown that for these non-additive systems, even two hydration shells of 50 explicit water molecules within  $7 \text{ \AA}$  of the solute complex may not be sufficient to ensure convergence of the solute–solute interactions towards the continuum limit, whereas for additive systems, only a single shell of 28 water molecules within  $3 \text{ \AA}$  is necessary for convergence.

In addition to developing approaches which may be leveraged to study non-covalent interactions in extended chemical systems and diverse chemical environments, this Thesis builds upon previous efforts to set forth the next generation of best-practices for all facets of the computational investigation of non-covalent interactions. From this foundation, a variety of new avenues forward are now emerging, with previously inaccessible chemical phenomena suddenly within reach of our methodologies. For example, the effect of protein environment on tuning active site binding activity and specificity can now be reliably and routinely quantified, opening the door to answering questions of the effects of distant point mutations or allosteric binding on enzyme function. While these and other applications are not explored here, it is my hope that the advances developed here may help further scientific progress throughout the broader computational molecular sciences community at large.

# CHAPTER 1

## INTRODUCTION

Since it is the goal of this Thesis to discuss the application of quantum mechanics to understand the fundamental nature of non-covalent interactions, it would seem natural to begin with an overview of the basics of quantum mechanics and a discussion of how it may be employed in this manner. This would be a fine approach, if not for the fact that quantum mechanics is perhaps the least intuitive scientific theory ever developed. However, the hallmark of a successful theory is not the intuitiveness of its predictions, but rather the precision and accuracy with which it may be experimentally validated. By this metric, quantum mechanics is also perhaps the most successful scientific theory ever developed. This Thesis is not, however, meant to be a treatise on the many successes of quantum mechanics; rather, this Thesis tells the story of the work undertaken by myself and my colleagues over the last five years, which itself is merely a continuation of the same journey undertaken by so many scientists who preceded us. We will begin, therefore, at a much more appropriate place: the beginning.

### 1.1 Historical Perspective on the Development of Quantum Mechanics

Central to the human experience is our desire to observe the world around ourselves and wonder about how we fit into the picture. From a young age, questions like “why does the sunshine make my skin feel warm?”, “where does the rain come from?”, and “how do water striders keep from sinking?” inspire our journey through life. Not only do these questions guide our individual growth, but they have also shaped our cultural heritage. While countless individuals have asked these questions in pursuit of understanding our place in the universe, perhaps the most relevant example of such a question to this Thesis is, “what exactly makes up all of the *stuff* around us?”

Of course, this is not the first time this question has been stated; indeed, this was one of the major questions pondered by ancient Greek philosophers, including Plato, Aristotle, Socrates, and many others. Two in particular, however, were remarkably prescient in their worldview: Democritus and Leucippus,<sup>1</sup> who more than two millennia ago (ca. 400 B.C.E) proposed that everything in the material world — what we have come to call *matter* — is made up of tiny, indestructible, eternal particles that they called *ατομος* (*atomos*), meaning “indivisible.” The two “atomists” believed these particles were inertly solid, with some jagged and sharp and others smooth and slippery, and that they interacted with one another mechanically. Furthermore, they reasoned that the properties they observed for macroscopic materials were a direct result of these interactions. Unfortunately, questions which we might consider scientific today were viewed primarily through a philosophical lens during antiquity.

This “natural philosophy,” where existential questions were rejected and accepted based purely on academic grounds rather than experimental validation, allowed a competing viewpoint where earth, water, air, and fire were the fundamental elements of nature to become the prevailing view for nearly two millennia. Natural philosophy was not only concerned with what matter *is*, but also the manner in which it transforms from one form to another. Since our modern conception of chemistry is fundamentally concerned with the study of such changes, chemistry remained more mythical than scientific for much of the intervening millennia from when atomism was first introduced.\* Fortunately for us, however, chemistry’s mystical shroud began to lift when it was first supposed that matter cannot be created or destroyed, only transfigured from one form to another.

### 1.1.1 Revolutionary Chemistry and Steady-State Physics

As early as 1630, Jean Rey first implicitly assumed that matter cannot be created or destroyed, only transformed — which we now know as the law of conservation of mass —

---

\*I suppose that’s why it was named *chemistry*. . .

setting in motion the advancements which would eventually lead to the evolution of chemistry into a fully-fledged natural science. More than a century later, the chemical revolution truly began in earnest in 1789 when Antoine Lavoisier's *Traité Élémentaire de Chimie*, "Elements of Chemistry," enshrined conservation of mass as empirical law and refuted the phlogiston theory of combustion with one based on the consumption of oxygen. This revolution continued in 1807, when John Dalton proposed that chemical *elements* — pure substances which cannot be broken down further into constituent parts by chemical means — owed their purity of composition to the fact that they were made up of a collection of identical, indivisible particles which, inspired by Leucippus and Democritus, Dalton referred to as atoms. Dalton's atomic theory postulated that chemical substances were formed by combining atoms in defined, whole-number ratios, and even though atoms were exchanged in chemical reactions, the properties of the atoms themselves were unchanged. Together with the discovery of an array of new chemical elements by Humphry Davy and others, and their subsequent organization into the periodic table by Demitri Mendeleev, atomic theory proved enormously successful for rationalizing, understanding, and predicting chemical phenomena. These advancements, combined with the rigor and meticulousness with which Lavoisier and his contemporaries performed their investigations, freed chemistry of its philosophical and mythical shroud by the mid-19th century.\*

The field of physics, on the other hand, was significantly less impeded by natural philosophy than was chemistry, thanks in large part to the fact that early physicists (who in truth were essentially pragmatic mathematicians) were more concerned with predicting macroscopic phenomena than addressing existential questions of reality.† From Newton and Gauss to Maxwell, Coulomb, and Faraday, physics steadily marched forward, producing a number of theories which provided remarkably accurate predictions of natural phenomena. Newtonian mechanics predicted the motions of the stars and planets, as well as the

---

\*Evidently, chemistry just needed to reflux for a while before it was ready to crystallize.

†Of course, physics has its own storied history of existentialism, including the radical ideas of Galileo and others who proposed that the Earth orbited the Sun, and not the other way around.



manner in which objects behave here on Earth. Thermodynamics explained the behavior of gases and the transfer of heat between objects, and its principles led to the invention of the steam locomotive and other major advances throughout the industrial revolution. Electromagnetism described the behaviors of light waves, including the diffraction of sunlight into a rainbow of colors by a prism. These theories were so successful, in fact, that by the middle of the 19th century it was widely believed that they were sufficient to describe all observable macroscopic phenomena. All that remained in order to label physics a solved science — which would surely occur by the turn of the 20th century — was to address the handful of unanswered questions yet to be resolved.\*†

### 1.1.2 Remaining Challenges in Physics

In 1859 and 1860, physicists Balfour Stewart and Robert Kirchhoff independently identified one such question while they were studying the thermal radiation properties of so-called “black bodies,” objects which perfectly absorb all incoming radiation. To avoid becoming indefinitely hot, these objects also emit radiation to achieve thermal equilibrium with their surroundings. At a given temperature, therefore, a black body will emit a characteristic spectrum of radiation, with certain wavelengths more intense than others, in order to maintain this thermal equilibrium.‡ Collectively, all of the wavelengths emitted are referred to as the object’s *emission spectrum*, and their relative strength (and consequently the “color” the object appears) is referred to as its *spectral intensity*. Despite our collec-

---

\*For a discussion and critique of whether we may be approaching the limit from another angle, namely the lack of testable hypotheses provided by modern advances in particle physics, see Ref. 2.

†As this Thesis is being written in 2020 — more than a century after physics was supposed to have been solved, however — it should come as no surprise that something big was about to happen. Double, double toil and trouble...

‡While this phenomenon may seem difficult to conceptualize, it can be illustrated with a simple thought experiment: imagine a lump of charcoal. At ambient temperature, the briquette is black, but as it is lit and starts to burn, it begins to appear red, orange, and eventually yellow-white as it gets progressively hotter. Then, as it cools, the briquette appears to regress in color in the reverse direction. Evidently, the apparent color of the charcoal depends only on its temperature. The color that the charcoal appears is due to the particular wavelengths of light being emitted from the charcoal with the highest intensity. Therefore, when the briquette appears red, it is because the wavelengths emitted from the charcoal corresponding to red light (approximately 650-750 nm) are the most intense. Similarly, when the briquette appears yellow, it is because the wavelengths being emitted with the highest intensity are primarily yellow, etc.

tive familiarity with this concept — after all, “red hot” is a common phrase in the English language for a reason — predicting the spectral intensity of emitted radiation from an ideal black body as a function of its temperature proved to be quite a bit more challenging than was originally anticipated.

Indeed, ever since Stewart and Kirchhoff originally introduced their proofs of this phenomenon, physicists had struggled to explain the spectrum and intensity of blackbody radiation. Several attempts to do so based on classical electromagnetism could reasonably predict the spectral intensity for large wavelengths (visible and infrared), but would fail miserably for shorter wavelengths (ultraviolet). Different approaches, relying either on fitting the spectral distribution against empirical measurements or attempting to derive it from macroscopic formulations of the second law of thermodynamics, however, agreed reasonably well with experiment for short wavelengths but diverge for longer ones. Regardless of approach, all such attempts produced equivalent results: the characteristic spectrum of blackbody radiation could not be predicted by theoretical arguments based on any existing physical theories.

At the same time, Kirchhoff and fellow German chemist R. W. Bunsen developed the foundation of elemental spectrochemical analysis. Built upon the earlier work of physicists A. J. Ångström and J. B. L. Foucault, who as early as 1849 had independently measured that elemental Hydrogen emitted at four characteristic wavelengths, spectrochemical analysis leveraged a flame source and state-of-the-art optics to measure the emission (or equivalently, absorption) lines for a given substance to identify its elemental composition. While electromagnetism was sufficient for describing the behavior of the light once it was emitted, it provided neither an explanation for why only certain wavelengths were emitted from each element nor could it predict further spectral lines for any elements given the first few. On the other hand, Sir W. N. Hartley, J. J. Balmer, and finally J. R. Rydberg independently observed and developed expressions connecting whole-number ratios to the visible spectral lines of the Hydrogen atom. While seeming to have no rigorous basis in elec-

tromagnetism, Rydberg's expression nevertheless predicted — and subsequently provided excellent agreement with — further series of Hydrogen atom spectral lines observed by Lyman in the ultraviolet and infrared regions of the electromagnetic spectrum. Unfortunately, however, the Rydberg formula could not be applied to predict the spectra of any element other than Hydrogen. Like blackbody radiation, it seemed, atomic spectra were among the phenomena not sufficiently described by existing physical theories.

As the second half of the 19th century continued onward, it became increasingly difficult to ignore that for phenomena that occurred at the micro- and sub-microscopic scale, especially in the case of the interaction of light and matter, classical physical theories were insufficient. Many physicists resisted acknowledging this fact, instead laying blame at the feet of the instrumentation (and sometimes, even their colleagues) measuring the phenomena for which these macroscopic theories broke down. A bold few, on the other hand, took inspiration from the world of chemistry, where Dalton's atomic theory had proven to be a source of clarity in the field. One of these luminaries was Ludwig Boltzmann, who had already pioneered the kinetic theory of gases (now widely referred to as kinetic molecular theory, KMT) during his doctoral studies. Boltzmann's greatest contribution, however, was his development of a statistical formulation of the Second Law of Thermodynamics, inspired by the atomic viewpoint of Dalton which reconciled atomic theory with the predictions of macroscopic thermodynamics.

In his formulation, the universe's ever-growing entropy was not an ethereal law, but rather a probabilistic byproduct of the fact it is relatively more likely to find particles in a statistically disordered state than an ordered one. Therefore, Boltzmann argued, it is also more likely for a system to tend towards disorder rather than order, thereby increasing the total entropy or "disorder" in the universe. Unfortunately, as this perspective was both consistent with and indeed inspired by Dalton's atomic theory, Boltzmann's formulation of entropy was rather controversial and drew widespread disapproval from many of his contemporaries. This disapproval became so pointed that it has been cited as a pos-

sible influence for Boltzmann's unexpected and tragic death by suicide in 1906.\* Despite the rejection of his approach, Boltzmann's formulation of entropy accomplished what his contemporaries failed to do: reconcile atomic theory with the predictions of macroscopic thermodynamics. As it would happen, it was precisely this success which provided the foundation for the next great scientific revolution.

### 1.1.3 The Quantization Revolution

In December of 1900, more than forty years since its original introduction, the problem of blackbody radiation was finally solved when Max Planck became the first physicist to present an expression for the spectral intensity of blackbody radiation which matched experimental observations over the entire spectral range. Planck derived his expression by imagining that the radiation emitted from a blackbody was produced by a finite number of oscillators, amongst which the energy of the radiation emitted from the body was distributed evenly in finite "energy elements" called *quanta*,<sup>†</sup> an idea which was inspired by Boltzmann's atomic formulation of entropy. The significance of Planck's division of the energy distribution into quanta was largely overlooked until Albert Einstein applied the same idea to describe light energy in 1905 with his work on the photoelectric effect, and again in 1908 with his treatment of the heat capacity of solids. All of the sudden, the cat was out of the bag<sup>‡</sup> and the significance of quantization, as well as its discrepancy with the predictions of classical electromagnetism and thermodynamics, could no longer be ignored.

Over the next three decades, an explosion of new physics emerged which revolutionized our understanding of both the world around us and our place in it. What began as a necessary assumption to solve outstanding physical problems was soon enshrined as its own theory, called *quantum mechanics*, which describes the behavior of quantized versions

---

\*Even though he had been prone to experiencing depressive episodes throughout his life, it has been widely suspected that his deteriorating mental health in the period leading up to his suicide was influenced by the rejection of his atomistic approach to thermodynamics by his contemporaries in the theoretical physics community.

<sup>†</sup>Something wicked this way comes!

<sup>‡</sup>Or, perhaps more appropriately, the box!

of classically continuous quantities, like energy, at very small length scales. The successes of Planck's expression for the spectral intensity of blackbody radiation and Einstein's work on the photoelectric effect and the heat capacity of solids were major victories for the burgeoning field of quantum mechanics. Furthermore, as they were partially based on Boltzmann's formulation of entropy, they vindicated his controversial "atomistic" approach to thermodynamics that reconciled John Dalton's atomic theory with that of macroscopic thermodynamics. Even though quantum mechanics (and Boltzmann's atomism) enjoyed early victories, the broader acceptance of these new ideas would be predicated on the development of a complete atomic theory, an effort which had by this time already been nearly 100 years in the making.

#### 1.1.4 Atomic Theory Redux

Despite the success of Dalton's atomic theory for describing and predicting chemical phenomena, many physicists resisted "atomism" on the basis that it seemed to disagree with classical thermodynamics and Newtonian mechanics, where all matter behaves in exactly the same way regardless of size, and quantities like energy were continuous. Even when Boltzmann reconciled these two perspectives by developing a microscopic formulation of entropy, many leading physicists resisted (and even openly ridiculed) atomism because it seemed to imply that events could be probabilistic, rather than the perfectly deterministic picture provided by Newtonian mechanics. Aside from hubris and a fundamental dissatisfaction with early statistical mechanics, several discrepancies existed between Dalton's atomic theory and experiments which fueled the predominant anti-atomic view within the physics community. Chief among these was that Dalton's theory rested on the hypothesis that no particle can be smaller than an atom. Only a few years earlier, however, in 1897, J. J. Thompson had measured the mass of the electron — the negatively charged particle responsible for electric current — to be nearly 1,000 times smaller than that of the lightest atom, Hydrogen. Additionally, in 1900, Henri Becquerel showed that the particles emitted

from Radium atoms (which had been labelled  $\beta$  particles upon their discovery by Ernest Rutherford in 1899) were of identical mass and charge to the electron, which seemed to indicate that electrons somehow *came from atoms*. Even with these discrepancies, however, the vindication of Boltzmann's atomistic formulation of entropy by Planck's law for blackbody radiation and by Einstein's work on the photoelectric effect made it clear that atomism must contain at least some kernel of truth, even if it was not perfect. This only further begged the question: *what exactly do atoms look like?*

In the wake of Bequerel's finding that electrons seemed to originate from atoms, Thompson proposed that atoms — which were known to be electrically neutral — are made up of electrons resting in a uniform cloud of positive charge. This model of atomic structure is referred to as the “plum pudding” model, as the distribution of negatively charged electrons throughout the positive charge cloud resembles chunks of fruit distributed throughout a traditional English Christmas pudding.\* According to the plum pudding model, atoms should be largely porous, as the only truly “solid” parts of the atom were the tiny electrons distributed at random throughout the interior of the charge cloud. In 1909, Ernest Rutherford, together with his students Hans Geiger and Ernest Marsden, tested this hypothesis by firing positively charged  $\alpha$  particles (discovered by Rutherford at the same time in 1899 as  $\beta$  particles) at a thin piece of gold foil. If the plum pudding model were correct, the  $\alpha$  particles should be uniformly deflected from their path by only a small angle, which should depend only on the total charge of the positive atomic cloud.

What the now famous gold-foil experiment revealed, however, was quite astonishing: instead of a nearly uniform, small angle of deflection, the team observed most of the  $\alpha$  particles fired at the gold foil passed straight through the material, as if there was no obstacle to their passage whatsoever. For a small fraction of particles, however, a very large deflection angle was measured; the deflection of some particles was so large, in fact, some

---

\*The English are a very strange people. To them, “pudding” refers to a dense, moist cake with fruit and nuts inside, rather than the gelatinous chocolate, vanilla, or butterscotch flavored concoction sold as part of Snack Packs© in the United States.

particles were even deflected *backwards* towards the particle source, as if the particles had been baseballs thrown at a massive boulder. The distribution of deflection angles observed led Rutherford to conclude first that atoms were comprised largely of empty space (allowing the majority of particles to pass through undeflected), and furthermore that all of the positive charge was concentrated at a single, dense point in the center of the atom called the *nucleus*, rather than distributed in a diffuse cloud. This model, published in 1911, hypothesized that electrons orbited the nucleus in much a similar manner to which the planets in our solar system orbit the Sun, earning it the label the “planetary model” of atomic structure.

As with every other model for atomic structure proposed to that point, Rutherford’s planetary model was not without its drawbacks. In particular, classical mechanics predicted electrons in the planetary model should lose energy by emitting light radiation as they orbit the nucleus, causing them to lose energy and collapse into the nucleus like a satellite crashing back into the Earth after running out of fuel. To address this deficiency, Danish physicist Niels Bohr proposed in 1913 that electrons revolve around the nucleus in certain stable orbits, with radius determined by the electron’s angular momentum. Bohr derived this result by assuming that the electron’s angular momentum could only take on integer multiples of Planck’s constant — in other words, by assuming electronic angular momentum was *quantized*. The results of this quantization were threefold: first, energy could only be gained or lost by an electron in discrete chunks; second, electrons could only move between orbits by leaping directly from one to the other, without moving through the space in between; and third, the energy gained or lost by an electron is in the form of light energy emitted or absorbed by the electron, the magnitude of which corresponds exactly to the difference in energy between the two orbits.

At the time, Bohr’s model was a triumph: not only did it corroborate Rutherford’s model of the atom, but it also was able to exactly reproduce — and provide a theoretical argument for — the emission spectrum of the Hydrogen atom, which as discussed above had been a source of consternation since even before blackbody radiation. Furthermore, Bohr’s

model provided the means by which to rigorously define the value of the Rydberg constant in terms of more fundamental quantities like the charge of the electron and Planck's constant, which beforehand had only been known empirically. Apparently, however, there was once again more physics to be discovered, as even Bohr's model failed to predict the emission spectra for any atom with two or more electrons.

#### 1.1.5 Towards a General Formulation of Quantum Mechanics

Through the first two decades of the twentieth century, all applications of “quantum mechanics” had been phenomenological (i.e., applied specifically to address particular unsolved problems) rather than having been developed into a theory which is generally applicable. The key to making this leap was published in 1924 as part of likely the most influential doctoral thesis of all time: *Researches on the quantum theory* by Louis de Broglie. Starting from Einstein's special relativity, which posited that light experiences *wave-particle duality* where photons can act both as a wave and as a particle (albeit one without mass), de Broglie showed electrons also experience wave-particle duality, and can behave both as a point particle and as a “matter wave.”\* Furthermore, de Broglie postulated that the wavelength for a particle with mass  $m$  travelling with velocity  $v$  was determined completely by these quantities.† This discovery led to the independent development in 1926 of two comprehensive formulations of a general quantum theory: *matrix mechanics*, proposed by Werner Heisenberg, Max Born, and Pascual Jordan, and *wave mechanics*, proposed by Erwin Schrödinger. These two formulations of quantum mechanics, though initially at odds, were later shown to be equivalent by John von Neumann and Marshall

---

\*The first observation of the wave nature of matter was made in 1927 by Clinton Davisson and Lester Germer, when they measured that the diffraction pattern of a beam of electrons incident upon a nickel surface was identical to that of X-ray radiation. Thanks to their confirmation of wave-particle duality, Louis de Broglie won the Nobel Prize in Physics in 1929; Davisson then went on to share part of the 1937 Nobel Prize in Physics for the 1927 experiment.

†The *de Broglie wavelength* provides a heuristic for evaluating when it is necessary to apply quantum, rather than classical, mechanics to describe the dynamics of a particle: if the de Broglie wavelength for a particle is within approximately three orders of magnitude ( $\sim 1,000\times$ ) of the particle's diameter, then quantum mechanics is more appropriate than classical mechanics.



Stone in 1931.\*

As stated above, chemistry is the study of the transfiguration of matter from one form into another. Referred to as *chemical reactions*, these processes typically involve the transfer of one or more atoms or electrons from one molecule to another. Based on the de Broglie wavelength, typical chemical processes should be expected to be governed by quantum mechanics;† in order to develop a complete understanding of chemical phenomena, therefore, we must build from a foundation of quantum mechanics. This task was first begun in 1927, within a single year of the introduction of wave- and matrix mechanics, when Walter Heitler and Fritz London‡ used the Schrödinger picture of the Hydrogen atom to describe the interatomic interaction in molecular Hydrogen, H<sub>2</sub>. Heitler referred to this interaction as a *covalent bond*, which he characterized by the sharing of the two electrons brought by each Hydrogen atom between the two nuclei in the molecule.

The next year, in 1928, this idea was generalized into the now-famous *valence bond theory* by Linus Pauling, in which a bond is comprised of one or more pairs of electrons shared between adjacent atomic nuclei in a molecule. At the same time, Robert Mulliken and Friedrich Hund developed a rival bonding theory known as *molecular orbital theory*, in which the local “bonds” between adjacent atoms from valence bond theory were replaced by molecular orbitals (MOs) which extended spatially over the entire molecule. Finally, then, in 1929, John Lennard-Jones introduced the notion that the molecular orbitals of Mulliken and Hund could be approximately constructed by taking a linear combination of Hydrogen-like “atomic orbitals” (AOs). This approach, referred to as LCAO-MO, has become the *de facto* bonding theory leveraged by chemists the world over to rationalize and predict chemical reactivity, geometry, and stability, among a host of other behaviors. It must be noted, however, that the LCAO-MO approach is not — for all its popularity —

---

\*Due only to the convenience of representation and more straightforward formulation, we will develop the basic framework of quantum mechanics in Chapter 2 using the Schrödinger formulation.

†For example, a water molecule at room temperature in the liquid phase will have de Broglie wavelength approximately 25% of the largest dimension of the water molecule itself. Since this is well within the heuristic limit of 1000×, quantum mechanics is the proper governing theory.

‡London, of “dispersion” fame, will turn out to be of central historical importance to this Thesis.

the truth, but rather a convenient construction allowing for the very abstract to be made (at least somewhat) tangible, a concept which will be explored more fully in Chapter 2. Let us therefore abandon our purely historical discussion in favor of considering the particular chemical phenomena of interest to this Thesis: non-covalent interactions.

## 1.2 Non-Covalent Interactions in Chemistry, Biology, and Physics

Based on our prior understanding of what constitutes a covalent bond, a *non-covalent interaction* (NCI) is any interaction between nearby atoms and/or molecules which does not expressly involve the sharing of a pair of electrons between the participating species. Included under the umbrella of NCI, therefore, are the attraction or repulsion felt between charged (ionic) or even neutral species with uneven distribution of electrons (i.e., dipolar, quadrupolar, etc. molecules), as well as a host of other intermolecular forces with a variety of underlying physical causes. In the earliest conceptualization of an atomic theory, Democritus and Leucippus envisioned that atoms would come together and join mechanically to create new substances, whose properties were distinct from those of the constituent atoms; depending, rather, on the manner of interaction between the atoms themselves. As it turns out, this picture corresponds better to *molecules* than to atoms, where the mechanical interactions envisioned by Democritus and Leucippus are actually NCI. While molecules do not often mechanically interlock with one another in the manner imagined by early atomists, their intuition that the macroscopic properties of a substance are dependent upon the manner in which its particles interact was remarkably prescient. Indeed, NCI govern both a substance's phase diagram and a host of other physical properties, including its density, surface tension, enthalpy of vaporization, and solubility.

In addition to these physical properties, NCI also play a vital role in chemical reactivity and mechanism. This may seem counterintuitive,\* however an extremely diverse array of reactions are influenced by *steric effects*, which arise when atoms or molecules “crowd”

---

\*Didn't we just define NCI to be non-bonding interactions?

each other. This does not mean different species literally touch one another; rather, they repel each other because their proximity is energetically unfavorable.\* Finally, NCI play a critical role in nearly every biochemical process. This includes the molecular recognition critical for intra- and intercellular signalling, allosteric or inhibitory control of enzyme function, and even the translation and transcription processes which transforms genetic information (itself dependent on NCI for storage and replication!) into proteins. Due to their ubiquitous role in governing this myriad of biological, chemical, and physical properties and processes, understanding the fundamental nature of NCI and the manner in which NCI control these processes is of central importance to understanding these phenomena. So, how can we study NCI themselves?

NCI could be probed indirectly by observing the processes controlled by them; this is the case for inferring that water–water interactions are stronger than He–He interactions due to the drastic difference in these substances’ boiling points. While this may seem crude, many clever and sophisticated experiments have been devised which can offer a great deal of insight from even this type of indirect observation. Far from observing the effects of non-covalent interactions on macroscopic chemical or physical behavior, however, investigating these interactions within a particular chemical system of interest can be a significant challenge. To do this, competition experiments or mass spectrometry can be utilized. Despite the ability of these (and other) experimental approaches to directly or indirectly quantify the strength of non-covalent interactions in a particular chemical system, they unfortunately cannot inform the optimization of these interactions by providing further insight into the *cause* of the interactions themselves.

In contrast, investigations of non-covalent interactions undertaken from a strictly theoretical perspective — built upon the quantum mechanical description of the atoms and molecules participating in the interaction — can provide the insight necessary to rationally design optimally interacting systems by answering questions not only of *how* NCI influ-

---

\*We will discuss this effect in much greater detail in Chapter 3, but for now, think of it like someone on the bus invading your personal space.

ence a chemical or physical process, but also *why*. Furthermore, theoretical investigations are predictive without ever needing to step into the laboratory and before performing even a single experiment, making them a green approach to molecular design. For these and many other reasons, theoretical investigation of NCI (among a host of other observable properties) has become a routine and integral part of chemical scientific discovery, especially when combined with subsequent experimental investigation. Towards this end, it is the goal of this Thesis to contribute to the existing body of knowledge regarding the theoretical investigation of non-covalent interactions by electronic structure theory, and to extend the current state of the art to the investigation of NCI directly within extended and complex chemical environments.

### 1.3 Prospectus

This Thesis will be organized into four Parts, each concerned with telling a piece of the story. In Part I, Chapters 2 and 3 will present the background information needed to make the rest of this thesis accessible, aimed at the advanced undergraduate and junior graduate level. In Part II, Chapters 4 and 5 will describe efforts to benchmark non-covalent interactions and the geometries of non-bonded complexes using explicitly correlated variants of coupled cluster theory, density functional theory, and symmetry-adapted perturbation theory, in order to obtain “the right answer for the right reasons.”\* Next, in Part III, Chapter 6 will discuss efforts to develop an approximate, semi-empirical variant to symmetry-adapted perturbation theory which can be applied to large chemical systems while still achieving high accuracy relative to reference interaction energies. Finally, in Part IV, Chapter 7 will discuss the application of both existing as well as our developed methodologies to address interesting questions in diverse chemical systems, namely to understand how NCI between solute molecules are tuned by their solvent environment.

---

\*This quote is taken from reviewer comments on the publication reproduced in Chapter 4. Even though it was slightly annoying at the time, I have since come to agree with this reviewer: getting the right answer for the right reason, rather than dumb luck, is so much more satisfying.

## **PART I**

### **THEORETICAL BACKGROUND TO THE THESIS PROJECT**

## CHAPTER 2

### INTRODUCTION TO ELECTRONIC STRUCTURE THEORY

Almost as soon as general frameworks for quantum mechanics were introduced (wave mechanics by Schrödinger and matrix mechanics by Heisenberg, Born, and Jordan) in 1926, the race to apply these newly proposed theories to rationalize chemical phenomena was on. Heitler and London were the first to successfully do so, when in 1927 they described the covalent bond in molecular hydrogen using Schrödinger's wave mechanics. Application of quantum mechanics to more complex molecules, however, required that a more general approach be developed. Towards this end, Pauling proposed valence bond (VB) theory in 1928 and Mulliken & Hund proposed molecular orbital (MO) theory in 1929. While both of these approaches provided a general bonding theory for arbitrary molecules, it became clear to the community that they were each too complex to be applied exactly. Indeed, as it was famously stated by Dirac in 1929,<sup>3</sup>

“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.”

In this Chapter, we will first introduce both the “insoluble” working equations of non-relativistic quantum mechanics, as well as presenting their simplified forms for atoms and molecules in *electronic structure theory*. We will then take inspiration from Dirac by introducing several approximate, practical methods for applying this theory to understand chemical phenomena, developed in the intervening 91 years since his prophecy. Finally,

we will discuss several practical considerations which must be kept in mind when utilizing the approaches developed here, as well as their consequences for the results presented in the remainder of the Thesis.

## 2.1 Basic Formulation

Electronic structure theory is concerned with solving the non-relativistic, time-independent Schrödinger equation

$$\widehat{\mathcal{H}}\Psi = \mathcal{E}\Psi, \quad (2.1)$$

where the fundamental quantity of interest is the *wavefunction*,  $\Psi$ , which contains all information necessary to compute any observable property of the system,  $\mathcal{E}$  is the energy of the system, and  $\widehat{\mathcal{H}}$  is the non-relativistic, time independent Hamiltonian operator. For a molecule, the Hamiltonian  $\widehat{\mathcal{H}}_{\text{molec}}$  is given by

$$\widehat{\mathcal{H}}_{\text{molec}} = -\sum_{i=1}^N \frac{1}{2} \nabla_i^2 - \sum_{A=1}^M \frac{1}{2M_A} \nabla_A^2 - \sum_{i=1}^N \sum_{A=1}^M \frac{Z_A}{r_{iA}} + \sum_{i=1}^N \sum_{j>i}^N \frac{1}{r_{ij}} + \sum_{A=1}^M \sum_{B>A}^M \frac{Z_A Z_B}{R_{AB}} \quad (2.2)$$

$$\widehat{\mathcal{H}}_{\text{molec}} = -\hat{T}_e - \hat{T}_N - \hat{V}_{eN} + \hat{V}_{NN} + \hat{V}_{ee}, \quad (2.3)$$

whose terms arise from (i) the kinetic energy of the electrons, (ii) the kinetic energy of the nuclei, (iii) the electron–nuclear attraction, (iv) the electron–electron repulsion, and (v) the nuclear–nuclear repulsion. For atoms and molecules more complex than the Hydrogen molecular ion  $\text{H}_2^+$ , the analytic solution to this equation is not known; to describe atoms and molecules with more than one electron, and more than two nuclei, therefore, a hierarchy of approximations must be made and various approximate solution methods must be employed.

## 2.2 Approximate Solution Methods

For several famous equations in physics, no analytic or general solution is known. In this section, we will briefly introduce two approximate solution methods employed throughout the rest of this and the following Chapter, and indeed the entire Thesis: linear variation and perturbation theory.

### 2.2.1 The Method of Linear Variations

Before moving on to discuss the various approximations which must be invoked to make solving the Schrödinger equation possible, we must first address how we will know that any approximate wavefunction is *good enough*. For this, we introduce the *Variation Principle*, as stated (and proved) in Ref. 4:

**Theorem 2.2.1.1** (The Variation Principle). *Given an eigenvalue problem*

$$\hat{\mathcal{O}}\phi = \omega\phi,$$

where  $\hat{\mathcal{O}}$  is a Hermitian operator and a normalized trial function  $\psi$  with the appropriate boundary conditions, the expectation value of  $\hat{\mathcal{O}}$  by  $\psi$  will always be an upper bound to the exact value of the lowest eigenvalue,  $\omega_0$ :

$$\langle \psi | \hat{\mathcal{O}} | \psi \rangle \geq \omega_0$$

Therefore, a trial wavefunction  $\tilde{\Psi}$  which approximately solves the Schrödinger equation will be improved by adjusting the parameters (or functional form, etc.) in such a way that the expectation value  $\langle \tilde{\Psi} | \hat{\mathcal{H}} | \tilde{\Psi} \rangle$  is lowered. In general, then, this will be our strategy: to construct a trial wavefunction and iteratively adjust it so that its expectation value with the exact Hamiltonian operator is lowered.



## 2.2.2 Rayleigh-Schrödinger Perturbation Theory

Alternatively, let's assume that  $\widehat{\mathcal{H}}$  is separable into a piece which can be solved exactly, denoted  $\widehat{H}^{(0)}$  and referred to as the zeroth-order Hamiltonian, plus some small additional *perturbation*  $\widehat{V}$ ,

$$\widehat{\mathcal{H}} = \widehat{H}^{(0)} + \lambda \widehat{V}, \quad (2.4)$$

where  $\lambda \in [0, 1]$  is referred to as the *perturbation strength*. Considering the full Schrödinger equation defined by  $\widehat{\mathcal{H}}$ ,

$$\widehat{\mathcal{H}}|\Psi\rangle = \left(\widehat{H}^{(0)} + \lambda \widehat{V}\right)|\Psi\rangle = \mathcal{E}_n(\lambda)|\Psi_n(\lambda)\rangle, \quad (2.5)$$

where the exact eigenvectors and eigenvalues are now functions of the perturbation strength. Since it is not clear exactly *how* the energy and wavefunction depend on  $\lambda$ , we may expand them as Taylor series about  $\lambda = 0$ :

$$|\Psi_n(\lambda)\rangle = |\Psi_n\rangle|_{\lambda=0} + \left.\frac{\partial|\Psi_n\rangle}{\partial\lambda}\right|_{\lambda=0} \lambda + \left.\frac{\partial^2|\Psi_n\rangle}{\partial\lambda^2}\right|_{\lambda=0} \frac{\lambda^2}{2!} + \cdots + \left.\frac{\partial^k|\Psi_n\rangle}{\partial\lambda^k}\right|_{\lambda=0} \frac{\lambda^k}{k!} + \cdots \quad (2.6)$$

$$\mathcal{E}_n(\lambda) = \mathcal{E}_n|_{\lambda=0} + \left.\frac{\partial\mathcal{E}_n}{\partial\lambda}\right|_{\lambda=0} \lambda + \left.\frac{\partial^2\mathcal{E}_n}{\partial\lambda^2}\right|_{\lambda=0} \frac{\lambda^2}{2!} + \cdots + \left.\frac{\partial^k\mathcal{E}_n}{\partial\lambda^k}\right|_{\lambda=0} \frac{\lambda^k}{k!} + \cdots \quad (2.7)$$

For the sake of brevity, let the Taylor coefficients in these expansions be denoted

$$|\Psi_n^{(k)}\rangle = \frac{1}{k!} \left.\frac{\partial^k|\Psi_n\rangle}{\partial\lambda^k}\right|_{\lambda=0} \quad (2.8)$$

$$\mathcal{E}_n^{(k)} = \frac{1}{k!} \left.\frac{\partial^k\mathcal{E}_n}{\partial\lambda^k}\right|_{\lambda=0}; \quad (2.9)$$

then, the exact wavefunction and energy will be given by

$$|\Psi_n\rangle = |\Psi_n^{(0)}\rangle + |\Psi_n^{(1)}\rangle\lambda + |\Psi_n^{(2)}\rangle\lambda^2 + \dots + |\Psi_n^{(k)}\rangle\lambda^k \quad (2.10)$$

$$\mathcal{E}_n = \mathcal{E}_n^{(0)} + \mathcal{E}_n^{(1)}\lambda + \mathcal{E}_n^{(2)}\lambda^2 + \dots + \mathcal{E}_n^{(k)}\lambda^k. \quad (2.11)$$

Substituting these expansions into the full Schrödinger equation yields

$$\left(\widehat{H}^{(0)} + \lambda\widehat{V}\right) \left[ \sum_{k=0}^{\infty} \lambda^k |\Psi_n^{(k)}\rangle \right] = \left[ \sum_{k=0}^{\infty} \lambda^k \mathcal{E}_n^{(k)} \right] \left[ \sum_{k=0}^{\infty} \lambda^k |\Psi_n^{(k)}\rangle \right] \quad (2.12)$$

For this expression to be true for all  $\lambda \in [0, 1]$ , the terms on the left hand side must equal the terms on the right hand side for a given power of  $\lambda$ . By equating the terms which are zeroth-order in  $\lambda$ , we have

$$\widehat{H}^{(0)} |\Psi_n^{(0)}\rangle = \mathcal{E}_n^{(0)} |\Psi_n^{(0)}\rangle; . \quad (2.13)$$

Therefore, the “zeroth order correction” to the energy and wavefunction are the zeroth-order energy and wavefunction themselves. Continuing this process, it can be shown that

$$|\Psi_n^{(1)}\rangle = \sum_{m \neq n} \frac{\langle \Psi_m^{(0)} | \widehat{V} | \Psi_n^{(0)} \rangle}{\mathcal{E}_m^{(0)} - \mathcal{E}_n^{(0)}} |\Psi_m^{(0)}\rangle \quad (2.14)$$

$$\mathcal{E}_n^{(1)} = \langle \Psi_n^{(0)} | \widehat{V} | \Psi_n^{(0)} \rangle, \quad (2.15)$$

by collecting terms which are first-order in  $\lambda$ , and also that the  $k$ th order energy correction\* is given by

$$\mathcal{E}_n^{(k)} = \langle \Psi_n^{(0)} | \widehat{V} | \Psi_n^{(k-1)} \rangle. \quad (2.16)$$

Ideally, subsequent correction orders will provide corrections which are progressively smaller in magnitude, i.e., that the Taylor expansions of the exact energy and wavefunc-

---

\*The derivation for the expression of the  $k$ th order wavefunction is beyond the scope of this discussion, and is therefore omitted.

tion are *convergent* series. Unfortunately, this is not typically the case: not only are these series not guaranteed to be convergent, where the corrections monotonically decrease in magnitude, in practice they often actually *diverge*. Therefore, while RSPT was developed to provide systematically improvable approximate solutions to *any* Schrödinger equation which is not analytically solvable, it is best applied in low perturbation orders (i.e., for  $k \leq 4$ ) to problems for which the perturbation  $\hat{V}$  is small in magnitude. Otherwise, the expansions in  $\lambda$  can diverge rapidly, and even low-order corrections can be suspect. Cautions aside, RSPT is a magnificent tool when applied within its scope, and can be used to derive approximate solutions for everything from a particle in a “slanted” box (where  $\hat{V} = x$ ) to the anharmonic oscillator. For our purposes here, however, its most relevant application is to computing the *correlation energy*, to be defined below, which will be developed and discussed in greater detail in Section. 2.7.2.

### 2.3 The Born–Oppenheimer Approximation

Now that we have a general strategy for approximately solving the Schrödinger equation, we can move on to the actual attempt. The most common first approximation to this equation originally introduced by Max Born and J. Robert Oppenheimer in 1927<sup>5</sup> is to assume that the momenta of the nuclei are sufficiently small compared to that of the electrons, that the nuclei of a molecule would appear stationary from the electrons’ perspective. This is due to the nearly  $1,800\times$  larger mass of a proton than that of an electron. Referred to as the *Born–Oppenheimer* (BO) approximation, the nuclei can therefore be effectively “clamped” in place to reduce the total molecular Hamiltonian to one only explicitly dependent on electronic coordinates:

$$\hat{\mathcal{H}}_{\text{elec}} = - \sum_{i=1}^N \frac{1}{2} \nabla_i^2 - \sum_{i=1}^N \sum_{A=1}^M \frac{Z_A}{r_{iA}} + \sum_{i=1}^N \sum_{j>i}^N \frac{1}{r_{ij}} + \sum_{A=1}^M \sum_{B>A}^M \frac{Z_A Z_B}{R_{AB}} \quad (2.17)$$

$$\hat{\mathcal{H}}_{\text{elec}} = -\hat{T}_e(\mathbf{r}) - \hat{V}_{eN}(\mathbf{r}; \mathbf{R}) + \hat{V}_{NN}(\mathbf{R}) + \hat{V}_{ee}(\mathbf{r}). \quad (2.18)$$

In this new electronic Hamiltonian, since the nuclei are stationary, the nuclear kinetic energy is zero and the nuclear-nuclear repulsion is a constant. Furthermore, the electron–nuclear attraction is only *parametrically* dependent on the nuclear positions, denoted by separating the electronic coordinates ( $\mathbf{r}$ ) from the nuclear ones ( $\mathbf{R}$ ) by a semicolon. Under the Born–Oppenheimer approximation, the Schrödinger equation becomes

$$\widehat{\mathcal{H}}_{\text{elec}}\Phi_{\text{elec}}(\mathbf{r}; \mathbf{R}) = \mathcal{E}_{\text{elec}}(\mathbf{R})\Phi_{\text{elec}}(\mathbf{r}; \mathbf{R}) \quad (2.19)$$

where the electronic energies  $\mathcal{E}_{\text{elec}}$  now depend explicitly on the nuclear positions  $\mathbf{R}$ . Therefore, for a polyatomic molecule of  $M$  atoms, the nuclear positions  $\mathbf{R}$  define a  $3M - 6$  dimensional potential energy surface (PES) (or  $3M - 5$  dimensional if the molecule is linear) upon which the nuclei rest. In the field of electronic structure theory, the primary concern is with solving the electronic Schrödinger equation for a given set of nuclear positions; often, this is sufficient, however some dynamical properties which require nuclear dynamics can require that the full PES be constructed, upon which the nuclear trajectories can be propagated.

## 2.4 Independent Particle Models

Even after decoupling the electronic and nuclear degrees of freedom by invoking Born–Oppenheimer, unfortunately, the problem of solving the electronic Schrödinger equation does not become any more tractable: there is still the problem of multiple interacting electrons. The simplest approximation which could address this issue would be to assume that the electrons do not interact; under this assumption, the electron–electron repulsion would be zero, and the Hamiltonian would become

$$\widehat{\mathcal{H}} = \sum_{i=1}^N \hat{h}(i), \quad (2.20)$$

where each  $\hat{h}(i)$  consists of the kinetic and potential energy operators corresponding to electron  $i$ . These one-particle Hamiltonian operators define their own Schrödinger equations,

$$\hat{h}(i)\chi_j(\mathbf{x}_i) = \epsilon_j\chi_j(\mathbf{x}_i), \quad (2.21)$$

where  $\mathbf{x}_i = (\mathbf{r}_i, \sigma_i)$  collects both spatial ( $\mathbf{r}$ ) and spin ( $\sigma$ ) coordinates for electron  $i$ , and the set of eigenfunctions  $\{\chi_j\}$  referred to as “spin orbitals.” The independent particle Schrödinger equation is given by

$$\left[ \sum_{i=1}^N \hat{h}(i) \right] \Psi^{\text{HP}} = E\Psi^{\text{HP}}, \quad (2.22)$$

where the eigenfunctions are a simple product of the one-electron spin orbitals

$$\Psi^{\text{HP}} = \chi_i(\mathbf{x}_1)\chi_j(\mathbf{x}_2) \cdots \chi_n(\mathbf{x}_N). \quad (2.23)$$

For this wavefunction, referred to as a *Hartree product* (HP), the subscripts on the variables  $\mathbf{x}$  denote the electrons (which are different than the identities of the spin orbitals), and the energy eigenvalue will be a sum of the spin orbital energies

$$E = \epsilon_i + \epsilon_j + \dots + \epsilon_n. \quad (2.24)$$

As was assumed when constructing the independent-particle Hamiltonian, individual electrons described collectively by a HP wavefunction are fully uncorrelated, i.e., their motions do not influence one another. From our wavefunction, it is clear that we have completely distinguished each electron from all others, since electron 1 is in spin orbital  $\chi_i$ , electron 2 is in spin orbital  $\chi_j$ , and so on. As we will see in the next section, this unfortunately violates a postulate of quantum mechanics — the famous Pauli Antisymmetry Principle.

## 2.5 The Antisymmetry Principle and Slater Determinants

Since in this Thesis we are working within a non-relativistic formulation of quantum mechanics, our Hamiltonian operators do not depend on electron spin. Since the Hamiltonian is spin-independent, our wavefunctions could be defined either over spatial ( $\mathbf{r}$ ) or spin [ $\mathbf{x} = (\mathbf{r}, \omega)$ ] coordinates. When incorporating relativistic effects, however, particle spin leads to drastically different behavior for fermions and bosons. Bosonic wavefunctions, where the particles have integer spin, are said to be *symmetric* with respect to the interchange of particles, while fermionic wavefunctions, where the particles have half-integer spin, are *antisymmetric* with respect to particle exchange.\* To illustrate this concept, let us first define the permutation operator  $\widehat{\mathcal{P}}_{ij}$ , which permutes two particles  $i$  and  $j$  by exchanging their coordinates  $\mathbf{x}_i$  and  $\mathbf{x}_j$ . Next, consider a wavefunction  $\phi$  describing two  $\text{He}^{2+}$  particles (each of which has spin 0), specified with coordinates  $\mathbf{x}_1$  and  $\mathbf{x}_2$ . Since bosonic wavefunctions are symmetric, applying the permutation operator to the wavefunction produces a result which is equal to the original one:

$$\widehat{\mathcal{P}}_{12}\phi(\mathbf{x}_1, \mathbf{x}_2) = \phi(\mathbf{x}_2, \mathbf{x}_1) = \phi(\mathbf{x}_1, \mathbf{x}_2)$$

Fermionic wavefunctions, on the other hand, are *antisymmetric* with respect to particle exchange. Considering now a two-electron wavefunction  $\psi(\mathbf{x}_1, \mathbf{x}_2)$ , the exchange of particle coordinates results in the negative of the original wavefunction:

$$\widehat{\mathcal{P}}_{12}\psi(\mathbf{x}_1, \mathbf{x}_2) = -\psi(\mathbf{x}_2, \mathbf{x}_1)$$

---

\*This behavior is incorporated into fermionic wavefunctions even for the non-relativistic Schrödinger equation by accepting the antisymmetry principle as a *postulate*.

From this expression, it should be clear why our Hartree product wavefunction is unsatisfactory, since it fails to be antisymmetric with respect to electron exchange:

$$\begin{aligned}
 \widehat{\mathcal{P}}_{12}\Psi_{\text{HP}} &= \widehat{\mathcal{P}}_{12}\chi_i(\mathbf{x}_1)\chi_j(\mathbf{x}_2)\cdots\chi_n(\mathbf{x}_N) \\
 &= \chi_i(\mathbf{x}_2)\chi_j(\mathbf{x}_1)\cdots\chi_n(\mathbf{x}_N) \\
 &\neq -\chi_i(\mathbf{x}_2)\chi_j(\mathbf{x}_1)\cdots\chi_n(\mathbf{x}_N) \\
 &= -\Psi_{\text{HP}}
 \end{aligned}$$

Even though the Hartree Product is not a suitable wavefunction, however, all is not lost. We could obtain an appropriately antisymmetrized wavefunction by constructing the following linear combination of the original and permuted Hartree products:

$$\Psi(\mathbf{x}_1, \mathbf{x}_2) = \frac{1}{\sqrt{2}} [\chi_i(\mathbf{x}_1)\chi_j(\mathbf{x}_2) - \chi_j(\mathbf{x}_1)\chi_i(\mathbf{x}_2)] \quad (2.25)$$

It is trivial to verify that this is indeed an antisymmetric wavefunction by applying our permutation operator over electrons 1 and 2:

$$\begin{aligned}
 \widehat{\mathcal{P}}_{12}\Psi(\mathbf{x}_1, \mathbf{x}_2) &= \frac{1}{\sqrt{2}} [\chi_i(\mathbf{x}_2)\chi_j(\mathbf{x}_1) - \chi_j(\mathbf{x}_2)\chi_i(\mathbf{x}_1)] \\
 &= -\frac{1}{\sqrt{2}} [\chi_i(\mathbf{x}_1)\chi_j(\mathbf{x}_2) - \chi_j(\mathbf{x}_1)\chi_i(\mathbf{x}_2)] \\
 &= -\Psi(\mathbf{x}_1, \mathbf{x}_2)
 \end{aligned}$$

Now that we have constructed a suitably antisymmetric wavefunction, we are free to complain about more inconsequential details, like how annoying it will be to antisymmetrize an  $N$ -electron Hartree product. Fortunately for mine and every other quantum chemists' sanity, John Slater realized that the antisymmetrized Hartree product above could be more

compactly represented as the determinant of a  $2 \times 2$  matrix,

$$\Psi(\mathbf{x}_1, \mathbf{x}_2) = \frac{1}{\sqrt{2}} \begin{vmatrix} \chi_i(\mathbf{x}_1) & \chi_j(\mathbf{x}_1) \\ \chi_i(\mathbf{x}_2) & \chi_j(\mathbf{x}_2) \end{vmatrix}.$$

This *Slater determinant* can then be easily generalized to an  $N$ -electron wavefunction

$$\Psi(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N) = \frac{1}{\sqrt{N!}} \begin{vmatrix} \chi_i(\mathbf{x}_1) & \chi_j(\mathbf{x}_1) & \cdots & \chi_k(\mathbf{x}_1) \\ \chi_i(\mathbf{x}_2) & \chi_j(\mathbf{x}_2) & \cdots & \chi_k(\mathbf{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \chi_i(\mathbf{x}_N) & \chi_j(\mathbf{x}_N) & \cdots & \chi_k(\mathbf{x}_N) \end{vmatrix}, \quad (2.26)$$

which we will always assume to be written in this order and adopt the notation introduced in Ref. 4, where the normalized Slater determinant is represented by a *ket\** containing the spin orbitals (or equivalently, their indices),

$$\Psi(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N) = |\chi_i \chi_j \cdots \chi_k\rangle = |i j \cdots k\rangle. \quad (2.27)$$

From Coulomb's law, we know that electrons should repel one another; indeed, this electron-electron repulsion creates what is known as a *Fermi hole* around each electron where the probability of finding any other electron is zero. Since it is truly an independent particle model, the Hartree product wavefunction completely ignores this behavior: since the motions of any two electrons are uncorrelated, there is a finite probability of violating the Fermi hole. Unlike the Hartree product, however, the picture provided by a Slater determinant wavefunction is one where the movements of electrons with parallel spin (i.e., with all  $\alpha$  or all  $\beta$  spin) are correlated, while those with antiparallel spin (i.e.,  $\alpha - \beta$  pairs) are not. Therefore, electrons with parallel spin obey the Fermi hole, even though antiparallel spin electrons do not. This phenomenon is referred to as *exchange correlation*,

---

\*"Kets" and "bras" refer to Dirac notation for denoting elements of a Hilbert space. See Ref. 4 for an overview.



as it arises only from the antisymmetrization of the wavefunction with respect to particle exchange. Even though a Slater determinant is formally half-correlated, since parallel-spin electrons are correlated but antiparallel-spin electrons are not, it is often referred to as uncorrelated because of a concept more fully explored in Section 2.7.

## 2.6 Hartree–Fock Molecular Orbital Theory

Now that we have constructed the simplest antisymmetric wavefunction for an  $N$ -electron system, a Slater determinant of  $N$  spin orbitals, let us apply the Variation Principle to obtain an approximate solution to the electronic Schrödinger equation. Given such a Slater determinant  $|\Psi_0\rangle = |\chi_i\chi_j\cdots\chi_k\rangle$ , the Variation Principle states that the best  $|\Psi_0\rangle$  is the one for which the expectation value of the electronic Hamiltonian,

$$E_0 = \langle \Psi_0 | \widehat{\mathcal{H}}_{\text{elec}} | \Psi_0 \rangle, \quad (2.28)$$

is minimized; for our trial function  $|\Psi_0\rangle$ , the variational flexibility lies in the choice of the spin orbitals  $\{\chi_i\}$ .

### 2.6.1 The Hartree–Fock Equations

By minimizing  $E_0$  with respect to the spin orbitals, it is possible to derive\* the *Hartree–Fock equations*

$$\hat{f}(i)\chi(\mathbf{x}_i) = \epsilon\chi(\mathbf{x}_i), \quad (2.29)$$

which are a set of  $N$  coupled integro-differential equations defined by the *Fock operator*

$$\hat{f}(i) = -\frac{1}{2}\nabla_i^2 - \sum_{A=1}^M M \frac{Z_A}{r_{iA}} + v^{\text{HF}}(i), \quad (2.30)$$

---

\*For a complete derivation of the Hartree–Fock and Roothaan equations, as well as all operators, integrals, and basis functions, we refer the reader to Chapters 2 and 3 of Ref. 4.

where  $v^{\text{HF}}(i)$  is the average potential experienced by the  $i$ th electron in the field of all other electrons. The Hartree–Fock energy,  $E_0$ , is given by

$$E_0 = \sum_a [a | \hat{h} | a] + \frac{1}{2} \sum_a \sum_b [aa||bb], \quad (2.31)$$

where  $[aa||bb]$  refers to the antisymmetrized two-electron integral

$$[aa||bb] = [aa|bb] - [ab|ba], \quad (2.32)$$

between the “Coulomb” integral

$$[aa|bb] = \int d\mathbf{x}_1 d\mathbf{x}_2 \chi_a^*(\mathbf{x}_1)\chi_a(\mathbf{x}_1) \frac{1}{r_{12}} \chi_b^*(\mathbf{x}_2)\chi_b(\mathbf{x}_2) \quad (2.33)$$

representing the electron–electron repulsion between two electrons with coordinates  $\mathbf{x}_1$  and  $\mathbf{x}_2$  occupying spin orbitals  $\chi_a$  and  $\chi_b$ , and the “exchange” integral

$$[ab|ba] = \int d\mathbf{x}_1 d\mathbf{x}_2 \chi_a^*(\mathbf{x}_1)\chi_b(\mathbf{x}_1) \frac{1}{r_{12}} \chi_b^*(\mathbf{x}_2)\chi_a(\mathbf{x}_2) \quad (2.34)$$

(since the indices have been exchanged from the Coulomb integral) which, unfortunately, does not have such a convenient physical interpretation. Collectively, these two-electron integrals are referred to as electron repulsion integrals (ERIs), since they arise from the effective potential operator  $v^{\text{HF}}(i)$ ; commonly, the Coulomb integral is denoted  $J_{ab}$ , and the exchange integral is denoted  $K_{ab}$ . This set of  $N$  Hartree–Fock equations is coupled due to the potential operator, as every spin orbital feels the potential induced by all other spin orbitals. The Fock operator is, however, a *one electron operator* because it only operates on a single electron at a time. Therefore, we have removed the explicit dependence of the electronic Hamiltonian on the reciprocal of the interelectronic distances,  $1/r_{ij}$ , by replacing this potential with a mean field; this fact has earned Hartree–Fock theory the moniker of a

*mean field theory.*

### *Transitioning from Spin to Spatial Orbitals*

Even though we have recast our  $N$ -electron problem into a set of  $N$  coupled one-electron problems, the Hartree–Fock equations given above present a particular challenge to solve in practice. Not only must the set of integro-differential equations be solved simultaneously, but the Slater determinant wavefunction is comprised of spin orbitals, each involving some mysterious functions  $\alpha(\omega)$  and  $\beta(\omega)$  multiplied by different, one-electron spatial functions. If we assume, however, that electrons with opposite spin are *paired* and are restricted to share a single spatial orbital, it can be shown<sup>4</sup> that the spin coordinate  $\omega$  can be integrated away to yield the *restricted Hartree–Fock* (RHF) equations,

$$\hat{f}(i)\psi(\mathbf{r}_i) = \epsilon\psi(\mathbf{r}_i), \quad (2.35)$$

where our orbitals  $\psi(\mathbf{r}_i)$  are now exclusively spatial.\* The restricted Hartree–Fock energy for a closed-shell (where all electrons are paired) ground state is

$$E_0 = 2 \sum_a (a | \hat{h} | a) + \sum_{ab} 2 (aa|bb) - (ab|ba), \quad (2.36)$$

where the parentheses in the electron repulsion integrals denote Chemist’s notation over spatial orbitals, rather than the spin orbitals as defined above. Using  $J$  and  $K$  notation, and

---

\*At this point, I believe it is worth mentioning that several mathematical artifacts employed in the derivation of the RHF equations are nearly always enshrined in any General Chemistry course as chemical truth: that electrons of opposite spin are somehow paired, and that because of this they are somehow allowed to occupy the same region of space. As we have seen, however, this entire statement is a lie; not only do “orbitals” not exist, as they are just a convenient picture we invoke because we’ve no better ideas to represent an  $N$ -electron molecular wavefunction, but electrons of opposite spins are *not* paired and they cannot occupy the same region of space, as the Fermi hole formally prevents this. As a personal aside, I find the truth that it’s all an illusion we’ve created for ourselves much more comforting than the possibility that these phenomena just *happen*, and that we have no idea why.

denoting  $h_{aa} = (a | \hat{h} | a)$ , this can be rewritten more compactly as

$$E_0 = 2 \sum_a h_{aa} + \sum_{ab} 2J_{ab} - K_{ab} \quad (2.37)$$

### 2.6.2 The Introduction of a Basis Set: The Roothaan Equations

After eliminating spin, the task still remains to solve the coupled set of  $N$  spatial Fock equations, which is not possible by direct or numerical means for all but atomic systems. Instead, Roothaan proposed to introduce a set of  $K$  basis functions  $\{\phi_\mu(\mathbf{r})\}$  and expand the spatial orbitals  $\psi_i$  as a linear combination of these basis functions:

$$\psi_i = \sum_{\mu=1}^K C_{\mu i} \phi_\mu \quad (2.38)$$

Therefore, from a linear combination of “atomic orbitals,” Hartree–Fock (via the Roothaan equations) constructs molecular orbitals. From our chemical intuition, it seems sensible for these basis functions to resemble Hydrogen-like orbitals, since the H-atom wavefunction exactly describes the behavior of a single electron as it orbits the nucleus.\* Rather than use exact Hydrogen atom functions, however, Slater proposed that basis functions should take the form

$$\phi_{abc}^{\text{STO}}(x, y, z) = N x^a y^b z^c e^{-\zeta r}, \quad (2.39)$$

where  $N$  is a normalization constant,  $a$ ,  $b$ , and  $c$  control angular momentum ( $L = a+b+c$ ), and  $\zeta$  controls the width of the orbital, with large  $\zeta$  corresponding to a tight orbital and small  $\zeta$  corresponding to a diffuse orbital.

While these *Slater type orbitals* (STOs) are Hydrogen-like for  $L = 0$  (i.e., for  $s$  orbitals), they are not for other angular momenta since they do not contain exact spherical harmonics to describe the angular distributions of orbital amplitude. STOs are not typically

---

\*This is, of course, the same intuition that led Lennard-Jones in 1929 to develop the LCAO-MO approach for constructing the molecular orbitals described by Mulliken and Hund.

used in practice, however, as evaluating electron repulsion integrals using STOs is computationally intensive. Instead, it is much more common to utilize *Gaussian type orbitals* (GTOs),

$$\phi_{abc}^{\text{GTO}}(x, y, z) = N x^a y^b z^c e^{-\zeta r^2}, \quad (2.40)$$

in which the simple exponential is replaced by a Gaussian function, because the Gaussian product theorem allows for one- and two-electron integrals to be evaluated much more rapidly for GTOs than STOs. Of course, GTOs are no longer Hydrogen-like even for  $L = 0$ , and deviate significantly from the more exact STOs as  $r \rightarrow 0$  and  $r \rightarrow \infty$ . This deficiency for GTOs may be remedied by approximating a STO with a linear combination of GTOs:

$$\phi_{abc}^{\text{CGTO}}(x, y, z) = N \sum_{i=1}^n c_i x^a y^b z^c e^{-\zeta_i r^2} \quad (2.41)$$

As the number of GTOs in the expansion increases, the agreement between this “contracted GTO” (CGTO) and a STO improves, thereby increasing the accuracy of the computation. Due to the combination of convenience and increased accuracy, therefore, CGTOs are the standard choice of basis function among quantum chemists.\*

The simplest method for constructing a basis set — the set of basis functions utilized in Eqn. 2.38 — is to use a single basis function (CGTO) for each electron in a molecule, referred to as a *minimal basis*. The accuracy for energies or properties computed in a minimal basis is, unfortunately — you guessed it — minimal. Aside from either using more primitives in the construction of each CGTO or switching entirely to STOs (both of which significantly increase the cost of integral evaluation), the best method by which to improve the accuracy of the basis set (and thereby the computed energy or property) is to utilize more than one basis function for each electron. In this approach, the number of basis functions included for each electron is indicated by referring to the basis set as a double-

---

\*To distinguish between the GTOs which comprise a single CGTO and the set of CGTOs which together form the basis set leveraged in Eqn. 2.38, we will refer to the raw GTOs as “primitive basis functions” (or more simply, “primitives”), reserving the term “basis functions” for the CGTOs they comprise.

, triple-, or even quadruple- $\zeta$  basis set, with marked increases in accuracy afforded by moving to larger  $\zeta$ -levels without significantly increasing the computational effort required to evaluate integrals. While many different types of basis sets exist, and indeed the design of new basis sets are still an active area of research, we will only concern ourselves with the *choice* of basis set in this Thesis, rather than worrying about the details of their construction.

Returning to Roothaan and his introduction of a basis set, we must substitute Eqn. 2.38 into the Fock equations in order to proceed. Before we do, however, we must recognize that even if the basis functions Eqn. 2.38 could be guaranteed to be mutually orthogonal by construction if they are centered on the same atom, there can definitely be no guarantee that basis functions centered on different atoms of the same molecule will be similarly orthogonal. Therefore, we introduce the *overlap integral*,  $S_{\mu\nu}$ , between basis functions  $\phi_\mu$  and  $\phi_\nu$ :

$$S_{\mu\nu} = \int d\mathbf{r}_1 \phi_\mu(1)\phi_\nu(1) \quad (2.42)$$

Now substituting our basis set expansion into the RHF equations, we transform this set of coupled integro-differential equations into a linear algebra problem known as the Roothaan equation:

$$\mathbf{FC} = \mathbf{SC}\epsilon \quad (2.43)$$

$\mathbf{C}$  is the *orbital coefficient matrix*, collecting the  $C_{\mu i}$ 's from Eqn. 2.38 for all molecular orbitals,  $\epsilon$  is a diagonal matrix containing the orbital energies  $\{\epsilon_i\}$ , and where  $\mathbf{F}$  is the *Fock matrix*, with elements

$$F_{\mu\nu} = \int d\mathbf{r}_1 \phi_\mu(1)\hat{f}(1)\phi_\nu(1). \quad (2.44)$$

At this point, we will refer to the set of basis functions  $\{\phi_\mu\}$  as *atomic orbitals* (AOs), which will always be labeled with Greek letters  $\mu, \nu, \lambda, \sigma$ , etc., to differentiate them from the set of occupied one-electron *molecular orbitals* (MOs)  $\{\psi_i\}$  which are produced by the Hartree–Fock procedure, labeled with  $i, j, k, l$ , etc.

*Solving the Roothaan Equations via the Self-Consistent Field (SCF) Procedure*

The Fock matrix above is the representation of the Fock operator in the AO basis; by inserting Eqn. 2.30 into the expression above, we obtain

$$F_{\mu\nu} = H_{\mu\nu}^{\text{core}} + \sum_i^{N/2} 2(\mu\nu|ii) - (\mu i|i\nu), \quad (2.45)$$

where we have defined the *core Hamiltonian* matrix  $\mathbf{H}^{\text{core}}$  with elements

$$H_{\mu\nu}^{\text{core}} = \int d\mathbf{r}_1 \phi_\mu(1)\hat{h}(1)\phi_\nu(1) = (\mu|\hat{h}|\nu). \quad (2.46)$$

At this point, our expression for the Fock matrix elements is partially represented in terms of AOs  $\phi_\mu$ ,  $\phi_\nu$ , and the MO  $\psi_i$ . Since we do not know *a priori* the form of this MO, we must insert the basis set expansion (Eqn. 2.38) into the expression to yield

$$F_{\mu\nu} = H_{\mu\nu}^{\text{core}} + \sum_i^{N/2} \sum_{\lambda\sigma} C_{\lambda i} C_{\sigma i} [2(\mu\nu|\lambda\sigma) - (\mu\lambda|\nu\sigma)] \quad (2.47)$$

$$= H_{\mu\nu}^{\text{core}} + C_{\lambda i} C_{\sigma i} [2(\mu\nu|\lambda\sigma) - (\mu\lambda|\nu\sigma)], \quad (2.48)$$

where in the second line we have adopted the Einstein summation convention where any repeated index labels are summed over. By defining the *density matrix*,  $\mathbf{D}$ , to have elements

$$D_{\lambda\sigma} = C_{\lambda i} C_{\sigma i}, \quad (2.49)$$

the Fock matrix can be represented as

$$F_{\mu\nu} = H_{\mu\nu}^{\text{core}} + 2(\mu\nu|\lambda\sigma) D_{\lambda\sigma} - (\mu\lambda|\nu\sigma) D_{\lambda\sigma} \quad (2.50)$$

$$= H_{\mu\nu}^{\text{core}} + 2J[D_{\lambda\sigma}]_{\mu\nu} - K[D_{\lambda\sigma}]_{\mu\nu}, \quad (2.51)$$

where

$$J [D_{\lambda\sigma}]_{\mu\nu} = (\mu\nu|\lambda\sigma) D_{\lambda\sigma} \text{ and} \quad (2.52)$$

$$K [D_{\lambda\sigma}]_{\mu\nu} = (\mu\lambda|\nu\sigma) D_{\lambda\sigma} \quad (2.53)$$

are the elements of the Coulomb and exchange matrices, denoted  $\mathbf{J}$  and  $\mathbf{K}$ , respectively, and whereby the Fock matrix will be given by

$$\mathbf{F} = \mathbf{H}^{\text{core}} + 2\mathbf{J}[\mathbf{D}] - \mathbf{K}[\mathbf{D}]. \quad (2.54)$$

Similarly, the electronic RHF energy can be represented in the AO basis as

$$E_{\text{elec}}^{\text{RHF}} = (F_{\mu\nu} + H_{\mu\nu}^{\text{core}}) D_{\mu\nu}, \quad (2.55)$$

which when added to the nuclear repulsion energy under the Born–Oppenheimer approximation,  $E_{\text{nuc}}^{\text{BO}}$ , yields the total RHF energy:

$$E_{\text{tot}}^{\text{RHF}} = E_{\text{elec}}^{\text{RHF}} + E_{\text{nuc}}^{\text{BO}}. \quad (2.56)$$

Since the Fock matrix  $\mathbf{F}$  itself depends on the orbital coefficient matrix  $\mathbf{C}$  through the Coulomb and exchange matrices, the Roothaan equations must be solved iteratively. This process, known as the *self-consistent field* (SCF) procedure, begins by first building a guess for the Fock matrix before solving the Roothaan equation for that guess to obtain the orbital coefficients, and finally computing the total RHF energy. If, between two iterations, the change in the RHF energy is smaller than a particular tolerance (specified by the user before the start of the SCF procedure), then the procedure is said to be *converged*. The most expensive step in a given SCF iteration is the formation of the Fock matrix  $\mathbf{F}$ , or more specifically, the formation of the Coulomb and exchange matrices. This is because



the contraction of the density matrix  $D_{\lambda\sigma}$  with the two-electron integrals  $(\mu\nu|\lambda\sigma)$  (to form  $\mathbf{J}$ ) or  $(\mu\lambda|\nu\sigma)$  (to form  $\mathbf{K}$ ) scales algorithmically as  $\mathcal{O}(N^4)$ , where  $N$  is the total number of atomic orbitals. For further discussion of the SCF procedure, as well as for details of the implementation of RHF in the popular Python programming language, please refer to Tutorial 3a\* of the PSI4NUMPY Project.<sup>6</sup>

### *Reducing the Computational Expense for SCF via Density Fitting*

As we have seen above, the computational bottleneck for the SCF procedure is in the construction of the Coulomb and exchange matrices, scaling as  $\mathcal{O}(N^4)$  for both  $\mathbf{J}$  and  $\mathbf{K}$ . This means that, even for a simple twofold increase in the number of atomic orbitals (or atoms if the same basis set is used), a  $2^4 = 16$ -fold increase in computational expense is incurred. In order to combat this, the *density fitting* approach may be employed, whereby the four-index electron repulsion integrals over spatial orbitals,

$$(\mu\nu|\lambda\sigma) = \int d\mathbf{r}_1^3 d\mathbf{r}_2^3 \phi_\mu(\mathbf{r}_1)\phi_\nu(\mathbf{r}_1)\frac{1}{r_{12}}\phi_\lambda(\mathbf{r}_2)\phi_\sigma(\mathbf{r}_2),$$

is written instead as a product of two three-index integrals and a two-index quantity

$$(\mu\nu|\lambda\sigma) \approx (\mu\nu|P) [\mathbf{J}^{-1}]_{PQ} (Q|\lambda\sigma), \quad (2.57)$$

defined as

$$(Q|\lambda\sigma) = \int d\mathbf{r}_1^3 d\mathbf{r}_2^3 \chi_Q(\mathbf{r}_1)\frac{1}{r_{12}}\phi_\lambda(\mathbf{r}_2)\phi_\sigma(\mathbf{r}_2) \quad (2.58)$$

$$\mathbf{J}_{PQ} = \int d\mathbf{r}_1^3 d\mathbf{r}_2^3 \chi_P(\mathbf{r}_1)\frac{1}{r_{12}}\chi_Q(\mathbf{r}_2) \quad (2.59)$$

---

\*[https://github.com/psi4/psi4numpy/blob/master/Tutorials/03\\_Hartree-Fock/3a\\_restricted-hartree-fock.ipynb](https://github.com/psi4/psi4numpy/blob/master/Tutorials/03_Hartree-Fock/3a_restricted-hartree-fock.ipynb)

where  $\chi_P$  and  $\chi_Q$  are *auxiliary basis functions*. Commonly, the Coulomb metric  $[\mathbf{J}^{-1}]_{PQ}$  is split into the product  $[\mathbf{J}^{-1}]_{PQ} = [\mathbf{J}^{-\frac{1}{2}}]_{PQ} [\mathbf{J}^{-\frac{1}{2}}]_{PQ}$  and folded into the so-called “dressed” three index integrals  $\widetilde{(\mu\nu|P)}$  and  $\widetilde{(Q|\lambda\sigma)}$ :

$$(\mu\nu|\lambda\sigma) \approx (\mu\nu|P) [\mathbf{J}^{-1}]_{PQ} (Q|\lambda\sigma) \quad (2.60)$$

$$= (\mu\nu|P) [\mathbf{J}^{-\frac{1}{2}}]_{PQ} [\mathbf{J}^{-\frac{1}{2}}]_{PQ} (Q|\lambda\sigma) \quad (2.61)$$

$$= \widetilde{(\mu\nu|P)} \widetilde{(Q|\lambda\sigma)}. \quad (2.62)$$

When constructing the electron repulsion integrals in this manner, the scaling of constructing the Coulomb matrix  $\mathbf{J}$  is actually reduced from  $\mathcal{O}(N^4) \rightarrow \mathcal{O}(N^3)$ ; unfortunately, no such luck is found for the exchange matrix, whose scaling is still  $\mathcal{O}(N^4)$ . Aside from reducing the expense of constructing  $\mathbf{J}$ , however, the other major benefit of the density fitting scheme is the reduction in the amount of space required to store the electron repulsion integrals, either on disk or in memory. For the exact ERIs, even when exploiting their eight-fold permutational symmetry, they are still rank-4 tensors and as such require a non-trivial amount of storage space, and are in fact the storage bottleneck for any Hartree–Fock computation. When using density fitted integrals, on the other hand, not only is the storage requirement reduced by a full order of magnitude, but only one set of integrals [i.e., either  $(\mu\nu|P)$  or  $(Q|\lambda\sigma)$ ] must be computed and stored, since they are transposes of one another. For further discussion on the details and implementing of a density-fitted RHF (DF-RHF) code, please refer to the density fitting tutorial \* in the PSI4NUMPY Project.<sup>6</sup>

---

\*[https://github.com/psi4/psi4numpy/blob/master/Tutorials/03\\_Hartree-Fock/density-fitting.ipynb](https://github.com/psi4/psi4numpy/blob/master/Tutorials/03_Hartree-Fock/density-fitting.ipynb)

## 2.7 Capturing Dynamical Electron Correlation with post-Hartree–Fock Methodologies

We have developed the framework of restricted Hartree–Fock molecular orbital theory, in which a Slater determinant wavefunction comprised of two-electron spatial orbitals is optimized by variationally minimizing the electronic energy with respect to the choice of the orbitals. Hartree–Fock theory rests on the assumptions that the electronic wavefunction is well represented by a single Slater determinant, and that the average potential operator  $v^{\text{HF}}(i)$  sufficiently approximates the exact pairwise electron–electron repulsion. The first of these, equivalent to assuming that the exact electronic wavefunction is comprised of only one electron configuration (i.e., a single, unique set of occupied spin orbitals), is a good approximation for chemical phenomena where near-degeneracies do not occur, e.g., for diradicals, open-shell singlet excited states, conical intersections, etc. Since this Thesis is not concerned with examining those phenomena, we will proceed without concern for the possibility of violating this assumption.

The second assumption, however, is clearly lacking; after all, in Hartree–Fock theory, only the motions of parallel-spin electrons are correlated. Even though Hartree–Fock can typically recover upwards of 98% of the electronic energy for most molecules, it is in this last 2% where much of the physics that governs chemical behavior lies. Since it is such an important piece of the total, we will define the *correlation energy*,  $E_{\text{corr}}$ , to be the energy not recovered by Hartree–Fock, assuming a complete basis set:

$$E_{\text{corr}} = \mathcal{E}_{\text{exact}} - E_{\text{HF}}^{\infty}, \quad (2.63)$$

where  $\mathcal{E}_{\text{exact}}$  is the exact, non-relativistic electronic energy under the Born–Oppenheimer approximation. Usually this quantity is unknown, so it makes more sense to refer to the

correlation energy in a given basis set with

$$E_{\text{corr}}^{\text{basis}} = \mathcal{E}_{\text{exact}}^{\text{basis}} - E_{\text{HF}}^{\text{basis}}. \quad (2.64)$$

Since we have already justified that the error incurred by assuming a single Slater determinant should be negligible for our purposes here, the correlation energy must be due to the second major approximation made by Hartree–Fock: the average potential operator which fails to capture the instantaneous electron–electron repulsion between electrons of antiparallel spins. This is generally referred to as *dynamical correlation*, because it arises from electron dynamics.

If Hartree–Fock fails to capture dynamical correlation, which is responsible for governing much of chemical behavior, what is there to do? Of course, we could simply not do Hartree–Fock — this is the approach taken by Density Functional Theory, which will be discussed in Section. 2.9 — but it seems wasteful to have gone through all of the derivation above to just abandon the theory now. Instead, in this section we will introduce several approaches whose goal is to recover dynamical electron correlation by starting with a Hartree–Fock description of the system of interest, before correcting it by some additional means. These *post-Hartree–Fock* methods are widely utilized, and are capable of recovering the correlation energy with astounding accuracy, given enough computational power.

For a molecule with  $N$  electrons, Hartree–Fock theory yields a collection of  $N$  coupled one-electron equations which each produce optimal one-electron orbitals that are combined to form a properly antisymmetrized  $N$ -electron wavefunction, the Slater determinant. This set of coupled equations can be solved one at a time (and the process iteratively converged) thanks to the fact that the Fock operator  $\hat{f}(i)$  is effectively a one-electron operator: each of the  $\hat{h}$  and  $v^{\text{HF}}(i)$  only operate on one electron at a time, even though  $v^{\text{HF}}(i)$  is parametrically coupled to the other electrons. For the exact electronic Schrödinger equation, how-

ever, the each electron is “connected” to every other electron via Coulomb’s law through the repulsion operator  $\sum_{i<j} \frac{1}{r_{ij}}$ , which is an  $N$ -electron operator. Therefore, the exact electronic wavefunction must be a proper  $N$ -electron wavefunction, instead of a simple collection of one-electron orbitals. This electronic “connectedness” is precisely what we must seek to recover when accounting for dynamical electron correlation, and what is missed by Hartree–Fock theory.

### 2.7.1 Configuration Interaction & The Exact Electronic Wavefunction

To construct a suitable  $N$ -electron wavefunction which can account for dynamical electron correlation, we can take inspiration from an unlikely place: Hartree–Fock theory, or more specifically, the Roothaan equations. Roothaan introduced a basis set of one-electron functions in which to expand the one-electron orbitals which comprises the  $N$ -electron Slater determinant, which transformed the set of  $N$  coupled Fock equations into a single linear algebra equation. In the same way, let us expand the exact  $N$ -electron wavefunction in an  $N$  electron basis of Slater determinants:

$$\Phi_{\text{exact}}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N) = \sum_i c_i |\Psi_i\rangle. \quad (2.65)$$

In the limit of an infinite expansion, this expression would be exact; unfortunately, this is not possible with finite computational power. We could, however, truncate this expansion to be over a finite set of Slater determinants to approximate the exact wavefunction. This, then, will be our first strategy for recovering dynamical electron correlation: choose a finite set of Slater determinants and optimize the expansion coefficients  $c_i$  in order to obtain an approximation to  $\Phi_{\text{exact}}$  which is correlated. But how can we construct a set of Slater determinants which are suitable for the system of interest?

Since the Hartree–Fock procedure yields a Slater determinant wavefunction  $|\Psi_0\rangle$  with a set of optimal occupied and virtual (unoccupied) orbitals, we can use this as a “refer-

ence” from which to build the other  $N$ -electron basis functions. These new, “substituted” determinants — which are comprised of different electronic configurations of occupied orbitals, together with the reference determinant — form the  $N$ -electron basis within which we will approximately describe the exact electronic wavefunction, according to the expansion above. In other words, it is thanks to the interaction of these different configurations of occupied orbitals that dynamical electron correlation is recaptured, even when starting from an uncorrelated reference. We will therefore refer to this approach as *configuration interaction* (CI), and our approach will be to apply the Variational Principle to solve for the expansion coefficients  $c_i$ . At this point, it is convenient to introduce notation to describe the substitution of occupied orbitals from the reference determinant with virtual ones. Letting  $|\Psi_0\rangle$  represent the reference determinant, we will allow  $|\Psi_{ij\dots k}^{ab\dots c}\rangle$  to denote a determinant where virtual orbitals  $\chi_a, \chi_b, \dots, \chi_c$  are substituted for occupied orbitals  $\chi_i, \chi_j, \dots, \chi_k$ . Then, the expansion of the exact wavefunction above can be rewritten as

$$\Phi_{\text{exact}} = \underbrace{c_0|\Psi_0\rangle}_{\text{Reference}} + \underbrace{\sum c_i^a|\Psi_i^a\rangle}_{\text{Singles}} + \underbrace{\sum c_{ij}^{ab}|\Psi_{ij}^{ab}\rangle}_{\text{Doubles}} + \underbrace{\sum c_{ijk}^{abc}|\Psi_{ijk}^{abc}\rangle}_{\text{Triples}} + \dots, \quad (2.66)$$

where we have grouped the substituted determinants in terms of how many orbitals are substituted.

### *Full Configuration Interaction (FCI)*

By examining the Eqn. 2.66 above, it is clear that the only manner in which to construct an infinite expansion of substituted determinants from a single reference determinant is if the reference determinant is itself constructed from an infinite one-electron basis set. As mentioned above, this is impossible to achieve. We can, however, accept this fact and construct the reference determinant  $|\Psi_0\rangle$  within a finite one-electron basis. This will reduce

the expansion above from containing an infinite number of substituted determinants to

$$|\Psi_{\text{FCI}}\rangle = \underbrace{c_0|\Psi_0\rangle}_{\text{Reference}} + \underbrace{\sum c_i^a|\Psi_i^a\rangle}_{\text{Singles}} + \underbrace{\sum c_{ij}^{ab}|\Psi_{ij}^{ab}\rangle}_{\text{Doubles}} + \underbrace{\sum c_{ijk}^{abc}|\Psi_{ijk}^{abc}\rangle}_{\text{Triples}} + \dots + \underbrace{\sum c_{ijk\dots l}^{abc\dots d}|\Psi_{ijk\dots l}^{abc\dots d}\rangle}_{N\text{-tuples}}, \quad (2.67)$$

where we have included each possible permutation of orbital occupations from the reference determinant in our expansion; this approach is referred to as *full configuration interaction* (FCI), and represents the exact solution to the electronic Schrödinger equation within a given one-electron basis set. Even though FCI is a finite expansion, the number of determinants in the FCI wavefunction grows factorially with the number of orbitals according to

$$N_{\text{det}} = \binom{n}{N_\alpha} \binom{n}{N_\beta}, \quad (2.68)$$

where each term is a binomial coefficient and reads, e.g., “ $n$  choose  $N_\alpha$ ”,  $n$  is the number of one-electron basis functions, and  $N_\alpha, N_\beta$  are the number of  $\alpha$  and  $\beta$  electrons, respectively.\* This may not seem like that many, but even for a small molecule like methane in a minimal basis (9 basis functions and 10 electrons with 5 each of  $\alpha$  and  $\beta$  spin), this would be 15,876 determinants!

Now that we have what appears to be a reasonable representation for the exact wavefunction in Eqn. 2.67, we must once again concern ourselves with the mechanism by which to obtain the expansion coefficients. Since we wish for our CI wavefunction to variationally approximate the exact electronic wavefunction (and electronic energy), we can begin by substituting the CI expansion given in Eqn. 2.66 into the full electronic Schrödinger equation:

$$\widehat{\mathcal{H}}_{\text{elec}}|\Psi_{\text{FCI}}\rangle = E_{\text{FCI}}|\Psi_{\text{FCI}}\rangle \quad (2.69)$$

$$\langle\Psi_{\text{FCI}}|\widehat{\mathcal{H}}_{\text{elec}}|\Psi_{\text{FCI}}\rangle = E_{\text{FCI}}, \quad (2.70)$$

---

\*These binomial coefficients are so named thanks to their presence in the binomial power series expansion. Even so, an arbitrary binomial coefficient  $\binom{N}{R}$  is identical to the “combination”  ${}_N C_R$ .

since the CI coefficients which minimize the total energy are the same as the eigenvectors of the electronic Hamiltonian in the basis of Slater determinants, and where we assume  $|\Psi_{\text{FCI}}\rangle$  has been normalized. We can therefore once again cast the problem of solving our eigenvalue equation in a given basis set into a linear algebra expression. Letting  $|S\rangle$ ,  $|D\rangle$ ,  $|T\rangle$ , etc. represent Slater determinants with a single, double, and triple orbital substitutions and  $|0\rangle$  refer to our reference determinant, we can form the matrix representation of the electronic Hamiltonian in our  $N$ -electron basis as

$$\underline{\underline{\mathbf{H}}} = \begin{pmatrix} \langle 0|\hat{H}|0\rangle & 0 & \langle 0|\hat{H}|D\rangle & 0 & \dots \\ 0 & \langle S|\hat{H}|S\rangle & \langle S|\hat{H}|D\rangle & \langle S|\hat{H}|T\rangle & \dots \\ \langle D|\hat{H}|0\rangle & \langle D|\hat{H}|S\rangle & \langle D|\hat{H}|D\rangle & \langle D|\hat{H}|T\rangle & \dots \\ 0 & \langle T|\hat{H}|S\rangle & \langle T|\hat{H}|D\rangle & \langle T|\hat{H}|T\rangle & \dots \\ 0 & 0 & \langle Q|\hat{H}|D\rangle & \langle Q|\hat{H}|T\rangle & \dots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{pmatrix}, \quad (2.71)$$

where the matrix elements  $H_{IJ} = \langle \Psi_I | \hat{\mathcal{H}}_{\text{elec}} | \Psi_J \rangle$  can be evaluated quite easily according to *Slater's Rules*, which are given in Table 4.1 of Ref. 4. Once  $\underline{\underline{\mathbf{H}}}$  has been constructed, a simple matrix diagonalization will yield the eigenvectors (CI coefficients) and eigenvalues for the FCI Hamiltonian, with the leading eigenvalue being the FCI electronic energy. Even though the formulation and solution of the FCI equation seems straightforward, in practice it is anything but simple to evaluate. This is because the size of the Hamiltonian matrix grows factorially with the number of electrons and the number of one-electron basis functions, which makes the diagonalization of the Hamiltonian matrix factorially expensive to perform. This becomes routinely prohibitive for systems larger than approximately 18 electrons in 18 orbitals; indeed, the largest CI computation ever performed was for the 22  $\pi$  electrons in 22 orbitals of the pentacene molecule.<sup>7</sup> So, if FCI cannot be routinely applied to systems larger than these, how can we account for dynamical correlation?



### *Truncated Configuration Interaction*

In order to make CI extensible for routine application to chemical systems of even modest size, the standard approach is to truncate the FCI expansion given in Eqn. 2.67 at a desired excitation level, i.e., by including singly- and doubly-substituted determinants but not those with triple substitutions and above, etc. These truncations of the CI expansion are abbreviated by appending the letters corresponding to the excitation level to “CI”; e.g., “CISD” for the aforementioned CI truncation including single and double substitutions. In addition to truncating the CI space, a multitude of technical and algorithmic advancements have been developed which have made CI at all levels more extensible. These truncated CI expansions offer a considerable computational advantage over FCI, where instead of scaling factorially with the number of orbitals, they scale polynomially. For example, while an optimal FCI algorithm scales with number of determinants  $N_{\text{det}}$  and number of orbitals  $N$  as  $\mathcal{O}(N_{\text{det}}N^4)$ , CISD scales as  $\mathcal{O}(N^6)$  and CISDT as  $\mathcal{O}(N^8)$ . While this may not seem like an appreciable decrease in complexity, compared against a computation with 10 electrons in 10 orbitals, a computation with 20 electrons in 20 orbitals (a simple factor of 2 increase) should require a factor of  $2^6 = 64$ ,  $5^8 = 256$ , and  $(N_{\text{det}}N^4) = (34, 134, 779, 536 \cdot 20^4) = 5.5 \times 10^{15}$  times longer to complete for CISD, CISDT, and FCI respectively. Given the fact that truncated CI seems to enjoy such a significant advantage over FCI in terms of its drastically reduced computational expense, why bother performing FCI computations at all?

To answer this question, let us imagine that some poor graduate student somewhere has been asked by their PI to perform the aforementioned experiment:\* to compute the total energy for a series of H-atom chains of different length, using CISD, CISDT, CISDTQ, and FCI, and plot the trends in total energy as a function of the number of electrons. Cursing their luck, the student decides that this task is too tedious to be worth their time. Intuitively, the student expects that they will find that the energies computed for two H atoms will simply be two times the energy of a single H atom, and that similarly the energy of a

---

\*R.I.P.

chain of three, four, and five H atoms will just be three, four, and five times the energy of a single H atom. Annoyed at their advisor, they devise a clever shortcut: compute the energy of a single H atom with each of CISD, CISDT, CISDTQ, and FCI, and then simply plot a linear trend of these values by the number of H atoms in the chain to represent the chains' energies. Feeling triumphant at their ability to outfox their advisor, and having already completed the task they had been assigned, the student decides to take the rest of the week off before presenting their results to their PI during their next meeting. At their next meeting, however, the PI was not pleased, having accused the student of cheating the assignment. In fact, not only did the PI know that the student had cheated their assignment, but they seemingly knew exactly *how* the student had cheated. In a panic, the student admitted their falsehood before asking, how the PI possibly could have known? "You fool!" they replied, "You fell victim to one of the classic blunders — the most famous of which is 'Never get involved in a land war in Asia,' but only slightly less well-known is this: Truncated CI is not size extensive!"

So, what exactly went wrong for the student in our thought experiment? Their error, it turns out, was not in their physical intuition for the *truth* — but rather in their understanding of the approximations involved in truncated CI. Indeed, it is true that the exact electronic energy of a system *does* scale linearly with the number of particles. When a theoretical approach satisfies this criteria, it is said to be *size extensive*; FCI is an example of such a method. Unfortunately for our student, however, the truncation of the  $N$ -electron basis (upon which all truncated CI expansions are based) prevents the resulting method from being size extensive. This is because the total wavefunction no longer contains all possible substitutions of the reference determinant; the CISD energy for a system with four electrons, for example, will be missing the contributions from triple and quadruple excitations, making it fall short of reproducing the exact (FCI) energy. Aside from size extensivity, a property of the exact electronic energy is that it is *size consistent*. Size consistency means that as two particles (or, more generally, any two chemical species) become infinitely sep-

arated, the total energy of the pair converges to the simple sum of the isolated particles. As we will see in the next Chapter, the property of size consistency is of central importance to the study of non-covalent interactions; we will therefore abandon truncations of the CI expansion to focus our attention on other post-Hartree–Fock electron correlation methods which do satisfy size extensivity and size consistency.

### 2.7.2 Møller–Plesset Perturbation Theory

As discussed above, especially for well-behaved molecules at (or suitably close to) their equilibrium geometries, the correlation energy accounts for only about 2% of the total molecular energy. Therefore, instead of using a variational approach to construct the exact electronic wavefunction, we could equivalently apply Rayleigh–Schrödinger Perturbation Theory (RSPT) to the problem of dynamical electron correlation. This has been among the standard strategies in the field of quantum chemistry since the 1930s, when it was originally introduced by Møller and Plesset.<sup>8</sup> This approach is therefore most commonly referred to as Møller–Plesset Perturbation Theory (MPPT), but is occasionally referred to as many-body perturbation theory (MBPT).

#### *Formulation of MPPT*

In order to develop MPPT, let us first recall the basic takeaways from its parent, RSPT. In RSPT, the exact Hamiltonian is assumed to be separable into a piece which can be solved exactly (the zeroth-order Hamiltonian,  $\hat{H}^{(0)}$ ) and a piece which cannot (the perturbation,  $\hat{V}$ ),

$$\hat{\mathcal{H}} = \hat{H}^{(0)} + \lambda \hat{V}.$$

Subsequently expanding the exact wavefunction as a Taylor series in the perturbation strength,  $\lambda$ , yields a series of shifts for the zeroth-order energy ( $\mathcal{E}_n^{(0)}$ ) and zeroth-order wavefunction ( $|\Psi_n^{(0)}\rangle$ ) which aim to correct these quantities towards their exact values. It is our first task, therefore, to define the zeroth-order Hamiltonian  $\hat{H}^{(0)}$  and the perturbation  $\hat{V}$ . Since in

MPPT we are concerned with applying RSPT to the problem of dynamical electron correlation on top of an existing Hartree–Fock description of a molecular system, it would be natural to define the zeroth-order Hamiltonian to be

$$\hat{H}_{\text{MPPT}}^{(0)} = \sum_i \hat{f}(i) = \sum_i \hat{h}(i) + v^{\text{HF}}(i), \quad (2.72)$$

which means that the perturbation must be the difference between the Hartree–Fock potential operator and the exact electron–electron repulsion:

$$\hat{V} = \hat{\mathcal{H}}_{\text{elec}} - \hat{H}^{(0)} = \sum_{i<j} \frac{1}{r_{ij}} - \sum_i v^{\text{HF}}(i), \quad (2.73)$$

where here we have used indices  $i, j$  to refer to electrons.

Now that we have defined these quantities, we may proceed according to the standard RSPT approach to evaluate the first few terms in the energy series. Before we do so, however, we must first briefly mention that hereafter, indices  $i, j, k, l$  refer to occupied molecular orbitals,  $a, b, c, d$  refer to virtual molecular orbitals,  $m, n$  refer to the wavefunction states themselves, and  $\mu, \nu, \lambda, \sigma$  refer to occupied atomic orbitals (basis functions), according to standard practice in the field. First, we solve the zeroth-order problem,

$$\hat{H}^{(0)} | \Psi_n^{(0)} \rangle = \mathcal{E}_n^{(0)} | \Psi_n^{(0)} \rangle, \quad (2.74)$$

where  $\mathcal{E}_n^{(0)} = \sum_i \epsilon_i$  (which is *not* the Hartree–Fock energy!) and  $|\Psi_n^{(0)}\rangle = \Psi_0^{\text{HF}}$ , the

ground-state Hartree–Fock wavefunction. The first order energy correction is given by

$$\begin{aligned}
\mathcal{E}_n^{(1)} &= \langle \Psi_n^{(0)} | \hat{V} | \Psi_n^{(0)} \rangle \\
&= \langle \Psi_0^{\text{HF}} | \sum_{i<j} \frac{1}{r_{ij}} - \sum_i v^{\text{HF}}(i) | \Psi_0^{\text{HF}} \rangle \\
&= \langle \Psi_0^{\text{HF}} | \sum_{i<j} \frac{1}{r_{ij}} | \Psi_0^{\text{HF}} \rangle - \langle \Psi_0^{\text{HF}} | \sum_i v^{\text{HF}}(i) | \Psi_0^{\text{HF}} \rangle \\
&= \frac{1}{2} \sum_{ij} \langle ij || ij \rangle - \sum_i \langle i | v^{\text{HF}}(i) | i \rangle \\
&= \frac{1}{2} \sum_{ij} \langle ij || ij \rangle - \sum_{ij} \langle ij || ij \rangle \\
&= -\frac{1}{2} \sum_{ij} \langle ij || ij \rangle
\end{aligned} \tag{2.75}$$

Hence the total energy *through* first order is

$$\mathcal{E}_n^{(0)} + \mathcal{E}_n^{(1)} = \sum_i \epsilon_i - \frac{1}{2} \sum_{ij} \langle ij || ij \rangle, \tag{2.76}$$

which *is* the Hartree–Fock energy; therefore, the first correction to the Hartree–Fock energy which accounts for dynamical electron correlation is at second order.

The second-order energy correction,  $\mathcal{E}_n^{(2)}$ , can be found by substituting the expression for the first-order correction to the wavefunction,  $|\Psi_n^{(1)}\rangle$  (as given in Eqn. 2.14) into

$$\mathcal{E}_n^{(2)} = \langle \Psi_n^{(0)} | \hat{V} | \Psi_n^{(1)} \rangle \tag{2.77}$$

$$\therefore \mathcal{E}_n^{(2)} = \sum_{n \neq m} \frac{|\langle \Psi_n^{\text{HF}} | \hat{V} | \Psi_m^{\text{HF}} \rangle|^2}{\mathcal{E}_n^{(0)} - \mathcal{E}_m^{(0)}}. \tag{2.78}$$

To compute the second-order energy correction for the ground state wavefunction, therefore, the sum in Eqn. 2.78 must run over all states *other* than the ground state, i.e., over all excited states of the Hartree–Fock reference determinant. From the FCI expansion, we already know how large this sum can become if every other excited state is considered; let us

consider, therefore, if any terms should necessarily be zero and therefore not be computed. First, since according to Brillouin's theorem no singles determinants  $|S\rangle$  will interact directly with the Hartree–Fock reference  $|0\rangle$ , all  $\langle 0|\widehat{V}|S\rangle = 0$  and hence the singly excited states will not contribute to  $\mathcal{E}_n^{(2)}$ . Furthermore, since  $\widehat{V}$  is a two-electron operator, excitation levels greater than or equal to triples will also not contribute to  $\mathcal{E}_n^{(2)}$ . Hence the only excited states which contribute to  $\mathcal{E}_n^{(2)}$  are *doubly* excited configurations, which for the excitation of an electron from occupied orbital  $i$  ( $j$ ) into virtual orbital  $a$  ( $b$ ) we will denote  $|\Psi_{ij}^{ab}\rangle$ , with energy difference (from the ground state)  $\epsilon_a - \epsilon_i + \epsilon_b - \epsilon_j$ . Substituting these facts into Eqn. 2.78 yields the second-order, ground state energy correction  $\mathcal{E}_0^{(2)}$ :

$$\mathcal{E}_0^{(2)} = \sum_{i<j} \sum_{a<b} \frac{\langle \Psi_0^{\text{HF}} | \widehat{V} | \Psi_{ij}^{ab} \rangle}{\epsilon_i - \epsilon_a + \epsilon_j - \epsilon_b}, \quad (2.79)$$

which, according to Slater's rules, yields

$$\mathcal{E}_0^{(2)} = \sum_{i<j} \sum_{a<b} \frac{|[ia||jb]|^2}{\epsilon_i - \epsilon_a + \epsilon_j - \epsilon_b}, \quad (2.80)$$

which is formulated in terms of molecular spin orbitals. If an RHF reference is used, the spin variable  $\omega$  can once again be integrated away to yield two equations over spatial orbitals,

$$E_{0,\text{SS}}^{(2)} = \sum_{i<j} \sum_{a<b} \frac{(ia|jb) [(ia|jb) - (ib|ja)]}{\epsilon_i - \epsilon_a + \epsilon_j - \epsilon_b} \quad (2.81)$$

$$E_{0,\text{OS}}^{(2)} = \sum_{i<j} \sum_{a<b} \frac{(ia|jb) (ia|jb)}{\epsilon_i - \epsilon_a + \epsilon_j - \epsilon_b}, \quad (2.82)$$

where the round brackets in the two-electron integrals denotes that they are over spatial orbitals, and the subscripts SS, OS denote the contributions from electron pairs with the same spin (SS) or opposite spin (OS). It is worth mentioning that there is an extra two-electron exchange integral in the numerator of the SS energy expression, thanks to the

fact that electrons with parallel spin experience exchange effects due to the antisymmetry principle. Finally, then, the second-order MPPT correction to the ground state RHF energy is given by

$$\mathcal{E}_0^{(2)} = E_{0,SS}^{(2)} + E_{0,OS}^{(2)} \quad (2.83)$$

In general, a particular truncation level of MPPT is referred to as  $MPn$ , where  $n$  is a number denoting the order of the perturbation. The expression given above in Eqn. 2.83 is therefore referred to as the *MP2* energy correction, symbolically denoted  $E_{\text{corr}}^{\text{MP2}}$ , and the MP2 total ground-state energy is given by

$$E_{\text{tot}}^{\text{MP2}} = \mathcal{E}_0^{(0)} + \mathcal{E}_0^{(1)} + \mathcal{E}_0^{(2)} = E_0^{\text{HF}} + E_{\text{corr}}^{\text{MP2}} \quad (2.84)$$

#### *Computing the MP2 Energy Correction for an RHF Reference*

As noted above, the purpose of  $MPn$  is to recapture dynamical electron correlation missed by Hartree–Fock via an  $n$ th order perturbation expansion. Above, we formulated the corresponding second-order correction, MP2, for the ground state configuration of an RHF reference wavefunction. From Eqns. 2.81 and 2.82, it appears that everything necessary to compute the MP2 energy correction (and subsequently,  $E_{\text{tot}}^{\text{MP2}}$  itself) is already available from the completed SCF procedure for the RHF reference. If this were true, MP2 would be a *free* addition to RHF — so, why would anyone *not* do MP2? As it turns out, unfortunately, even though all of the quantities needed to compute the MP2 energy (two-electron integrals, orbital energies, HF energy, and nuclear repulsion energy) are constructed during the SCF procedure, not all of them are in a format which is ready to be immediately used. In particular, the two-electron integrals required for Eqns. 2.81 and 2.82 are represented in the molecular orbital basis, which should be clear from the use of orbital labels  $i, \dots, l$  in those expressions. During the SCF procedure, however, these integrals are generated and used in the *atomic orbital* basis, denoted by the use of Greek orbital labels  $\mu, \dots, \sigma$  in, e.g., Eqn. 2.51.

The fact that the two-electron integrals are available from the SCF procedure in the AO basis is not a significant problem, since the SCF procedure also offers the means by which to move from the AO basis to the MO basis: the orbital coefficients,  $C_{\lambda_i}$ . Transforming the integrals into the MO basis requires transforming each AO index into a corresponding MO index, which can be accomplished by contracting

$$(ia|jb) = C_{\mu i} C_{\nu a} (\mu\nu|\lambda\sigma) C_{\lambda j} C_{\sigma b}. \quad (2.85)$$

As written, however, this operation would scale as  $\mathcal{O}(N^8)$ , since there are 8 unique indices involved in the contraction. Luckily, this can be refactored from a single  $\mathcal{O}(N^8)$  step into four  $\mathcal{O}(N^5)$  steps by contracting over each AO index separately:

$$(ia|jb) = [C_{\mu i} [C_{\nu a} [C_{\lambda j} [C_{\sigma b} (\mu\nu|\lambda\sigma)]]]]. \quad (2.86)$$

Once the integrals have been transformed into the MO basis, evaluating Eqns. 2.81 and 2.82 actually only scales as  $\mathcal{O}(N^4)$ , since there are only four unique indices in these expressions which must be looped over. Therefore, the overall cost of MP2 is  $\mathcal{O}(N^5)$ , with the bottleneck arising from the AO→MO transformation of the two-electron integrals.\* For further details of the procedure for computing conventional and density-fitted MP2, as well as for details of their implementation starting from an RHF (or DF-RHF) reference in the popular Python programming language, please refer to Tutorials 5a and 5b<sup>†</sup> of the PSI4NUMPY Project.<sup>6</sup>

---

\*The cost of transforming the integrals to the MO basis can be reduced if three-index density-fitted integrals are used to compute the RHF reference (i.e., DF-RHF) to only  $\mathcal{O}(N_{aux}N^3)$ , where  $N_{aux}$  is the number of auxiliary functions and  $N$  is the number of atomic orbitals. This does not reduce the overall algorithmic scaling of DF-MP2, however, as computing the DF-MP2 energy correction scales as  $\mathcal{O}(N^5)$ . The real win for DF-MP2 is that only one  $\mathcal{O}(N^5)$  step is required (to compute the correlation energy), as opposed to four for conventional MP2, where each quarter-transform of the four-index integrals are each  $\mathcal{O}(N^5)$  but the final evaluation of the correlation energy in Eqn. 2.83 is only  $\mathcal{O}(N^4)$ .

<sup>†</sup>[https://github.com/psi4/psi4numpy/blob/master/Tutorials/05\\_Moller-Plesset](https://github.com/psi4/psi4numpy/blob/master/Tutorials/05_Moller-Plesset)



### 2.7.3 Coupled-Cluster Theory

Despite the utility of MPPT, particularly of its lowest-order truncation, MP2, the failure of the MP series to be convergent towards the exact electronic energy prevents it from being systematically improvable beyond convergence towards the complete one-particle basis set limit. Configuration interaction, on the other hand, is convergent — but for all but its exact formulation, FCI, it is neither size extensive nor size consistent, preventing it from being able to be reliably applied to compare the energies of chemical systems with different numbers of electrons. It would be desirable, therefore, to develop a method which combines the best of both of these post-Hartree–Fock methods: to be both convergent *and* size extensive. Fortunately, such a theory exists, known as coupled-cluster (CC) theory. Originally developed by the high-energy physics community as an approach to model nuclear structure, CC has become the *de facto* approach for quantum chemists interested in highly accurate energies, structures, and properties for molecular systems, and will be utilized heavily throughout this Thesis as the reference against which more approximate methods will be evaluated.

In order to develop the basic framework of CC theory, let us begin by rewriting our expression for the FCI wavefunction:

$$\begin{aligned} |\Psi_{\text{FCI}}\rangle &= c_0|\Psi_0\rangle + \sum c_i^a|\Psi_i^a\rangle + \sum c_{ij}^{ab}|\Psi_{ij}^{ab}\rangle + \sum c_{ijk}^{abc}|\Psi_{ijk}^{abc}\rangle + \dots + \sum c_{ijk\dots l}^{abc\dots d}|\Psi_{ijk\dots l}^{abc\dots d}\rangle \\ &= (1 + \hat{T}_1 + \hat{T}_2 + \hat{T}_3 + \dots)|\Psi_0\rangle \\ &= \hat{T}|\Psi_0\rangle, \end{aligned}$$

where  $\hat{T}$  is the total *excitation operator*, and each of the  $\hat{T}_1$ ,  $\hat{T}_2$ ,  $\hat{T}_3$ , etc. produce the singly, doubly, and triply etc. excited determinants in the expansion. Since in the FCI expansion all possible substitutions of the reference determinant  $|\Psi_0\rangle$  are included, it is equivalent to

write the series above with the *cluster operator*,  $e^{\hat{T}}$ , instead:

$$|\Psi_{\text{FCC}}\rangle = e^{\hat{T}}|\Psi_0\rangle = \left(1 + \hat{T} + \frac{1}{2!}\hat{T}^2 + \frac{1}{3!}\hat{T}^3 + \dots\right)|\Psi_0\rangle, \quad (2.87)$$

where we have denoted the wavefunction generated by the cluster operator as the full coupled-cluster (FCC) wavefunction,  $|\Psi_{\text{FCC}}\rangle$ , even though it is formally equivalent to the FCI wavefunction, and used the Taylor series expansion for the exponential function to expand the cluster operator. Since FCI and FCC are equivalent, it may not be obvious why we have gone to the extra trouble to rewrite the wavefunction thus; to illustrate the advantages of the cluster operator, let us consider a singles-and-doubles truncation of the excitation operator such that  $\hat{T} = \hat{T}_1 + \hat{T}_2$ .

In the FCI series, this truncation would produce the CISD wavefunction, which we know to lack size extensivity. A similar “coupled cluster with single and double substitutions” (CCSD) wavefunction may be constructed by substituting this truncated excitation operator into the cluster operator:

$$\begin{aligned} |\Psi_{\text{CCSD}}\rangle &= e^{\hat{T}_1 + \hat{T}_2}|\Psi_0\rangle \\ &= \left(1 + \hat{T}_1 + \hat{T}_2 + \frac{1}{2}\hat{T}_1^2 + \hat{T}_1\hat{T}_2 + \frac{1}{2}\hat{T}_2^2 + \dots\right)|\Psi_0\rangle \\ &= |\Psi_0\rangle + \sum_i^a t_i^a |\Psi_i^a\rangle + \sum_{ij}^{ab} t_{ij}^{ab} |\Psi_{ij}^{ab}\rangle + \frac{1}{2} \sum_i^a t_i^a \sum_j^b t_j^b |\Psi_{ij}^{ab}\rangle \\ &\quad + \sum_i^a t_i^a \sum_{jk}^{bc} t_{jk}^{bc} |\Psi_{ijk}^{abc}\rangle + \frac{1}{2} \sum_{ij}^{ab} t_{ij}^{ab} \sum_{kl}^{cd} t_{kl}^{cd} |\Psi_{ijkl}^{abcd}\rangle + \dots, \end{aligned}$$

where  $t_i^a$  are the *singles amplitudes* and  $t_{ij}^{ab}$  are the *doubles amplitudes*, which determine the magnitude to which singly- and doubly substituted determinants contribute to the overall CCSD wavefunction. It is worth noting that even though the excitation operator has been truncated at the level of double substitutions, the exponential form of the cluster operator produces so-called “disconnected” triply, quadruply, etc. substituted determinants

through infinite order, which while failing to recover every possible substitution at a given excitation level (and thus the full energetic contribution of, e.g.,  $\hat{T}_3$ ) do recover parts of this higher-order electron correlation which would be neglected by CISD alone. Furthermore, it is precisely these disconnected substitutions which allow for truncated CC methods to remain size-extensive. Progressively more complete truncations of the CC series, e.g., CCSD, CCSDT, CCSDTQ, etc. are therefore not only convergent towards the exact electronic energy, but can generate reasonable comparisons between systems of different size and composition. Like the levels of the truncated CI expansion, progressively more complete truncations of the CC expansion also increase in computational expense, with CCSD scaling as  $\mathcal{O}(N^6)$ , CCSDT scaling as  $\mathcal{O}(N^8)$ , etc. Unlike for CI, however, the amplitude equations for CC are not variational; therefore, even though CC is convergent, a given level of truncation in the CC series is not guaranteed to produce an energy which is above the exact electronic energy. In fact, the energies computed at different truncation levels of the CC series tend to oscillate about the FCC limit as they converge, with CCSD above the FCC limit, CCSDT below, CCSDTQ above, and so on.

By substituting the expression above into the electronic Schrödinger equation, it is possible to derive both the correlation energy and amplitude equations for CCSD; as this derivation is quite lengthy, however, we will direct the interested reader elsewhere for it.<sup>9</sup> The resulting CCSD correlation energy is often sufficiently accurate for many molecular properties, and total energies, especially given its modest computational expense. For some properties, however, where a more robust treatment of electron correlation is required for accurate predictions to be made (e.g., non-covalent interactions), CCSD can fall sufficiently short of the exact correlation energy that its application is inappropriate. Unfortunately, however, the sizeable increase in computational expense from  $\mathcal{O}(N^6) \rightarrow \mathcal{O}(N^8)$  when including triple excitations in the cluster operator with CCSDT is often prohibitive for application to relevant chemical systems. To address this challenge, a number of *approximate* approaches for including some (but not all) triple substitutions have been proposed. The

most successful, denoted CCSD(T),<sup>10</sup> perturbatively includes the most important triples amplitudes at the cost of a single  $\mathcal{O}(N^7)$  step. While the order of magnitude increase in expense of CCSD(T) compared to CCSD does make CCSD(T) prohibitively expensive for some systems, CCSD(T) happens to benefit from fortuitous error cancellation which actually make it a more reliable estimate of the FCI limit than either CCSD or CCSDT;\* CCSD(T) has therefore been touted as the “gold standard” of quantum chemistry, and is the standard choice for producing high-quality energies and properties. For reference implementations of each of the CCSD and CCSD(T) methods, please refer to the Psi4NUMPY project.<sup>†</sup>

## 2.8 Practical Considerations for Correlated Computations

We have thus far developed several methodologies leveraging perturbative or variational approaches to provide a description of dynamical electron correlation on top of a Hartree–Fock reference wavefunction. These methods range in computational expense from  $\mathcal{O}(N^5)$  for MP2 to  $\mathcal{O}(N!)$  for FCI with respect to the number of orbitals  $N$ , and (at least for CI and CC theories) are systematically improvable towards recovering the exact electronic energy within a given basis set. What remains to be seen, however, is how these methods perform for describing chemical properties which, in turn, depend on an accurate description of electron correlation to be described properly. In this section, we will address this and other questions related to the practical aspects of correlated computations, to offer a general overview of when applying these methods are appropriate (or even necessary) and the best manner in which to do so.

---

\*The convergence of CC truncations towards the FCI limit is discussed more thoroughly in Section 2.8.1.

<sup>†</sup>[https://github.com/psi4/psi4numpy/tree/master/Coupled-Cluster/Spin\\_Orbitals/CCSD](https://github.com/psi4/psi4numpy/tree/master/Coupled-Cluster/Spin_Orbitals/CCSD)

### 2.8.1 Convergent Methodologies for Electron Correlation

Since it is the goal of any post-Hartree–Fock methodology to recover the correlation energy neglected by Hartree–Fock, and indeed both full coupled-cluster (FCC) and full configuration-interaction (FCI) can do so exactly in a given one-electron basis, it is desirable also for successively more complete truncations to these expansions to provide progressively more accurate approximations to the FCC and FCI result. Similarly, it would be desirable for successively higher perturbation levels in the Møller–Plesset series to behave in the same manner. Such a theory, where successively less approximate formulations recover the exact result, is said to be *convergent*. Since both truncated CI and truncated CC formulations satisfy such a requirement, they are examples of convergent theories;  $MP_n$ , on the other hand, is not guaranteed to be convergent, and has even famously been demonstrated to be *divergent* in some cases.\* As truncated CC and CI approaches are tractable for application to larger chemical systems than could ever be treated with FCC or FCI, it is necessary to determine the accuracy of successively higher truncations of CC and CI for a given property with respect to the exact result computed for systems to which FCC and FCI *can* be applied. This process of benchmarking then allows for truncated CC and CI methods to be applied to new chemical systems with confidence. In this section, we will review some of the literature on the convergence of correlated approaches with respect to both truncation level and the completeness of the one-electron basis set, in order to set the stage for the benchmarking of even more approximate approaches for describing non-covalent interactions described in Part II of this Thesis.

---

\*The most well-known divergence of the  $MP_n$  series occurs for even mildly non-equilibrium molecular geometries<sup>11</sup> due to the fact that, especially for geometries with stretched bond distances, the presence of non-dynamical correlation can cause the magnitude of the total correlation energy to be such that it cannot be reasonably recovered by a perturbative approach.

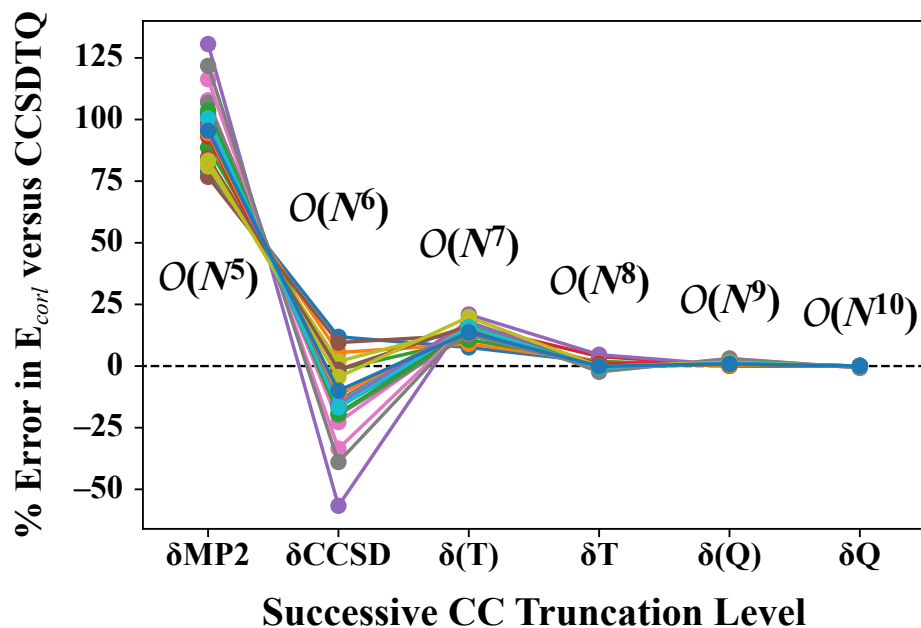


Figure 2.1: Convergence of correlation energy contribution to non-covalent interaction energies at various truncation orders for the CC series. Values taken from Table S-2 of Ref. 12.

#### *Convergence of the Correlation Energy with Truncation Level*

The hallmarks of truncated CI and CC theories is that more complete truncations of the excitation operator applied to the Hartree–Fock reference generates produces more exact wavefunctions with correspondingly more exact energies. This can be verified by examining the change in correlation energy with truncation level, which should exhibit the monotonic convergence of the absolute error from the FCI limit with increase in the truncation order. Visualized in Fig. 2.1 is precisely such an analysis, whereby the percent error in the correlation energy contribution to non-covalent interaction energies versus CCSDTQ is plotted for 21 small bimolecular complexes (data taken from Fig. S-2 of Ref. 12 ). Clearly, monotonic convergence of the absolute percent error (A%E) is present for the CC series; furthermore, oscillations across the zero of percent error is not unexpected, as truncated CC is not variational due to the similarity transformation of the Hamiltonian operator. An equivalent plot for truncations in the CI series would therefore be expected to converge exclusively from above, but exhibit much the same magnitudes of signed percent errors at

each truncation level. In practice, post-(T) excitations contribute very little to the total computed interaction energy, with typical contributions on the order of several hundredths of kcal mol<sup>-1</sup>,<sup>13-15</sup> which is on the same order of magnitude as core correlation effects.<sup>14,16-18</sup> For typical CCSD(T)/CBS computations employing the frozen-core approximation, therefore, the most fine-grained level to which comparisons with reference interaction energies can be made is 0.01 kcal mol<sup>-1</sup>.

### *Convergence of the Correlation Energy with Basis Set*

Truncated CI and CC are both convergent with respect to the truncation level for the number of excited configurations included in the wavefunction, in a given one-electron basis set. Since FCI and FCC are only exact in a given one-electron basis set, it is also worth examining how these theories (and their truncations) converge towards their respective complete basis set (CBS) limits. Furthermore, convergence towards the CBS limit is of interest even for particular levels of non-convergent theories, like Hartree–Fock and MP2.\* For Hartree–Fock, convergence of the total energy towards the CBS limit (also referred to as the Hartree–Fock limit) is exponential with respect to the  $\zeta$ -level of the basis set, with the convergence essentially achieved even at the quadruple- $\zeta$  level.

Unfortunately, correlated approaches do not converge towards their respective CBS limits nearly as rapidly as for Hartree–Fock; in fact, they converge quite slowly with basis set, showing a  $\ell^{-3}$  dependence on the order of the basis set,  $\ell$ . This is due to the fact that in order for the Fermi hole surrounding each electron to be modeled correctly, the wavefunction itself must exhibit a node at the electronic positions. Referred to as *interelectronic cusps* in the wavefunction, the behavior of the wavefunction in these cusp regions is well known thanks to the work of Kato,<sup>19</sup> and are known as the Kato cusp conditions. Unfor-

---

\*DFT, unfortunately, is neither systematically improvable with respect to functional construction on successive rungs of Jacob’s ladder or with respect to the completeness of the one-electron basis set. It is therefore more appropriate to choose a combination of functional/basis set whose accuracy has previously been assessed. The most common such choice is the classic combination of B3LYP/6-31G\*, whose accuracy for computing total energies and reaction barriers has led to its popularity and widespread adoption by the organic chemistry community.

tunately, however, the cusp conditions require substantial variational flexibility to recover exactly, thereby necessitating large basis sets to do so even approximately. Often, correlated methods are poorly converged at even the quadruple- or quintuple- $\zeta$  levels, making their systematic application to chemical systems of even modest size challenging. Furthermore, the slow convergence of the correlation energy presents significant challenges to the benchmarking of more approximate methods, due to the difficulty in obtaining adequate reference values for certain properties.\* To address both the basis set and truncation level challenges, several approaches have been developed which combine two or more computations with more affordable methods or basis sets to obtain a result which is much more accurate.

### 2.8.2 Accelerating Convergence for Correlated Approaches

In this section, we introduce several approaches for accelerating the convergence of correlated computations, both toward the FCI and CBS limits, which will be of critical importance later in the Thesis to the particular application of benchmarking non-covalent interaction energies.

#### *Basis Set Extrapolation*

As discussed previously, the correlation energy converges slowly towards the complete basis set limit because of the difficulty in describing the interelectronic cusp. Thanks to the work of Dunning, however, a basis set family exists which are “correlation-consistent” and converge sufficiently smoothly towards the CBS limit that computations performed at successive  $\zeta$ -levels can be extrapolated towards this asymptotic limit with a simple polynomial functions. The most successful such “basis set extrapolation” approach fits the correlation

---

\*Among the great ironies of quantum chemistry is the fact that Hartree–Fock — the most affordable electronic structure method which is convergent towards the CBS limit — does so exponentially, while correlated approaches whose computational scaling is at least an order of magnitude less favorable than Hartree–Fock converge much more slowly. It has been the bane of my own and many others’ existence.



energy in a basis set of cardinal number  $\ell$ ,  $E_\ell^{\text{corl}}$ , to the power law

$$E_\ell^{\text{corl}} = E_\infty^{\text{corl}} + A\ell^{-3}, \quad (2.88)$$

where  $E_\infty^{\text{corl}}$  is the correlation energy at the complete basis set limit. Helgaker and coworkers showed<sup>20</sup> that this expression could be rearranged to provide an approximation to the CBS-limit via a “two-point” extrapolation of correlation energies computed in two different basis sets with cardinal numbers  $X$  and  $Y$  as

$$E_\infty^{\text{corl}}(X, Y) = \frac{E_X^{\text{corl}}X^{-3} - E_Y^{\text{corl}}Y^{-3}}{X^3 - Y^3} \quad (2.89)$$

Typically, since the Hartree–Fock energy converges exponentially towards the CBS limit, only the correlation energy is extrapolated, with the Hartree–Fock contribution to the total energy simply being taken to be converged in the larger of the basis sets employed in the extrapolation.<sup>21</sup> Using the correlation consistent basis sets of Dunning,<sup>22</sup> extrapolations of the correlation energy by the Helgaker formula above are typically denoted CBS(XZ, YZ) or even XYZ, e.g., MP2/CBS(aTZ, aQZ) and MP2/aTQZ both refer to the same two-point Helgaker extrapolation of the MP2 correlation energy computed in the aug-cc-pVTZ and aug-cc-pVQZ basis sets, where the Hartree–Fock energy is taken to be converged in the aug-cc-pVQZ basis set.

### *Focal-Point Analysis & Composite Approaches*

While CBS extrapolations do address the challenge of converging the correlation energy in a given one-electron basis set towards the complete basis set limit, they do not address deficiencies in the method itself (i.e., errors with respect to FCI) nor reduce the computational expense of a given method itself. Thanks to the fact that higher-order correlation effects (arising from e.g., triple, quadruple, etc. substitutions) are fairly well-described by “small” basis sets, the difference between a higher-order correlation method [e.g., CCSD(T)] and

lower-order one (e.g., MP2) does not change significantly with basis set. In other words, the difference between high-level and low-level methods computed in a small basis set is a good approximation of the same difference at the CBS limit. The CBS limit of a lower-level method can therefore be “corrected” towards the CBS limit of a higher-order method:

$$E_{\infty}^{\text{high}} \approx E_{\infty}^{\text{low}} + \delta_{\text{low}}^{\text{high}} \quad (2.90)$$

$$\delta_{\text{low}}^{\text{high}} = E_{\text{small}}^{\text{high}} - E_{\text{small}}^{\text{low}} \quad (2.91)$$

For this Thesis, the most common such approach corrects the CBS limit of MP2 to that of CCSD(T) by

$$E_{\text{CBS}}^{\text{CCSD(T)}} \approx \text{MP2/CBS(aTZ, aQZ)} + \delta_{\text{MP2}}^{\text{CCSD(T)}/\text{aTZ}} \quad (2.92)$$

$$\delta_{\text{MP2}}^{\text{CCSD(T)}/\text{aTZ}} = E_{\text{aTZ}}^{\text{CCSD(T)}} - E_{\text{aTZ}}^{\text{MP2}}, \quad (2.93)$$

which is denoted CCSD(T)/[aTQZ;  $\delta$ :aTZ]. These “focal-point” (or “composite”) approaches have become a popular approach for addressing both the basis set and truncation challenges for computing high-quality correlation energies in a variety of contexts,<sup>23,24</sup> and indeed have been applied widely to compute properties for which electron correlation is important, particularly to describe non-covalent interactions.<sup>25–30</sup>

### *Explicitly Correlated Approaches*

Despite their success for computing high-accuracy energies and properties, basis set extrapolations and focal-point approaches do suffer from the fact that they require multiple, separate computations to be performed and combined to yield a single result. Not only does this necessitate additional computational strain than performing a single computation, but until recent advancements made in automating quantum chemistry workflows,<sup>31</sup> performing focal-point analyses or CBS extrapolations was often significantly challenging also for the user of quantum chemistry software. Since a major reason for the extrapolation of cor-

relation energies is that the interelectronic cusp requires significant variational flexibility to model correctly, an alternative approach would be to require that the wavefunction depend explicitly on the distance between any pair of electrons 1 and 2,  $r_{12}$ . The simplest manner in which to include this “explicit correlation” would be to force the wavefunction to depend explicitly on terms linear in  $r_{12}$ ; indeed, this is precisely the approach (together with terms depending on  $r_1 + r_2$  and  $r_1 - r_2$ ) taken by Hylleraas in his highly accurate treatment of the Helium atom.<sup>32</sup> Methods which depend only linearly on the interelectronic separation are often referred to as R12 methods; more recently, generalizations of R12 theory in which a function of  $r_{12}$ ,  $f(r_{12})$ , is used have been developed. These methods, referred to as F12 methods, typically leverage a simple Slater geminal function  $f_{12}(r_{12}) = e^{-\beta r_{12}}$ ,<sup>33</sup> and have been shown to greatly accelerate the convergence of the correlation energy in a given basis set, where, e.g., MP2-F12/aTZ may be just as accurate as MP2/a5Z versus MP2/CBS.

## 2.9 Density Functional Theory

In this section, we present an alternative approach to the quantum mechanical description of a molecular system than the one taken by the wavefunction-based methods discussed above, or, indeed, even the electronic Schrödinger equation itself. Taking the electron density  $n(\mathbf{r})$  to be the fundamental quantity defining a quantum mechanical system, rather than the electronic wavefunction  $\Psi_{\text{elec}}$ , *density functional theory* (DFT) has become the *de facto* approach for applying quantum mechanics to predict both molecular and solid-state properties in a great number of contexts, the development of which even earned a share of the 1998 Nobel Prize in Chemistry. While the application of DFT plays a relatively minor role in this Thesis, it represents such a significant fraction of computational chemistry that it deserves specific mention in this Chapter. Therefore, we will introduce in this section the basics of DFT and its most successful molecular formulation, Kohn-Sham DFT (KS-DFT).

### 2.9.1 The Hohenberg-Kohn Theorems

In order to build density functional theory, let us first consider a general electronic Hamiltonian operator for a many-electron system under an external potential  $v(\mathbf{r}_i)$ , which includes any effects external to the electrons (including Coulombic nuclear-electron attraction, etc.):

$$\widehat{\mathcal{H}}_{\text{elec}} = -\frac{1}{2} \sum_i \nabla_i^2 + \frac{1}{2} \sum_{i>j} \frac{1}{|\mathbf{r}_i - \mathbf{r}_j|} + \sum_i v(\mathbf{r}_i) = \hat{T}_e + \hat{V}_{ee} + \hat{v}_{\text{ext}} \quad (2.94)$$

Beginning with this Hamiltonian, the two Hohenberg-Kohn theorems can be stated as follows.

**Theorem 2.9.1.1** (The First Hohenberg-Kohn Theorem). *The external potential  $v(\mathbf{r}_i)$  is a unique functional of the electron density  $n(\mathbf{r})$  in the ground state, and therefore the total ground-state energy (and, by extension, all observable properties) are uniquely determined by the ground state electron density.*

**Theorem 2.9.1.2** (The Second Hohenberg-Kohn Theorem). *The total energy of a many-electron system is minimized for the correct ground-state electron density.*

The power of the Hohenberg-Kohn theorems lies in the fact that they establish a bijective map between the exact ground state energy and the ground state electron density, via a *functional* of the electron density, i.e., as  $E[n(\mathbf{r})] = E[n]$ , where the brackets indicate the functional dependence of the energy  $E$  on the density function  $n$ . Therefore, in order to determine either the ground state energy or the exact ground state density, all that is required is to variationally minimize the ground state energy with respect to changes in the electron density. Unfortunately, the Hohenberg-Kohn theorems do not provide any clues about *how* this may be done in practice, rather only establishing the possibility of doing so.

### 2.9.2 The Universal Functional

In order to move closer to a practical formulation of DFT, let us define the *universal functional*,  $F[n]$ , such that  $F$  minimizes the ground-state energy with respect only to electronic

physics (i.e., neglecting  $\hat{v}_{\text{ext}}$  for the moment):

$$F[n] = \min_{\Psi \rightarrow \rho} \langle \Psi | \hat{T}_e + \hat{V}_{ee} | \Psi \rangle = T[n] + J[n] + E_{QM}[n], \quad (2.95)$$

where  $T[n]$  is the kinetic energy functional,  $J[n]$  is the classical Coulomb interaction of electron density with itself, and  $E[n]$  encapsulates the energetic dependence on non-classical electronic behavior. While the form of  $J[n]$  is known exactly and the form of  $E[n]$  can be approximated reasonably well, the kinetic energy functional is nearly impossible to define. Fortunately, however, Kohn and Sham introduced a formulation of DFT in which  $T[n]$  can be approximated. This advancement allowed DFT to actually be applied in practice, rather than being a purely academic exercise.

### 2.9.3 Kohn-Sham DFT

In the Kohn-Sham formulation of DFT (KS-DFT), the electron density is assumed to be generated by a set of one-electron, non-interacting orbitals  $\{\phi_i(\mathbf{r}_i)\}$ :

$$n(\mathbf{r}) = \sum_i |\phi_i(\mathbf{r}_i)|^2 \quad (2.96)$$

In this picture, the kinetic energy functional is simply

$$T_s[n] = -\frac{1}{2} \sum_i \langle \phi_i | \nabla_i^2 | \phi_i \rangle \quad (2.97)$$

Of course,  $T_s[n]$  is not exactly representative of the true kinetic energy functional  $T[n]$ , but for cases when an independent particle approximation is appropriate,  $T_s[n]$  will be good enough.\* Based on this form of the electron density and effective kinetic energy functional,

---

\*As it turns out, KS-DFT is based on the same mean-field approximation as Hartree–Fock, so it will be vulnerable in the same instances in which a single Slater determinant is not an appropriate representation of the full  $N$ -electron wavefunction.

the total energy functional is given by

$$E_{\text{KS}}[n] = F_{\text{KS}}[n] + v_{\text{ext}}[n] = T_{\text{s}}[n] + J[n] + v_{\text{ext}}[n] + E_{XC}[n], \quad (2.98)$$

where  $v_{\text{ext}}[n]$  is (typically) the electron-nuclear attraction and  $v_{XC}$  is the *exchange-correlation functional*, which encapsulates the many-body non-classical electronic interactions arising from Pauli exchange and dynamical electron correlation.\* This energy functional  $E_{\text{KS}}[n]$  and its associated Hamiltonian operator define a set of one-particle equations which, after spin-integration, yield

$$\left\{ -\frac{1}{2}\nabla_i^2 - \sum_A \frac{Z_A}{r_{iA}} + \int \frac{n(\mathbf{r}_2)}{r_{12}} d\mathbf{r}_2 + V_{xc}(\mathbf{r}_1) \right\} \phi_i = \epsilon_i \phi_i, \quad (2.99)$$

where  $V_{xc}(\mathbf{r}_1)$  is the *exchange-correlation potential*, is given by

$$V_{xc}(\mathbf{r}_1) = \frac{\delta E_{xc}}{\delta n}.$$

These *restricted Kohn-Sham* equations, after the introduction of a basis and recasting as a linear algebra problem, can be solved self-consistently in much the same way that the RHF equations are. By doing so, we would *exactly* solve the electronic Schrödinger equation by uniquely determining the ground-state electron density, with the caveat that we have invoked an independent particle model to do so. Fortunately, however, the exchange-correlation functional nominally contains all of the information necessary to correct for the independent-particle formulation of Kohn-Sham DFT — but what exactly is the form of the exchange-correlation functional?

---

\*For more details on electron correlation, refer to Section 2.7.

#### 2.9.4 The Exchange-Correlation Functional

The exchange-correlation functional,  $E_{xc}[n]$ , formally encompasses all of the non-classical physical interaction between electrons, as well as correcting for the presence of the independent-particle kinetic energy functional:

$$E_{xc}[n] = (T[n] - T_s[n]) + (E_{ee}[n] - J[n]) = E_{\Delta}[n] + E_x[n] + E_c[n], \quad (2.100)$$

where  $E_x[n]$  is the *exchange functional*, which attempts to recover the exchange behavior resulting from the indistinguishability of electrons (i.e., Pauli exchange) and  $E_c[n]$  is the *correlation functional*, which attempts to recover the dynamical electron correlation neglected by Hartree–Fock. Often, the kinetic energy correction functional  $E_{\Delta}[n]$  is neglected, as it is assumed to be negligible. Therefore, if both  $E_c[n]$  and  $E_x[n]$  were known, then KS-DFT would be an exact theory. Of course, we do know the exact form of  $E_x[n]$ , as it is identical to the exchange contribution to the Fock matrix. Unfortunately, however, the exact form of  $E_c[n]$  is not known; furthermore, it has been speculated that even if it were known, it would be so complex that it would render KS-DFT intractable. For practical KS-DFT, therefore, approximate correlation functionals must be developed. Furthermore, to avoid the computational expense of computing the exact Hartree–Fock exchange, approximate exchange functionals are also often employed.

In formulating approximate exchange and correlation functionals, a hierarchy of approximations may be invoked, with successively more complex formulations of these quantities nominally providing\* higher accuracy in computed energies and properties. This hierarchy is typically referred to as Jacob’s ladder for density functional approximations (DFAs), with “less approximate” functionals occupying higher rungs on the ladder. On the lowest rung are placed the most approximate functionals, which depend only on the

---

\*Unlike wavefunction methods, DFT is convergent neither in the completeness of the one-electron basis set nor in the rung of Jacob’s ladder for a given functional. The accuracy of DFT is highly system-dependent, and must be assessed before application of a particular combination of density functional and basis set may be chosen for a given problem of interest.

local density at given spatial coordinates. Referred to as the local density approximation (LDA), functionals in this family scale as  $\mathcal{O}(N^3)$  with basis set size, but can sometimes fail spectacularly. For no increase in algorithmic complexity, functionals occupying the next rung on Jacob’s ladder depend on both the local density and its gradient, referred to as generalized gradient approximation functionals (GGAs). Even though these functionals offer higher accuracy for energies and properties, they are still fundamentally local; therefore, they may not perform well for non-local properties like non-covalent interactions or extended molecular geometries. The next logical step up the ladder is to force the dependence of exchange-correlation functional on local density, its gradient, and its Hessian — the *second* derivative of the density with respect to spatial coordinates. These functionals are referred to as meta-GGAs, and while offering superior performance to GGAs, can suffer from numerical instabilities due to the order of derivative required to build the exchange correlation potential. To this point, we have employed only approximate exchange-correlation functionals; why not use exact (Hartree–Fock) exchange instead? For the additional cost of one order of magnitude higher scaling [i.e.,  $\mathcal{O}(N^3) \rightarrow \mathcal{O}(N^4)$ ], the exchange functional can be corrected by mixing in some fraction of exact, Hartree–Fock exchange. These are referred to as *hybrid* functionals, with numerous examples of hybrid-GGAs and hybrid-meta-GGAs throughout the literature. A similar approach could be taken to correct the correlation functional by mixing in some fraction of MP2 correlation, with the caveat that these “double-hybrid” functionals exhibit the same scaling as MP2 [ $\mathcal{O}(N^5)$ ].

### 2.9.5 Non-Local Corrections to DFT

As discussed above, KS-DFT contains only *local* information until exact exchange is incorporated by hybrid functionals. As such, non-local properties like non-covalent interactions, charge migration, etc. will not be well described by KS-DFT alone. Furthermore, mixing in a part of the exact exchange, either for all orbital pairs or (as in the range-separated hybrid approach) switching it on at long range, cannot solve all of the locality issues with



KS-DFT, as exact exchange only correlates parallel-spin electrons. Furthermore, punting to a double-hybrid functional results in a significant increase in computational expense. It would therefore be desirable to correct KS-DFT for non-local electron correlation without computing it at the MP2 level. Towards this end, several such schemes have been introduced (see Ref. 34 for a very thorough review); the most relevant of which to this Thesis is the third-generation dispersion correction of Grimme,<sup>35</sup> denoted by appending “-D3” to the functional abbreviation, e.g., B3LYP-D3. This “dispersion-corrected DFT” (DFT-D) approach seeks to correct the total molecular energy for the presence of long-range electronic correlation (for particles separated by more than  $\approx 3.5 \text{ \AA}$ ) by adding to the DFT energy the semi-empirical, atom-pairwise correction

$$E_{\text{disp}}^{\text{molec}} = - \sum_{AB} E_{\text{disp}}^{(n)} = - \sum_{AB} \sum_{n=6,8,10,\dots} \frac{C_n^{AB}}{R^n} f_{\text{damp}}^{(n)}, \quad (2.101)$$

where the  $C_n^{AB}$  coefficients are based on the Casimir-Polder expression for the dispersion interaction between two spherically-symmetric electron densities oscillating in imaginary frequency,

$$E_{\text{disp}}^{(6)} = - \frac{3}{\pi R^6} \int_0^\infty d\omega \alpha_A(i\omega) \alpha_B(i\omega) = - \frac{C_6^{AB}}{R^6}, \quad (2.102)$$

and the damping function  $f_{\text{damp}}^{(n)}$  tends to zero as  $R_{AB} \rightarrow 0$  to ensure that the multipole expansion of the dispersion energy remains well-defined, since it exhibits singularities at  $R_{AB} = 0$ .

Since the  $C_6$ ,  $C_8$ , etc. coefficients are fixed for a given pair of atoms, the flexibility in the -D correction arises from the choice of damping function and its parameterization. Two of the most widely used damping functions are the Becke-Johnson damping function<sup>36,37</sup> [denoted (BJ)] and the “zero damping” scheme of Chai and Head-Gordon,<sup>35,38</sup> [denoted

(0)] which are combined with the -D3 correction to yield

$$E_{\text{disp}}^{-\text{D3(BJ)}} = -\frac{1}{2} \sum_{A<B} \sum_{n=6,8} s_n \frac{C_n^{AB}}{r_{AB}^n + (\alpha_1 R_0^{AB} + \alpha_2)^n} \quad (2.103)$$

$$E_{\text{disp}}^{-\text{D3(0)}} = -\frac{1}{2} \sum_{A<B} \sum_{n=6,8} s_n \frac{C_n^{AB}}{r_{AB}^n} \frac{1}{1 + 6(r_{AB}/(s_{r,n} R_0^{AB}) + R_0^{AB} \beta)^{-\alpha_n}} \quad (2.104)$$

These damping functions were originally parameterized to reproduce reaction energies and barrier heights, together with non-covalent interactions, computed using CCSD(T).<sup>35,36</sup> While this initial parameterization did afford additional accuracy for DFT-D3 versus uncorrected DFT, very high-quality reference energies have become significantly more widely available since these parameterizations were introduced. Recently, Smith *et al.* revised the damping parameters for both BJ and zero-damping functions for a wide variety of density functionals by training them against a very large set of CCSD(T)/CBS-quality energy points comprised of approximately 1,600 interaction energies and potential energy curves.<sup>39</sup> The result of this revision was that the DFT-D errors versus CCSD(T)/CBS were reduced significantly, while also reducing the variability in performance with the choice of density functional.

## CHAPTER 3

### THEORETICAL APPROACHES FOR NON-COVALENT INTERACTIONS

In this Chapter, we will briefly discuss the manner in which the electronic structure methods introduced above are applied to study non-covalent interactions (NCI). First, however, we will introduce the quantities of interest when studying NCI, the interaction energy and its physically and chemically meaningful components, before finally introducing a family of electronic structure methods specifically designed to be applied to NCI. There are a number of excellent reviews and book chapters discussing this material in significantly more detail than possible here; therefore, we will only present what discussion is absolutely necessary to understand the following Chapters. For further detail, we refer the interested reader to those references where relevant.

#### 3.1 Defining the “Interaction Energy”

When two particles are infinitely separated, the energy of the pair is simply the sum of the energies of the individual particles. For fermionic particles, this size-consistency is not quite obeyed when the particles are separated by finite distances, however, as their motions (and, consequently, their energies) are correlated. As such, fermionic particles begin to *interact* when they become close to one another. Since atoms and molecules can be considered to be simply clouds of fermions, as under the Born-Oppenheimer approximation their nuclei (which may be either bosons or fermions depending on nuclear spin) are clamped, atoms and molecules will also interact at finite distances. This begs the question of by *how much* two or more atoms/molecules interact, which could also be asked “how much do the motions of the fermions belonging to one species impact the motions of the fermions in the opposite species?”, or equivalently, “by how much does the energy change in the pair when they become close?”.

We may define the *interaction energy* between two molecules  $\mathcal{A}$  and  $\mathcal{B}$  (the *monomers*) as the difference between the total energy of the pair  $\mathcal{AB}$  (the *dimer*) and the sum of the differences of each monomer in isolation. This may be written as

$$\text{IE}_{\mathcal{AB}} = E_{\mathcal{AB}} - E_{\mathcal{A}} - E_{\mathcal{B}}, \quad (3.1)$$

where  $\text{IE}_{\mathcal{AB}}$  is the interaction energy of the dimer  $\mathcal{AB}$ , and  $E_{\mathcal{AB}}$ ,  $E_{\mathcal{A}}$ , and  $E_{\mathcal{B}}$  are the energies of the dimer and each monomer, respectively. This expression can be generalized to describe the interaction energy for an arbitrary number of associating atoms or molecules  $\mathcal{A}, \mathcal{B}, \dots \mathcal{K}$  by replacing the energy of the dimer with the energy of the “supersystem”  $\mathcal{S} = \mathcal{AB} \cdots \mathcal{K}$  and subtracting off the energy of each monomer in isolation:

$$\text{IE}_{\mathcal{S}} = E_{\mathcal{S}} - \sum_{\mathcal{I} \in \mathcal{S}} E_{\mathcal{I}} \quad (3.2)$$

It is worth noting that we have not yet specified the level of theory at which these dimer and monomer energies are computed.\* For the sake of argument, let us assume that it is desired to compute the interaction energy of a particular configuration of the benzene pentamer at the CCSD(T)/aug-cc-pVTZ level of theory. This quantity, according to Eqn. 3.2, would naïvely require separate computations to be performed for the total energy of each of the five benzene monomers in the complex, together with a single computation on the full pentamer. The monomer computations are routine enough; the pentamer computation, however, would likely not be possible at all, thanks to the steep  $\mathcal{O}(N^7)$  algorithmic scaling of CCSD(T). Indeed, given that the monomer energies would take roughly two hours each to compute, the pentamer computation would require roughly 18 years!

---

\*In general, we will use calligraphic capital letters ( $\mathcal{A}, \mathcal{B}, \mathcal{C}$ , etc.) to represent molecules (or collections thereof) as abstract chemical entities, while we will denote chemical entities described at a particular level of theory with italic capital letters ( $A, B, C$ , etc.).

### 3.2 The Many-Body Expansion

Such a time commitment for a single result would seem to render accurate investigations of such systems intractable. This can be avoided, however, by leveraging the many body expansion (MBE), whereby the total energy of a collection of  $N$  molecules can be exactly written as

$$E^{(n)} = \sum_{\alpha=1}^N E_{\alpha}^{(1)} + \sum_{\alpha=1}^{N C_2} \Delta E_{\alpha}^{(2)} + \sum_{\alpha=1}^{N C_3} \Delta E_{\alpha}^{(3)} + \dots, \quad (3.3)$$

where each  $\Delta E_{\alpha}^{(i)}$  is the  $i$ th-order energy correction, defined recursively by

$$\Delta E_{\alpha}^{(n)} = E_{\alpha}^{(n)} - \sum_{\beta=1}^{n C_{n-1}} \Delta E_{\beta}^{(n-1)} - \sum_{\gamma=1}^{n C_{n-2}} \Delta E_{\gamma}^{(n-2)} - \dots - \sum_{\omega=1}^n E_{\omega}^{(1)}, \quad (3.4)$$

and the summations over unified indices  $\alpha, \beta, \gamma, \dots$  for each  $i$ th-order correction are over all  $N C_i$  unique  $i$ -mers in the supersystem. Under the MBE, an *approximate* total energy can be computed for the supersystem by truncating the expansion at a desired level, e.g., by including all trimers but no tetramers. This can drastically reduce the overall expense of computing the total energy of the supersystem by constructing it from computations on smaller subsystems, each of which can furthermore be performed in a pleasantly parallelizable fashion across many compute nodes. If the interactions between monomers are largely non-cooperative (i.e., no substantial mutual polarization or dispersion effects among trimers, tetramers, etc.), the MBE can be truncated at fairly low order while still retaining accuracy.

The MBE also provides a convenient expression for the interaction energy of a large assembly of non-bonded fragments, provided by simply subtracting away the energies of all monomers from that of the supersystem. This yields the same expression as the one defined above for the  $i$ th order energy correction, where  $i$  is simply taken to be the total

number of monomers:

$$\text{IE}^{(n)} = E_{\alpha}^{(n)} - \sum_{\beta=1}^{nC_{n-1}} \Delta E_{\beta}^{(n-1)} - \sum_{\gamma=1}^{nC_{n-2}} \Delta E_{\gamma}^{(n-2)} - \dots - \sum_{\omega=1}^n E_{\omega}^{(1)}. \quad (3.5)$$

It is clear that this expression could also be truncated at a chosen order to approximate the total IE by neglecting higher-body contributions. Just like for the total supersystem energy, it has been shown that the truncation level for this expression necessary to achieve a particular level of accuracy system-dependent; for assemblies bound largely by dispersion, truncation at three- and even two-body terms yields a reasonable estimate of total interaction energy,<sup>40-43</sup> while for hydrogen bonded complexes, four- and five-body contributions to the interaction energy are still significant.<sup>44-46</sup>

### 3.3 Types of Non-Covalent Interactions

We have already seen what non-covalent interactions *do* — they change the total energy of a collection of molecules from what they would be if they were infinitely isolated to what they are when they become close to one another — but we have yet to address the question of exactly *what* causes this phenomenon and *why* it does so. Fortunately, we are already equipped with all of the knowledge necessary to answer these questions, at least at the conceptual level. For the rigorous mathematical details, we refer the interested reader to the proverbial bible of NCI, Anthony Stone’s *The Theory of Intermolecular Forces*,<sup>47</sup> as well as several books, collections, and review articles written by leaders in the field.<sup>48-50</sup>

Nearly every particular type of non-covalent interaction, of which there are many specific examples — C-H- $\pi$ , hydrogen bonding, halogen bonding,  $\pi - \pi$ , and hydrophobic interactions, among others — can be thought of as the interplay between one or more of four basic forces: electrostatics, exchange, induction, and dispersion. For the sake of simplicity, therefore, we will provide here a birds-eye view of what each of these forces are and how they arise, before later discussing the particular computational and/or theoretical

challenges that each of these forces present.

**Electrostatics** Electrostatic forces are felt between all charged particles according to Coulomb's Law. For molecular interactions, therefore, electrostatics encompasses the attraction between the electrons of one monomer to the nuclei of the other and the repulsion between the nuclei of nearby monomers. Furthermore, a variety of electrostatic interactions can arise from the interaction of electron densities. When one or more molecule with a net electronic dipole moment (or even local dipole moments on particular functional groups) interact, the attraction or repulsion between the partial positive and partial negative "ends" of the dipole(s) is also an electrostatic force.

**Exchange** Exchange forces occur when electron densities overlap, and are always repulsive. Therefore, this is often referred to as *exchange-repulsion* or *steric repulsion*. Exchange, as well as its connection to the antisymmetry principle, will be discussed in more detail in Section 3.5.

**Induction** Induction forces, often referred to as *polarization*, occur when a permanent full or partial charge on one species *induces* a separation of charge on another species, which then interact electrostatically. This effect can be either attractive or repulsive, but is typically attractive.

**Dispersion** Dispersion forces, originally introduced by London<sup>51,52</sup> (and thereby typically referred to as London dispersion), arise when instantaneous fluctuations of electrons in one species form a temporary dipole moment, thereby inducing an instantaneous, complementary dipole moment in another species. Always attractive, these mutual, instantaneous fluctuations are due exclusively to electron correlation, and as such cannot be recovered by a mean-field theory like Hartree–Fock. Furthermore, the strength of the dispersion interactions are proportional to the total number of electrons.

### 3.4 Supramolecular Approaches for Non-Covalent Interactions

As for any chemical property, many practical considerations must be taken into account before submission to a computational approach. First and foremost, since electron correlation plays a large role in determining the strength of NCI, many of the same considerations which were discussed in Section 2.8 are relevant when computing NCI — including the question of size consistency for a given method, the truncation level for the correlated method employed, and the choice and completeness of the one-electron basis set. Since the interaction energy is defined as a difference of total energies, it would seemingly follow that a better description of these total energies should also lead to a better description of the interaction energy. In some cases, however, this may not actually be the case due to factors unique to interaction energy computations. It is with these additional factors which will dominate much of the discussion in this section, however let us begin by addressing factors which we already know to be important considerations for any correlated computation: truncation level and basis set completeness.

#### 3.4.1 Correlation Effects and Size Consistency

We have defined the interaction energy to be the difference between the energies of a dimer and the energies of its isolated constituent monomers. In doing so, we have relied on the fact that at an infinite separation distance the two monomers do not interact, and thus the energy of the full system is simply the sum of the monomer energies. Recalling the previous discussion in Section 2.7.1, this limiting behavior is exactly what we had defined as *size consistency*, a well-known shortcoming of truncated CI variants. Therefore, it is not advisable to apply anything short of FCI to supermolecular IE computations, which of course is impossible for all but the smallest interacting systems, e.g., He-H<sub>2</sub>, etc. Since MP<sub>n</sub> and CC are size-consistent for all truncation levels of either the perturbation series or the cluster operator, these approaches are preferred for application to supermolecular



computations of the interaction energy. Once either  $MP_n$  or CC has been chosen, however, the appropriate level of truncation remains to be determined. The appropriate choice thereof must be determined on a case-by-case basis, driven by the interplay between the desired level of accuracy and the computational expense which can be afforded. Since the cost of a given computation — as well as its accuracy — also depends significantly on the completeness of the one-electron basis set, however, we must first address concerns related to the basis set which are unique to supermolecular IE computations before determining a preferred *level of theory* for a particular computation, which we define to be a combination of method and basis set.

### 3.4.2 Basis Set Incompleteness and Basis Set Superposition Errors

As with any property, the completeness of the one-electron basis set can significantly affect the quality of computed supermolecular interaction energies. This is especially true for NCI, however, since so much of this property is determined by dynamical electron correlation, which as we have already seen is particularly slow to converge towards the complete basis set (CBS) limit. Therefore, the incompleteness in the basis set can cause significant errors to occur; this is referred to as the *basis set incompleteness error* (BSIE), and is formally defined to be the difference between the value computed in a finite basis from that of the CBS limit:

$$BSIE = IE^\infty - IE^{\text{finite}} \quad (3.6)$$

Fortunately, we may employ the same tactics to reduce the BSIE in an IE computation as were introduced in Section 2.8.2, including CBS extrapolation, focal point approaches, and even explicit correlation; indeed, the application of these convergence acceleration methods to NCI has received significant attention, including in the work reproduced in Chapter 4. The largest effect of BSIE on a supermolecular IE computation is that an insufficient amount of “cross talk” is possible between the two monomers to allow for the stabilization of the dimer relative to the isolated monomers, causing the computed IE to be under bound

relative to the true value.

Unfortunately, however, even if the individual energies of the dimer and each monomer are computed at the CBS limit using, e.g., a two-point extrapolation of energies computed in the cc-pVTZ and cc-pVQZ basis sets, there still may be significant errors in the total interaction energy. This is due to the fact that, in order to properly describe the interaction of the two molecular species, some spatial overlap of basis functions from each monomer must be present. For typical non-bonded contact distances of between 3.5-5 Å, the spatial extent of the cc-pVXZ family of basis functions is often insufficiently diffuse to allow for overlap to occur. The effect of this lack of overlap is identical to that of BSIE, namely, that cross talk between the monomers cannot occur, leading to an under-binding of the complex. While some have advocated adding “midbond” functions centered between the two monomers as a method to allow for crosstalk,<sup>16</sup> the simplest (and most even-handed) approach is to simply augment the one-electron basis with a set of more spatially diffuse functions that would allow for overlap. Such basis sets are prefixed with “aug-,” i.e., aug-cc-pVTZ (abbreviated aTZ), and are the standard choice for NCI computations.

Addressing BSIE by performing computations directly at the CBS limit is not always practical, or indeed even possible; in such cases, when finite basis sets must be used, BSIE is not the only basis set error to which IE computations are susceptible. If the dimer component of a supermolecular IE is performed in a so-called “dimer-centered” basis (where the basis functions from both monomers are included) but the monomers are each only described by their own monomer-centered basis sets, an artificial stabilization of the dimer will occur thanks to the increased variational flexibility of the monomers within the dimer wavefunction relative to their isolated wavefunctions. This over-stabilization of the dimer relative to the monomers causes supermolecular IEs computed in this manner to bound artificially strongly (i.e., too negative an IE) compared to the value computed in a complete basis set. This effect is referred to as *basis set superposition error* (BSSE), and the best method for addressing BSSE (or whether to address it at all) has been the center of much

debate in the community (see, e.g., Refs. 53 and the references therein). The most common approach to removing BSSE from a two-body supermolecular IE computation is to simply perform each of the computations in the dimer basis set; this way, the monomers feel the same variational flexibility on their own as they do in the dimer, which will cancel when the energies are subtracted to compute the IE.\* This approach, termed the *counterpoise correction* (CP), was introduced independently in two separate works but is largely attributed to Boys and Bernardi.<sup>54</sup>

The main argument for *not* addressing BSSE is that the over-binding of the supermolecular IE due to BSSE can fortuitously cancel with the under-binding due to BSIE. The result is one of the many instances of a contradiction in quantum chemistry: a more theoretically rigorous result which is, in fact, less accurate than the less rigorous one. This is especially true for CP vs. non-CP corrected IEs computed at the MP2/6-31G\* level of theory: when these IEs are CP-corrected, they can exhibit errors greater than 1 kcal mol<sup>-1</sup> from reference IEs, but when they are uncorrected for BSSE, they can be accurate to within 0.05 kcal mol<sup>-1</sup> of reference IEs.<sup>55,56</sup> The question of correcting or not correcting for BSSE — or indeed taking an average approach — was recently examined by Burns *et al.*, who found that for basis sets larger than aQZ (including CBS extrapolations and focal-point methods) CP-corrected IEs agreed better with high-quality reference data than did uncorrected IEs.<sup>57</sup> In general, it is our goal in this Thesis to apply rigorous theoretical methodologies to address chemical phenomena whenever possible; therefore, it is our belief that the CP correction should be combined with CBS extrapolation whenever tractable.

### 3.4.3 Medal Winners for Non-Covalent Interactions

Reviewed recently in an excellent book chapter by Sherrill,<sup>59</sup> the application of wavefunction theory methods to the supermolecular computation of NCI is a well-established practice which has been explored thoroughly. Furthermore, the application of XDM, wdW-DF,

---

\*For the sake of simplicity, we will limit our discussion of BSSE to two-body IE computations; for a treatise on BSSE in many-body IE computations, please refer to Ref. 53 and the references therein.

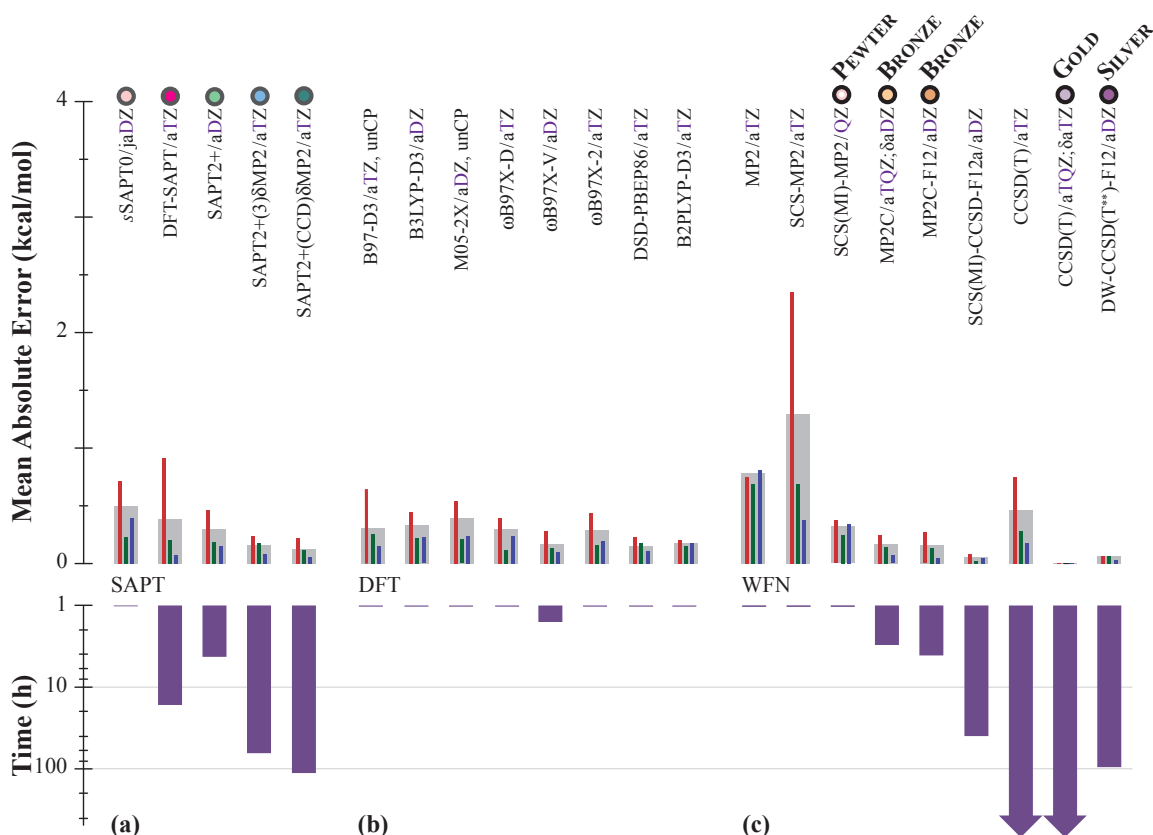


Figure 3.1: Comparison of computational classes for NCI. (a) The recommended SAPT model chemistries from Ref. 58 are compared to (b) common DFT approaches and (c) common or recommended wavefunction techniques from Ref. 55 according to both efficiency (purple; time required for adenine-thymine) and accuracy (grey; MAE averaged over S22, HBC6, HSG, and NBC10 databases) metrics. Subset MAE values are shown as inset bars for hydrogen-bonding (red), mixed-influence (green), and dispersion-dominated (blue) NCI motifs. Figure reproduced from Ref. 55.

and other DFT and DFT-D approaches to NCI is similarly well studied, with several chapters in that same volume dedicated to their discussion,<sup>60–62</sup> together with excellent collections in the literature.<sup>63</sup> Much of these efforts have gone to establishing best-practices for properly choosing a level of theory which balance the desired level of accuracy with computational expense, a non-trivial fraction of which has been put forth by my current and former co-workers in the Sherrill group.\* NCI benchmarking, whereby the accuracy of IEs computed by a particular combination of method/basis set are assessed against high-level reference energies, has and continues to be an active area of research, with a significant portion of this Thesis so dedicated.

Before introducing those results in the following Chapters, however, it is worth assessing the status of the field of NCI benchmarking before the contributions set forth in this Thesis. To this end, summarized in Fig. 3.1 is the interplay between accuracy and computational expense for several of the most popular model chemistries for computing NCI from among the DFT, wavefunction, and SAPT (introduced below) families of methods, reproduced with permission from Ref. 55 (ca. 2014). Also visualized are “NCI medalists,” as appointed in Refs. 55, 58, 64, which designate particular levels of theory as best balancing a particular level of accuracy against the associated computational expense. The most notable takeaways for DFT and wavefunction methods are that (a) for DFT methods, no systematic improvement is seen when either moving to a higher rung of Jacob’s ladder or increasing the size of the one-electron basis set, and (b) for wavefunction methods, the more accurate a level of theory, the more expensive it is.†

### 3.5 Symmetry-Adapted Perturbation Theory

Even though we have seen that the supermolecular approaches can directly quantify the strength of interactions between two or more chemical species — a considerable advan-

---

\*Indeed, Dr. Lori Burns herself is a veritable sage of NCI benchmarking, and it is under her tutelage that a long line of Sherrill group graduate students and postdocs (with myself just being the latest) have contributed in this area.

†As we will see below, this is also true for levels of SAPT.

tage over indirect experimental observation — supermolecular IE computations still only provide a single scalar value by which to do so. Even though it is illustrative, the strength of the IE cannot offer any information to rationalize *why* a particular value has been computed, or to support any chemically or physically meaningful arguments thereof. Even for molecular aggregates whose total IE may be decomposed into contributions from two-body, three-body, etc. interactions via the MBE, there is still no information provided to justify why even this decomposition is observed. Instead, it would be desirable for an approach to be able to decompose the interaction energy into contributions from chemically meaningful forces, such as electrostatic attraction or steric repulsion. Fortunately, various energy decomposition analysis (EDA) approaches have been developed which offer precisely this information. Some EDA schemes partition a supermolecular IE according to charge or population analysis, etc., however these approaches can suffer from the fact that any attempt to localize the wavefunction is arbitrary, and can therefore lead to non-unique partitions of the IE.

Rather than by partitioning a supermolecular IE into semi-arbitrary contributions from different physical forces, a more rigorous approach is obtained by treating the interaction between exactly two molecules as a perturbation of their isolated monomer wavefunctions. By further ensuring that the resulting fully interacting wavefunction is antisymmetric with respect to interchange of electrons between monomers, it is possible to develop the symmetry-adapted perturbation theory (SAPT) approach. The major advantage of SAPT over other EDA schemes is that the contributions to the total interaction energy from underlying physical forces are well-defined, arising from the perturbation operator acting on the molecular orbitals of each zeroth-order monomer wavefunction, rather than as an *a posteriori* partitioning of the IE itself. As SAPT will play a significant role in the body of this Thesis, we will devote the remainder of this Chapter to its discussion; thanks to a wonderful review recently written by Patkowski,<sup>65</sup> however, we will concern ourselves here with only the details of SAPT and its applications which are most relevant to this Thesis, and refer

the interested reader to that wonderful review for all other details.

### 3.5.1 Basic Formulation of SAPT

In SAPT, the total Hamiltonian operator is written as

$$\widehat{H}_{AB}^{\text{total}} = \widehat{F}_A + \widehat{F}_B + \lambda \widehat{V}_{AB} + \zeta \widehat{W}_A + \xi \widehat{W}_B, \quad (3.7)$$

where  $\widehat{F}_A$ ,  $\widehat{F}_B$  are Fock operators corresponding to monomers  $A$  and  $B$ ,  $\widehat{V}_{AB}$  is the perturbation corresponding to intermolecular interactions, and  $\widehat{W}_A$ ,  $\widehat{W}_B$  are the perturbations corresponding to intramolecular electron correlation. This triple perturbation series can be truncated for particular values of  $\lambda$ ,  $\zeta$ , and  $\xi$ ; terms arising from such a truncation are denoted  $E^{(\lambda, \zeta + \xi)}$ . The most basic such truncation, where the intramolecular correlation is neglected but the intermolecular perturbation is included through second order, yields the so-called SAPT0 expression for the interaction energy,

$$E_{int}^{\text{SAPT0}} = E_{elst}^{(10)} + E_{exch}^{(10)} + E_{ind,r}^{(20)} + E_{exch-ind,r}^{(20)} + E_{disp}^{(20)} + E_{exch-disp}^{(20)} + \delta_{\text{HF}}^{(2)} \quad (3.8)$$

where the  $\delta_{\text{HF}}^{(2)}$  term corrects for the presence of higher order induction and exchange-induction, and is defined as

$$\delta_{\text{HF}}^{(2)} = E_{\text{int}}^{\text{HF}} - \left( E_{elst}^{(10)} + E_{exch}^{(10)} + E_{ind,r}^{(20)} + E_{exch-ind,r}^{(20)} \right) \quad (3.9)$$

Unlike supermolecular approaches for computing the interaction energy, SAPT is considered to be formally BSSE-free, which is one of the major advantages of the approach. To ensure this desirable property of the SAPT interaction energy is maintained, therefore, it is necessary to ensure that the  $\delta_{\text{HF}}^{(2)}$  correction term remain similarly BSSE-free. Therefore,  $\delta_{\text{HF}}^{(2)}$  is typically counterpoise-corrected according to the CP scheme of Boys and Bernardi.<sup>54</sup> The form of the second-order dispersion (and exchange-dispersion) terms is analogous to

the expression for the MP2 correlation energy, and also, therefore, is their expense. While the other terms in SAPT0 require at most for the coupled-perturbed Hartree–Fock (CPHF) equations to be solved, thereby scaling no worse than  $\mathcal{O}(N^4)$ ,  $E_{disp}^{(20)}$  and  $E_{exch-disp}^{(20)}$  scale with total number of occupied orbitals  $o$  and virtual orbitals  $v$  as  $\mathcal{O}(o^3v^2)$  and  $\mathcal{O}(o^2v^3)$ , respectively. Similarly to MP2, however, SAPT benefits greatly from the density fitting approach introduced in Section 2.6.2, with DF-SAPT applicable to systems as large as  $\sim 200$  atoms. Unfortunately, the similarities of  $E_{disp}^{(20)}$  and  $E_{exch-disp}^{(20)}$  to MP2 also extend to their physical accuracy; since MP2 (and the second-order dispersion) are both formulated in terms of excitations within Hartree–Fock monomer densities that are uncoupled from each other, the effects of orbital relaxation on the dispersion energy is neglected. This often leads to marked over-binding in complexes for which dispersion is significant, particularly for  $\pi - \pi$  interactions.

### 3.5.2 Levels of SAPT

In order to further improve the interaction energy from the description afforded by SAPT0, the perturbations corresponding to *intramolecular* correlation,  $\widehat{W}_A$  and  $\widehat{W}_B$ , may be included at second (or higher) order in addition to  $\widehat{V}_{AB}$  to yield various higher-order truncations of SAPT. The performance of these various truncation levels, together with their basis set dependence, was examined extensively by Parker *et al.* for computing total interaction energies.<sup>58</sup> In that work, medalists for SAPT were assigned which best combined accuracy and computational expense in much the same manner as was discussed above for wavefunction methods. In particular, the jun-cc-pVDZ basis set — where the diffuse functions on H atoms and diffuse  $d$  functions on second-row elements are neglected — was determined to be the most favorable basis set for use with the SAPT0 truncation level, due to the fact that the limited diffuse space in the basis set (which would normally under-bind the IE) caused a fortuitous cancellation of errors with the uncoupled dispersion present in SAPT0. This level of theory, SAPT0/jun-cc-pVDZ, was set forth by Parker *et al.* as



the *bronze standard* for SAPT theory, has been widely applied in a variety of chemical contexts ranging from drug intercalation in DNA<sup>66</sup> to understanding the enantioselectivity of organocatalyzed reactions<sup>67,68</sup> and the differential binding of substituted Factor Xa inhibitor drugs to its binding target.<sup>69</sup> Each of these applications, however, leveraged additional partitions of the SAPT energy into contributions from particular pairs of atomic or functional-group contacts, which we will leverage later in this Thesis to our own work on extended chemical systems.

### 3.5.3 Additional SAPT Partitions

While SAPT offers increased chemical intuition into the nature of a particular non-covalent interaction when compared to a simple supermolecular IE, it still offers only total interaction quantities, e.g., the electrostatic attraction between two entire monomers. For interacting molecules which each contain multiple functional groups or residues, however, it would be desirable also to quantify the interaction between a particular pair of atoms or functional groups on opposite monomers. Fortunately, exactly this type of fine-grained partitioning of SAPT energies was recently developed by Parrish *et al.* for SAPT0, termed the atomic<sup>70</sup> and functional-group<sup>66</sup> partitions of SAPT (ASAPT and F-SAPT, respectively). After the introduction of F-SAPT, Parrish and Gonthier introduced a practical formulation of intramolecular SAPT (ISAPT)<sup>71</sup> based on Hartree–Fock embedding, building upon Gonthier’s earlier efforts towards generating appropriate zeroth-order wavefunctions for intramolecular SAPT.<sup>72</sup> These three methods, ASAPT, F-SAPT, and ISAPT,\* together greatly expand the quantum mechanical toolkit for providing detailed insight into the fundamental nature of non-covalent interactions for a variety of applications, and will be used exten-

---

\*The semantic reason behind the inclusion of a hyphen in the abbreviations for F-SAPT and ASAPT/ISAPT is that, in the cases of ASAPT and ISAPT, the differences with conventional SAPT0 are more fundamental, i.e., they appear at the algorithmic level. F-SAPT, on the other hand, is based on the functional group accumulation of atomic-pairwise contacts, so it is rather an *a posteriori* partition of what are essentially ASAPT terms. Of course, the formulation and implementation details of both ASAPT and F/I-SAPT are actually much more complex than this simplified explanation; nevertheless, their abbreviations are determined based on this principle.

sively throughout the rest of this Thesis.

### *Atomic and Functional-Group Partitions of SAPT: ASAPT & F-SAPT*

Both ASAPT and F-SAPT partition the full SAPT0 interaction energy and components computed between monomers  $\mathcal{A}$  and  $\mathcal{B}$  (an order-0 partition) into (i) contributions from fragments of, e.g., monomer  $\mathcal{A}$  interacting with the entirety of monomer  $\mathcal{B}$  (an order-1 partition) and (ii) contributions from unique pairs of a fragment from monomer  $\mathcal{A}$  interacting with a fragment from monomer  $\mathcal{B}$  (an order-2 partition). Within this scheme, each of the total electrostatics, exchange, induction, and dispersion are partitioned by leveraging the iterative stockholder analysis (ISA) localization procedure and the intrinsic locality of the density-fitted two electron integrals. Unfortunately,  $\delta_{\text{HF}}^{(2)}$  has no convenient local representation, and thus remains a correction to the total (order-0) interaction energy. Even though the total component energies are conserved under localization, and the total  $\delta_{\text{HF}}^{(2)}$  correction is included at the order-0 level, both order-1 and order-2 partitions of the SAPT0 energy and components only incorporate an “approximate”  $\delta_{\text{HF}}^{(2)}$  partition based on scaling the total  $\delta_{\text{HF}}^{(2)}$  correction in an even-handed manner. Therefore, order-1 and order-2 partitions of the SAPT0 energy are generally considered semi-quantitative, since no *exact* partition of  $\delta_{\text{HF}}^{(2)}$  is included to correct fragment interactions for the presence of higher-order induction and exchange-induction effects.

In spite of the semi-quantitative nature of the fragment interactions, the analysis provided by both ASAPT and F-SAPT offer significant chemical insight that is lacking even with the conventional formulation of SAPT. These partitions are not without their shortcomings, however. In ASAPT, for example, the electrostatic terms tend to exhibit wild oscillations between interactions of adjacent atoms, due to the proximity of partial atomic charges from the assignment of the partitioned molecular electron density to atoms. F-SAPT, on the other hand, largely removes these inconsistencies by coarse-graining the partition of electron density to functional groups comprised of several atoms. In this method,

however, functional groups must not be defined such that anything but single  $\sigma$ -bonds are “cut” between adjacent fragments,\* as this will result in spurious multipoles.

### *The Intramolecular Formulation of SAPT: ISAPT*

From its earliest inception, the main drawback of SAPT has been that it is formulated to compute the interaction between *exactly* two monomers. In order to study the interactions between two different functional groups belonging to the same molecule (let’s label them  $\mathcal{A}$  and  $\mathcal{B}$ ) — an *intramolecular* non-covalent interaction — a cut-and-cap approach was therefore necessary to fragment the molecule into two separate monomers before SAPT could be applied, even if the interacting moieties were many bonds removed from one another. In the cut-and-cap approach, a  $\sigma$  bond between the fragment(s) of interest and the rest of the molecule (to which we will refer as  $\mathcal{C}$ ) would first be cut, before removing all atoms in  $\mathcal{C}$  and replacing them with a single hydrogen atom. This approach has been justified by arguing that, especially if  $\mathcal{C}$  is a saturated aliphatic chain, the electron densities of the interacting fragments  $\mathcal{A}$ ,  $\mathcal{B}$  are not significantly changed by the presence (or lack thereof) of  $\mathcal{C}$ , and therefore neither is the  $\mathcal{A} \cdots \mathcal{B}$  interaction. For some systems, the cut-and-cap approach may indeed provide a realistic representation of the interactions between fragments of interest,<sup>73</sup> however the appropriateness of this approach in general is an open question in the field and will likely vary on a system-by-system basis.

Recently, however, Parrish and Gonthier developed an intramolecular formulation of SAPT (ISAPT) in which the presence of the connecting backbone  $\mathcal{C}$  is effectively incorporated into the  $\mathcal{A} \cdots \mathcal{B}$  interaction via a HF-in-HF embedding procedure. In their approach, the zeroth-order wavefunctions for  $\mathcal{A}$  and  $\mathcal{B}$  are prepared by relaxing the density of  $\mathcal{A}$  ( $\mathcal{B}$ ) in the presence of  $\mathcal{C}$  so that they each remain fully orthogonal to  $\mathcal{C}$  but no longer orthogonal to each other. A standard SAPT computation performed using these zeroth-order

---

\*For example, the oxygen atom in an epoxide should not be defined as an independent fragment, even though only  $\sigma_{C-O}$  bonds are cut by the fragmentation, since two bonds are severed. Similarly, even aliphatic rings should not be partitioned into different fragments.

wavefunctions for  $\mathcal{A}$  and  $\mathcal{B}$  then *effectively* incorporates the electronic deformation of each fragment by the presence of the linker  $\mathcal{C}$ , thereby quantifying by how much  $\mathcal{C}$  tunes the  $\mathcal{A} \cdots \mathcal{B}$  interaction.

## **PART II**

### **BENCHMARKING NON-COVALENT INTERACTIONS: TOWARDS THE “RIGHT ANSWER FOR THE RIGHT REASONS”**

# CHAPTER 4

## COMPARISON OF EXPLICITLY CORRELATED METHODS FOR COMPUTING HIGH-ACCURACY BENCHMARK ENERGIES FOR NONCOVALENT INTERACTIONS

### 4.1 Abstract

The reliability of explicitly correlated methods for providing benchmark-quality noncovalent interaction energies was tested at various levels of theory and compared to estimates of the complete basis set (CBS) limit. For all systems of the A24 test set, computations were performed using both aug-cc-pVXZ (aXZ; X = D, T, Q, 5) basis sets and specialized cc-pVXZ-F12 (XZ-F12; X = D, T, Q, 5) basis sets paired with explicitly correlated coupled cluster singles and doubles [CCSD-F12n (n = a, b, c)] with triple excitations treated by the canonical perturbative method and scaled to compensate for their lack of explicit correlation [(T<sup>\*\*</sup>)]. Results show that aXZ basis sets produce smaller errors versus the CBS limit than XZ-F12 basis sets. The F12b *ansatz* results in the lowest average errors for aTZ and larger basis sets, while F12a is best for double- $\zeta$  basis sets. When using aXZ basis sets (X  $\geq$  3), convergence is achieved from above for F12b and F12c *ansatz* and from below for F12a. The CCSD(T<sup>\*\*</sup>)-F12b/aXZ approach converges quicker with respect to basis than any other combination, although the performance of CCSD(T<sup>\*\*</sup>)-F12c/aXZ is very similar. Both CCSD(T<sup>\*\*</sup>)-F12b/aTZ and focal point schemes employing density-fitted, frozen natural orbital [DF-FNO] CCSD(T)/aTZ exhibit similar accuracy and computational cost, and both are much more computationally efficient than large-basis conventional CCSD(T) computations of similar accuracy.<sup>†</sup>

---

<sup>†</sup>This Chapter reproduces the work in Ref. 74.

## 4.2 Introduction

Recent improvements have made density functional theory (DFT) and lower-level post-Hartree–Fock wavefunction methods much more accurate for noncovalent interactions (NCI), with mean absolute errors as low as 0.1–0.3 kcal mol<sup>-1</sup> for some of the popular benchmark test sets of small van der Waals dimers.<sup>75–77</sup> Thus, further improvements will require more accurate benchmark values for parameterization and testing. For small molecules (up to around 15 heavy atoms), it seems desirable for benchmark interaction energies (IEs) to be within about 0.1 kcal mol<sup>-1</sup> of the “exact” result (for present purposes defined as the non-relativistic Born–Oppenheimer limit). Complete-basis-set (CBS) extrapolations of coupled-cluster theory through perturbative triple substitutions [CCSD(T)]<sup>10</sup> have been widely regarded as the most appropriate method to estimate this benchmark limit. Residual errors from neglected higher-order electron correlation effects (e.g., quadruple excitations) appear to be very small (on the order of hundredths of one kcal mol<sup>-1</sup> for small systems),<sup>13,14</sup> and core correlation effects can also be neglected as being on this same order.<sup>14,16</sup> On the other hand, incomplete basis sets can easily yield errors of more than 0.1 kcal mol<sup>-1</sup> for CCSD(T) IEs.<sup>57</sup> Hence, for “gold standard” level computations benchmarking noncovalent interactions, it is important to obtain good estimates of the CBS limit.

Unfortunately, direct application of conventional CCSD(T) to NCI is difficult for large molecules and/or large basis sets because of the method’s steep  $\mathcal{O}(o^3v^4)$  computational scaling, where  $o$  and  $v$  are the number of occupied and virtual (unoccupied) molecular orbitals, respectively. One approach to reducing the steep cost of this method is the introduction of focal-point schemes,<sup>23,24</sup> which obtain the CBS limit for a less expensive method like second-order Møller–Plesset perturbation theory (MP2) and then apply a correction for electron correlation neglected by MP2, e.g., by adding  $\delta_{\text{MP2}}^{\text{CCSD(T)}}$ , the difference between MP2 and CCSD(T) as computed in a smaller basis set. This focal-point approach has been widely applied and generally yields at least a “free”  $\zeta$ -level of accuracy in terms

of accuracy vs. computational cost.<sup>25–28,30,78–86</sup> Due to the difference between CCSD(T) and MP2 energies being already relatively constant across basis sets, however, these corrections are not necessarily systematically improvable; indeed they often exhibit small oscillations and lack of definite convergence for basis sets beyond aug-cc-pVTZ (aTZ).<sup>57</sup>

By including explicit dependence upon the interelectronic distance, “explicitly correlated” methods (see Refs. 87–89 and references therein) can greatly accelerate convergence towards the CBS limit. In particular, the explicitly correlated CCSD-F12 methods<sup>90</sup> and the associated F12a,<sup>91</sup> F12b,<sup>92</sup> and (F12\*)<sup>93</sup> [also referred to as F12c]<sup>94</sup> approximations, have gained significant traction. These popular F12 approximations have not been extended to triple excitations; that is, CCSD(T)-F12a, etc., apply explicit correlation only to the single and double substitutions. As a practical way to remedy this deficiency, Werner and co-workers<sup>92,95</sup> introduced a scaling approach that assumes explicit correlation should magnify the triples contribution to the correlation energy by the same ratio as the magnification of the MP2 doubles contribution:

$$E_{\text{corr}}^{(\text{T}^*)} = E_{\text{corr}}^{(\text{T})} \frac{E_{\text{corr}}^{\text{MP2-F12}}}{E_{\text{corr}}^{\text{MP2}}}. \quad (4.1)$$

In order to preserve size-consistency for noncovalent interaction energies in weakly bound dimers, it is helpful to use the same scaling ratio for each contribution (monomer A, monomer B, and the dimer);<sup>95</sup> we denote this use of the dimer ratio for each component by (T\*\*).

Marchetti and Werner<sup>95</sup> showed that CCSD(T\*\*)-F12a/aug-cc-pVDZ achieves errors of less than 0.2 kcal mol<sup>-1</sup> vs. CCSD(T)/CBS for all members of the popular S22 test set,<sup>79</sup> despite the small basis set used. Later comparisons against revised values<sup>30</sup> for the S22 test set showed mean absolute errors of about 0.1 kcal mol<sup>-1</sup> for both F12a and F12b variants of CCSD(T\*\*) in the modest aug-cc-pVDZ basis.<sup>30</sup> The small remaining errors in CCSD(T\*\*)-F12a and CCSD(T\*\*)-F12b exhibit varying behavior for different systems; in



an aug-cc-pVDZ basis set, CCSD(T<sup>\*\*</sup>)-F12a tends to do best for hydrogen-bonded systems, while CCSD(T<sup>\*\*</sup>)-F12b tends to do best for dispersion-bound or mixed systems. This led to the introduction of a “dispersion weighted” (DW) CCSD(T<sup>\*\*</sup>)-F12 that mixes F12a and F12b results to best effect<sup>96</sup> (analogous to the dispersion-weighted MP2-F12 method of Marchetti and Werner,<sup>95</sup> which mixes MP2-F12 with spin-component-scaled MP2-F12). DW-CCSD(T<sup>\*\*</sup>)-F12/aug-cc-pVDZ exhibits a mean absolute error of only 0.05 kcal mol<sup>-1</sup> for the S22B test set<sup>96</sup> and may be useful as a “silver standard” in benchmarking NCI,<sup>77</sup> with errors nearly as small as CBS extrapolations of conventional CCSD(T) (or focal-point estimates of the same), but with much smaller basis set requirements.

Recently, the question of applying explicitly correlated methods to noncovalent interactions has received renewed attention.<sup>97–99</sup> In particular, Patkowski<sup>100</sup> compared several approximate CCSD(T)-F12 methods and basis sets (with and without midbond functions) for computing the IE of a few rare gas, water-methane, and water dimers; Martin and coworkers<sup>101</sup> sought to investigate the role of basis set superposition error within explicitly correlated methods; and Schmitz *et al.*<sup>102</sup> examined the accuracy of a novel pair natural orbital (PNO) formulation of the MP2-F12, CCSD(2)<sub>F12</sub>,<sup>103</sup> and CCSD[F12]<sup>93</sup> methods. Additionally, efforts to revise the reference values in a number of popular test sets for noncovalent interactions or extend them to include more systems have been undertaken utilizing such explicitly correlated methods. Most notably, this includes revisions of the S66 test set<sup>76,104,105</sup> (partially by Werner and coworkers<sup>91</sup> and then more completely by Tew and coworkers<sup>102</sup>) and the S66×8 test set<sup>105,106</sup> by Martin and coworkers,<sup>Brauer:2016:xxx</sup> as well as the extension of the S22×5 test set<sup>107</sup> by Smith *et al.*<sup>39</sup> to include two additional points at 0.7 and 0.8· $R_e$ , dubbed S22×7. The very high quality of approximate CCSD(T)-F12 methods shown in those studies for describing noncovalent interactions raises the question of whether these approaches, even with modest basis sets, might be suitable replacements for current “gold standard” benchmark procedures based on conventional CCSD(T). Hence, an improved understanding of the convergence behavior of these CCSD(T)-F12 approxi-

mations for NCI, as well as their performance at given basis set levels, would be helpful.

To this end, we seek here to present a complete, systematic study of the basis set convergence of both CCSD(T) and explicitly correlated CCSD(T)-F12 with the F12a, F12b, and F12c approximations for the IEs of several small bimolecular complexes. In order to make generalized recommendations about the application of these methods for benchmarking NCI in arbitrary systems, we have chosen to examine the A24 test set,<sup>14</sup> which provides a diverse set of systems still small enough to perform benchmark computations in very large basis sets, and the S22 test set<sup>79</sup> with revised reference values<sup>30</sup> which tests the efficacy of our recommendations for systems of larger size and stronger interaction. We employ the aug-cc-pVXZ (X = D, T, Q, 5)<sup>22</sup> and cc-pVXZ-F12 (X = D, T, Q, 5)<sup>108,109</sup> basis sets. The latter basis set family was specialized for use with explicitly correlated methods; X = D, T, Q were optimized for use with MP2-F12(3C),<sup>108</sup> while X = 5 was optimized for the CCSD(T)-F12b method.<sup>109</sup> Additionally, we will compare both the accuracy and computational expense of the best of these explicitly correlated methods, as well as existing approaches, in order to determine which (if any) should be considered as viable alternatives to the currently recommended focal-point prescriptions for benchmark-quality interaction energies.

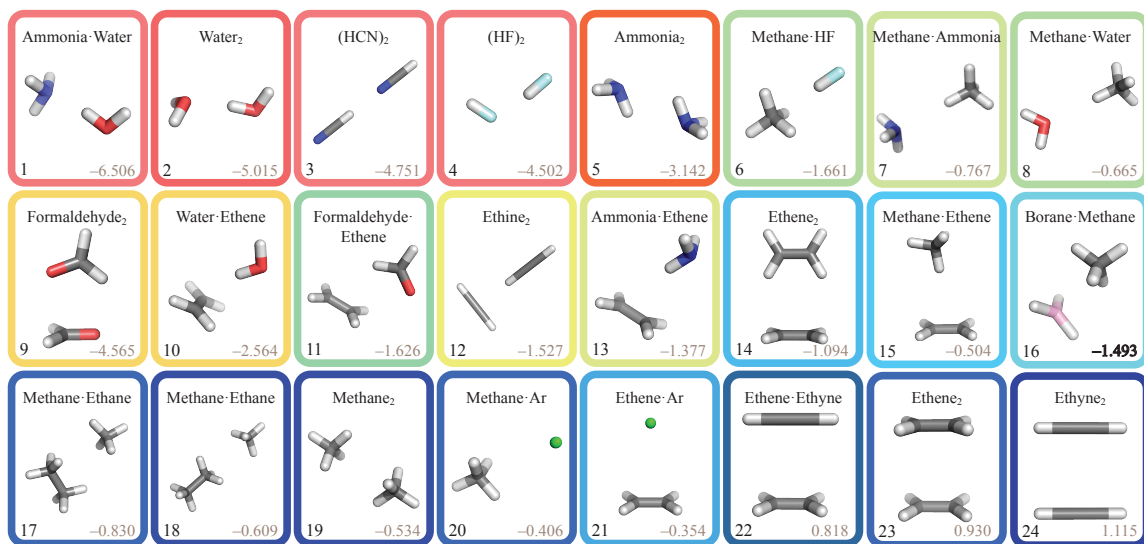
### 4.3 Theoretical & Computational Methods

#### 4.3.1 Overview of Approximate CCSD(T)-F12 Methods.

Here, we present an overview of the approximate CCSD(T)-F12n methods examined in this work; we direct the interested reader to Refs. 87–95 and the references therein for a more thorough discussion. We will adopt the Einstein convention for which any repeated indices are summed over. The form of the CCSD-F12 wavefunction is given by

$$|\Psi\rangle_{\text{CCSD-F12}} = e^{\hat{T}_1 + \hat{T}_2} |\Phi\rangle_{\text{HF}} \quad (4.2)$$

### (a) A24



### (b) S22

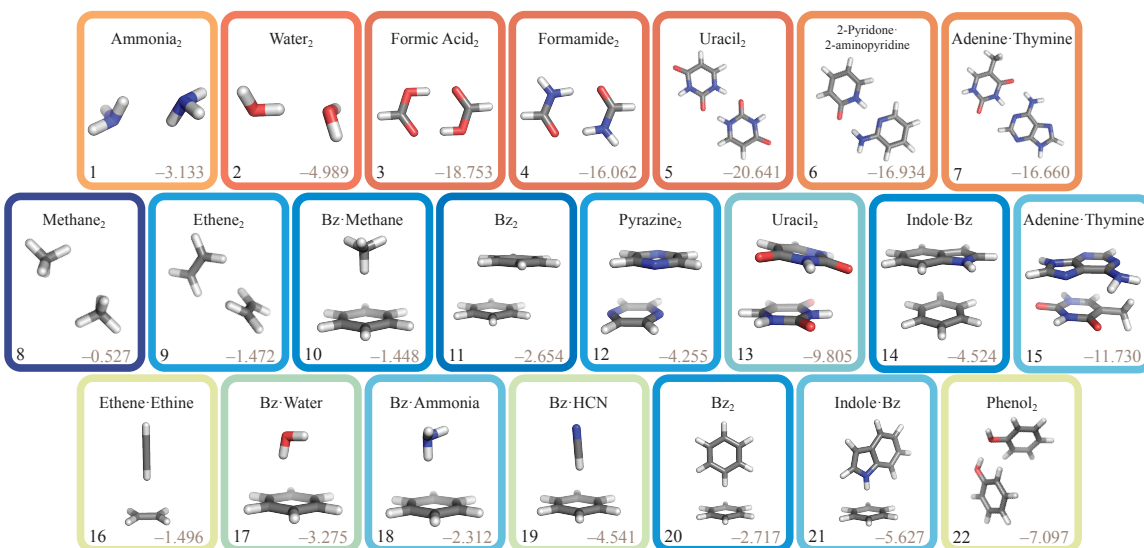


Figure 4.1: Bimolecular complexes included in (a) the A24 and (b) S22 test sets of Hobza and coworkers<sup>14,79</sup> with revised A24B and S22B<sup>30</sup> reference energies, in kcal mol<sup>-1</sup>. Coloring is based on SAPT2+(3)/aTZ results reported previously by Burns *et al.*<sup>57,77</sup> and indicates interaction type: red for electrostatically dominated interactions (typically hydrogen bonding), blue for dispersion dominated interactions, and yellow-green for interactions of mixed character.

with the cluster operators  $\hat{T}_1$  and  $\hat{T}_2$  defined by

$$\hat{T}_1 = t_a^i \hat{E}_i^a, \quad (4.3)$$

$$\hat{T}_2 = T_{ab}^{ij} \hat{E}_{ij}^{ab} + \mathcal{T}_{\alpha\beta}^{ij} \hat{E}_{ij}^{\alpha\beta}, \quad (4.4)$$

for occupied orbital indices  $i, j$ , and virtual orbital indices  $a, b$ .  $\alpha, \beta$  denote virtual orbitals in a complete basis set. In these expressions,  $\hat{E}_i^a$  and  $\hat{E}_{ij}^{ab} = \hat{E}_i^a \hat{E}_j^b$  are spin-free one and two electron excitation operators and  $t_a^i$  and  $T_{ab}^{ij}$  are the canonical singles and doubles amplitudes from CCSD theory.

The amplitudes  $\mathcal{T}_{\alpha\beta}^{ij}$  represent additional configurations not spanned by the primary orbital basis, and they are given by

$$\mathcal{T}_{\alpha\beta}^{ij} = T_{mn}^{ij} \mathcal{F}_{\alpha\beta}^{mn}, \quad (4.5)$$

$$\mathcal{F}_{\alpha\beta}^{mn} = \langle mn | F_{12} \hat{Q}_{12} | \alpha\beta \rangle. \quad (4.6)$$

The short-range correlation factor,  $F_{12} = f(r_{12})$ , is commonly given by the Slater function

$$f(r_{12}) = e^{-\beta r_{12}}, \quad (4.7)$$

which accurately describes the nature of the wavefunction near the interelectronic cusp, i.e., when  $r_{12} \rightarrow 0$ . The optimal value of the geminal parameter,  $\beta$ , will typically depend on the system and basis set;  $\beta$  can therefore be optimized for each individual computation, but more commonly it is set *a priori* to a fixed value (the approach taken in this work). The F12 amplitudes,  $T_{mn}^{ij}$ , are often taken to be diagonal, such that  $T_{ii}^{ii} = t_s$ ,  $T_{ij}^{ij} = \frac{1}{2}[t_s + t_t]$  (for  $i = j$ ),  $T_{ji}^{ij} = \frac{1}{2}[t_s - t_t]$  (for  $i \neq j$ ), and all other amplitudes  $T_{kl}^{ij} = 0$ . Additionally, the values of the parameters  $t_s$  and  $t_t$  are commonly taken to be fixed:

$$t_s = -\frac{1}{2\beta}, \quad t_t = -\frac{1}{4\beta}, \quad (4.8)$$

with  $\beta$  the geminal parameter from Eqn. 4.7; these values are from the cusp conditions for  $s$  and  $p$  functions. This diagonal, fixed *ansatz* [referred to as 3C(FIX)] is both size consistent and orbital invariant, and is employed by each of the CCSD-F12a, CCSD-F12b, and CCSD(F12\*) methods. Finally, the projector  $\widehat{Q}_{12}$ , given by

$$\widehat{Q}_{12} = (1 - \hat{o}_1)(1 - \hat{o}_2)(1 - \hat{v}_1\hat{v}_2) \quad (4.9)$$

where  $\hat{o}$  and  $\hat{v}$  project onto the occupied and virtual spaces, respectively, ensures strong orthogonality between the F12 states  $|\Phi_{ij}^{mn}\rangle = \mathcal{F}_{\alpha\beta}^{mn} \widehat{E}_{ij}^{\alpha\beta} |\Phi\rangle_{\text{HF}}$  and the configurations spanned by the primary orbital basis.

The CCSD-F12a and CCSD-F12b methods are very similar, approximating the projector  $\widehat{Q}_{12}$  in equation 4.9 as  $\widehat{Q}_{12} = 1 - |rs\rangle\langle rs|$  (with  $r, s$  full MO basis indices). The only difference between the two is that the CCSD-F12b method contains an extra energy correction within the coupled cluster residual that increases the coupling between conventional and explicitly correlated parts of the wavefunction; this means that CCSD-F12b more effectively constrains the orthogonality between F12 states and orbital states. Therefore, in the F12a approximation, increased variational flexibility from this reduced orthogonality (relative to F12b) may be free to cancel with basis set incompleteness errors (BSIE) in small basis sets to yield artificially accurate results. On the other hand, the F12b approximation yields no such artificial flexibility; therefore, CCSD-F12b may be more susceptible to BSIE and exhibit less accuracy in small basis sets than CCSD-F12a.

The last of these approximate methods, namely CCSD(F12\*) (a.k.a. F12c) is instead based on the CCSD(F12) method,<sup>110</sup> which neglects higher order terms involving the F12 amplitudes  $\mathcal{T}_{\alpha\beta}^{ij}$  within the fixed-amplitude *ansatz*. (F12\*) further approximates the amplitude equations by first ignoring any contributions that are fourth order (or higher) in the MP2-F12 treatment, then adding the most important higher-order coupling contributions

between orbital and F12 states resulting from the less approximate form of the projector,

$$\widehat{Q}_{12}^{(1)} = 1 - \widehat{P}_1 \widehat{P}_2 - \widehat{P}_1 \widehat{P}'_2 - \widehat{P}'_1 \widehat{P}_2,$$

where  $\widehat{P} = \widehat{o} + \widehat{v}$  with  $\widehat{o}$  and  $\widehat{v}$  projecting onto the occupied and virtual spaces, respectively, and  $\widehat{P}'$  projecting onto the complimentary virtual space. Due to the less approximate projector employed, the CCSD(F12\*) method is more theoretically rigorous than either CCSD-F12a or CCSD-F12b methods; it is, however, more computationally intensive due to the need to evaluate additional intermediates.<sup>93</sup> Again, due to the stronger orthogonality between F12 and orbital states, BSIE may be present in small basis sets, but this method should yield high accuracy with increasing basis set size.

#### 4.3.2 Computational Details

In this study, we compute the interaction energies (IEs) of all complexes in A24<sup>14</sup> and S22<sup>79</sup> test sets via the supermolecular approach, whereby the total energy for each monomer is subtracted from the total energy of the dimer to form the total IE. In order to correct for basis set superposition error (BSSE), we employ the counterpoise correction scheme of Boys and Bernardi.<sup>54</sup> Computations were performed using each combination of CCSD(T\*\*) - F12n (abbreviated F12n; n = a, b, c) method paired with either aug-cc-pVXZ (abbreviated aXZ; X = D, T, Q, 5)<sup>22</sup> or cc-pVXZ-F12 (abbreviated XZ-F12; X = D, T, Q, 5)<sup>108,109</sup> basis sets. Triples contributions were scaled according to equation (4.1), with the scale factor determined for the dimer also used for the monomer computations to preserve size consistency,<sup>95</sup> denoted (T\*\*). All computations using these explicitly correlated methods were performed using the MOLPRO 2010.1 suite of *ab initio* quantum chemistry programs,<sup>111</sup> and employed the complete auxiliary basis set (CABS)<sup>112</sup> singles correction.<sup>91</sup> Both the density fitting (DF) and resolution of the identity (RI) basis sets were kept at their MOLPRO default values: aug-cc-pVXZ/MP2FIT was used for the overall density fitting ba-

sis (keyword `DF_BASIS`) and `cc-pVXZ/JKFIT` was used for the computing the exchange and Fock operators (keyword `DF_BASIS_EXCH`) for both `aXZ` and `XZ-F12` orbital basis sets, while `cc-pVXZ/JKFIT` or `cc-pVXZ/OPTRI` were used for the RI basis sets (keyword `RI_BASIS`) for `aXZ` and `XZ-F12` orbital basis sets, respectively. According to the recommendation of Patkowski,<sup>100</sup> the F12 geminal parameter  $\beta$  (given above in equation 4.7) was kept at its MOLPRO default value of  $\beta = 1.0 a_0^{-1}$ .

As discussed above, in Section 4.3.1, the approximate F12a, F12b, and F12c methods become increasingly more physically justified in the sequence  $a < b < c$ . Each of the basis set types examined here have been employed previously to compute noncovalent interactions (see Ref. 77 and the references therein, as well as Refs. 101, 102, 113 for more details); these basis sets are not equivalent, however, and should therefore not be expected to attain equivalent accuracy. For instance, the specialized `cc-pVXZ-F12` basis sets are larger than their canonical `aug-cc-pVXZ` counterparts by virtue of containing several more *s* and *p* functions, while remaining as diffuse in these low angular momentum basis functions. However, the `aug-cc-pVXZ` basis sets are more diffuse in the higher angular momentum functions (*d*, *f*, ...) than `cc-pVXZ-F12`.

To assess the various model chemistries (combinations of method and basis set)<sup>114</sup> considered here, we obtained “best estimate” reference IEs using CBS extrapolations of conventional CCSD(T) correlation energies  $E_{\text{corr}}$  according to the popular two-point extrapolation formula of Helgaker and co-workers.<sup>20</sup> As suggested previously,<sup>21</sup> we take the total Hartree-Fock energy computed in the larger of these two basis sets as being converged to the CBS limit.

For this work, we use the `aQZ` and `a5Z` basis sets in our extrapolation procedure for constructing our reference IEs, which we denote `CCSD(T)/CBS(aQZ, a5Z)`. These data were previously reported<sup>57</sup> for the A24 test set,<sup>14</sup> along with `CCSD(T)/CBS(a5Z, a6Z)` extrapolations for a few of the smaller members of the test set, namely the  $\text{H}_2\text{O}\cdot\text{NH}_3$ ,  $\text{H}_2\text{O}\cdot\text{H}_2\text{O}$ ,  $\text{HCN}\cdot\text{HCN}$ ,  $\text{HF}\cdot\text{HF}$ ,  $\text{NH}_3\cdot\text{NH}_3$ ,  $\text{CH}_4\cdot\text{HF}$ ,  $\text{CH}_4\cdot\text{H}_2\text{O}$ ,  $\text{CH}_4\cdot\text{Ar}$ , and  $\text{CH}_2\text{CH}_2\cdot\text{Ar}$  complexes.

In that prior study, various counterpoise correction schemes were compared against a weighted average of counterpoise-corrected and uncorrected energies; this reference was denoted A24A to distinguish these values from the originally published benchmarks.<sup>14</sup> Because we apply counterpoise correction universally to all explicitly correlated methods in the present study, for consistency we also use counterpoise corrected CCSD(T)/CBS benchmark values, which we denote here as A24B (although our previous study<sup>57</sup> indicates that the counterpoise treatment matters very little once basis sets as large as a5Z and a6Z are employed). Recently, Hobza and co-workers have presented<sup>115</sup> even more accurate interaction energies for A24, by using larger basis sets for estimating the importance of quadruple excitations, which might be denoted as A24C. In this study, however, we limit ourselves to consideration only of the CCSD(T)/CBS limit, without quadruples corrections; hence, we have adopted the A24B energies as reference values.

The accuracy of each model chemistry is characterized by mean absolute errors (MAE) and mean absolute percent errors (MA%E) across all systems in the A24 and S22 test sets, and also across each of three subsets of systems grouped by interaction type: hydrogen bonding, dispersion dominated, and mixed interaction (whereby the interactions have both electrostatics and dispersion character). Additionally, we have reported values for a full set of summary statistics, including each of the minimal, maximal, and mean (both signed and unsigned) errors and percent errors for each model chemistry over the A24 and S22 test sets, in the Supplemental Information (see Section ??). For reference, these databases are visualized in Fig. 4.1.

#### 4.4 Results and Discussion

Before comparing with existing high-quality canonical schemes, we first examine in Section 4.4.1 the convergence of each explicitly correlated CCSD(T<sup>\*\*</sup>)-F12n method towards the complete basis set limit. Interaction energies computed using aXZ basis sets tend to converge to reference energies differently than those computed with an XZ-F12 basis



set; results will therefore be presented separately for each basis set type in Sections 4.4.1 and 4.4.1, respectively. Next, in Section 4.4.2, the performance of the various model chemistries for the A24 test set will be compared directly to determine which are superior for NCI at each basis set level. Then, in Section 4.4.3, the explicitly correlated CCSD(T<sup>\*\*</sup>)-F12n methods will be applied to the S22 test set to determine if our findings for the small systems in A24 can be generalized to slightly larger systems and interaction strength. Finally, in Section 4.4.4, the accuracy and computational expense of several approaches will be compared in order to assess their potential candidacy for benchmark quality computations, and final recommendations will be given.

#### 4.4.1 Basis Set Convergence for A24 Systems

Figs. 4.2 and 4.3 show the convergence of interaction energies computed with various CCSD(T<sup>\*\*</sup>)-F12n methods [along with conventional CCSD(T)] with respect to aXZ and XZ-F12 basis sets, for the ammonia-water complex, formaldehyde-ethylene complex, ethylene dimer, and the methane-argon complex (additionally, Fig. 4.3.e shows the methane dimer). MAEs and MA%Es for all computations are presented in Table 4.1. For a complete set of figures detailing method convergence for all members of the A24 test set, as well as raw interaction energies and summary statistics for every combination of method and basis set considered, refer to the Supplementary Information (see Section ??).

##### *F12n/aXZ Convergence*

A notable difference in the convergence behavior of A24 systems with different binding motifs exists for model chemistries involving aXZ-type basis sets. For hydrogen bonded (HB) systems (e.g., the ammonia-water complex, Fig. 4.2.a), F12b and F12c methods converge to the best-estimate energies from above. In contrast, F12a methods, although always above the best-estimate IE for aDZ, converge to the reference energy from below for aTZ and larger. Because the F12 methods are so rapidly convergent to the CBS limit, for

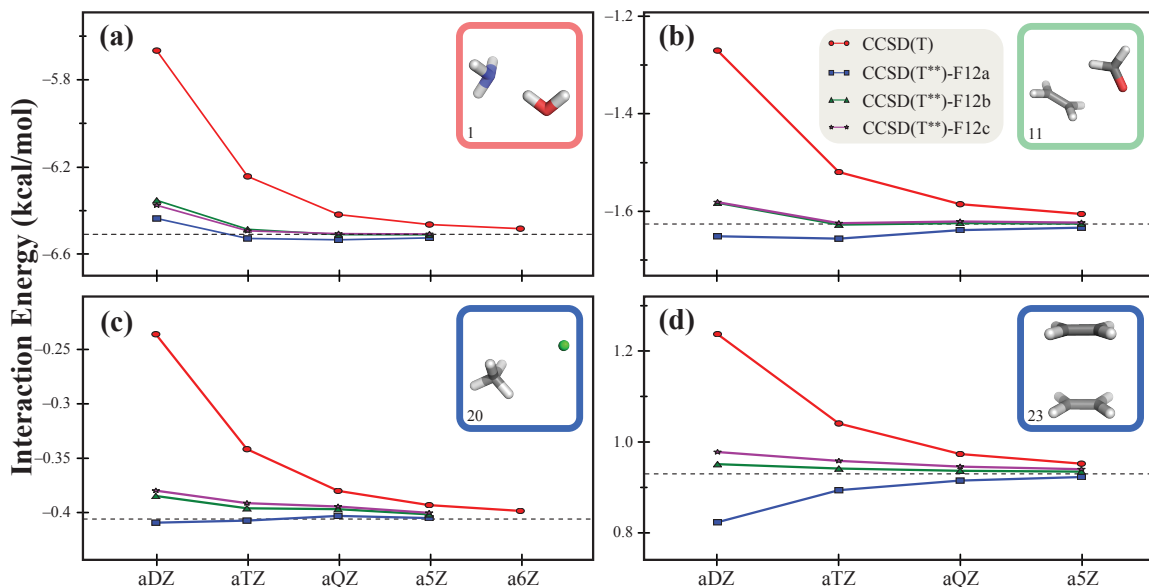


Figure 4.2: Convergence of CCSD(T<sup>\*\*</sup>)-F12n/aXZ ( $n = a, b, c$ ;  $X = D, T, Q, 5, 6$ ) IEs for (a) the ammonia-water complex, (b) the formaldehyde-ethylene complex, (c) the methane-Ar complex, and (d) the ethylene dimer in forced  $\pi$ -stacking geometry. Also plotted are canonical CCSD(T)/aXZ IEs and our revised A24B reference energies (dotted line) obtained at the CCSD(T)/CBS(aQZ, a5Z) [(a) & (c)] or CCSD(T)/CBS(a5Z, a6Z) [(b) & (d)] levels of theory (see text).

these small systems it is feasible to ask what basis set is required to achieve nearly exact convergence to the CBS limit (within, say, 0.01 kcal/mol). We will denote this level of agreement “benchmark convergence” for this paper; however, it is worth bearing in mind that other sources of error (quadruple excitations, core correlation) are several hundredths of one kcal/mol for systems of this size, and hence for “benchmark quality” results we could tolerate somewhat larger errors in the CBS convergence. Nevertheless, for the A24 hydrogen-bonding systems, this strict level of convergence is achieved with just the modest aTZ basis set for both F12a and F12c, while F12b is nearly as good in this basis (MAE 0.02 kcal/mol).

For systems of mixed (MX) interaction type (e.g., the formaldehyde-ethane complex, Fig 4.2.b), convergence patterns for F12b/aXZ and F12c/aXZ model chemistries show similar behavior to hydrogen bonded systems. Unlike for HB systems, F12a/aDZ is sometimes above the reference interaction energy (methane-HF complex, water-ethylene complex,

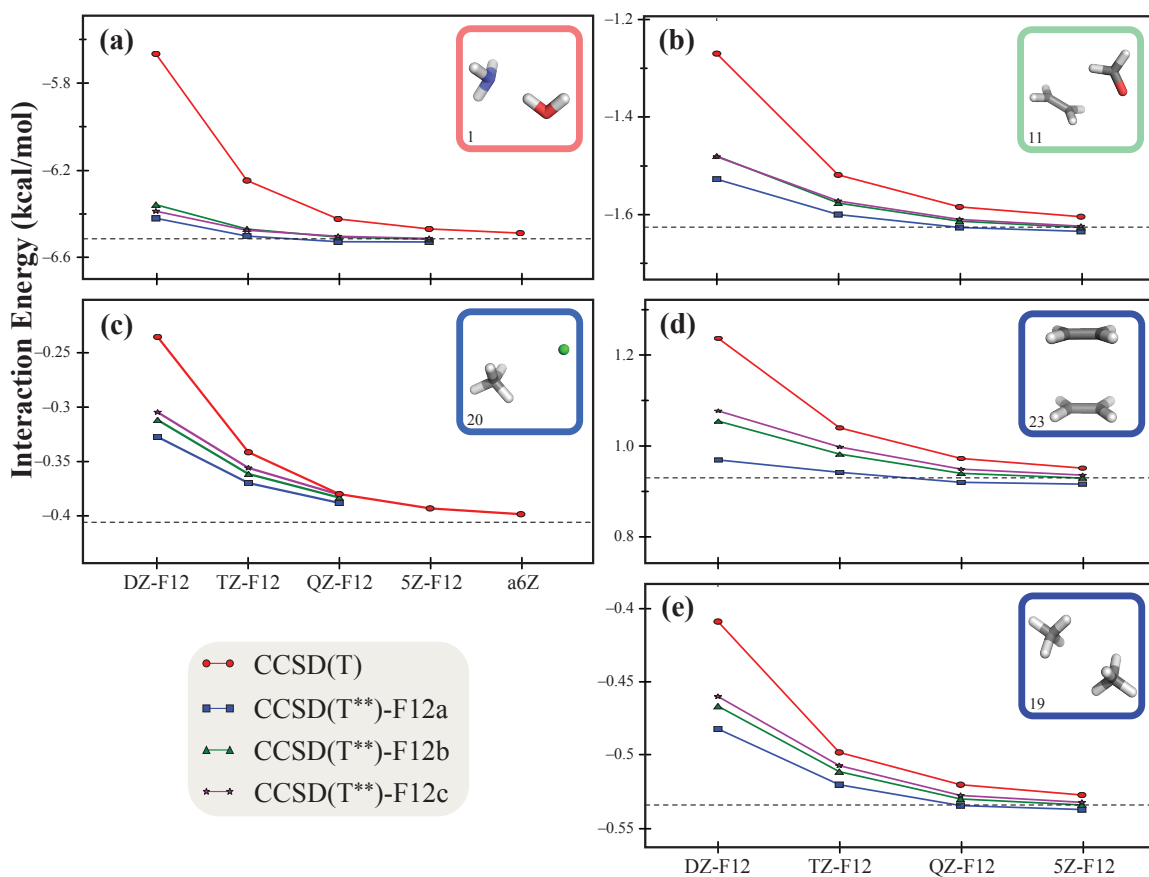


Figure 4.3: Convergence of CCSD(T<sup>\*\*</sup>)-F12n/XZ-F12 (n = a, b, c; X = D, T, Q, 5) IEs for the (a) ammonia-water complex, (b) formaldehyde-ethylene complex, (c) methane-Ar complex, (d) ethylene dimer in forced  $\pi$ -stacking geometry, and (e) methane dimer. Also plotted are canonical CCSD(T)/aXZ IEs and our revised A24B reference energies (dotted line) obtained at the CCSD(T)/CBS(aQZ, a5Z) [a & c] or CCSD(T)/CBS(a5Z, a6Z) [b, d, & e] levels of theory (see text).

ethylene dimer) and sometimes below (methane-ammonia complex, methane-water complex, formaldehyde dimer, formaldehyde-ethylene complex, ammonia-ethylene complex). For mixed systems, F12a/aXZ and F12c/aXZ converge similarly as for hydrogen bonded systems, while F12b is now a little more tightly converged in the aTZ basis than F12a.

For dispersion-dominated (DD) systems of A24 (e.g., the  $\pi$ -stacked ethylene dimer and methane-argon complex, Figs. 4.2c,d respectively), the convergence behavior of F12a is notably different than for HB and MX systems, while F12b and F12c behave similarly. F12a typically converges toward the reference IE monotonically from below, shown in Fig. 4.2.d. Two notable exceptions to this behavior are those complexes containing a rare gas: the methane-argon complex (Fig. 4.2.c) and ethylene-argon complex (Fig. S-30 in the SI, see Section ??) These interaction energies oscillate across the reference IE, while still converging towards it with increasing  $\zeta$ . For DD systems, all F12 methods converge very rapidly with basis set, and achieve MAE of 0.01-0.02 kcal/mol by the aTZ basis set.

Table 4.1: Interaction energy (kcal/mol) error statistics vs CCSD(T)/CBS for F12n/aXZ and F12n/XZ-F12 methods, gathered for all A24 systems, as well as for the hydrogen bonding (HB), mixed interaction (MX) and dispersion dominated (DD) subsets.

Method & Basis Set	A24		HB		MX		DD	
	MAE	MA %E	MAE	MA %E	MAE	MA %E	MAE	MA %E
<b>aug-cc-pVDZ</b>								
CCSD(T)	0.38	24.5	0.61	12.8	0.43	24.6	0.21	30.1
DW-CCSD(T**)-F12	0.05	3.8	0.06	1.1	0.05	2.9	0.04	5.9
CCSD(T**)-F12a	0.04	3.0	0.05	1.1	0.03	2.3	0.04	4.5
CCSD(T**)-F12b	0.06	3.6	0.11	2.2	0.07	3.9	0.02	4.1
CCSD(T**)-F12c	0.06	4.1	0.08	1.7	0.07	3.9	0.04	5.6
<b>aug-cc-pVTZ</b>								
CCSD(T)	0.12	7.5	0.22	4.6	0.12	6.4	0.07	9.9
DW-CCSD(T**)-F12	0.01	1.1	0.01	0.2	0.01	0.8	0.01	1.7
CCSD(T**)-F12a	0.02	1.4	0.01	0.3	0.02	1.2	0.02	2.1
CCSD(T**)-F12b	0.01	0.6	0.02	0.4	0.01	0.3	0.01	1.1
CCSD(T**)-F12c	0.01	1.1	0.01	0.3	0.01	0.5	0.01	2.2
<b>aug-cc-pVQZ</b>								
CCSD(T)	0.04	2.8	0.07	1.5	0.04	2.3	0.03	4.0
CCSD(T**)-F12a	0.01	0.7	0.01	0.3	0.01	0.7	0.01	1.0
CCSD(T**)-F12b	0.00	0.4	0.00	0.0	0.00	0.1	0.00	0.8
CCSD(T**)-F12c	0.01	0.7	0.00	0.1	0.01	0.3	0.01	1.3
<b>aug-cc-pV5Z</b>								
CCSD(T)	0.02	1.4	0.03	0.7	0.02	1.2	0.01	2.0
CCSD(T**)-F12a	0.01	0.4	0.01	0.2	0.01	0.4	0.00	0.5
CCSD(T**)-F12b	0.00	0.2	0.00	0.0	0.00	0.1	0.00	0.4
CCSD(T**)-F12c	0.00	0.4	0.00	0.0	0.00	0.2	0.00	0.7
<b>cc-pVDZ-F12</b>								
CCSD(T**)-F12a	0.07	6.2	0.09	1.9	0.08	4.9	0.06	9.7
CCSD(T**)-F12b	0.12	9.5	0.13	2.7	0.13	7.2	0.10	15.1
CCSD(T**)-F12c	0.11	10.0	0.11	2.3	0.12	7.1	0.11	16.4
<b>cc-pVTZ-F12</b>								
CCSD(T**)-F12a	0.02	1.8	0.02	0.4	0.02	0.9	0.02	3.3
CCSD(T**)-F12b	0.04	3.5	0.04	0.9	0.04	2.1	0.04	6.0
CCSD(T**)-F12c	0.04	4.0	0.04	0.8	0.04	2.3	0.05	7.1
<b>cc-pVQZ-F12</b>								
CCSD(T**)-F12a	0.01	0.7	0.01	0.2	0.00	0.2	0.01	1.5
CCSD(T**)-F12b	0.01	1.0	0.01	0.2	0.01	0.5	0.01	1.9
CCSD(T**)-F12c	0.01	1.4	0.01	0.2	0.01	0.8	0.01	2.6
<b>cc-pV5Z-F12</b>								
CCSD(T**)-F12a	0.01 <sup>a,b,c</sup>	0.6	0.01 <sup>a</sup>	0.1	0.01 <sup>b</sup>	0.5	0.01 <sup>c</sup>	1.0
CCSD(T**)-F12b	0.00 <sup>a,b,c</sup>	0.1	0.00 <sup>a</sup>	0.0	0.00 <sup>c</sup>	0.1	0.00 <sup>c</sup>	0.2
CCSD(T**)-F12c	0.00 <sup>a,b,c</sup>	0.2	0.00 <sup>a</sup>	0.0	0.00 <sup>b</sup>	0.1	0.00 <sup>c</sup>	0.3

<sup>a</sup> Missing IE for ammonia dimer, <sup>5zf12\_LinDep</sup> <sup>b</sup> Missing IE for ethylene-water complex, <sup>5zf12\_LinDep</sup>

<sup>c</sup> Missing IEs for methane-Ar and ethylene-Ar complex, due to no 5Z-F12 basis set existing for Ar.

Considering MAE for the entire A24 test set, all of the explicitly correlated CCSD(T<sup>\*\*</sup>)-F12 methods converge within 0.02 kcal/mol of the benchmark values by the aTZ basis set. Overall, F12b and F12c always converge to the CBS limit IE from above, while F12a generally converges from below. While the (counterpoise-corrected) canonical CCSD(T)/aXZ IEs steadily converge from above, they can still fall short of the CBS limit even with the a5Z basis set [with differences between (0.006, 0.069) kcal mol<sup>-1</sup>, MAE = 0.021]. Indeed, even the CCSD(T)/a6Z values can fall short of the reference CCSD(T)/CBS(a5Z, a6Z) values by between (0.004, 0.026) kcal mol<sup>-1</sup>, with MAE = 0.014. In contrast, the F12n results are much more rapidly convergent. Results in the aTZ basis for F12a are within (-0.039, -0.001) kcal mol<sup>-1</sup> of the reference energy, with MAE = 0.017, which is superior than the canonical CCSD(T)/a5Z results and nearly as good as CCSD(T)/a6Z. F12c/aTZ is on par with CCSD(T)/a6Z, with energies falling between (0.001, 0.029) kcal mol<sup>-1</sup> of the reference with MAE = 0.012; remarkably, while F12b/aTZ exhibits a comparable error range to CCSD(T)/a6Z [energies within (0.000, 0.029) kcal mol<sup>-1</sup> of the reference], it *exceeds* the convergence of canonical CCSD(T)/a6Z on average, with MAE = 0.009.

Moreover, the aDZ results are not much worse. Energies computed in the aDZ basis exhibit error ranges and MAEs of [(-0.107, 0.002); 0.040], [(0.001, 0.156); 0.059], and [(0.004, 0.154); 0.059] kcal mol<sup>-1</sup> for F12a, F12b, and F12c respectively. These tests suggest that, at least for NCI, one obtains on average two additional  $\zeta$ -levels of accuracy when moving from conventional CCSD(T) to CCSD(T<sup>\*\*</sup>)-F12a, and *three* additional  $\zeta$ -levels when moving to CCSD(T<sup>\*\*</sup>)-F12b and CCSD(T<sup>\*\*</sup>)-F12c. This is a significant advantage, given that the computational cost of CCSD(T)-F12n methods is typically not much more (on average, within a factor of 1.3) than that of conventional CCSD(T).<sup>TimingsNotes</sup>

### *F12n/XZ-F12 Convergence*

Again, different convergence behavior between F12n methods is exhibited for A24 systems of different binding motif when paired with XZ-F12 basis sets. For electrostatics-

dominated HB systems (e.g., ammonia-water complex, Fig. 4.3.a), F12a remains the most converged of the three *ansatz* for the double- $\zeta$  DZ-F12 basis set, although the improvement over F12b and F12c is very small. F12a remains very slightly better than F12b and F12c on average for the TZ-F12 basis, and results are essentially identical for the QZ-F12 basis. These methods converge toward the reference IE from above. Systems of mixed character (e.g., the formaldehyde-ethene complex, Fig. 4.3.b), exhibit similar convergence as HB systems for F12n/XZ-F12 IEs. Again, F12a exhibits slightly smaller MAE than F12b or F12c for the DZ-F12 and TZ-F12 basis sets, with nearly exact agreement achieved by the QZ-F12 basis.

For dispersion-bound systems, IEs computed with F12n/XZ-F12 display different convergence behavior depending on the complex. For DD systems containing rare gases [e.g., the methane-Ar complex (Fig. 4.3.c) and ethylene-Ar complex, Fig. S-54 in the Supplemental Information, (see Section ??)], while all explicitly correlated *ansatz* are significantly more tightly converged than canonical CCSD(T) for DZ-F12, canonical CCSD(T) converges toward the reference energy quickly enough to nearly overtake them by QZ-F12. Unfortunately, no interaction energy for these complexes could be computed for 5Z-F12, because this basis set is as yet unavailable for third row atoms. DD systems involving neither rare gases nor  $\pi$ -stacking, e.g., the methane dimer (Fig. 4.3.e), behave somewhat more normally; steady convergence toward the reference energy is exhibited while consistently outperforming the canonical CCSD(T).

Each of the F12n/XZ-F12 model chemistries, when compared across all systems from the A24 test set, exhibit similar convergence speed toward the CBS limit. For DZ-F12, F12a lies within  $0.1 \text{ kcal mol}^{-1}$  of the reference, while both F12b and F12c lie just outside this bound. All three *ansatz* are converged to within  $0.05 \text{ kcal mol}^{-1}$  of the reference for TZ-F12, with the MAE for F12a about half that for F12b and F12c. For QZ- and 5Z-F12, all F12n exhibit MAE of about  $0.01 \text{ kcal mol}^{-1}$  or less.

#### 4.4.2 aXZ vs. XZ-F12: Accuracy Comparison over A24 and Subsets

Individual errors for every bimolecular complex in A24 and average error metrics over the entire test set, computed versus our A24B reference energies, are visualized in Fig. 4.4 for each of the F12n/aXZ and F12n/XZ-F12 model chemistries examined above. For comparison, the dispersion-weighted DW-CCSD(T<sup>\*\*</sup>)-F12 method (abbreviated DW-F12) is also included in Table 4.1 and Fig. 4.4. Interaction energies for A24 systems computed using F12a/aDZ span the spectrum of slightly over- to slightly under-bound, while every complex of the A24 test set is underbound for the F12a/DZ-F12 model chemistry, as shown in Fig. 4.4.a. Typically, F12a/aXZ overbinds complexes for  $\zeta \geq 3$ ; in contrast, all complexes are underbound by the F12a/TZ-F12 model chemistry, but are more likely to be overbound by F12a/QZ-F12. F12b/aXZ (X = D-5), unlike F12a/aXZ, typically underbinds A24 complexes, as is shown in Fig. 4.4.b. Analogous to F12a, F12b/DZ-F12 and F12b/TZ-F12 underbind all A24 complexes. Distinct from F12a, however, F12b/QZ-F12 also underbinds every complex of A24. The F12c *ansatz* behaves qualitatively the same as F12b, as can be seen in Fig. 4.4.c. While F12a/aDZ is slightly more converged than F12a/DZ-F12, MAE values are nearly identical between F12a/aXZ and F12a/XZ-F12 for larger basis sets. The F12b/aXZ and F12c/aXZ model chemistries exhibit lower MAE and MA%E than F12b,c/XZ-F12 methods for double and triple- $\zeta$  basis sets before achieving benchmark convergence in the quadruple- $\zeta$  basis sets, as seen in Figs. 4.4.b&c and Table 4.1.

For the A24 test set, it is clear that F12a/aDZ performs the best out of all examined combinations for double- $\zeta$  basis sets, including silver-standard DW-F12/aDZ. For  $\zeta \geq 3$ , however, the F12b method exhibits the lowest MAE, achieving equivalent accuracy to DW-F12/aTZ over all A24 as well as the MX and DD subsets; for HB systems, however, the DW-F12 method achieves slightly higher accuracy due to the admixture of F12a/aTZ. For double- and triple- $\zeta$  basis sets in particular, computations utilizing the Dunning-style aXZ basis sets produce more accurate results than those done using Peterson and coworkers'



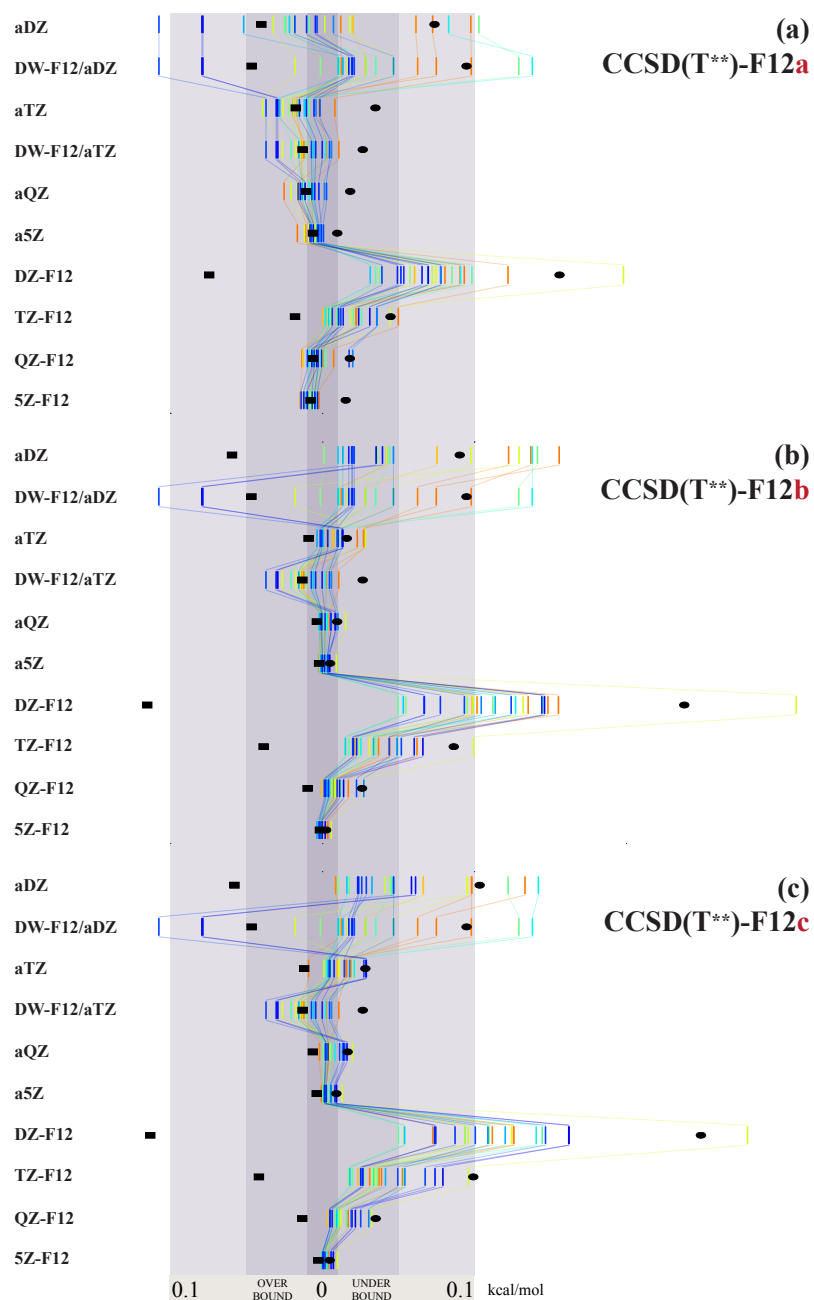


Figure 4.4: Error in IE computed for all bimolecular complexes in the A24 test set relative to A24B reference energies for the (a) CCSD(T<sup>\*\*</sup>)-F12a, (b) CCSD(T<sup>\*\*</sup>)-F12b, and (c) CCSD(T<sup>\*\*</sup>)-F12c methods for both aXZ and XZ-F12 (X = D, T, Q, 5) basis sets. Vertical lines represent individual members of A24, color-coded by interaction type (red = hydrogen bonding, blue = dispersion dominated, yellow/green = mixed interaction). For each level of theory, MAE (black rectangles, given on the left) and MA%E (black ovals, given on the right) are presented. Three shaded error regions are shown: the lightest encompasses  $\pm 0.1$  kcal mol<sup>-1</sup> &  $\sim 4\%$ , next lightest region  $\pm 0.05$  kcal mol<sup>-1</sup> &  $\sim 2\%$ , and darkest region  $\pm 0.01$  kcal mol<sup>-1</sup> &  $\sim 1\%$ . For comparison, errors computed using the current silver-standard DW-CCSD(T<sup>\*\*</sup>)-F12 method, paired with aDZ and aTZ basis sets, are also presented.

specialized XZ-F12 basis sets for *each* F12a,b,c *ansatz*. For quadruple- and quintuple- $\zeta$ , however, the two basis set types give approximately equal MAE. One could imagine that adapting the dispersion weighting scheme may remedy the deficiencies exhibited by the DZ- and TZ-F12 basis sets for NCI relative to the aXZ basis sets. The DW-CCSD(T<sup>\*\*</sup>)-F12 method was introduced to combine the good performance of F12a/aXZ for describing electrostatic interactions and F12b/aXZ for describing dispersion interactions in small basis sets. We have found, however, that the F12a *ansatz* yields the highest accuracy across the entire A24 test set and *each* of the hydrogen bonding, dispersion dominated, and mixed interaction subsets when employing XZ-F12 basis sets. Therefore, the benefit of mixing F12a/XZ-F12 and F12b/XZ-F12 will be lost; as such, we have chosen not to extend the dispersion weighting scheme to the DZ- and TZ-F12 basis sets.

Given below are total orderings for the accuracy of all model chemistries at each  $\zeta$ -level versus the CCSD(T)/CBS reference; methods are listed from left to right in order from most to least accurate:

D $\zeta$ : F12n/aDZ (n = a>b $\gtrsim$ c) > F12n/DZ-F12 (n = a>c>b)

T $\zeta$ : F12n/aTZ (n = b $\gtrsim$ c>a)  $\sim$  F12n/TZ-F12 (n = a>b $\gtrsim$ c)

Q $\zeta$ : F12b/aQZ>F12a/QZ-F12 $\sim$ F12c/aQZ $\sim$ F12a/aQZ $\sim$ F12a/QZ-F12 $\gtrsim$ F12c/QZ-F12

5 $\zeta$ : F12b/5Z-F12 $\sim$ F12b/a5Z $\sim$ F12c/5Z-F12  $\sim$ F12c/a5Z>F12a/a5Z $\sim$ F12a/5Z-F12

These orderings were constructed to a MAE resolution of  $\pm 0.01$  kcal mol<sup>-1</sup> and MA%E resolution of  $\pm 0.5\%$ ; between two adjacent methods, > indicates that the corresponding MAE is distinguishable to this resolution,  $\gtrsim$  indicates indistinguishable MAE but distinguishable MA%E, and  $\sim$  indicates indistinguishable MAE and MA%E.

The most surprising result of this work is the comparatively poor performance of F12n/XZ-F12 model chemistries for A24 systems at the double- and triple- $\zeta$  levels. This finding, however, confirms the same conclusion by Patkowski.<sup>100</sup> One potential cause for this be-

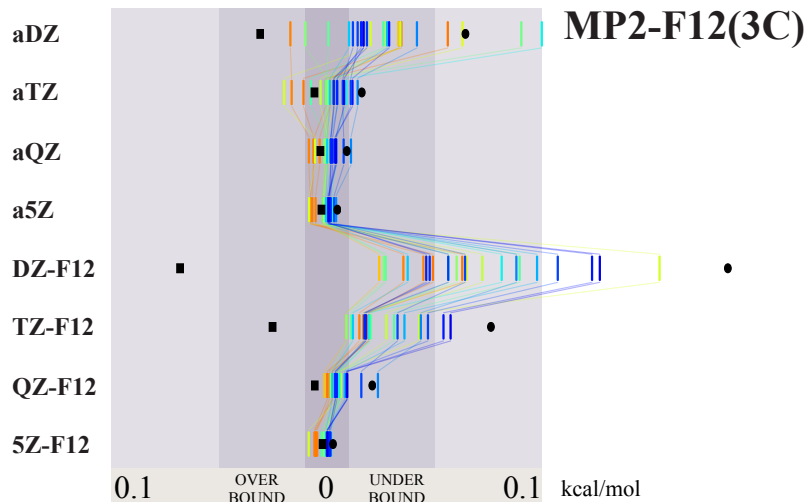


Figure 4.5: Error in IE computed for all bimolecular complexes in the A24 test set relative to Helgaker-extrapolated<sup>20</sup> MP2/CBS(aXZ, a(X+1)Z) reference IEs (analogous  $\zeta$ -levels to A24B), computed using (a) the MP2-F12(3C)/aXZ (X = D, T, Q, 5) and (b) MP2-F12(3C)/XZ-F12 (X = D, T, Q) model chemistries. For each level of theory, MAE (black rectangles, given on the left) and MA%E (black ovals, given on the right) are presented. Shaded error regions are given with identical ranges as in Fig.4.4.

havior is that the XZ-F12 basis sets were optimized for use with the explicitly correlated MP2-F12(3C) method and *not* for use with explicitly correlated coupled-cluster methods (CCSD-F12n). In their study introducing the XZ-F12 (X = D, T, Q) basis sets, Peterson and coworkers showed that MP2-F12(3C)/XZ-F12 outperforms both MP2-F12(3C)/aXZ and MP2-F12(3C)/aug-cc-pV(X+d)Z. For the sake of completeness, we have also examined the performance of the MP2-F12(3C)/XZ-F12 and MP2-F12(3C)/aXZ methods for interaction energies of systems in the A24 test set, compared to reference IEs computed at the canonical MP2/CBS level.

MAE and MA%E error metrics for the MP2-F12(3C)/XZ-F12 and MP2-F12(3C)/aXZ methods computed over the entire A24 test set, as well as each subset, are presented in Table S-6 in the Supplemental Information (see Section ??) and visualized in Fig. 4.5. For A24 systems, the MP2-F12(3C)/aXZ method produces somewhat tighter convergence to MP2/CBS reference energies than MP2-F12(3C)/XZ-F12. This result seemingly contradicts the findings by Peterson;<sup>116</sup> however, the XZ-F12 basis sets were optimized for

computations of *molecular correlation energy*, not noncovalent *interaction energies*. Our results for correlation energies computed for the dimer and both monomers of all A24 systems using both MP2-F12(3C)/aXZ and MP2-F12(3C)/XZ-F12 methods show that XZ-F12 basis sets produce both faster convergence and better agreement to MP2/CBS limits for these energies, in agreement with Peterson. Hence, it appears that the XZ-F12 basis sets are more slowly convergent toward the CBS limit for overall noncovalent interaction energies than the aXZ basis sets, despite being more rapidly convergent for the individual molecular correlation energies. This can most likely be attributed to the comparative compactness of XZ-F12 basis sets versus aXZ basis sets. Indeed, while XZ-F12 contains *s* and *p* functions that are as diffuse or even slightly more diffuse than the corresponding aXZ basis, the functions of higher angular momentum are significantly more compact in XZ-F12 than aXZ; the worse performance of modest DZ- and TZ-F12 basis sets for NCI vs. the aDZ and aTZ basis sets is therefore understandable. In contrast, both basis set types achieve comparable levels of accuracy for quadruple- and quintuple- $\zeta$  levels, where the XZ-F12 basis sets reach sufficient levels of diffusivity to accurately describe weak intermolecular interactions.

#### 4.4.3 aXZ vs. XZ-F12: Extension to the S22 Test Set

In order to validate and generalize the above analysis regarding the performance of CCSD-T(\*\*)-F12n methods to larger and more strongly bound bimolecular complexes, we have applied these methods to the S22 test set,<sup>79</sup> using the revised S22B interaction energies.<sup>30</sup> Previously, Burns *et al.*<sup>77</sup> reported interaction energies and statistics for the F12n/aDZ (n = a, b) and DW-F12/aDZ levels of theory. We have extended that investigation to include F12c/aDZ and each of the F12n/DZ-F12 (n = a, b, c) levels of theory; MAE and MA%E error statistics for each of these model chemistries are given in Table 4.2, and visualized in Fig. 4.6. Again, for comparison, statistics for DW-F12/aDZ are also in Table 4.2 and Fig. 4.6.

The previously reported interaction energies for S22 systems computed using the DW-

Table 4.2: Interaction energy (kcal/mol) error statistics for F12n/aDZ and F12n/DZ-F12 methods, applied to the S22B test set, as well as for the hydrogen bonding (HB), mixed interaction (MX) and dispersion dominated (DD) subsets. Included for reference are error statistics computed with the DW-CCSD(T<sup>\*\*</sup>)-F12/aDZ method.

Method & Basis Set	S22		HB		MX		DD	
	MAE	MA%E	MAE	MA%E	MAE	MA%E	MAE	MA%E
<b>aug-cc-pVDZ</b>								
DW-CCSD(T <sup>**</sup> )-F12 <sup>a</sup>	0.05	1.3	0.06	0.6	0.07	1.4	0.03	1.8
CCSD(T <sup>**</sup> )-F12a <sup>a</sup>	0.12	2.3	0.06	0.6	0.15	1.8	0.15	4.6
CCSD(T <sup>**</sup> )-F12b <sup>a</sup>	0.10	1.8	0.18	1.6	0.09	2.1	0.03	1.8
CCSD(T <sup>**</sup> )-F12c	0.08	1.9	0.12	1.2	0.07	1.9	0.05	2.7
<b>cc-pVDZ-F12</b>								
CCSD(T <sup>**</sup> )-F12a	0.17	3.4	0.21	1.6	0.18	3.4	0.11	5.2
CCSD(T <sup>**</sup> )-F12b	0.28	5.7	0.31	2.4	0.30	5.5	0.23	9.3
CCSD(T <sup>**</sup> )-F12c	0.27	5.8	0.26	2.0	0.31	5.4	0.25	10.1

<sup>a</sup> Values from Ref 77.

F12/aDZ, F12a/aDZ, and F12b/aDZ model chemistries all straddle the best-estimate IE, while F12c/aDZ and each F12n/DZ-F12 underbind complexes in S22 (with the lone exception of the adenine-thymine complex, which is slightly overbound by F12c/aDZ). This underbinding of complexes with F12c/aDZ and F12n/DZ-F12 was expected, as it is identical to the behavior of these model chemistries for the A24 test set. F12a/aDZ and F12b/aDZ, however, were also expected to underbind these complexes; apparently, when applied to larger systems, these methods no longer reliably provide an upper bound for the CBS limit IE.

From Table 4.2, it is clear that F12n/aDZ model chemistries are more converged than F12n/DZ-F12 for these larger systems, just as was seen for A24. Across the entire S22 test set, as well as each subset, the best F12n/DZ-F12 methods exhibit MAEs a factor of 2-4 times larger than the best F12n/aDZ methods. Curiously, in contrast to the A24 test set, F12c/aDZ is the most converged *ansatz* overall, having the smallest MAE versus the S22B reference values. Again, F12a is the best choice of *ansatz* for modeling electrostatic interactions and F12b is best for dispersion. For the F12n/DZ-F12 methods an identical trend is observed for S22 as was for A24: F12a is superior to both F12b and F12c over all S22 systems and in each subset by a non-trivial factor. These results indicate the following

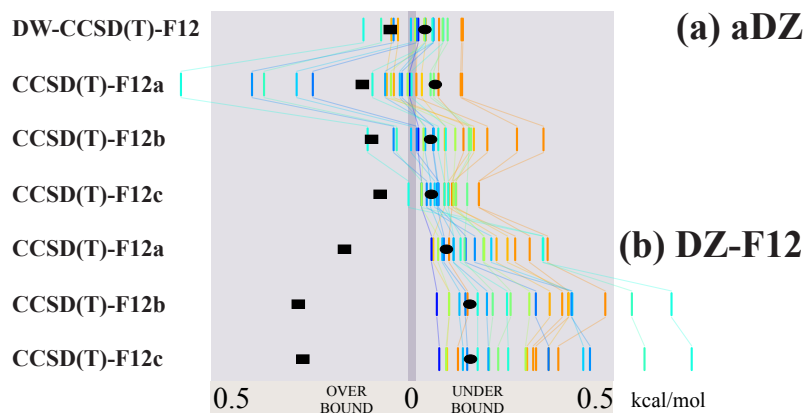


Figure 4.6: Error in IE computed for all bimolecular complexes in the S22 test set,<sup>79</sup> relative to revised S22B reference energies,<sup>30</sup> computed using (a) the CCSD(T<sup>\*\*</sup>)-F12n/aDZ (n = a, b, c) and (b) CCSD(T<sup>\*\*</sup>)-F12n/DZ-F12 (n = a, b, c) model chemistries. Included for reference are IEs computed using the silver-standard DW-CCSD(T<sup>\*\*</sup>)-F12/aDZ method. For each model chemistry, MAE (black rectangles, given on the left) and MA%E (black ovals, given on the right) are presented. The outer, lightly shaded region encompasses  $\pm 0.5$  kcal mol<sup>-1</sup> &  $\sim 23\%$ , and the inner, darkly shaded region is given to indicate the location of zero error.

total ordering for the examined explicitly correlated methods:

$$\text{F12n/aDZ (n = c > b > a)} > \text{F12n/DZ-F12 (n = a > c > b)}$$

While this manuscript was in preparation, Martin and coworkers<sup>MartinS66x8</sup> published revised benchmark energies for the S66 $\times$ 8 test set of Hobza and coworkers,<sup>105,106</sup> computed using the same XZ-F12 basis sets of Peterson considered here. Their study examined the most accurate manner in which to scale the (T) correction to the CCSD-F12n energy, what explicitly correlated method to utilize within the focal-point correction approach for reaching the complete basis set limit, and what counterpoise correction scheme is most appropriate. For the double- $\zeta$  DZ-F12 basis set, CCSD-F12c was superior to CCSD-F12b in the context of the  $\delta_{\text{MP2-F12}}^{\text{CCSD-F12n}}$  correction for dimers at their equilibrium geometries ( $r = r_e$ ). In their study, only the XZ-F12 basis set was used to compute S66 $\times$ 8 benchmark energies; aXZ was not examined, nor was the F12a *ansatz* considered. In light of our findings for the noncovalently bound complexes within the A24 and S22 test sets, perhaps these new benchmarks should be revisited with computations making use of aXZ basis sets rather

than XZ-F12.

#### 4.4.4 Benchmark Procedures for NCI: Combining Accuracy and Computational Cost

The vastly superior performance of explicitly correlated CCSD(T<sup>\*\*</sup>)-F12n approaches for NCI over canonical CCSD(T) raises the question of whether or not one of these methods should replace the focal point scheme we previously recommended<sup>30</sup> for NCI benchmarking, namely MP2/CBS(aTZ, aQZ) +  $\delta_{\text{MP2}}^{\text{CCSD(T)}}$ /aTZ, denoted more compactly as CCSD(T)/[aTQZ;  $\delta$ :aTZ]. In this section, we therefore compare the explicitly correlated CCSD(T<sup>\*\*</sup>)-F12b/aTZ level of theory, canonical CCSD(T) with CBS(aDZ, aTZ) and CBS(aTZ, aQZ) Helgaker extrapolations, and the CCSD(T)/[aTQZ;  $\delta$ :aTZ] focal-point procedure. The error incurred by employing density fitting in MP2/CBS computations of noncovalent interactions is as little as 0.001 kcal mol<sup>-1</sup> in the aTZ basis;<sup>117,118</sup> we have therefore also examined the focal point CCSD(T)/[DF-aTQZ;  $\delta$ :aTZ] scheme, where DF-aTQZ indicates using DF-MP2/CBS(aTZ, aQZ) in the focal-point procedure.

In their extensive study of wavefunction approaches for computing noncovalent interaction energies, Burns *et al.*<sup>77</sup> touted DW-F12/aDZ as the “silver standard” method for NCI, indicating its ability to produce near-benchmark accuracy [MAE to within 0.06 kcal mol<sup>-1</sup> of estimated CCSD(T)/CBS reference values] for an extensive test set of noncovalently bound complexes with significantly lower computational cost than expected for such accuracy. Additionally, for the 10 smallest members of the S22 test set where the highest-quality benchmark IE’s are available, F12b/aTZ and DW-F12/aTZ exhibited mean errors of only 0.028 and 0.008 kcal mol<sup>-1</sup>, respectively, versus CCSD(T)/CBS.<sup>96</sup> We will therefore also consider the DW-CCSD(T<sup>\*\*</sup>)-F12/aTZ model chemistry in our comparison. The double- $\zeta$  counterparts of these methods, namely DW-F12/aDZ, F12b/aDZ, and CCSD(T)/[aTQZ;  $\delta$ :aDZ] (using both MP2/ and DF-MP2/CBS) methods are also included in this analysis. We have decided to consider neither F12a nor F12c *ansatz* in this examination: first, even while F12a/aDZ is the best converged method with double- $\zeta$  basis sets for A24 systems, its

poor performance for the S22 test set indicates that this combination does not reliably yield well-converged results for a variety of chemical systems, and second, as was discussed in Section 4.3.1, F12c is more computationally intensive than F12b, while the two methods achieve similar accuracy.

Recently, the combination of frozen natural orbital coupled-cluster [FNO-CCSD(T)]<sup>119–123</sup> and density-fitted coupled-cluster<sup>124–126</sup> has been shown to be accurate and efficient for reaction energies and NCI.<sup>124</sup> Hence, we also considered focal-point approaches similar to those mentioned above, but with DF-FNO-CCSD(T) used to compute the coupled-cluster correction. In particular, we considered DF-FNO-CCSD(T)/[DF-aTQZ;  $\delta$ :aDZ] and DF-FNO-CCSD(T)/[DF-aTQZ;  $\delta$ :aTZ] methods. Interaction energies involving the DF-FNO-CCSD(T) method<sup>124</sup> were computed with PSI4,<sup>127</sup> a freely available, open-source suite of *ab initio* quantum chemistry programs. The conservative, PSI4 default frozen natural orbital cutoff value of  $1.0 \times 10^{-6}$  was used for all DF-FNO-CCSD(T) computations.

To assess the combined computational expense and accuracy of these methods, interaction energy computations were timed for a single representative system, the methane-ethane complex (A24-17), which has 700 basis functions for the aTZ basis set. All timings were performed on a workstation computer with an Intel Core i7 3930K CPU with 6-cores at 3.2 GHz, 64GB of memory, and a RAID0 array of 3×3TB hard disks for scratch space. Both serial (single core) and parallel (six cores) timings were obtained; we report total (“wall”) times for computation of the counterpoise-corrected interaction energy (consisting of total energy computations for the dimer and both monomers). Point-group symmetry was not utilized for these timings.

As shown in Table 4.3, explicitly correlated F12b/aTZ and DW-F12/aTZ are as accurate as the currently recommended method from the literature, CCSD(T)/[DF-aTQZ;  $\delta$ :aTZ].<sup>30</sup> Additionally, DW-F12/aTZ is much less costly than any canonical (non-DF-FNO) method considered here that achieves near-benchmark convergence. In fact, this explicitly correlated combination is converged as well as (non-DF) CCSD(T)/[aTQZ;  $\delta$ :aTZ] and



Table 4.3: Interaction energy (kcal/mol) error statistics, error distributions, and timing summaries for benchmark procedure candidates, along with their double- $\zeta$  counterparts, computed over the A24 test set and each subset versus A24B reference energies.

Method & Basis Set	A24	HB	MX	DD	Error Distribution <sup>a</sup>	Time <sup>b</sup>	
						Serial <sup>c</sup>	Parallel <sup>d</sup>
<b>CCSD(T**) - F12b</b>							
aug-cc-pVDZ	0.06	0.11	0.07	0.02	■     ■■■■■     ■■■■     ■■■■	7.6	2.5
aug-cc-pVTZ	0.01	0.02	0.01	0.01	■■■■     ■■■■     ■■■■	161.8	28.7
<b>DW-CCSD(T**) - F12</b>							
aug-cc-pVDZ	0.05	0.06	0.05	0.04	■     ■■■■     ■■■■     ■■■■	7.7	2.5
aug-cc-pVTZ	0.01	0.01	0.01	0.01	■■■■     ■■■■     ■■■■	165.9	46.0
<b>CCSD(T)</b>							
aug-cc-pVDTZ <sup>e</sup>	0.03	0.06	0.02	0.02	■■■■     ■■■■     ■■■■     ■■■■	158.0	44.1
aug-cc-pVTQZ <sup>e</sup>	0.01	0.02	0.01	0.00	■■■■     ■■■■     ■■■■	2614.5	1150.7
[aTQZ; $\delta$ :aDZ] <sup>e</sup>	0.03	0.07	0.02	0.03	■■■■     ■■■■     ■■■■     ■■■■	589.3	481.3
[aTQZ; $\delta$ :aTZ] <sup>e</sup>	0.02	0.01	0.01	0.02	■■■■     ■■■■     ■■■■	706.8	510.6
[DF-aTQZ; $\delta$ :aDZ]	0.03	0.07	0.02	0.03	■■■■     ■■■■     ■■■■     ■■■■	11.8	5.1
[DF-aTQZ; $\delta$ :aTZ]	0.01	0.01	0.01	0.02	■■■■     ■■■■     ■■■■	156.7	45.2
<b>DF-FNO-CCSD(T)<sup>f</sup></b>							
[DF-aTQZ; $\delta$ :aDZ]	0.03	0.07	0.02	0.03	■■■■     ■■■■     ■■■■     ■■■■	13.4	4.1
[DF-aTQZ; $\delta$ :aTZ]	0.01	0.01	0.01	0.01	■■■■     ■■■■     ■■■■	126.0	28.7

<sup>a</sup> A24 Errors with respect to A24B reference within  $\pm 0.15$  kcal/mol. Guide lines are at 0.0, 0.01, 0.05, and 0.1 kcal mol<sup>-1</sup> overbound (−) and underbound (+).

<sup>b</sup> Times given are the total walltime (min) necessary to compute the counterpoise-corrected interaction energy for the methane-ethane complex (A24-17) without use of spatial symmetry.

<sup>c</sup> Computations performed on a single core.

<sup>d</sup> Computations performed with threading over 6 cores.

<sup>e</sup> Values from Ref. 77.

<sup>f</sup> Computed using PSI4 with default natural orbital cutoff value of  $1.0 \times 10^{-6}$ .

CCSD(T)/CBS(aTZ, aQZ), while taking only 23% and 6% (8% and 4%) of the computational time used by these composite methods, respectively, when running in serial (parallel). The speedup relative to CCSD(T)/CBS(aTZ, aQZ) is not surprising given that procedure's need to compute the expensive CCSD(T)/aQZ IE. The speedup versus CCSD(T)/[aTQZ;  $\delta$ :aTZ] may be more surprising because the CCSD(T) computations are being done in the same aTZ basis. However, the composite method's required MP2/aQZ computations can be time consuming using conventional (non-DF) codes. When using CCSD(T)/[DF-aTQZ;  $\delta$ :aTZ], the timings are essentially the same as for DW-F12/aTZ [as is the MAE vs. the CCSD(T)/CBS limit]. Hence, these methods appear to be interchangeable for benchmarking purposes, at least as far as can be discerned from the A24 test set. Remarkably, the F12b/aTZ combination achieves this level of accuracy while enjoying a 1.6x speedup over the currently recommended CCSD(T)/[DF-aTQZ;  $\delta$ :aTZ] method when using 6 cores.

Considering now the triple- $\zeta$  DF-FNO-based composite approach, we see no additional error incurred by either DF or FNO approximations compared to conventional CCSD(T)/CBS[DF-aTQZ;  $\delta$ :aTZ]. With respect to FNO, this can be attributed to the quite conservative default  $1.0 \times 10^{-6}$  natural orbital occupation number cutoff in PSI4. However, the synergistic DF and FNO approximations afford the possibility of helpful speedups.<sup>124</sup> From Table 4.3 we see a 1.6x speedup of DF-FNO-CCSD(T)/CBS[DF-aTQZ;  $\delta$ :aTZ] over CCSD(T)/CBS[DF-aTQZ;  $\delta$ :aTZ] when running on 6 cores. The combination of high accuracy and lower computational expense is very promising for the application of DF-FNO-CCSD(T) and F12b to larger noncovalently bound systems. Indeed, the DF-FNO scheme sped up CCSD(T) computations of the three-body contribution to the interaction energy of benzene trimer by about a factor of four, while incurring an error of only 0.002 kcal mol<sup>-1</sup>.<sup>124</sup>

Overall, these results suggest that the DW-F12/aTZ combination should be considered to be equivalently accurate and cost-effective as the currently recommended method from the literature, namely, the CCSD(T)/CBS[DF-aTQZ;  $\delta$ :aTZ] focal-point approach. When density fitting and a truncated frozen natural orbital space are employed for the CCSD(T)

procedure, no significant errors are added, but the computations of counterpoise-corrected interaction energies are sped up significantly; additionally, an identical combination of speed and accuracy is achieved by the F12b/aTZ model chemistry.

#### 4.5 Summary and Conclusions

In an effort to understand and evaluate the performance of several popular approximations to the explicitly correlated CCSD(T)-F12 method, as well as to determine whether these approaches could be suitable for computing benchmark-quality noncovalent interaction energies (IEs), we have examined the convergence behavior and accuracy for each of the CCSD(T<sup>\*\*</sup>)-F12n (n = a, b, c; abbreviated as F12n) methods over a variety of bimolecular complexes of diverse binding motifs and interaction strengths. These methods were paired with the correlation-consistent polarized valence basis sets of Dunning, augmented with diffuse functions (aug-cc-pVXZ, X = D, T, Q, 5; abbreviated aXZ)<sup>22</sup> and the specialized explicitly correlated correlation-consistent polarized valence basis sets of Peterson (cc-pVXZ-F12, X = D, T, Q, 5; abbreviated XZ-F12),<sup>108,109</sup> which were designed for use with the explicitly correlated F12 methods. The accuracy of these methods at each basis set level was assessed for the A24 test set of Hobza and coworkers<sup>14</sup> and each of the hydrogen bonding, mixed interaction, and dispersion dominated subsets against our revised best estimate reference values (denoted A24B); all IE computations were counterpoise corrected to account for basis set superposition errors (BSSE),<sup>54</sup> and findings were validated by examining the S22 test set,<sup>79</sup> with previously revised reference values (S22B).<sup>96</sup> Although we focused on counterpoise-corrected values, we partially examined uncorrected interaction energies for the A24 test set and found them to yield substantially larger mean absolute errors [MAE = 0.13 kcal/mol for each of F12a/aTZ and F12b/aTZ, compared to counterpoise-corrected values of 0.02 and 0.01 kcal/mol, respectively]. Additional comparisons of counterpoise-corrected vs uncorrected values are available in Ref. 101.

For all members of A24, the F12b/aXZ and F12c/aXZ model chemistries converge to-

wards reference interaction energies from above, while F12a/aXZ typically converges from below towards the reference interaction energy for  $\zeta > 3$ , after being an upper bound with the double- $\zeta$  aDZ basis. When paired with the explicitly correlated XZ-F12 basis sets, however, each of the F12n *ansatz* converge monotonically from above towards the reference IE for all  $\zeta$ . F12a very slightly overshoots the best-estimate IE for the majority of A24 systems with the 5Z-F12 basis set, converging towards an overbound interaction energy at its complete basis set limit. F12n/aXZ model chemistries converge to the best-estimate IE more quickly than their F12n/XZ-F12 counterparts for the A24 test set. The F12b/aXZ model chemistry converges the most rapidly out of any method examined, affording about three free  $\zeta$ -levels of accuracy over canonical CCSD(T) with only a slight increase in computational cost. The performance of F12c/aXZ is extremely similar.

We also considered the S22 test set, which contains bimolecular complexes involving slightly larger molecules than A24. Here we only tested F12n methods in conjunction with double- $\zeta$  quality basis sets, because the CCSD(T)/CBS limits are not as precisely known for many of the complexes in S22 as for A24, leading to the possibility that CCSD(T)-F12n computations using aTZ or aQZ basis sets might actually be closer to the CBS limit than currently available reference values. For A24, the best F12n *ansatz* paired with an aDZ basis set was F12a, but for S22, it becomes the worst, with accuracy decreasing in the order DW-F12 > F12c > F12b > F12a.

Considering computations in the aDZ basis specifically, hydrogen-bonded complexes of A24 and S22 are most accurately computed by F12a and DW-CCSD(T)-F12, while F12b is best for dispersion-bound complexes. For mixed interaction types, F12a and DW-CCSD(T)-F12 are best on average for A24, but F12a is the worst *ansatz* for S22, with an MAE of 0.15 kcal mol<sup>-1</sup> (all other F12 *ansatz* considered have MAE < 0.1 kcal mol<sup>-1</sup> across the mixed interaction subsets of A24 and S22). Considering larger aXZ basis sets, there is little to distinguish any of the explicitly correlated approaches considered, with MAE vs. the best estimates of 0.02 kcal mol<sup>-1</sup> or less for A24 or any of its interaction type

subsets. Nevertheless, F12b exhibits slightly lower mean absolute percent errors than the other explicitly correlated procedures considered.

Perhaps surprisingly, the DZ and TZ-F12 basis sets do not perform as well as aDZ and aTZ for CCSD(T)-F12n computations of noncovalent interaction energies, either for A24 or S22. This is not a result of these basis sets being optimized for use with the MP2-F12(3C) method instead of CCSD(T)-F12n; rather it appears that their optimization for *molecular correlation energies* and relative compactness in high angular momentum  $d$ ,  $f$ , ... basis functions does not translate well to noncovalent *interaction energies*.

Finally, prompted by the outstanding performance of the explicitly correlated F12n and DW-F12 methods compared to canonical CCSD(T), we have examined the F12b/aTZ and DW-F12/aTZ combinations, as well as composite focal-point estimates of the CCSD(T)/CBS limit combining CCSD(T) and MP2 computations, with consideration of each method's computational expense and convergence to the CBS limit. Explicitly correlated DW-F12/aTZ and F12b/aTZ both yielded benchmark-quality interaction energies for systems in the A24 test set; the former attained such accuracy with nearly identical computational expense as the currently recommended benchmark procedure, DF-MP2/CBS(aTZ, aQZ) +  $\delta_{\text{MP2}}^{\text{CCSD(T)}}$ /aTZ, also denoted CCSD(T)/[DF-aTQZ;  $\delta$ :aTZ], while the latter achieves a 1.6x speedup over the focal-point procedure when using 6 cores (timings for a single representative test case). Replacing CCSD(T) energies with their density-fitted, frozen natural orbital truncated counterparts via DF-FNO-CCSD(T) does not lead to an increase in the mean absolute errors for the A24 test set, but also yields a 1.6x speedup over the non-DF, non-FNO focal-point approach when running on 6 cores.

## CHAPTER 5

### ASSESSMENT OF DENSITY FUNCTIONAL METHODS FOR GEOMETRY OPTIMIZATION OF BIMOLECULAR VAN DER WAALS COMPLEXES

#### 5.1 Abstract

We explore the suitability of three popular density functionals (B97-D3, B3LYP-D3, M05-2X) for producing accurate equilibrium geometries of van der Waals (vdW) complexes with diverse binding motifs. For these functionals, optimizations using Dunning's aug-cc-pVDZ basis set best combine accuracy and a reasonable computational expense. Each DFT/aug-cc-pVDZ combination produces optimized equilibrium geometries for 21 small vdW complexes of organic molecules (up to four non-hydrogen atoms total) that agree with high-level CCSD(T)/CBS reference geometries to within  $\pm 0.1$  Å for the averages of the center-of-mass displacement and the mean least root-mean-squared displacement. The DFT/aug-cc-pVDZ combinations are also able to reproduce the optimal center-of-mass displacements interpolated from CCSD(T)/CBS radial potential energy surfaces in both NBC7x and HBC6 test sets to within  $\pm 0.1$  Å. We therefore conclude that each of these density functional methods, together with the aug-cc-pVDZ basis set, are suitable for producing equilibrium geometries of generic non-bonded complexes.<sup>†</sup>

#### 5.2 Introduction

Structure-based computer-aided drug design (SB-CADD) has emerged as a valued approach in the development of novel pharmaceutical compounds. Optimization of binding affinity of a lead chemical series is an iterative process, often requiring a detailed understanding of existing host-guest interactions and the ability to successfully predict new

---

<sup>†</sup>This Chapter reproduces the work in Ref. 128.

ones. These assessments are typically performed with molecular mechanics forcefields in combination with structural models based on crystallographic structures of closely related molecules. Such methods, however, are not always sufficiently fine-grained to accurately describe or quantify the guest-host interactions of interest. Capable of augmenting the existing SB-CADD paradigm by providing the information necessary to allow for the rational refinement of the resulting drug candidates, *ab initio* quantum-chemical methods offer a first-principles description of the non-covalent interactions (NCI) which govern host-guest binding. In particular, energy decomposition analysis (EDA) schemes such as the absolutely localized molecular orbital EDA (ALMO-EDA)<sup>129-132</sup> and symmetry-adapted perturbation theory (SAPT)<sup>133-136</sup> approaches offer a physically meaningful breakdown of interaction energies into contributions from more fundamental components, such as electrostatics, induction, dispersion, and exchange-repulsion. Furthermore, the atomic<sup>137</sup> and functional-group<sup>138</sup> partitionings of SAPT (A-SAPT and F-SAPT, respectively) offer an additional layer of interaction energy decomposition into the specific interactions between pairs of atoms or functional groups on each interacting species. Indeed, these methods have already provided insight into the relative stability of chlorinated vs. methylated factor Xa inhibitors<sup>139</sup> and the role of NCI on transition-state stabilization in organocatalyzed aldol addition.<sup>68</sup>

Before performing any of these quantum-chemical computations within SB-CADD applications, a model system must first be constructed which mimics the NCI of interest, and a geometry of suitable quality obtained, as resolution of tenths of kcal/mol or less may be necessary to distinguish between relevant binding configurations (e.g., the sandwich and T-shaped configurations of the pyridine dimer differ in interaction energy by only 0.1 kcal mol<sup>-1</sup>!<sup>39,140</sup>) or seemingly minor chemical modifications. While significant attention has been paid to the optimization of individual molecules, a general protocol that is capable of generating accurate geometries of supermolecular assemblies is conspicuously lacking in the literature. Even among benchmark sets of non-covalently bound complexes, sig-

nificantly more consideration is given to the computation of interaction energies than to the geometry optimization of the complexes themselves.<sup>15,30,76,79,141–145</sup> Hence, despite the availability of high-quality interaction energies approaching the coupled-cluster through perturbative triples [CCSD(T)]<sup>10</sup> complete-basis set (CBS) limit, there are very limited data on high-quality geometries of van der Waals dimers that might be used to assess various approximate methods for geometry optimization.

Shown previously to be capable of reproducing benchmark quality IEs to within sub-kcal mol<sup>-1</sup> accuracy<sup>75–77</sup> while maintaining low execution time relative to post-Hartree–Fock electron correlation methods like second-order Møller–Plesset perturbation theory (MP2) or CCSD(T), density functional theory approaches that include a treatment of London dispersion forces seem like natural candidates for routine application to the geometry optimization of such complexes. Indeed, some of the best such approaches yield MAE of only a few tenths of one kcal mol<sup>-1</sup> for the S22 test set.<sup>75</sup> We therefore explore here the suitability of three of the best of these methods (B97-D3,<sup>146</sup> B3LYP-D3,<sup>147,148</sup> and M05-2X<sup>149</sup>), where -D3 denotes the third-generation dispersion correction of Grimme;<sup>35</sup> these functionals have exhibited MAD = 0.48, 0.79, and 0.36 kcal mol<sup>-1</sup>, respectively, for non-counterpoise-corrected interaction energies versus benchmarks for complexes in the S22 test set.<sup>75</sup> The performance of each functional is assessed first by comparing DFT-optimized geometries for the 21 minimum-energy complexes in the A24 test set of Hobza and co-workers,<sup>14</sup> for which CCSD(T)/CBS-quality geometries are available due to the small size of these complexes (A24 systems contain up to four non-hydrogen atoms). Fully-optimized CCSD(T)/CBS geometries of larger systems would be computationally difficult to obtain; however, here we also present CCSD(T)/CBS potential energies vs intermolecular separation for 13 systems with up to twelve non-hydrogen atoms. These one-dimensional potential energy curves allow us to assess density functionals for their ability to reproduce the optimal intermolecular separation in these larger complexes, in order to validate the conclusions drawn from the optimization of the A24 systems.



### 5.3 Computational Methods

Throughout this study, we employ three density functionals which have become routinely applied to NCI:<sup>75–77</sup> B97-D3 (generalized gradient approximation, GGA),<sup>146</sup> B3LYP-D3 (hybrid-GGA),<sup>147,148</sup> and M05-2X (hybrid-meta-GGA);<sup>149</sup> each of these is paired with the popular correlation-consistent basis sets of Dunning both with and without augmentation by diffuse functions [(aug-)cc-pVXZ, X = D, T; abbreviated throughout as aXZ and XZ, respectively]. The relative computational expense of each of these functionals increases with each successive rung, from B97  $\rightarrow$  B3LYP  $\rightarrow$  M05-2X, with B97 exhibiting an order of magnitude smaller overall algorithmic scaling compared to both B3LYP and M05-2X when density fitting is employed [ $\mathcal{O}(N^3)$  versus  $\mathcal{O}(N^4)$ , respectively, with  $N$  proportional to overall system size]. This increase in computational scaling results from an increase in the amount of physics recovered by each successive functional: GGAs, which depend only on the gradient of the density, incorporate only local correlation; hybrid-GGAs incorporate a percentage of Hartree–Fock exchange, recovering some nonlocal correlation; and hybrid-meta-GGAs incorporate exact exchange in addition to a functional dependence on both the gradient and Laplacian of the local density, recovering both nonlocal correlation and a more correct description of the topological dependence of the energy on the overall electron density. For the interested reader, we have included in Section I C in the supplementary information a brief summary of timings for the construction of these gradients with the methods examined here.

We have applied the -D3 dispersion correction of Grimme<sup>35</sup> to B97 and B3LYP, as correction for missing dispersion in these functionals has been shown to be necessary for a high-quality description of NCI.<sup>75</sup> We have chosen this pairwise dispersion treatment as opposed to a many-body,<sup>150–154</sup> exchange-dipole moment,<sup>155–157</sup> or non-local<sup>158–161</sup> approaches due to the availability of low-cost analytical gradients for the -D family of corrections, allowing for minimal additional expense when incorporated into the geometry

optimization procedure. The -D3 correction should give correct long-range behavior for the London dispersion interactions, while also accounting for the local chemical environment around each atom.<sup>35</sup> M05-2X, on the other hand, can describe London dispersion interactions at short to intermediate distances (up to  $\sim 5 \text{ \AA}$ ),<sup>162</sup> but fails to have correct long-range behavior. This deficiency should not be significant for the smaller molecular systems examined here, but can become a problem for large systems with many long-range contacts.<sup>163,164</sup> For a thorough discussion of the ladder of approximations within density functional theory, a recent review of dispersion corrections in DFT and other mean-field electronic structure methods, and the application of density functional theory to study non-covalent interactions, we refer the interested reader to Refs. 165, 166, and 75, respectively, and the references therein.

### 5.3.1 Optimized Geometries for A24 systems

Geometries were optimized for the 21 minimum-energy complexes in the A24 test set of Hobza and co-workers<sup>14</sup> (complexes 1–21, denoted A21; visualized in Fig. 5.1.a) using the dispersion corrected functionals described above. Optimizations with B3LYP-D3 and B97-D3 used a development version of the open-source PSI4 electronic structure program,<sup>31</sup> and optimizations with M05-2X used the Q-Chem 4 program package.<sup>167</sup> In both cases, full optimizations were performed, allowing for monomer relaxation. For the optimization of these systems, we did not attempt to correct for basis set superposition error (BSSE), as this would add computational expense and would also be more difficult to automate with standard geometry optimizers. Fortunately, our results indicate that BSSE correction is not necessary for reliable geometries. Convergence criteria used for optimizations in this work were a) the energy difference between successive optimization steps below  $1 \times 10^{-6} E_h$ , b) the maximum component of the gradient below  $1.5 \times 10^{-5} E_h/a_0$ , c) the root-mean-square of the elements of the gradient below  $1.0 \times 10^{-5}$ , d) the maximum atomic displacement between successive optimization steps below  $6.0 \times 10^{-4} a_0$ , and e) the root-

mean-square of the atomic displacements between successive optimization steps below  $4.0 \times 10^{-4}$ . For optimizations performed using PSI4, each of these five criteria must be satisfied for convergence to be achieved; for those performed using Q-Chem, however, convergence was achieved when criterion (b) and either of (a) or (d) were satisfied. These thresholds were chosen such that the difference between optimized geometries were less on average than the differences between the average errors of the DFT methods chosen.

As meta-GGAs have been shown to exhibit oscillations in intermolecular potential energy surfaces of dispersion-bound complexes,<sup>168</sup> we have adopted a dense integration grid (150 radial points, 434 spherical points) for all density functional computations. To reduce the computational expense of DFT computations incurred by employing dense integration grids, the density-fitting approximation was applied to the electron repulsion integrals for computations performed using PSI4.<sup>169–176</sup> To examine the effect of basis set on the quality of the optimized geometries, we have employed Dunning’s correlation consistent polarized valence basis sets,<sup>22</sup> both with (aug-cc-pVXZ; X = D, T) and without (cc-pVXZ; X = D, T) augmentation by diffuse functions. For the convenience of the reader, these basis sets will be abbreviated as aXZ and XZ, respectively.

For this work, we take the originally published geometries for each A21 complex as benchmarks.<sup>177</sup> These were optimized by minimizing the numerical gradient of the counterpoise-corrected (CP) CCSD(T) interaction energy for each complex, at the complete basis set (CBS) limit. These interaction energies were estimated using the popular focal-point composite approach,<sup>23,24</sup> whereby the CBS limit estimate of the total MP2 energy [computed using the two-point extrapolation scheme of Helgaker,<sup>20</sup> denoted as MP2/CBS(aXZ, a[X+1]Z)], is corrected for higher-order correlation effects by adding the difference between CCSD(T) and MP2 as computed in a smaller basis set (denoted  $\delta_{\text{MP2}}^{\text{CCSD(T)}}$ ). In particular, these benchmarks were computed at the MP2/CBS(aTZ, aQZ) +  $\delta_{\text{MP2}}^{\text{CCSD(T)}}$ /aDZ level; this treatment will be denoted here as CCSD(T)/[aTQZ;  $\delta$ :aDZ]. This approach has been widely applied to estimate the CCSD(T) complete basis set limit for interaction ener-

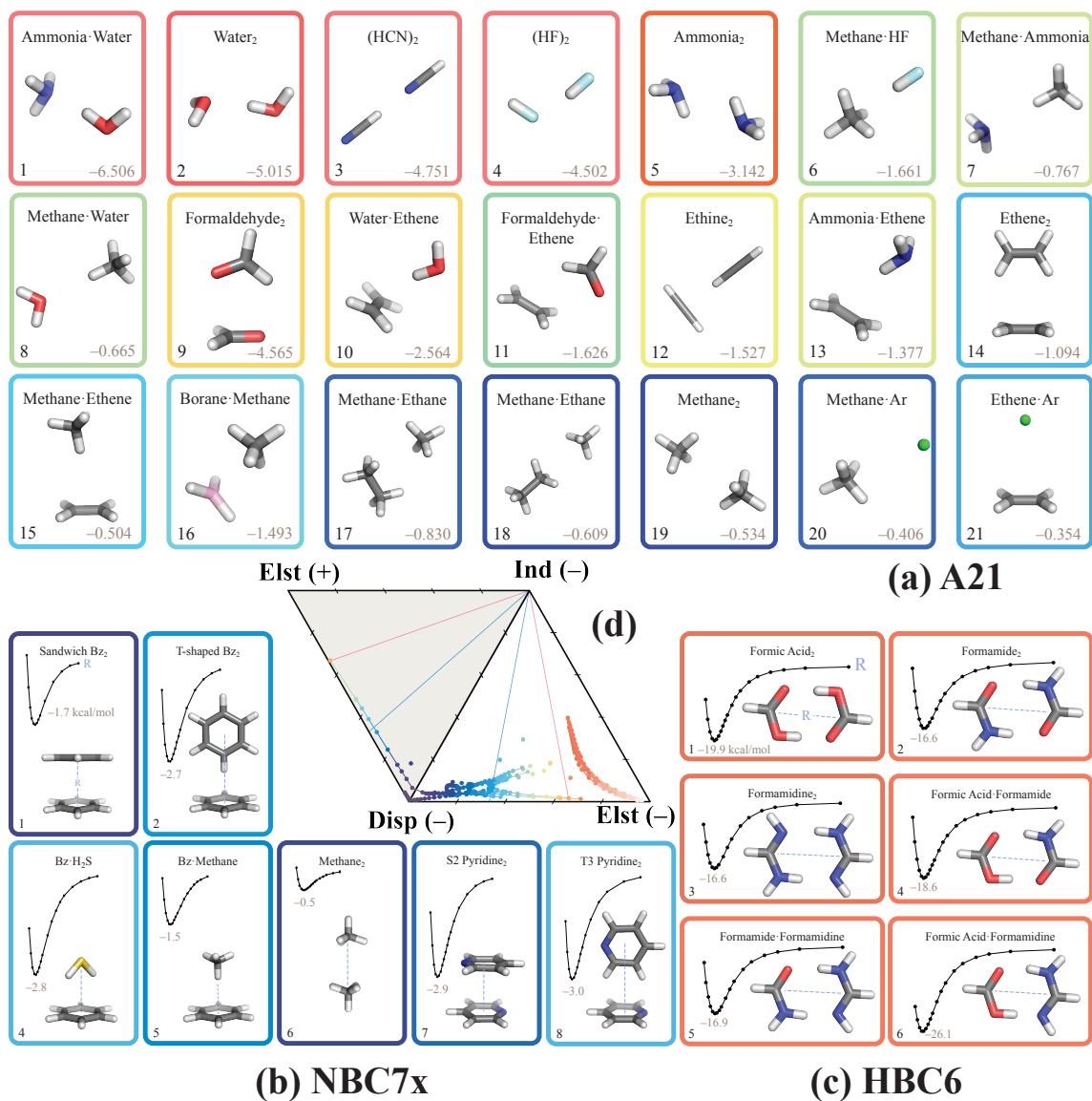


Figure 5.1: Test sets of bimolecular complexes examined here. (a) A21: 21 bound complexes contained in the A24 test set of Hobza and co-workers,<sup>177</sup> (b) NBC7x: seven (recently extended<sup>39</sup>) radial potential scans from the NBC10 test set<sup>140</sup> and (c) HBC6:<sup>30,145</sup> radial potential scans for six doubly hydrogen bonded complexes. Indicated by box coloring [(a)–(c)] or by dot color [(d)] are the noncovalent interaction type for each complex, reported previously:<sup>57,77</sup> red for electrostatic interactions, blue for dispersion interactions, and yellow/green for mixed electrostatic and dispersion interactions. The ternary diagram (d) further indicates the relative magnitude of the interaction energy components for these complexes,<sup>145,178</sup> by placing a colored dot according to the ratios of attractive dispersion/induction and attractive/repulsive electrostatic contributions to the total interaction energy. Proximity to each labeled vertex indicates an increasing fraction of the attraction (repulsion) arising from that component.

The quality of optimized supermolecular geometries computed using each model chemistry examined here is assessed according to the following metrics:

- (i) the center-of-mass displacement ( $\Delta\text{COM}$ ) between the monomers comprising each complex, compared to the  $\Delta\text{COM}$  in the corresponding benchmark geometry, and
- (ii) the least root mean square deviation (LRMSD) between the optimized geometry and the accompanying benchmark geometry.

For the purposes of this work, we consider both  $\Delta\text{COM}$  and LRMSD metrics with values less than 0.1 Å to correspond to “satisfactory” optimizations. However, we expect that for in many applications, larger errors of LRMSD  $\approx$  0.15–0.2 Å and  $\Delta\text{COM} \approx \pm$  0.1–0.15 Å could remain acceptable. A more detailed discussion regarding the assessment of optimization quality using the metrics listed above and the choice of optimization thresholds is presented in Sections I A & B of the Supporting Information.

### 5.3.2 Radial Potential Surface Scans of HBC6 and NBC10x systems

In order to assess the generality of the conclusions drawn for the geometry optimizations of the A21 test set, the ability of DFT to reproduce optimal intermolecular separation distances between monomers in complexes from the NBC7x<sup>39,140</sup> and HBC6<sup>30,145</sup> test sets (visualized in Fig. 5.1b & c, respectively) was examined. To do this, radial potential energy surface scans of the selected bimolecular complexes were constructed from both counterpoise-corrected<sup>54</sup> (CP) and uncorrected (unCP) interaction energies computed with each combination of density functional and basis set examined above, using a 0.1 Å step size. In order to estimate the optimal intermolecular separation for each curve, a second-degree polynomial was fit to the three ( $R, IE$ ) points straddling the well minima using the Numerical Python (NumPy) library,<sup>179</sup> which was subsequently used to interpolate the optimal center-of-mass displacement ( $R_{eq}$ ) for each complex. We use here the label  $R_{eq}$  to distinguish these interpolated intermolecular separation distances from the center-of-mass displacements ( $\Delta\text{COM}$ ) reported for optimized A21 complexes, since (i) the separation co-

ordinate used to construct curves in the HBC6 and NBC7x test sets does not necessarily coincide with the vector connecting monomer centers-of-mass, and (ii) to further differentiate these interpolated distances computed from interaction energy curves from the monomer center-of-mass displacements within fully optimized structures.

Benchmark values for the  $R_{eq}$  corresponding to each curve were determined by applying this same curve-fitting procedure to the HBC6 revision A<sup>30</sup> and NBC10 revision B<sup>39</sup> reference curves, respectively, each constructed from energies computed at the CCSD(T)/CBS limit. Information on these interaction energy benchmarks and revisions for HBC6 and NBC7x can be found in Table S-1 of the Supporting Information. For the formic acid dimer (HBC6-1) with all DFT model chemistries, and for the formic acid–formamidinium complex (HBC6-6, see Fig. 5.1) with all M05-2X model chemistries, radial curves do not exhibit clear potential wells within which  $R_{eq}$  could be interpolated; we have therefore removed these curves from the statistical analysis visualized in Fig. 5.5 and discussed in Section 5.4.2. We have, however, included these curves in the SI (see Section ??).

## 5.4 Results and Discussion

We first assess each model chemistry for its ability to produce optimal geometries for A21 complexes, in Section 5.4.1, before examining the generality of these conclusions in Section 5.4.2 by computing the optimal center-of-mass displacements for radial potential energy curves of larger bimolecular complexes in the HBC6 and NBC7x test sets.

### 5.4.1 Optimization of A21 Systems

As can be seen in Fig. 5.2.a & b for the B3LYP-D3 and B97-D3 density functionals, respectively, the aDZ basis set yields the best results for LRMSD values (light blue box-and-whisker plots) for A21 complexes. Indeed, for every density functional paired with aDZ, the average value of LRMSD,  $\mu$ LRMSD, is  $\leq 0.05 \text{ \AA}$  and, for all A21 complexes except a single outlier (LRMSD =  $0.15 \text{ \AA}$  for the water–ethene complex with M05-2X/aDZ; see

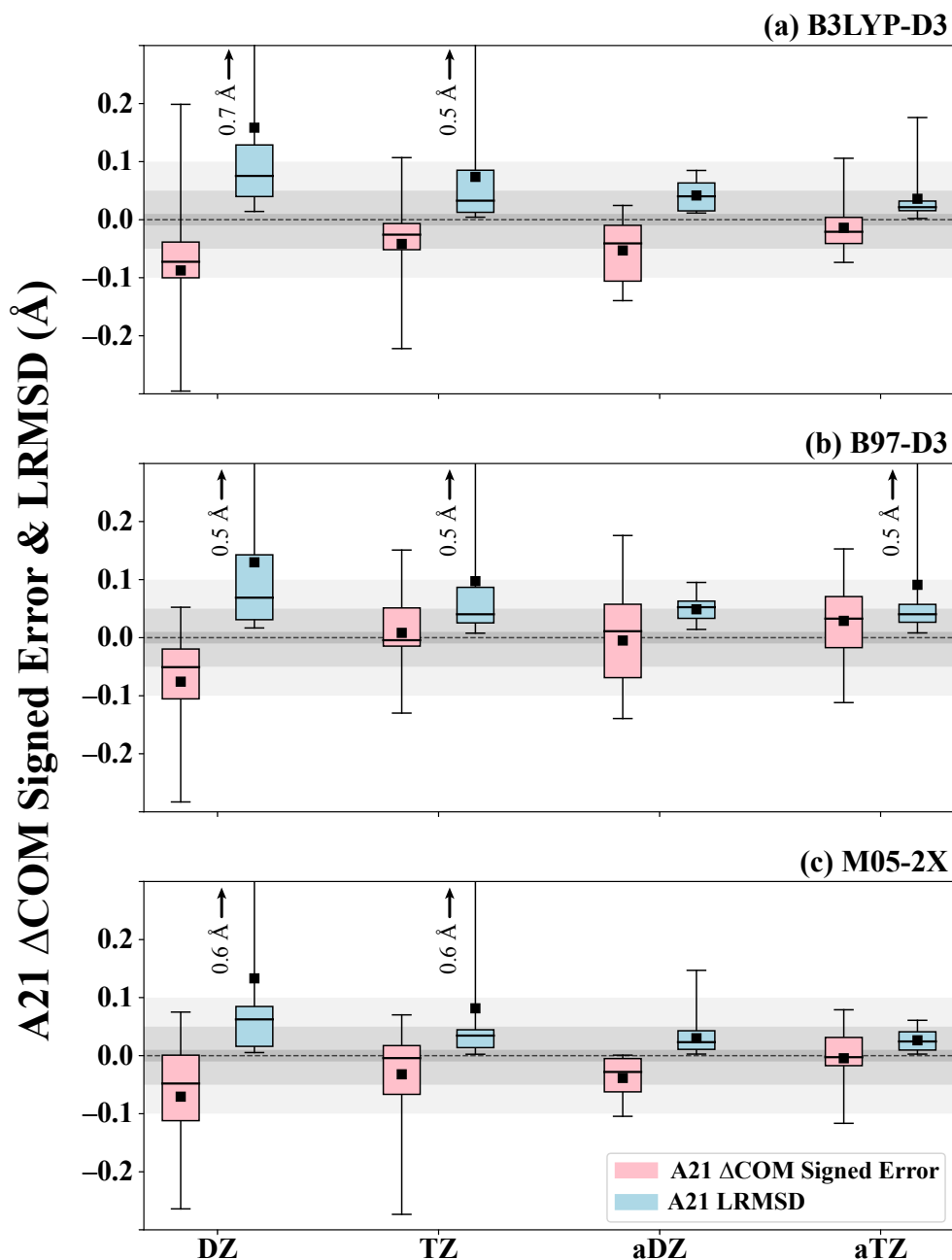


Figure 5.2: Box-and-whisker plots representing both  $\Delta\text{COM}$  signed errors (boxes shaded pink) and LRMSD values (boxes shaded blue) for systems in the A21 test set, optimized using the (a) B3LYP-D3, (b) B97-D3, and (c) M05-2X density functionals together with the DZ, TZ, aDZ, and aTZ basis sets. Boxes encompass the first (Q1) through third (Q3) quartiles of each data set, with values corresponding to the median (Q2) and mean LRMSD and  $\Delta\text{COM}$  signed error indicated as a solid black bar and black square, respectively. Whiskers encompass the full range of LRMSD values and  $\Delta\text{COM}$  signed errors; maximum values are indicated when whiskers surpass the area shown. For reference, a dotted line indicates  $0.0 \text{ \AA}$ , and three levels of shading are provided: light grey encompassing  $\pm 0.1 \text{ \AA}$ , medium grey encompassing  $\pm 0.05 \text{ \AA}$ , and medium-dark grey encompassing  $\pm 0.01 \text{ \AA}$ .

Fig. 5.2.c),  $\text{LRMSD} \leq 0.1 \text{ \AA}$ . The aTZ basis set also exhibits good performance, with inner quartile ranges (IQRs) of less than  $0.03 \text{ \AA}$  for each density functional. Despite this good performance for the majority of A21 systems, four complexes optimized with B97-D3/aTZ exhibit  $\text{LRMSD} \geq 0.1 \text{ \AA}$  [methane–water (A24-8;  $\text{LRMSD} = 0.5 \text{ \AA}$ ), ammonia–ethene (A24-13;  $\text{LRMSD} = 0.1 \text{ \AA}$ ),  $C_s$  methane–ethane (A24-15;  $\text{LRMSD} = 0.2 \text{ \AA}$ ), and borane–methane (A24-16;  $\text{LRMSD} = 0.5 \text{ \AA}$ )], as opposed to just a single complex ( $C_s$  methane–ethane, A24-15;  $\text{LRMSD} = 0.2 \text{ \AA}$ ) for B3LYP-D3/aTZ, and no such complexes for M05-2X/aTZ. B3LYP-D3/aTZ and M05-2X/aTZ have slightly smaller IQRs and  $\mu\text{LRMSD}$  values than for aDZ, but the improvement is quite small. For M05-2X, a more noticeable improvement is observed in the overall range of LRMSD values, which decreases from  $0.14 \text{ \AA}$  for aDZ to only  $0.06 \text{ \AA}$  for aTZ.

For the  $\Delta\text{COM}$  metric, visualized in Fig. 5.2 with light red box-and-whisker plots, signed errors (SE) for individual model geometries and mean signed errors (MSEs) over all A21 complexes are not so clearly superior for the aDZ basis as was observed for the LRMSD metric, and  $\Delta\text{COM}$  IQRs seem to be largely comparable between basis sets for each density functional. In fact, for both B3LYP-D3 and M05-2X, while the overall ranges of  $\Delta\text{COM}$  signed errors are slightly smaller when using the aDZ basis set, the MSEs for these functionals are smallest with the aTZ basis set. For B97-D3, each of the aDZ, TZ, and aTZ basis sets yield error ranges which are nearly identical, with maximum and minimum signed errors lying slightly outside the target range of  $\pm 0.1 \text{ \AA}$ . The mean signed errors for these combinations benefit from this nearly symmetric distribution; both B97-D3/aDZ and B97-D3/TZ exhibit  $\text{MSE} \leq \pm 0.01 \text{ \AA}$ , and B97-D3/aTZ is not much worse, with  $\text{MSE} = 0.03 \text{ \AA}$ .

Generally, however, the relative quality of model geometries with respect to the  $\Delta\text{COM}$  metric is again not *significantly* improved when moving from aDZ to aTZ basis sets. IQRs improve slightly for each functional; however, the overall range of  $\Delta\text{COM}$  values increases for both B3LYP-D3 and M05-2X functionals. Regardless of this increase in the total range



for these model chemistries, the highly symmetric distribution of  $\Delta\text{COM}$  values about  $0.0 \text{ \AA}$  signed error yields very small MSEs, with  $\text{MSE} = -0.01, 0.00 \text{ \AA}$  for B3LYP-D3/aTZ and M05-2X/aTZ; the double- $\zeta$  counterparts are not much worse, however, with  $\text{MSE} = -0.05, -0.04$  for B3LYP-D3/aDZ and M05-2X/aDZ. B97-D3, on the other hand, exhibits the opposite trend when moving from aDZ to aTZ: while the total range of  $\Delta\text{COM}$  values improves from  $0.32 \text{ \AA}$  to  $0.26 \text{ \AA}$ , MSE increases slightly, from  $\text{MSE} = 0.00 \text{ \AA}$  to  $\text{MSE} = 0.03 \text{ \AA}$ . The similar or only marginally improved performance of optimizations utilizing aTZ over aDZ for both LRMSD and  $\Delta\text{COM}$  metrics, together with the increased computational cost associated with the increase in  $\zeta$ -level, implies that the smaller aDZ basis set is generally preferable for optimizations of nonbonded complexes similar to those within the A21 test set.

To more closely examine the performance of each density functional for optimizing A21 complexes, we next consider the quality of individual equilibrium geometries optimized using the recommended basis set, aDZ. Visualized in Fig. 5.3 are values corresponding to each A21 complex, and box-and-whisker plots describing these values' distributions for (a) LRMSD and (b)  $\Delta\text{COM}$  metrics. Within each density functional, electrostatically bound complexes (HB subset, red circles) exhibit LRMSD values and  $\Delta\text{COM}$  signed errors that are more clustered than for the other two A21 subsets [dispersion-dominated (DD) subset (blue circles) and mixed-interaction (MX) systems (green circles)]. Among these density functionals, M05-2X/aDZ generates model geometries that are notably superior for HB and DD complexes, with respect to both LRMSD and  $\Delta\text{COM}$  metrics; this model chemistry is also slightly superior to B3LYP-D3/aDZ and B97-D3/aDZ for MX systems, with the lone exception being the water–ethene complex, exhibiting  $\text{LRMSD} = 0.15 \text{ \AA}$ . For B3LYP-D3/aDZ,  $R_{eq}$  is underestimated in nearly all A21 model geometries (negative signed error for  $\Delta\text{COM}$ ), while M05-2X/aDZ generally underestimates  $R_{eq}$  for A21 complexes. B97-D3/aDZ model geometries, on the other hand, exhibit different behavior for the  $\Delta\text{COM}$  metric depending on the interaction type of the complex;  $R_{eq}$  is typically over-

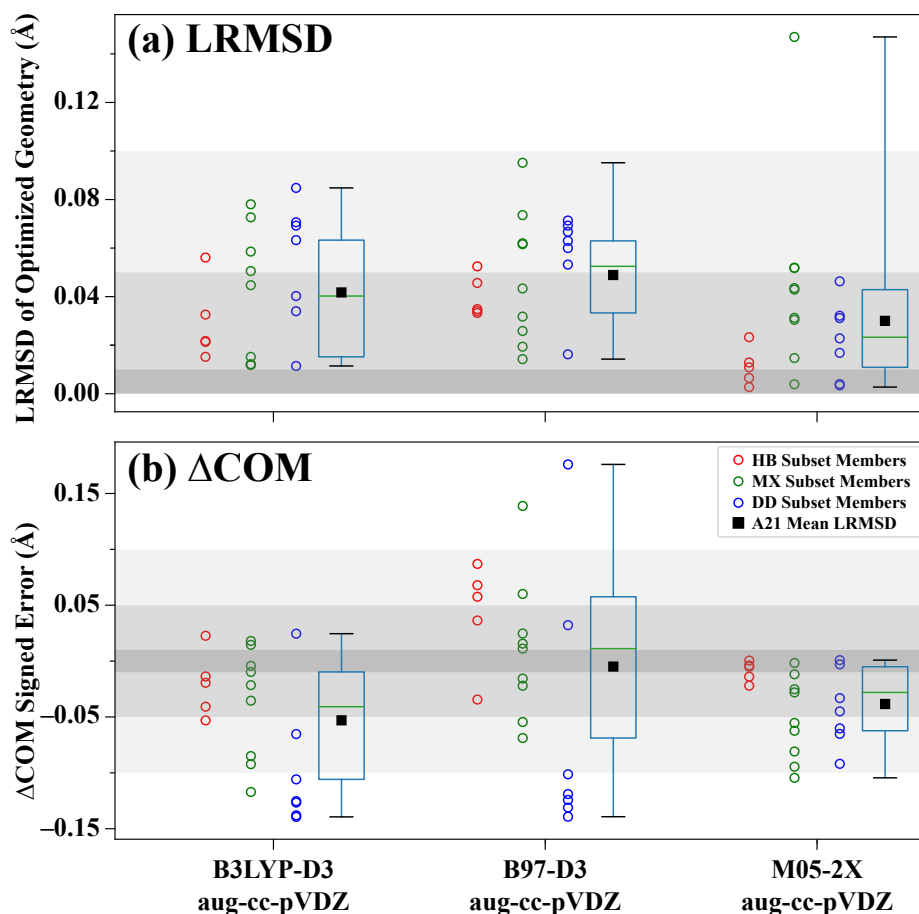


Figure 5.3: Values corresponding to individual A21 complexes and box-and-whisker plots detailing test-set-wide distributions for (a) least root mean square displacements (LRMSD) and (b) signed errors in center-of-mass distance ( $\Delta\text{COM}$ ), computed with each density functional using the aug-cc-pVDZ basis set. Values for individual A21 systems, shown with empty circle markers, are grouped and colored according to interaction motif:<sup>57,77</sup> red for electrostatically bound complexes (HB subset), blue for dispersion bound complexes (DD subset), and green for mixed interaction complexes (MX subset). For box-and-whisker plots of each model chemistry, boxes encompass the first (Q1) through third (Q3) quartiles of each data set, with values corresponding to the median (Q2) and mean LRMSD and  $\Delta\text{COM}$  signed error indicated as a solid green bar and black square, respectively. Whiskers encompass the full range of LRMSD values and  $\Delta\text{COM}$  signed errors; maximum values are indicated when whiskers surpass the area shown. For reference, a dotted line indicates 0.0 Å, and three levels of shading are provided: light grey encompassing  $\pm 0.1$  Å, medium grey encompassing  $\pm 0.05$  Å, and medium-dark grey encompassing  $\pm 0.01$  Å.

estimated in HB systems, underestimated in DD systems, and no trend is observed for the MX subset.

Based on the performance of these model chemistries for producing optimal geometries of A21 test set, a total ordering which ranks the performance of the best such model chemistries can be constructed. Here, we consider first the  $\mu$ LRMSD and  $\Delta$ COM MSE statistics for each set of values, then the sample inner quartile and total ranges; these considerations produce the following ordering:

$$\text{M05-2X/aTZ} \sim \text{B3LYP-D3/aTZ} \gtrsim \text{M05-2X/aDZ} \gtrsim \text{B3LYP-D3/aDZ} \succ \text{B97-D3/aDZ}$$

where “ $\sim$ ” indicates roughly equivalent performance of sample means and IQR, “ $\gtrsim$ ” indicates superior performance with respect to sample mean, but roughly similar performance in IQR, and “ $\succ$ ” indicates superior performance in both sample means and IQR. The reason for the classification of B97-D3/aDZ as inferior to B3LYP-D3/aDZ and M05-2X/aDZ, despite a seemingly excellent sample mean for  $\Delta$ COM signed errors, is the larger spread of the errors for B97-D3/aDZ; indeed, the range in the  $\Delta$ COM signed error is 0.32 Å for B97-D3/aDZ versus 0.16 Å for B3LYP-D3/aDZ and 0.11 Å for M05-2X/aDZ. While this wider distribution of errors cancels fortuitously for B97-D3/aDZ to produce a very low MSE (0.00 Å), the mean absolute error (MAE) for this model chemistry is nearly double that of B3LYP-D3/aDZ, with MAE = 0.06, 0.03 Å, respectively. Despite the presence of some cases with errors slightly larger than the target value for B97-D3/aDZ, each density functional is able to produce equilibrium geometries of the desired accuracy level for a significant percentage of the A21 complexes when paired with the aDZ basis set.

#### 5.4.2 Prediction of Optimal Intermolecular Separation in NBC7x and HBC6 Interaction

##### Energy Scans

As illustrated in Fig. 5.4 for the formamidine dimer, optimal intermolecular separations ( $R_{eq}$ ) were interpolated from radial potential scans constructed for the 13 complexes in the HBC6 and NBC7x test sets using both counterpoise-corrected (CP) and uncorrected

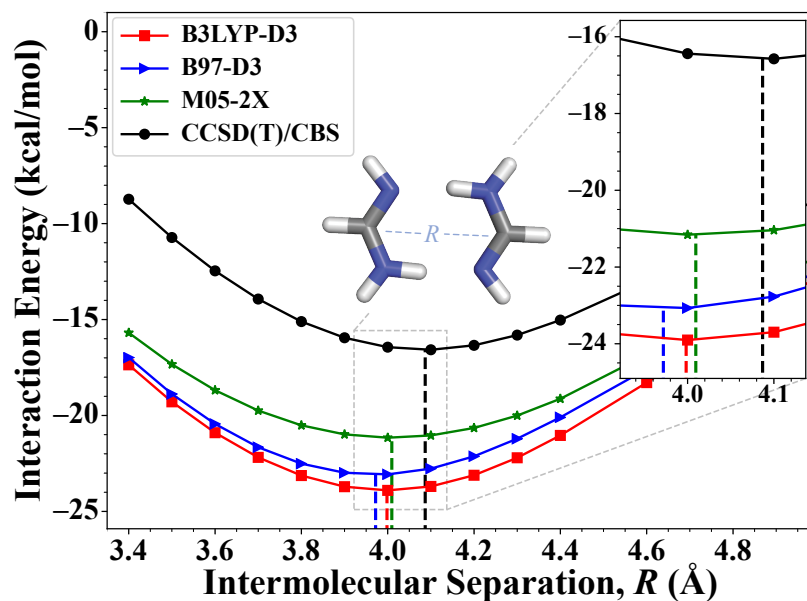


Figure 5.4: Scans of the non-counterpoise-corrected interaction energy (unCP IE) along the radial separation coordinate  $R$  in the formamidine dimer (HBC6-3; inset shown) computed with B3LYP-D3 (red), B97-D3 (blue) and M05-2X (green) using the cc-pVDZ basis set. The interpolated optimal intermolecular separation for each curve is indicated with a vertical dotted line in the same colors. For reference, a curve constructed from the CCSD(T)/CBS benchmark IEs at each value of  $R$  is presented in black.

(unCP) interaction energies, computed using each combination of density functional and basis set examined above. The CCSD(T)/CBS reference curve, from which the reference  $R_{eq}$  value is interpolated, is also shown; for a complete set of equivalent figures (108 total), please refer to the Supporting Information. Provided in Fig. 5.5 are box-and-whisker plots describing the distribution of signed errors of these interpolated minima for each DFT model chemistry, as compared to the minima interpolated from reference curves. Regardless of the choice of BSSE treatment (either CP or unCP), interpolated minima for curves in the NBC7x test set exhibit slightly larger signed errors than those for HBC6 curves; while all model chemistries produce  $MSE \leq \pm 0.06 \text{ \AA}$  for both CP and unCP curves within the HBC6 test set (and twelve model chemistries with  $MSE \leq \pm 0.01 \text{ \AA}$ !), several model chemistries for NBC7x complexes yield MSE slightly outside this range, albeit still within the target of  $R_{eq} \leq \pm 0.1 \text{ \AA}$ . Errors in interpolated  $R_{eq}$  for CP-curves are largely independent of basis set size or augmentation for both NBC7x and HBC6, with the lone exception

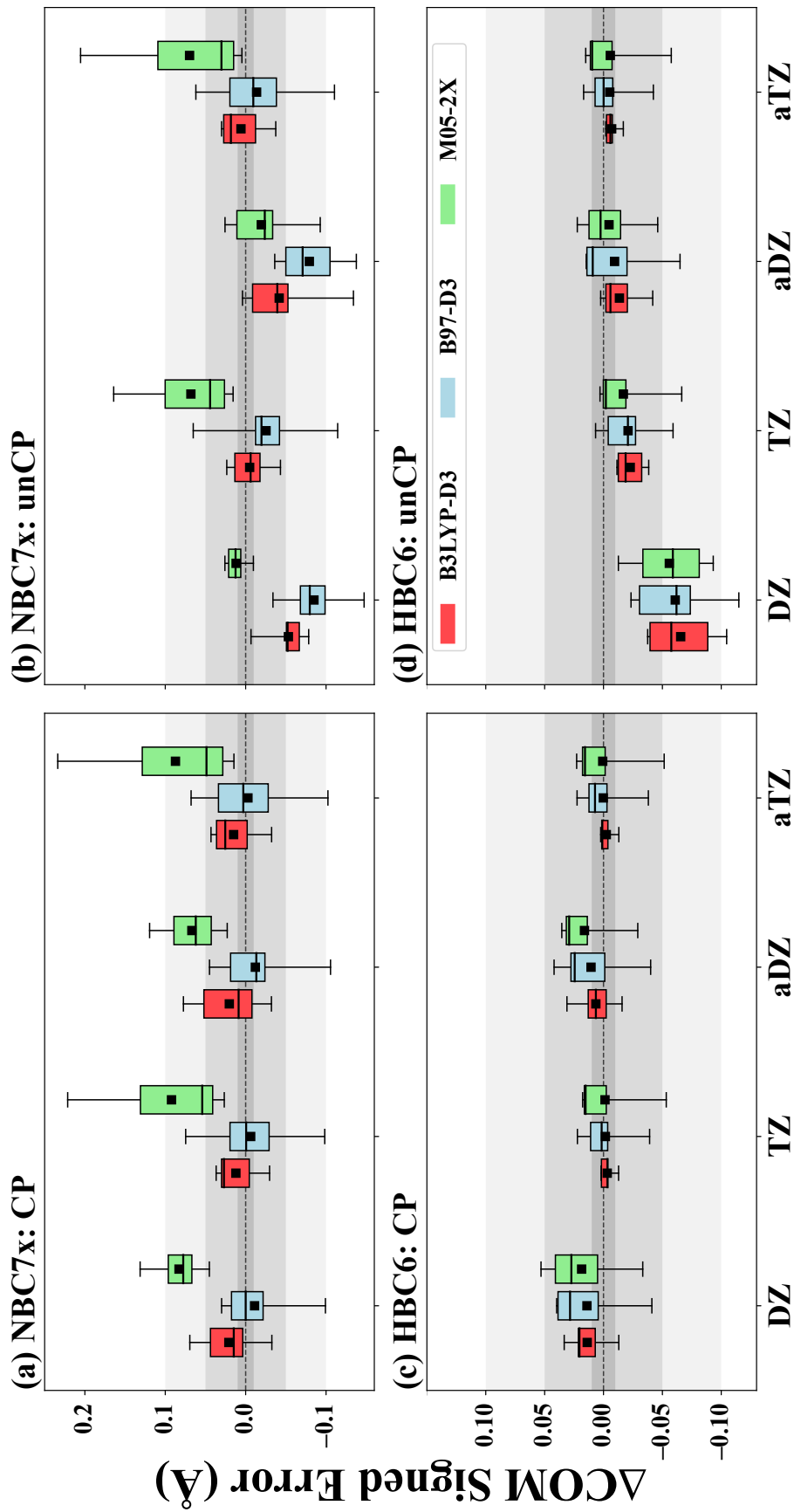


Figure 5.5: Box-and-whisker plots depicting the distribution of signed error in interpolated optimal center-of-mass displacement ( $\Delta\text{COM}$ ) for radial interaction energy curves in the NBC7x (a & b) and HBC6 (c & d) test sets. For both test sets, box-and-whisker plots representing curves constructed from both counterpoise-corrected (CP; left panels a & c) and uncorrected (unCP right panels, b & d) IEs are given. Whiskers encompass the full range of  $\Delta\text{COM}$  signed errors for the indicated test set, correction scheme, and model chemistry, and boxes illustrate the first (Q1), second (median, black bar), and third quartiles (Q3) of these data; additionally, the mean signed error for each data set is indicated with a black square. For reference, a dotted line indicates 0.0 Å, and three levels of shading are provided: light grey encompassing  $\pm 0.1$  Å, medium grey encompassing  $\pm 0.05$  Å, and medium-dark grey encompassing  $\pm 0.01$  Å.

of HBC6 curves with TZ exhibiting improved MSE and IQR over DZ. For unCP curves, however, the quality of interpolated  $R_{eq}$  is sensitive to both  $\zeta$ -level and augmentation.

Considering next the performance of individual model chemistries, B3LYP-D3 and B97-D3 perform about the same (CP curves) or better (unCP curves) for NBC7x complexes with a triple- $\zeta$  basis set, regardless of augmentation, while for M05-2X, the double- $\zeta$  basis sets are better. For HBC6, interpolated minima are quite good for all model chemistries considered, although errors are slightly larger for unCP curves using the DZ basis set. For both test sets, B3LYP-D3 and B97-D3 exhibit excellent performance for  $R_{eq}$ ; among the 32 total combinations of these two functionals, the four examined basis sets, and two possible BSSE treatments,  $MSE \leq 0.01 \text{ \AA}$  for  $R_{eq}$  of 16 model chemistries,  $0.01 \leq MSE \leq 0.05 \text{ \AA}$  for  $R_{eq}$  of 11 model chemistries, and the remaining five model chemistries all produce interpolated  $R_{eq}$  values with  $0.05 \leq MSE \leq 0.1 \text{ \AA}$ . While M05-2X yields high quality interpolated  $R_{eq}$  for HBC6 systems, this functional tends on average to underbind for CP curves of NBC7x systems, as well as for unCP curves with triple- $\zeta$  basis sets (see, e.g., Figs. S-76–S-82 and S-111–117 in the Supplemental Information), leading to a slight overestimation of the optimal intermolecular separation distance  $R_{eq}$ ; the MSE for these cases ranges between 0.07–0.1  $\text{\AA}$ .

Interestingly, NBC7x systems with  $\pi - \pi$  stacking (NBC7x-1 & 7; sandwich benzene and pyridine dimers, respectively) seem particularly susceptible to this drastic underbinding by M05-2X, exhibiting deviations from reference interaction energies as large as +2  $\text{kcal mol}^{-1}$  in the neighborhood of the minima of CP curves, a full factor of four larger than deviations exhibited by either B97-D3 or B3LYP-D3 (see, e.g., Figs. S-62 and S-67 in the Supplemental Information). The T-shaped counterparts to these complexes (NBC7x-2 & 8 for benzene and pyridine dimers, respectively), while slightly less underbound, still exhibit significant deviations from CCSD(T)/CBS IEs of nearly 1  $\text{kcal mol}^{-1}$  in the neighborhood of the curve minima (Figs. S-63 and S-68). While these T-shaped complexes are much more realistically described by M05-2X with double- $\zeta$  basis sets when not employing the

counterpoise-correction procedure, the sandwich configurations are similarly underbound for unCP curves even with double- $\zeta$ , and both sandwich and T-shaped configurations are severely underbound in unCP curves with triple- $\zeta$  basis sets. Strangely, this underbinding of systems involving  $\pi$ -stacking or CH- $\pi$  interactions by M05-2X is not present in other dispersion-dominated NBC7x systems, e.g., methane dimer, where M05-2X underbinds either by only tenths of kcal mol<sup>-1</sup> for CP curves (e.g., Fig. S-66) or nearly perfectly reproduces CCSD(T)/CBS IEs for unCP curves (e.g., Fig. S-73); these results are consistent with values reported previously in Ref. 75.

Finally, we examine in particular the performance of aDZ model chemistries for unCP curves, as our optimizations of A21 complexes did not employ BSSE corrections, and aDZ was found to perform similarly to aTZ. For these larger systems, MSE  $\leq$  0.02 Å for HBC6 curves, and MSE  $\leq$  0.05 Å for NBC7x curves constructed with B3LYP-D3 and M05-2X. B97-D3/aDZ performs only slightly worse for NBC7x curves, with MSE = -0.08 Å; each of these model chemistries, however, produce  $R_{eq}$  MSEs within the target range of accuracy. This indicates that the high quality of DFT/aDZ geometries observed in Section 5.4.1 for the optimization of A21 systems is likely to generalize to the optimization of systems as large as those in HBC6 or NBC7x, and perhaps somewhat larger. We must note, however, that in systems that are much larger, long-range effects neglected by M05-2X and many-body dispersion interactions neglected by all of the approaches tested here may become significant.<sup>150–154</sup>

## 5.5 Summary & Conclusions

We have shown that each of B3LYP-D3, B97-D3, and M05-2X density functionals paired with Dunning’s aug-cc-pVDZ (aDZ) basis set combine accuracy and reasonable computational expense for producing equilibrium geometries of 21 small bimolecular van der Waals complexes from the A24 test set. Each DFT/aug-cc-pVDZ level of theory performs well compared to CCSD(T)/CBS references: both B3LYP-D3/aDZ and M05-2X/aDZ consis-

tently yield equilibrium geometries with very small least root-mean-square displacement (LRMSD) and center-of-mass displacement signed error ( $\Delta\text{COM SE}$ ), each within 0.05 Å on average. B97-D3/aDZ nearly as good, but exhibits slightly larger range of  $\Delta\text{COM SE}$ . We have also shown that these DFT/aDZ combinations are capable of reproducing optimal intermolecular separation distances ( $R_{eq}$ ) interpolated from radial interaction energy scans of 13 larger complexes in the HBC6 and NBC7x test sets. Minima interpolated from curves in both HBC6 and NBC7x test sets constructed from non-counterpoise-corrected (unCP) interaction energies computed using DFT/aDZ are of similarly high quality, with B3LYP-D3/aDZ and M05-2X/aDZ yielding minima within 0.05 Å of CCSD(T)/CBS, while B97-D3/aDZ is again nearly as good but slightly less reliable for unCP curves in NBC7x. Overall the analysis of optimized A21 systems, together with the quality of interpolated NBC7x and HBC6 minima, indicate that each of these DFT/aDZ combinations are well suited to produce equilibrium geometries of given conformations of bimolecular van der Waals complexes of diverse binding motif.

In the course of this work, we developed a software tool to maximally align the approximate geometries optimized with DFT against the CCSD(T) benchmark geometries, to obtain the LRMSD metric for each A21 complex analyzed above. This tool, consisting of Python implementations of two general algorithms solving the maximal alignment problem, together with all data presented here and all Python code necessary to perform the above data analysis and visualization, are available free of charge via an open-source GitHub repository at [www.github.com/cdsgroup/dftoptbench-si](http://www.github.com/cdsgroup/dftoptbench-si). All software contained in this repository can be executed without local installation via the Jupyter Hub cloud server, or cloned to be used offline.



## **PART III**

# **DEVELOPING APPROXIMATE PERTURBATIVE METHODS FOR NON-COVALENT INTERACTIONS**

## CHAPTER 6

### OPTIMIZED DAMPING PARAMETERS FOR EMPIRICAL DISPERSION CORRECTIONS TO SYMMETRY-ADAPTED PERTURBATION THEORY

#### 6.1 Abstract

Symmetry adapted perturbation theory (SAPT) has become a valuable computational tool offering physical insight into the fundamental nature of non-covalent interactions in diverse chemical systems by directly computing the electrostatics, exchange (steric) repulsion, induction (polarization), and London dispersion contributions to the interaction energy using quantum mechanics. Further application of SAPT to novel chemical problems is limited primarily by its computational expense, where even for its most affordable variant, SAPT0, computing the London dispersion contribution to the interaction energy (IE) scales as the fifth power of system size, [ $\mathcal{O}(N^5)$ ]. Here we optimize damping parameters for the semi-empirical -D3 dispersion correction of Grimme and co-workers, so that they are suitable for use as replacements of the computationally expensive dispersion term in SAPT0. Parameters are obtained by fitting to a large set of 2295 interaction energies computed at the CCSD(T)/CBS level of theory. This reduces the algorithmic scaling of SAPT0 from  $\mathcal{O}(N^5) \rightarrow \mathcal{O}(N^4)$  while retaining the physically meaningful interpretation of IE components characteristic of all SAPT methods. This scaling reduction translates into a nearly  $2.5\times$  speedup over conventional SAPT0 for systems with  $\sim 300$ -atoms. Furthermore, this allows for SAPT-D computations to be performed on systems with over 450 atoms, while offering nearly equivalent accuracy to SAPT0 when compared against reference IEs for a diverse set of approximately 8,100 bimolecular complexes. We have extended our formulation of SAPT-D to be consistent with the functional group partition (F-SAPT-D) and applied this method to conclude that the difference in binding affinity for partial agonist

salbutamol to the G-protein coupled  $\beta_1$ -adrenergic receptor between active and inactive forms is due to the cooperative effects of both peptide bonds and residues outside the immediate binding pocket, indicating that a local contact model for protein–ligand binding is insufficient to discriminate between binding conformations which possess similar activities.\*

## 6.2 Introduction

Symmetry-Adapted Perturbation Theory (SAPT) has proven useful in computing the strength and character of interactions between molecules.<sup>133–136</sup> SAPT computes the physical components of the interaction (electrostatics, induction/polarization, London dispersion forces, and exchange repulsion) directly, not as a decomposition or difference of total energies. SAPT has been formulated in terms of many-body perturbation theory,<sup>133,134,181</sup> coupled-cluster theory,<sup>182,183</sup> and density functional theory.<sup>184,185</sup> The highest orders of SAPT include terms analogous to the perturbative triples correction in the popular CCSD(T) method,<sup>10</sup> and exhibit similar accuracy (and computational expense).<sup>58,134,135</sup> Fortunately, even the simplest SAPT treatments can be reasonably accurate and can provide insight into the nature of intermolecular interactions. The lowest-order truncation of the perturbation series, sometimes referred to as SAPT0, uses a Hartree–Fock treatment of the monomers and treats the intermolecular perturbation through second order.<sup>133</sup> When paired with the jun-cc-pVDZ basis set, which is a truncation of the aug-cc-pVDZ basis in which diffuse functions on H atoms and diffuse  $d$  functions on heavy atoms are neglected, SAPT0 is reasonably accurate (mean absolute error of 0.86 kcal mol<sup>-1</sup> over four high-quality benchmark test sets, or only 0.49 kcal mol<sup>-1</sup> using an exchange-scaled variant labeled  $s$ SAPT0).<sup>58</sup> SAPT0, formally scaling as  $\mathcal{O}(N^5)$ , where  $N$  is proportional to the size of the molecular system, is applicable to systems of a few hundred atoms when density-fitting techniques are employed.<sup>186,187</sup> However, computations of this size can be time consuming, and for

---

\*This Chapter reproduces the work in Ref. 180.

applications to protein-ligand interactions, molecular crystals, solvated species, etc., one may wish to go to even larger systems.

Parrish *et al.* recently introduced a very efficient implementation of SAPT0 utilizing graphics processing units (GPUs),<sup>188</sup> which they demonstrated on the entire indinavir/HIV-II protease complex. To reach systems of this size, they dropped all diffuse functions from the orbital basis, formulated the SAPT energy in terms of only potential integral primitives which were computed with a highly specialized, mixed-precision algorithm leveraging both CPU and GPU architectures, and avoided the  $\mathcal{O}(N^5)$  scaling of the dispersion part of the computation, replacing the dispersion term with force-field-type pairwise-atomic  $-C_{6,ij}/R_{ij}^6$  terms. Such force-field dispersion models had been previously utilized in SAPT<sup>189–191</sup> and in Hartree–Fock treatments of intermolecular interactions.<sup>192–194</sup> Parrish *et al.* used -D corrections from the very popular DFT-D3 approach of Grimme and co-workers,<sup>35</sup> which provides  $C_6$  coefficients and parameters for functions to damp the corrections at short range. Although the -D3 corrections were formulated for use with DFT, the DFT-D3 program<sup>195</sup> from the Grimme group includes damping parameters for Hartree–Fock, and these are the damping parameters that were used by Parrish *et al.*<sup>188</sup>

We have previously found that the damping parameters used in the -D3 method do not perform as well for short-range contacts, and we recommended modified parameters based on reparameterization using a much larger training set (1526 data points vs. an original 130).<sup>39</sup> Here, we perform an analogous investigation of optimal damping parameters for -D3 corrections as replacements for the dispersion energies of SAPT0, which are the rate-determining step in SAPT0 computations. We reparameterize the damping functions by training against a large set of 2295 interaction energies estimated at the CCSD(T) complete-basis-set limit, and the updated parameters provide around a XXX reduction in error compared to the original Hartree–Fock damping parameters. We also illustrate the savings in computer time, storage and required memory that are possible for larger systems by replacing the SAPT0 dispersion terms with -D corrections, using the implementation

in the Psi4 program package.<sup>31</sup> We investigate this through calculations on dimers of regularly increasing size obtained by fragmenting a simplified, publicly available cocrystal of C(30) carotenoid dehydrosqualene synthase in complex with the BPH-673 ligand (PDB ID 3ACX).<sup>196</sup>

Finally, as an illustration of the utility of our approach, we have adapted our implementation to be compatible with the functional group partition of SAPT (F-SAPT),<sup>66</sup> (denoted F-SAPT-D), and we have used it to address the question of the differential binding of ligands to the activated and inactive states of G-protein coupled receptors (GPCRs), a subject of biological relevance but heretofore inaccessible to this method due to the associated computational requirements. GPCRs can exist in an ensemble of conformations, but they generally bind with higher affinity for agonists when they are in an active (bound to G proteins) versus inactive state. Recently, Warne *et al.*<sup>197</sup> provided a potential means to understanding this difference by solving the crystal structures for the  $\beta_1$ -adrenergic receptor ( $\beta_1$ AR) bound to various agonists and partial agonists, while stabilized in an active versus inactive state by combining the GPCRs with conformation-specific nanobodies. A detailed but qualitative comparison of the structures led to the conclusion that an overall contraction of the binding site, as well as changes in contact distances between ligands and particular binding-site residues were responsible for the increased potency of active  $\beta_1$ AR for full and partial agonists. Here, we apply F-SAPT-D to one example discussed at length by Warne *et al.*,<sup>197</sup>  $\beta_1$ AR bound to the partial agonist salbutamol, for which the experimentally obtained increase in affinity for the ligand by active versus inactive state of GPCR is 1.88 log units (76-fold, or 2.6 kcal mol<sup>-1</sup>). We first verify that the method qualitatively matches the experimentally observed increase in affinity, and then obtain a quantum-mechanically-based assessment of which interactions contribute primarily to the GPCR-conformation-specific preference.

Table 6.1: Datasets utilized in the training and validation sets. All benchmark datasets are of MP2/CBS +  $\Delta$ CCSD(T)/aDZ quality or better. For further details of reference levels of theory for each dataset, please refer to Table S8 in the Supplementary Materials of Ref. 39.

Database	Points	Curves	Largest*	Reference	Description
<b>Training</b>	<b>2295</b>	<b>Curves</b>	<b>104</b>		
ACHC	25	6	19	198	rise, twist, slide, shift, roll, and tilt of adenine:cytosine nucleobase step
HBC6	118	6	6	145, 30	dissociation curves of doubly hydrogen-bonded (HB) complexes
NBC10x <sup>34</sup>	192	10	12	30, 75, 39	dissociation curves of dispersion-bound (DD) complexes
S22x7 <sup>34</sup>	149	22	19	79, 107, 39	dissociation curves for a mix of HB and DD complexes
S66x10 <sup>3</sup>	658	66	16	76, 39	dissociation curves for a balanced mix of biomolecule NCI bonding motifs
BBI	56	- <sup>2</sup>	10	199	peptide backbone-backbone complexes
SSI	750	- <sup>2</sup>	20	199	a representative subset of 500 structures from SSI
Water2510	347	- <sup>2</sup>	2	200-202	water dimer PES
<b>Validation</b>	<b>5789</b>	<b>Curves</b>	<b>158</b>		
ACHC	29	6	19	198	rise, twist, slide, shift, roll, and tilt of adenine:cytosine nucleobase step
C <sub>2</sub> H <sub>4</sub> · NT	75	15	26	203	ethene with curved coronene
CH <sub>4</sub> · PAH <sup>3</sup>	405	45	25	204, 118, 39	methane with PAHs the size of benzene through coronene and curved coronene
CO <sub>2</sub> · NPHAC	96	16	27	205	CO <sub>2</sub> with nitrogen-doped polyheterocyclic aromatic compounds
CO <sub>2</sub> · PAH	249	45	27	206	CO <sub>2</sub> with PAHs the size of benzene through coronene and curved coronene
Water2510	2142	- <sup>2</sup>	2	200-202	water dimer PES
X31x10 <sup>5</sup>	310	31	18	144	dissociation curves of organic halide, halohydrate, and halogen complexes
BBI	44 <sup>6</sup>	- <sup>2</sup>	10	199	peptide backbone-backbone complexes
SSI	2439 <sup>7</sup>	- <sup>2</sup>	21	199	peptide sidechain-sidechain complexes

<sup>1</sup> The largest number of heavy atoms in the dataset.

<sup>2</sup> Database does not contain curves.

<sup>3</sup> Database was extended to shorter ranges.

<sup>4</sup> Database was recomputed at a higher level of theory; see Supplemental Information, Table Sxx.

<sup>5</sup> The X40x10 database with complexes containing iodine removed.

<sup>6</sup> SSI contains 3380 complexes. The stated figure is less 750 from the SSI fitting subset.

<sup>7</sup> BBI contains 100 complexes. The stated figure is less 56 from the BBI fitting subset.

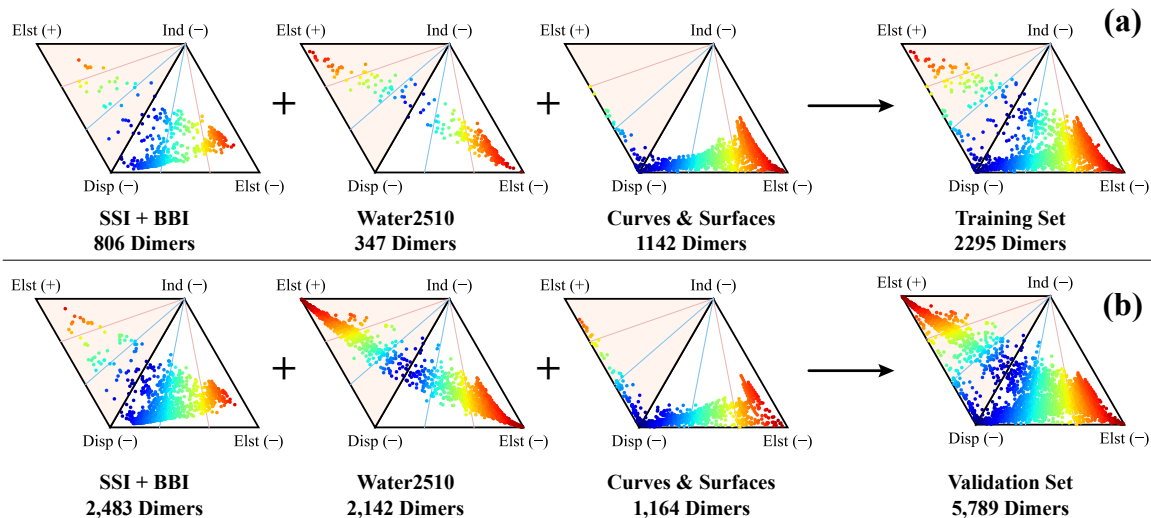


Figure 6.1: Ternary diagrams visualizing relative contributions of attractive (-) and/or repulsive (+) electrostatics, induction, and dispersion interaction energy components based on the SAPT0/jun-cc-pVDZ description of IEs for all systems comprising the (a) training and (b) validation sets for SAPT-D parameter training. Each system is represented by a single dot, colored according to the most dominant contribution to the overall SAPT0 IE: red indicating an electrostatically dominated interaction, blue indicating a dispersion dominated interaction, and yellow-green indicating an interaction for which neither electrostatics nor dispersion components dominate.

### 6.3 Theoretical & Computational Methods

#### 6.3.1 Formulation of SAPT0-D

The simplest truncation of the SAPT perturbation series, SAPT0, treats the intermolecular perturbation  $\hat{V}$  through second order and the intramolecular correlation energy through 0th order (i.e., uses a Hartree-Fock description of the monomers). The SAPT0 interaction energy is given by<sup>58</sup>

$$\begin{aligned}
 \text{IE}^{\text{SAPT0}} &= \text{IE}^{\text{HF}} + \left[ E_{\text{disp}}^{(20)} + E_{\text{exch-disp}}^{(20)} \right]_{\text{disp}} \\
 &= \left[ E_{\text{elst}}^{(10)} \right]_{\text{elst}} + \left[ E_{\text{exch}}^{(10)} \right]_{\text{exch}} + \left[ E_{\text{ind,r}}^{(20)} + E_{\text{exch-ind,r}}^{(20)} + \delta E_{\text{HF}}^{(2)} \right]_{\text{ind}} + \left[ E_{\text{disp}}^{(20)} + E_{\text{exch-disp}}^{(20)} \right]_{\text{disp}},
 \end{aligned} \tag{6.1}$$

where the two superscripts in parentheses on the interaction energy components denote the order of perturbation theory with respect to the intermolecular perturbation and intramolec-

ular electron correlation, respectively. The Hartree–Fock correction  $\delta E_{\text{HF}}^{(2)}$  is defined such that Eqns. (6.1) and (6.2) are equal:

$$\delta E_{\text{HF}}^{(2)} \equiv \text{IE}^{\text{HF}} - \left( \left[ E_{\text{elst}}^{(10)} \right]_{\text{elst}} + \left[ E_{\text{exch}}^{(10)} \right]_{\text{exch}} + \left[ E_{\text{ind,r}}^{(20)} + E_{\text{exch-ind,r}}^{(20)} \right]_{\text{ind}} \right). \quad (6.2)$$

SAPT0 exhibits similar computational expense to MP2, due to the terms  $E_{\text{disp}}^{(20)}$  and  $E_{\text{exch-disp}}^{(20)}$  scaling with number of occupied orbitals  $o$  and virtual orbitals  $v$  as  $\mathcal{O}(o^3v^2)$  and  $\mathcal{O}(o^2v^3)$ , respectively. More generically, these terms scale as  $\mathcal{O}(N^5)$ , where  $N$  is proportional to the size of the dimer, whereas the remaining terms scale as  $\mathcal{O}(N^4)$  or better.

SAPT0–D replaces the computationally expensive second-order dispersion  $\left[ E_{\text{disp}}^{(20)} + E_{\text{exch-disp}}^{(20)} \right]_{\text{disp}}$  with the third-generation semiclassical “–D3” dispersion correction of Grimme and coworkers.<sup>35</sup> In the –D3 approach, the dispersion correction to the total *molecular* energy is given by

$$E_{\text{disp}}^{-\text{D3}} = - \sum_{A < B} \sum_{n=6,8}^{\text{atoms}} s_n \frac{C_n^{AB}}{R_{AB}^n} f_{\text{damp}}^{(n)}(R_{AB}), \quad (6.3)$$

where  $f_{\text{damp}}^{(n)}(R_{AB})$  is a damping function which decays as the interatomic distance  $R_{AB} \rightarrow 0$ , and atom-pairwise  $C_6$  and  $C_8$  coefficients are obtained from tabulated values available in Grimme’s DFT-D3 program,<sup>195</sup> as distributed with PSI4.<sup>31</sup> Since SAPT0–D seeks to directly compute an *interaction* energy, the SAPT0–D dispersion interaction is computed according to the supermolecular approach:

$$\text{IE}_{\text{disp}}^{-\text{D3}} = E_{\text{disp}}^{-\text{D3}}(AB) - E_{\text{disp}}^{-\text{D3}}(A) - E_{\text{disp}}^{-\text{D3}}(B), \quad (6.4)$$

where  $(AB)$  denotes that the computation is performed on the dimer and  $(A)$ ,  $(B)$  denote that the computation is performed on either monomer  $A$  or  $B$ , respectively. The SAPT0–D



interaction energy is given therefore by

$$\text{IE}_{\text{SAPT0-D}} = \text{IE}^{\text{HF}} + [\text{IE}_{\text{disp}}^{-\text{D3}}]_{\text{disp}} \quad (6.5)$$

$$= \left[ E_{\text{elst}}^{(10)} \right]_{\text{elst}} + \left[ E_{\text{exch}}^{(10)} \right]_{\text{exch}} + \left[ E_{\text{ind,r}}^{(20)} + E_{\text{exch-ind,r}}^{(20)} + \delta E_{\text{HF}}^{(2)} \right]_{\text{ind}} + [\text{IE}_{\text{disp}}^{-\text{D3}}]_{\text{disp}} \cdot \quad (6.6)$$

As in the DFT-D3 approach, different damping functions  $f_{\text{damp}}^{(n)}(R)$  can be employed in SAPT0-D. In an analogous manner to Ref. 39, we here consider two forms of this damping function. The damping function of Becke and Johnson (BJ)<sup>36,37</sup> expresses the dispersion energy as

$$E_{\text{disp}}^{-\text{D3(BJ)}} = -\frac{1}{2} \sum_{A<B} \sum_{n=6,8} s_n \frac{C_n^{AB}}{r_{AB}^n + (\alpha_1 R_0^{AB} + \alpha_2)^n}, \quad (6.7)$$

with global parameters  $s_8$ ,  $\alpha_1$ , and  $\alpha_2$ . The “zero-damping” function of Chai and Head-Gordon (CHG),<sup>35,38</sup> writes the total dispersion energy as

$$E_{\text{disp}}^{-\text{D3(0)}} = -\frac{1}{2} \sum_{A<B} \sum_{n=6,8} s_n \frac{C_n^{AB}}{r_{AB}^n} \frac{1}{1 + 6(r_{AB}/(s_{r,n} R_0^{AB}) + R_0^{AB} \beta)^{-\alpha_n}} \quad (6.8)$$

and global parameters  $s_8$ ,  $s_{r,6}$ , and  $\beta$  (the  $s_{r,8}$  parameter is fixed at 1, as is the parameter  $s_6$ , while  $\alpha_6$  and  $\alpha_8$  are fixed to 14 and 16, respectively). As in Ref. 39, we have introduced an additional parameter,  $\beta$ , to the original CHG damping function to give the same number of parameters as BJ damping. For clarity, we will hereafter utilize the “-D3(0)” suffix to the SAPT0 method abbreviation to denote the use of CHG damping function and “-D3(BJ)” to denote the use of the BJ damping function.

### 6.3.2 Refitting Damping Parameters for SAPT0-D

From Eqn. 6.5, it is clear that optimal damping parameters for SAPT0-D will be identical to the optimal damping parameters for a supermolecular Hartree-Fock IE. While damping parameters already exist for each of HF-D3(BJ) and HF-D3(0) (these are the parameters

employed by Parrish *et al.*<sup>188</sup>), these parameters were optimized over a relatively small set of training data; indeed, out of the 130 energy points in the original training set, only 72 were interaction energies.<sup>35,36</sup> Furthermore, these IEs sampled a relatively small space of possible noncovalent interactions, making their general application in the context of SAPT0–D an open question. Therefore, optimal parameters for CHG and BJ damping functions must be obtained for a broader, more diverse set of training data than was available when these damping functions were first parameterized. To do this, we have revised these damping parameters in the fashion of Ref. 39, which provided damping parameters for a wide array of different density functionals. As a part of that work, Smith *et al.* also computed SAPT interaction energies (at various levels) for all systems in their test set to classify their interaction motifs.

As our first proof of principle for developing SAPT0–D, we will herein optimize damping parameters for “dispersion-less” SAPT based on those previously reported data, by the nonlinear least-squares minimization of the mean capped unsigned relative error (MCURE) between IEs computed with SAPT and CCSD(T)/CBS taken from Ref. 39. The MCURE was developed previously<sup>39</sup> as a balanced error metric which avoids singularities present in, e.g., mean unsigned error (MUE) for potential energy curves when they cross the zero of interaction energy. The MCURE is given by

$$\text{MCURE} = \left\langle \left| \frac{E - E_{\text{ref}}}{E_{\text{weight}}} \right| \right\rangle \cdot 100\%, \quad E_{\text{weight}} = \max \left\{ |E_{\text{ref}}|, \frac{\xi |E_{\text{ref-eq}}|}{z^3} \right\}, \quad (6.9)$$

where  $\xi$  is a flexible dimensionless parameter that determines the severity of the capping. As in Ref. 39, we have chosen  $\xi = 0.2$ ; additionally, for systems included in our fitting set which are potential energy scans, we have set the cap to be 0.5 kcal mol<sup>-1</sup>. Provided in Table 6.1 are details of the NCI datasets which comprise both the training and validation sets for this parameterization. The space of NCI spanned by these datasets is further visualized in Fig. 6.1, containing “ternary diagrams,” which plot the relative contribution

of the IE components provided by the SAPT0 truncation which can be attractive: electrostatics, induction, and dispersion. It is worth noting that the particular NCI datasets included in the training versus validation sets is different in this work than in our previous revision of DFT-D3 damping parameters: most notably, inclusion of BBI and Water2510 systems in the training set markedly improves the performance of the resulting SAPT0-D methods on this dataset, without negatively impacting accuracy for other datasets. Finally, we distinguish SAPT0-D variants which make use of our optimized damping parameters as SAPT0-D3M(0) and SAPT0-D3M(BJ), where the “M” refers to our newly “modified” parameters.

### 6.3.3 Preparation of Test Systems to Evaluate Scaling of Computational Cost

To investigate the dependence of the computational costs of our SAPT-D implementation on system size, dimer structures composed of increasing numbers of atoms were constructed from the publicly available cocrystal structure 3ACX<sup>196</sup> (resolution: 1.31 Å). This structure was selected due to its high resolution, and the stretch of continuous helical protein near the ligand, which led to easy deconstruction into ligand/protein dimer subsets of increasing size. The input files were created as follows: The 3ACX structure was prepared in Maestro v2019-1<sup>207</sup> using the Protein Preparation utility<sup>208</sup> with default parameters, then the BPH-673 ligand minus its terminal isopropylamino group, along with the continuous stretch of residues 119-181, were extracted from the structure. The protein ends were capped with an acyl (N-terminus) and N-methylacetamide (C-terminus) group. All non-glycine residues in this subset were mutated to alanines, and nonpolar hydrogens were minimized, to form the source of the dimer structures. The smallest dimer example was composed of the reduced ligand (monomer A) plus the closest protein residues, Ala141 and Ala157 (monomer B), along with their neighboring residue caps, resulting in a total of 83 atoms. To create the larger dimer examples, the two nearest additional residues from the source structure were added, with their caps, in turn. This resulted in dimer structures

of sizes increasing by 17-20 atoms, depending on whether the additional residues were glycines or alanines. As a final step, nonpolar hydrogens were minimized to ensure minimal system strain. Files containing all structures examined here are available in plain text format among the Supplementary Materials. Timings were performed on an isolated server equipped with 3.0 TB of RAM and 4.0 TB of scratch space, parallelized over an 8-core Intel Xeon Gold processor running at 2.3 GHz.

#### 6.3.4 Preparation of Structures for Application to GPCR Binding

The crystal structures of turkey  $\beta_1$ AR bound to salbutamol in the active state (PDB ID 6H7M, stabilized by the nanobody NB6B9 at 2.76 Å resolution)<sup>197</sup> and in the inactive state (PDB ID 2Y04, minus stabilizing nanobody at 3.05 Å resolution)<sup>209</sup> were prepared in Maestro v2019-13 with the Protein Preparation utility<sup>208</sup> using default parameters. The input structures for F-SAPT-D calculations were generated by extracting out the ligand along with surrounding residues, with an approximate radius of 7 Å. Residues nearby were included in full; in the case of more distant residues, sometimes only side chains or backbone atoms were retained, according to distance from the ligand. Capping groups from the backbones of neighboring residues were invariably maintained. The resulting salbutamol/reduced- $\beta_1$ AR structures include the full ligand with surrounding residue environment to at least a first shell (structure files in PDB format are included in the Supplementary Materials). The identical set of protein atoms was extracted from the active and inactive cocrystal structures, enabling the comparison of interaction energies. As a final step, nonpolar hydrogens were minimized to ensure minimal system strain. In all, each system is comprised of 459 atoms.

Table 6.2: Summary statistics for signed and unsigned interaction energy error distributions computed using SAPT0/jaDZ and variants of SAPT0–D/jaDZ over each of the training and validation sets, as well as over the full set of all systems examined. Values provided for first and third quartiles correspond to the box borders in a box-and-whisker plot, and the second quartile corresponds to the median value for each SE distribution.

	SAPT0		SAPT0–D3M(0)		SAPT0–D3M(BJ)	
	Signed	Unsigned	Signed	Unsigned	Signed	Unsigned
<b>Fitting Set</b>						
<b>Mean</b>	-0.01	0.51	0.18	0.75	0.30	0.69
<b>StDev</b>	1.15	1.03	1.29	1.07	1.21	1.04
<b>Minimum</b>	-20.06	0.00	-2.77	0.00	-2.84	0.00
<b>1<sup>st</sup> Quartile</b>	-0.15	0.06	-0.40	0.15	-0.26	0.09
<b>2<sup>nd</sup> Quartile</b>	0.00	0.19	-0.12	0.37	-0.04	0.31
<b>3<sup>rd</sup> Quartile</b>	0.25	0.52	0.31	0.86	0.42	0.79
<b>Maximum</b>	3.28	20.06	8.26	8.26	7.24	7.24
<b>Validation Set</b>						
<b>Mean</b>	0.22	0.40	0.06	0.51	0.22	0.52
<b>StDev</b>	0.63	0.54	0.88	0.72	0.94	0.81
<b>Minimum</b>	-4.64	0.00	-2.92	0.00	-4.11	0.00
<b>1<sup>st</sup> Quartile</b>	-0.06	0.05	-0.30	0.12	-0.17	0.07
<b>2<sup>nd</sup> Quartile</b>	0.04	0.16	-0.12	0.26	-0.02	0.21
<b>3<sup>rd</sup> Quartile</b>	0.36	0.50	0.14	0.63	0.31	0.59
<b>Maximum</b>	4.45	4.64	8.27	8.27	7.63	7.63
<b>Full Set</b>						
<b>Mean</b>	0.16	0.43	0.10	0.58	0.25	0.57
<b>StDev</b>	0.82	0.71	1.02	0.84	1.02	0.88
<b>Minimum</b>	-20.06	0.00	-2.92	0.00	-4.11	0.00
<b>1<sup>st</sup> Quartile</b>	-0.09	0.05	-0.32	0.13	-0.19	0.08
<b>2<sup>nd</sup> Quartile</b>	0.03	0.16	-0.12	0.28	-0.02	0.23
<b>3<sup>rd</sup> Quartile</b>	0.32	0.51	0.17	0.70	0.34	0.63
<b>Maximum</b>	4.45	20.06	8.27	8.27	7.63	7.63

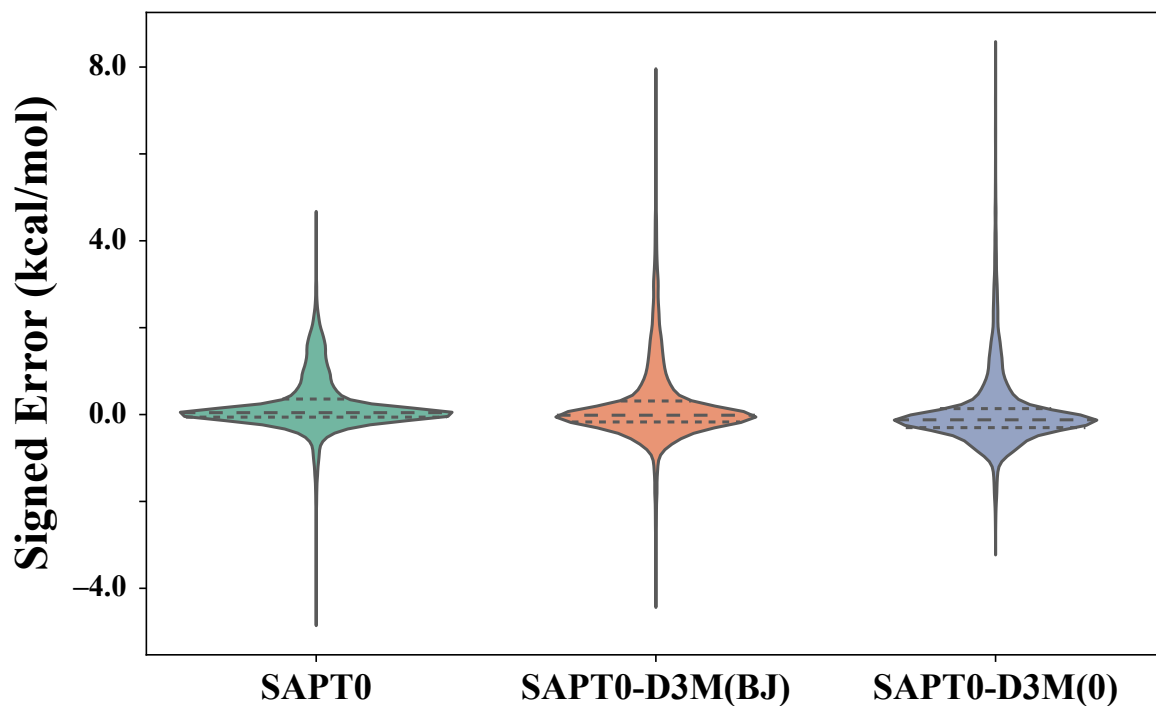


Figure 6.2: Violin plots visualizing the distribution of signed errors (SE) of IEs computed for complexes in the validation set with SAPT0/jun-cc-pVDZ (green), SAPT0-D3M(BJ)/jun-cc-pVDZ (orange), and SAPT0-D3M(0)/jun-cc-pVDZ (purple) as compared to CCSD(T)/CBS reference IEs. Violin widths at a given SE correspond to the relative frequency of complexes exhibiting that value of SE. Also provided for convenience in horizontal dotted lines are the first (Q1), second (Q2), and third (Q3) quartiles for each distribution of SE values.

## 6.4 Results & Discussion

### 6.4.1 Accuracy of SAPT0–D Variants

#### *Statistics over Validation Set*

Visualized in Fig. 6.2 are signed errors for interaction energies of complexes in the validation set computed at the SAPT0/jaDZ, SAPT0–D3M(0)/jaDZ, and SAPT0–D3M(BJ)/jaDZ levels of theory. Summary statistics for these distributions, as well as for the systems in the fitting set, are further provided in Table 6.2. Over these systems included in the validation set, each of SAPT0–D3M(0) and SAPT0–D3M(BJ) exhibits a mean signed error (MSE) whose magnitude is smaller than or equal to SAPT0, with MSE = -0.15, -0.22, and -0.22 kcal mol<sup>-1</sup>, respectively. Furthermore, mean unsigned errors (MUE) for both SAPT0–D3M(0) (MUE = 0.51 kcal mol<sup>-1</sup>) and SAPT0–D3M(BJ) (MUE = 0.52 kcal mol<sup>-1</sup>) are only slightly worse than SAPT0 (MUE = 0.40). The middle 99% of SE distributions are quite similar between SAPT0–D3M(0) and SAPT0–D3M(BJ); SAPT0–D3M(BJ) is, however, slightly more balanced around the mean signed error than SAPT0–D3M(0), where relatively more systems are overbound than for SAPT0–D3M(BJ). For systems in our validation set, therefore, the performance of SAPT0–D3M(0) and SAPT0–D3M(BJ) are effectively equivalent, and nearly as accurate as SAPT0.

#### *Analysis of SAPT–D Over All Systems Considered*

In addition to the signed errors visualized in Fig. 6.2 for the validation set, we provide in Fig. 6.3 and Table 6.2 an analogous analysis of signed errors for interaction energies of all complexes considered here, combining the fitting and validation sets. Summary statistics for these distributions are further provided in Table 6.2. When considering all systems, significant outliers exist for signed errors of interaction energies computed by SAPT0. Indeed, the maximal signed error for SAPT0 IEs is 20.06 kcal mol<sup>-1</sup>, for the closest intermolecular separation of the H<sub>2</sub>S···Benzene complex from the NBC10x data set.<sup>210</sup> This error

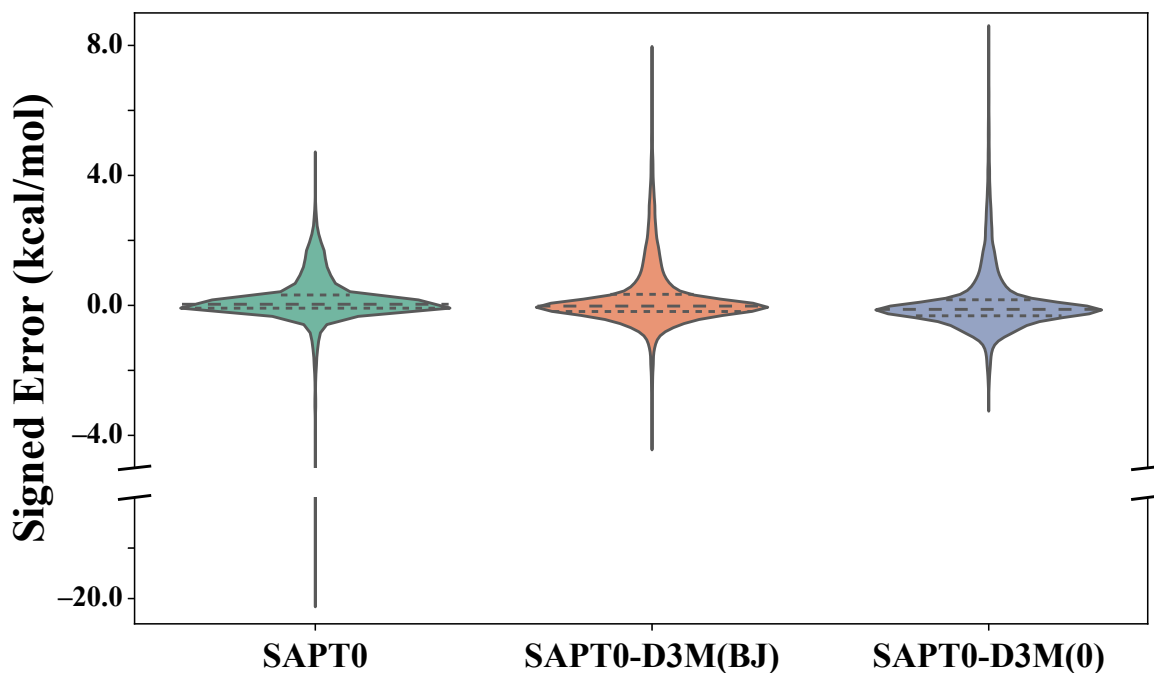


Figure 6.3: Violin plots visualizing the distribution of signed errors (SE) of IEs computed for complexes in the full data set (training and validation sets) with SAPT0/jun-cc-pVDZ (green), SAPT0-D3M(BJ)/jun-cc-pVDZ (orange), and SAPT0-D3M(0)/jun-cc-pVDZ (purple) as compared to CCSD(T)/CBS reference IEs. Violin widths at a given SE correspond to the relative frequency of complexes exhibiting that value of SE. Also provided for convenience in horizontal dotted lines are the first (Q1), second (Q2), and third (Q3) quartiles for each distribution of SE values.



in particular is troubling, especially for potential pharmacological applications of SAPT0, because there may exist close SH- $\pi$  contacts between cysteine and aromatic sidechains in biologically relevant systems. Indeed, in their recent work cataloging nearly 11,000 high-quality X-ray crystal structures taken from the Protein Data Bank, Qi and Kulik<sup>211</sup> found several examples of closer-than-expected contact distances between cysteine or methionine and an aromatic sidechain.

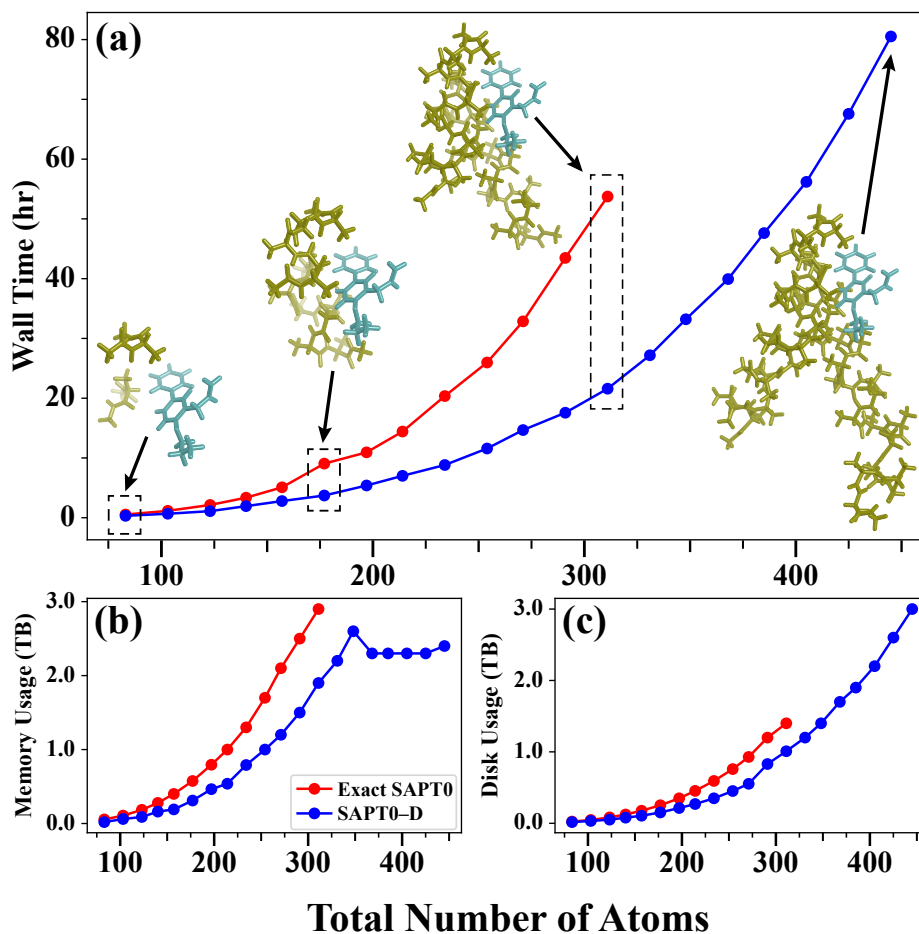


Figure 6.4: Metrics describing the computational expense of SAPT0-D (blue) as compared to exact SAPT0 (red) for (a) the total wall time for each computation, (b) the total disk space utilized by each computation, and (c) the total memory utilized by each computation for scaling tests using progressively larger subsystems of the 3ACX co-crystal structure.<sup>196</sup> SAPT0 computations on more than 311 atoms failed; however, computations successfully completed for SAPT0-D on 3ACX subsystems with up to 445 atoms. **Insets:** Structures for selected 3ACX subsystems with 83, 177, 311, and 445 atoms.

## 6.4.2 Computational Scaling of SAPT0–D

Visualized in Fig. 6.4 are metrics to assess the relative computational expense of SAPT0–D as compared to exact SAPT0, including (a) the total wall time, (b) the total disk space utilized, and (c) the total memory utilized by computations on progressively larger subsystems of mutant 3ACX. Even for systems with as few as 83 atoms, a nontrivial  $1.7\times$  speedup is observed for SAPT0–D versus exact SAPT0; observed speedups increase with system size, up to  $2.5\times$  at 311 atoms, the largest system we were able to run on the test hardware with the conventional approach (see Fig. 6.4.a). The other major improvement for SAPT0–D over SAPT0 is in the total memory utilized by the computations (Fig. 6.4.b), where SAPT–D enjoys between a  $1.5\times$ – $2.9\times$  reduction in the total memory consumed, due to the fact that the storage of intermediate tensors used to compute both  $E_{\text{disp}}^{(20)}$  and  $E_{\text{exch–disp}}^{(20)}$  is required for SAPT0, while no such storage is required by SAPT0–D. Furthermore, while both SAPT–D and SAPT0 are capable of leveraging both traditional, out-of-core and in-core algorithms in PSI4, the reduction in computational expense afforded by SAPT–D enables application to large enough systems (beginning at 331 atoms for 3ACX) where even our test hardware — equipped with an astonishing 3 TB of memory — is forced to switch to the out-of-core (disk-based) algorithm as the maximum amount of physical RAM on the node is reached.

In addition to information regarding the computational efficiency of the SAPT0–D approaches relative to SAPT0, we may also examine their interaction energies and components for progressively larger subsystems of 3ACX. Visualized in Fig. 6.5.a are total SAPT0 and SAPT0–D interaction energies and components for progressively larger subsystems of 3ACX (see, e.g., inset structures in Fig. 6.4.a). As expected, both the total IE and SAPT0(–D) components asymptotically converge as the number of 3ACX residues increases; interestingly, however, the rate of decay is not uniform for all components, as can be seen in Fig. 6.5.b, visualizing differences in interaction energy computed between successively larger 3ACX subsystems. Electrostatics and induction, for example, exhibit differences

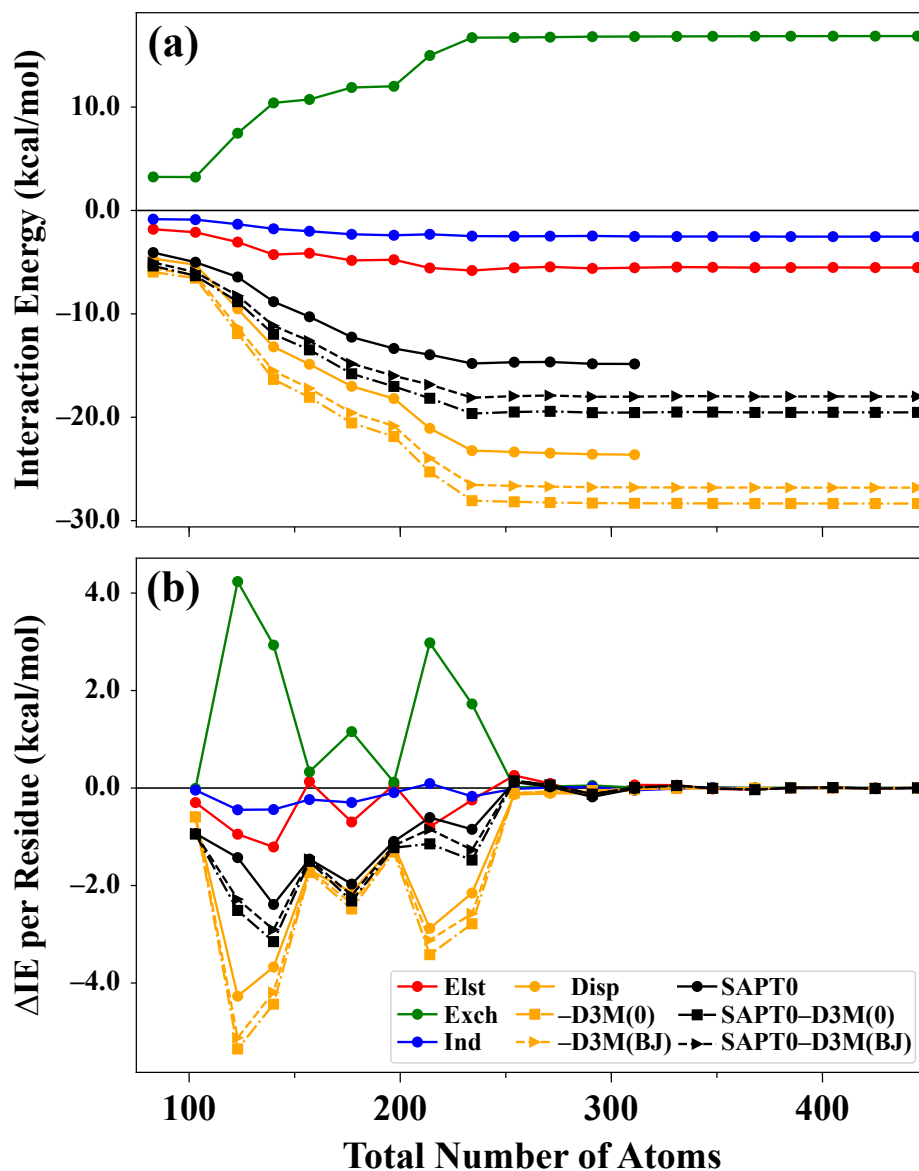


Figure 6.5: (a) Total interaction energies and components ( $\text{kcal mol}^{-1}$ ) for 3ACX subsystems computed with SAPT0, SAPT0-D3M(0), and SAPT0-D3M(BJ) in the jun-cc-pVDZ basis set. (b) Differences between interaction energies of subsequent 3ACX subsystems ( $\text{kcal mol}^{-1}$ ), illustrating the convergence of SAPT0 and SAPT0-D components and total interaction energies.

in IE between successive computations no larger in absolute magnitude than 1.2 and 0.4 kcal mol<sup>-1</sup>, respectively, while exchange and each flavor of dispersion [exact, -D3M(0), -D3M(BJ)] exhibit successive absolute differences as large as 4.2, 4.3, 5.4, and 5.3 kcal mol<sup>-1</sup>, respectively. Convergence for both components and total IE appears to be reached for this system with the inclusion of X residues (254 atoms), indicating that in order to properly describe these interactions, a sufficiently large number of atoms must be included which may make application of exact SAPT0 intractable.

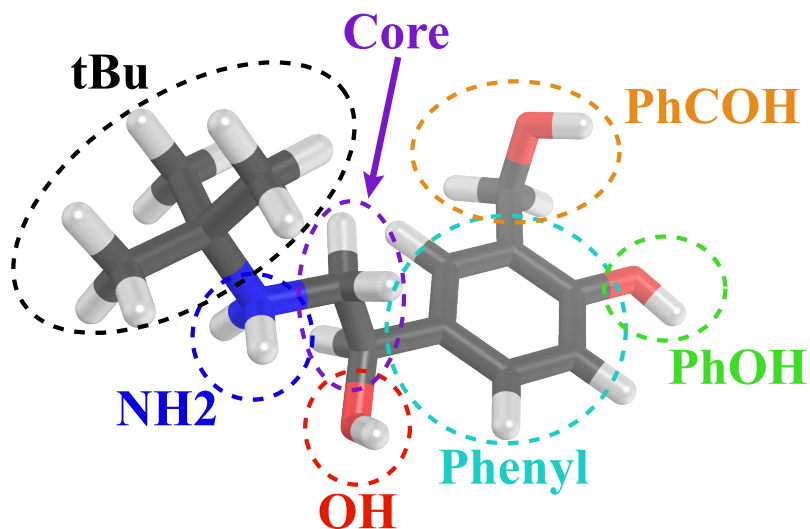


Figure 6.6: Fragmentation scheme for salbutamol monomer in F-SAPT analysis of interaction energies of  $\beta_1$ AR–salbutamol complexes examined here.

#### 6.4.3 Differential Binding of Salbutamol to Active vs. Inactive $\beta_1$ AR

The differential binding strength of salbutamol to the active vs inactive forms of  $\beta_1$ AR was determined by applying the F-SAPT–D approaches presented here to compute the interaction energy between the ligand and model binding pocket of  $\beta_1$ AR, at the F-SAPT–D3M(BJ)/jun-cc-pVDZ and F-SAPT–D3M(0)/jun-cc-pVDZ levels of theory. Each computation was performed on a server equipped with 3.0 TB of RAM and 4.0 TB of local scratch space, parallelized over a 16-core Intel Xeon Gold 2.3 GHz processor. To our knowledge, the computations performed here are the largest symmetry-adapted perturba-

tion theory computations performed that do not rely on software written specifically for hybrid CPU/graphical processing unit (GPU) architectures. F-SAPT post-processing and fragmentation analysis was then performed, allowing for the differences in binding strength to be assessed in terms of individual functional group contacts between ligand and protein side chain/backbone moieties. To do this, fragments were defined for each of the  $\beta_1$ AR binding pocket and salbutamol ligand. For the truncated structure of  $\beta_1$ AR, amino acid sidechains and peptide bonds between residues were grouped into separate fragments, while for salbutamol, functional groups were chosen to maximize the chemical information which could be extracted from the F-SAPT analysis (see Fig. 6.6). For further details on the fragmentation procedure and fragmentation scheme for  $\beta_1$ AR binding pocket, refer to the Supplementary Information.

Table 6.3:  $\Delta\Delta E_{\text{int}}$  values (kcal mol<sup>-1</sup>) computed between active and inactive forms of the  $\beta_1$ AR–salbutamol complex by F-SAPT0–D3M(BJ) and F-SAPT0–D3M(0) in the juncc-pVDZ basis set, decomposed into functional group contacts between the full binding pocket of  $\beta_1$ AR and fragments of salbutamol. Fragment labels are consistent with those shown in Fig. 6.6, with the row labeled “All” corresponding to the total interaction energy of the  $\beta_1$ AR–salbutamol complex.

	Elst	Exch	Ind	D3M(BJ)	D3M(0)	F-SAPT0-D3M(BJ)	F-SAPT0-D3M(0)
<b>PhOH</b>	-4.04	4.22	-1.27	-0.75	-0.87	-1.84	-1.96
<b>PhCOH</b>	-1.59	2.16	0.20	-1.76	-1.82	-0.98	-1.05
<b>Phenyl</b>	-7.31	3.58	0.17	-3.16	-3.30	-6.72	-6.86
<b>OH</b>	-3.93	3.35	-0.99	0.32	0.77	-1.25	-0.80
<b>Core</b>	0.54	-0.29	0.38	-0.74	-0.76	-0.11	-0.14
<b>NH<sub>2</sub></b>	0.27	0.93	-2.29	-0.76	-0.81	-1.85	-1.90
<b>tBu</b>	1.04	1.65	-1.51	-2.25	-2.34	-1.07	-1.17
<b>All</b>	-15.01	15.59	-5.31	-9.09	-9.14	-13.83	-13.87

Presented in Table 6.3 are differences between the interaction energy of the binding pocket of  $\beta_1$ AR in its active and inactive forms, denoted  $\Delta\Delta E_{\text{int}}$ , with the functional groups of partial agonist salbutamol (see Fig. 6.6). Contrary to the experimentally measured  $\Delta\Delta G_{\text{bind}}$  of -2.6 kcal mol<sup>-1</sup> stabilization of the  $\beta_1$ AR—salbutamol complex in the active versus inactive forms, we have computed the difference in interaction energy between the two states to be  $\Delta\Delta E_{\text{int}} = -13.8$  and  $-13.9$  kcal mol<sup>-1</sup> with F-SAPT0–D3M(BJ)

and F-SAPT0–D3M(0), respectively. This deviation from experiment is expected, as our computed  $\Delta\Delta E_{\text{int}}$  represents a 0K energy difference and neglects zero-point energies, finite-temperature contributions to enthalpy, and entropic terms, in addition to neglecting differential solvation effects which would likely dampen the difference in interaction. Despite the difference in magnitude, we nevertheless expect that the  $\Delta\Delta E_{\text{int}}$  computed here is likely to yield meaningful, semi-quantitative insight into how the change in geometry of the binding pocket affects the strength of the interactions between the protein residues and the ligand.

Table 6.4: Order-1 and order-2 F-SAPT–D analysis quantifying the contributions of contacts predicted by Warne *et al.* to be important for explaining difference in binding affinity for salbutamol to active vs. inactive states of  $\beta_1$ AR ( $\Delta\Delta E_{\text{int}}$ , computed by F-SAPT–D3M(BJ) and F-SAPT0–D3M(0) in the jun-cc-pVDZ basis set. Fragment labels for salbutamol are consistent with those shown in Fig. 6.6. (top) Contributions of residue–functional group pairs hypothesized by Warne *et al.* to make polar or hydrogen-bonding contacts, as identified in Figure 3 of Ref. 197; (bottom) Contributions of total contact strength between all amino acid sidechains identified in Figure 3 of Ref. 197 and the full salbutamol ligand (labeled “All”).

$\beta_1$ AR	Salbutamol	Elst	Exch	Ind	D3M(BJ)	D3M(0)	F-SAPT0-D3M(BJ)	F-SAPT0-D3M(0)
<i>Predicted H-Bonds &amp; Polar Contacts from Fig. 3 of Ref. 197</i>								
<b>D121</b>	<b>NH<sub>2</sub></b>	-2.92	1.59	-0.63	-0.67	-0.72	-2.62	-2.67
	<b>OH</b>	-6.43	5.51	-1.37	0.54	1.10	-1.75	-1.19
<b>S211</b>	<b>PhOH</b>	5.75	-0.85	0.53	0.15	0.27	5.57	5.70
	<b>PhCOH</b>	-2.40	2.07	-0.35	0.01	-0.13	-0.68	-0.81
<b>S215</b>	<b>PhOH</b>	-8.24	5.23	-1.71	-0.73	-1.00	-5.46	-5.72
<b>N329</b>	<b>NH<sub>2</sub></b>	-0.09	-0.65	-0.11	0.01	0.00	-0.84	-0.86
	<b>OH</b>	1.71	-2.53	0.61	0.11	0.05	-0.10	-0.15
<i>Contacts with all residues listed in Fig. 3 of Ref. 197</i>								
<b>W117</b>	<b>All</b>	-0.36	0.66	-0.09	-0.29	-0.30	-0.08	-0.09
<b>D121</b>	<b>All</b>	-8.47	7.97	-1.72	-1.08	-0.66	-3.29	-2.87
<b>V122</b>	<b>All</b>	0.12	0.33	-0.13	-0.40	-0.39	-0.08	-0.08
<b>V125</b>	<b>All</b>	-0.41	2.26	-0.19	-2.22	-2.36	-0.56	-0.70
<b>F201</b>	<b>All</b>	-4.07	1.84	-0.01	-1.24	-1.26	-3.47	-3.50
<b>Y207</b>	<b>All</b>	-0.68	0.24	0.06	-0.09	-0.13	-0.47	-0.51
<b>S211</b>	<b>All</b>	-1.85	1.20	0.45	0.32	0.32	0.12	0.11
<b>S215</b>	<b>All</b>	-10.46	5.50	-2.05	-1.46	-1.78	-8.48	-8.79
<b>F306</b>	<b>All</b>	0.47	1.08	-0.13	-0.71	-0.75	0.71	0.67
<b>F307</b>	<b>All</b>	-0.04	-0.48	0.03	-0.10	-0.11	-0.60	-0.61
<b>N310</b>	<b>All</b>	-0.32	-2.00	0.71	-0.27	-0.04	-1.87	-1.64
<b>N329</b>	<b>All</b>	2.11	-2.94	0.45	0.17	0.09	-0.22	-0.30
<b>Y333</b>	<b>All</b>	-1.31	-0.76	-0.12	0.41	0.50	-1.77	-1.68

Turning our attention to the role of the particular functional group contacts between

residues in the binding pocket of  $\beta_1$ AR and salbutamol which were previously hypothesized by Warne *et al.* as being important contributors to the  $\Delta\Delta G_{\text{bind}}$ , F-SAPT-D analysis of contacts identified in Fig. 3 of Ref. 197 are provided in Table 6.4. Despite the significant  $-1.2\text{\AA}$  movement by residue D121 of  $\beta_1$ AR upon activation, the interaction strength between D121 with salbutamol  $\text{NH}_2$  and OH do not change significantly according to F-SAPT-D. This is likely due to the fact that even though the aspartate sidechain moves markedly closer to the ligand, the D121-OH contact distance only contracts by  $0.3\text{\AA}$ , while the D121-NH<sub>2</sub> contact distance actually *lengthens* by  $0.5\text{\AA}$ . Even smaller changes in direct contact distance are present for the N329 sidechain, resulting therefore in smaller contributions to the total  $\Delta\Delta E_{\text{int}}$  than for D121. Among all polar and hydrogen bonding contacts predicted by Warne *et al.* to be important contributors to observed  $\Delta\Delta G_{\text{bind}}$ , only S211-PhOH and S215-PhOH are significant contributors to  $\Delta\Delta E_{\text{int}}$ ; interestingly, however, these contacts largely cancel. All told, the predicted H-bonding and polar contacts contribute only 40% (44%) of the total  $\Delta\Delta E_{\text{int}}$  as computed by F-SAPT0-D3M(BJ) (F-SAPT0-D3M(0)). Unfortunately, the situation is not improved when expanding the interactions considered to include residue interactions with the full salbutamol ligand (as opposed to only suspected H-bonds and polar contacts) and also including other nearby residues noted in Fig. 3 of Ref. 197, as the total contribution of these contacts is 144% (145%) of the overall F-SAPT0-D  $\Delta\Delta E_{\text{int}}$ . If the total difference in interaction strength between active and inactive states of the  $\beta_1$ AR-salbutamol complex are so incompletely described by these residue sidechains, where exactly does the rest of the energy difference come from?

Provided in Table 6.5 are the relative order-1 F-SAPT contributions for the ten fragments of  $\beta_1$ AR which contribute most significantly to the total  $\Delta\Delta E_{\text{int}}$  for the  $\beta_1$ AR-salbutamol complex, as computed by F-SAPT0-D3M(BJ) and F-SAPT0-D3M(0) with the jun-cc-pVDZ basis set. Quite surprisingly, out of the ten most important fragments of  $\beta_1$ AR contributing to the total  $\Delta\Delta E_{\text{int}}$ , only *half* were among those identified by Warne

Table 6.5: Relative contributions of order-1 [i.e.,  $\beta_1$ AR(fragment)–salbutamol] contacts to the total  $\Delta\Delta E_{\text{int}}$  for the full  $\beta_1$ AR–salbutamol complex, computed with F-SAPT0–D3M(BJ) and F-SAPT0–D3M(0) in the jun-cc-pVDZ basis set. Also provided for reference are whether or not the particular residue sidechain or peptide bond was hypothesized to be important by Warne *et al.* in Ref. 197. Fragment labels for  $\beta_1$ AR are consistent with the fragmentation procedure described in the Supplementary Information.

$\beta_1$ AR Fragment	% SAPT0-D3M(BJ)	% SAPT0-D3M(0)	Warne <i>et al.</i> ?
<b>200p201</b>	-71.00	-71.02	No
<b>S215</b>	61.30	63.37	Yes
<b>F201</b>	25.11	25.20	Yes
<b>D121</b>	23.77	20.67	Yes
<b>211p212</b>	15.03	15.27	No
<b>201p202</b>	-13.90	-13.86	No
<b>N310</b>	13.51	11.84	Yes
<b>Y333</b>	12.78	12.10	Yes
<b>208p209</b>	10.42	10.39	No
<b>T126</b>	9.98	10.05	No

*et al.* either in Fig. 3 of Ref. 197 or the main text of that work. Furthermore, four out of the five fragments neglected by the analysis of Warne *et al.* are peptide bonds, which thanks to their non-trivial dipole moments (on average, approximately 2.5 D) can interact strongly with a binding ligand. We find it important to stress that, instead of an indictment of the methodology or chemical intuition of Warne *et al.*, this finding is rather an indication that a local contact model — whereby the total protein–ligand interactions are assumed to be well captured by the sum of nearest-neighbor contacts between sidechains in the binding pocket and the ligand itself — is an insufficient picture to justify the computed  $\Delta\Delta E_{\text{int}}$  (and therefore the experimentally determined  $\Delta\Delta G_{\text{bind}}$ !). Indeed, similar conclusions have been drawn in the case of the differential activity of chloro- versus methyl-aryl substituted factor Xa inhibitor drugs.<sup>69</sup>

## 6.5 Conclusions

Thanks to its intuitive decomposition of the interaction energy between two chemical species into well-defined, chemically meaningful contributions from electrostatics, exchange (steric) repulsion, induction (polarization), and London dispersion, symmetry-adapted



perturbation theory (SAPT) has become an invaluable computational tool for understanding the physical basis for non-covalent interactions. Furthermore, its atomic and functional group partitions (ASAPT and F-SAPT, respectively) have already been used to investigate the difference in binding strength between functional isomers of factor Xa inhibitor drugs<sup>69</sup> and the stereoselectivity of reactions whose transition states are preferentially stabilized by non-covalent interactions.<sup>67,68</sup> Unfortunately, further application of SAPT and its partitions to larger chemical systems has been limited by its computational expense, with its lowest-order truncation, SAPT0, scaling as  $\mathcal{O}(N^5)$  with  $N$  proportional to system size. Inspired by the recent developments of Parrish *et al.*, we have presented here a reduced-scaling approach combining the empirical dispersion correction of Grimme<sup>35</sup> with SAPT, denoted SAPT-D, whereby we have obtained optimal damping parameters for the -D3 dispersion interaction component of our method over a diverse training and validation set of nearly 8,100 bimolecular complexes. Over the total set of bimolecular complexes, our two SAPT0-D variants (leveraging different damping functions) achieve equivalent accuracy to SAPT0 when compared against silver-standard<sup>77</sup> reference interaction energies computed at the DW-CCSD(T<sup>\*\*</sup>)-F12/aug-cc-pVDZ level of theory, while removing several egregious outliers present for SAPT0.

We have also examined the computational expense of our formulation of SAPT-D by comparing interaction energy computations performed using SAPT0-D and exact SAPT0 on increasingly large truncations of the publicly available 3ACX cocrystal structure<sup>196</sup> with a computer equipped with an 8-core Intel Xeon Gold processor at 2.4 GHz and 3.0 TB of available RAM. While exact SAPT0 computations failed to complete within the queue limit of 80 hours of wall time for all systems larger than 311 atoms, SAPT0-D completed successfully for systems with up to 445 atoms in that time. We have also shown the convergence of total interaction energies and SAPT0(-D) components towards some macroscopic limit, indicating that the relative contribution by successive residues to the strength of the ligand binding to the head group of the truncated 3ACX helix decays as distance to the

binding site increases. This result implies that for systems with residues which are very distant from the binding site, these residues may be able to be removed from the computation without significantly effecting the validity of conclusions drawn using our approach.

Finally, we have expanded our SAPT-D formulation to be consistent with the atomic and functional-group partitions of SAPT, and applied F-SAPT-D to investigate the difference in binding strength between 460-atom subsystems of active and inactive forms of the G-protein coupled receptor  $\beta_1$ AR complexed with the partial agonist salbutamol, whose crystal structures were recently published by Warne *et al.*<sup>197</sup> In that work, the difference in binding affinity between the receptor coupled (the active state) or not coupled (the inactive state) to G-protein was rationalized by appealing to the significant decrease in binding-site volume upon G-protein coupling, which consequently decreased contact distances between salbutamol and several amino acid residues in the binding pocket which were hypothesized to play a large role in the binding affinity difference. In addition to determining that the  $\Delta\Delta E_{\text{int}}$  computed by our F-SAPT-D approach provides a reasonable approximation to the  $\Delta\Delta G_{\text{bind}}$  observed experimentally by Warne *et al.*, we have presented here a functional-group partition analysis of the computed  $\Delta\Delta E_{\text{int}}$  which indicates that, contrary to the justification of Warne *et al.* for the difference in activity based on only local sidechain–ligand contacts, peptide bonds and even more distant cooperative effects play a major role in determining the difference in binding affinity upon G-protein coupling. This finding indicates that a local contact model is likely insufficient to justify small differences in binding affinity based on conformational differences in protein environment.

## **PART IV**

### **APPLICATION TO INTERESTING CHEMICAL SYSTEMS**

**CHAPTER 7**  
**THE INFLUENCE OF SOLVATION ON NON-COVALENT INTERACTIONS IN**  
**BIMOLECULAR COMPLEXES: AN INTRAMOLECULAR**  
**SYMMETRY-ADAPTED PERTURBATION STUDY**

**7.1 Abstract**

High-level quantum chemical computations have provided significant insight into the fundamental physical nature of non-covalent interactions (NCI). To date, these studies have focused primarily on gas-phase computations of small van der Waals dimers; however, these interactions are frequently taking place in complex chemical environments such as proteins, solutions, or solids. In order to better understand how chemical environment affects non-covalent interactions, we have undertaken a quantum chemical study of  $\pi$ - $\pi$  interactions in aqueous solution, as exemplified by T-shaped benzene dimers surrounded by 28 or 50 explicit water molecules. We report interaction energies using second-order Møller-Plesset perturbation theory, and we also apply the intramolecular and functional-group partitioning extensions of symmetry-adapted perturbation theory (ISAPT and F-SAPT, respectively) to analyze how the solvent molecules tune the  $\pi$ - $\pi$  interactions of the solute. For complexes containing neutral monomers, even 50 explicit waters change total SAPT interaction energies (IEs) between the two solute molecules by only tenths of a kcal mol<sup>-1</sup>, while significant screening of up to 3 kcal mol<sup>-1</sup> of the electrostatic component is seen for the cationic pyridinium–benzene dimer. These differences between solvation levels are attributed to large non-additive interactions within solvated ion-containing complexes of  $\sim 40\%$  the gas-phase IE on average, an order of magnitude larger a fraction than for neutral complexes where the extent of solvation is significantly less influential.

## 7.2 Introduction

Non-covalent interactions (NCI) continue to receive significant attention in the computational chemistry literature, due in large part to their fundamental importance to governing important chemical and physical phenomena, such as the relative stability of crystal polymorphs and host-guest binding in drug design. Reliable prediction of these phenomena relies therefore on an understanding of NCI within the system of interest. To do this, the typical approach has been to construct gas-phase model systems comprised of two interacting molecules which mimic the interactions present in the full system of interest, e.g., investigating the  $\pi - \pi$  interaction in the stacked benzene dimer as a model for the interaction of aromatic side chains in a protein. Quantum chemical methods can then be applied to compute the interaction energy (IE) present in the complex, by subtracting the energy of the isolated monomers from the energy of the dimer:

$$\Delta E_{AB}^{\text{int}} = E_{AB} - E_A - E_B,$$

whereby  $\mathcal{A}$  and  $\mathcal{B}$  denote the monomers and  $\mathcal{AB}$  denotes the dimer. This “supermolecular” approach to quantifying NCI has been employed to study a variety of interaction motifs, including  $\pi - \pi$ , cation- $\pi$ , and halogen bonding, among others. Furthermore, this approach has also been leveraged to establish very high-quality IE benchmarks using the “gold-standard” method in quantum chemistry, coupled cluster through single, double, and perturbative triple substitutions [CCSD(T)]. With a significant computational expense of  $\mathcal{O}(o^3v^4)$ , where  $o$  is the total number of occupied and  $v$  is the total number of virtual (unoccupied) orbitals, however, CCSD(T) can only be applied to small systems, with no more than approximately 30 non-Hydrogen atoms in the total dimer complex. The accuracy of a host of less expensive approaches [including density functional theory (DFT) and second-order Møller–Plesset perturbation theory (MP2)] has furthermore been assessed against CCSD(T) over test sets of bimolecular complexes for which these high-level reference IEs

exist, such that now NCI in arbitrary chemical systems can now be reliably studied.

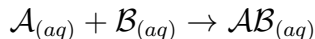
Despite its utility for quantifying the strength of NCI, the supermolecular approach offers only a single scalar quantity by which to do so, namely, the IE. Of course, the underlying reasons for the different behavior of, e.g., hydrogen bonds versus  $\pi - \pi$  stacking is not due only to their relative strengths. For further details into the nature of NCI, therefore, it is desirable to quantify the contributions to the total IE arising from different physical forces, such as electrostatics or steric repulsion. Fortunately, the importance of these types of interactions and their relative contribution to an IE may be quantified either by the post-hoc partitioning of the total IE (an energy decomposition analysis, EDA) or by computing them directly via symmetry-adapted perturbation theory (SAPT). SAPT has been particularly successful at unambiguously providing a detailed description of the physics governing NCI by directly computing each of the electrostatic, exchange (steric) repulsion, induction (polarization), and London dispersion components contributing to the IE of a bimolecular complex, and has been repeatedly applied to analyze and classify interaction motifs in a wide range of chemical systems. Thanks to the formulation of the lowest-order truncation of SAPT (SAPT0) to leverage density-fitted two-electron integrals in the atomic-orbital basis (DF-SAPT0), SAPT has recently become routinely applicable to systems as large as  $\sim 300$  atoms. Furthermore, IEs computed with SAPT0 in the jun-cc-pVDZ basis set — where the diffuse space is truncated by neglecting diffuse functions on H atoms and diffuse  $d$  functions on second-row atoms — have been shown to exhibit mean absolute errors of  $0.9 \text{ kcal mol}^{-1}$  relative to CCSD(T) benchmarks, which can be improved for no additional cost to  $0.5 \text{ kcal mol}^{-1}$  by leveraging exchange scaling (sSAPT0).<sup>58</sup>

The distinct advantage offered by SAPT over supermolecular approaches is that it provides a chemically-intuitive decomposition of interaction energies into their physical components; unfortunately, however, both supermolecular approaches and SAPT share the shortcoming that these methods are applicable only to describing the interactions between exactly two chemical species at a time. Of course, a “many-body” interaction energy has

been previously defined in an analogous supermolecular manner to the two-body case, and through the many-body expansion (MBE) has been used to successfully predict the lattice energies of several organic crystals with very high accuracy.<sup>CrystalLatte</sup> Furthermore, a “three-body” formulation of SAPT has been developed, but due to its computational expense is intractable for systems much larger than the water trimer. To remedy this, Herbert and coworkers developed the XSAPT approach based on the self-consistent XPol treatment of many-body induction. What these many-body approaches do not provide, however, is the ability to investigate how a particular non-covalent interaction is *tuned* by its chemical environment, i.e., how it differs *in situ* from what it would be in the gas-phase. This question is most notably relevant in the context of rational pharmaceutical design, as the binding strength of drug to protein target has recently been shown to be modulated by the complete protein environment, not only local contacts.<sup>69,180</sup> An ideal approach to addressing this fact would be to simply perform quantum mechanics on the entire protein;<sup>212</sup> this solution is, unfortunately, intractable for all practical purposes. It seems desirable, therefore, to develop approaches which may access this information without resorting to a full-system QM solution.

One might consider at least three mechanisms by which a chemical environment can affect individual NCI: (1) *direct electronic modification* of an NCI due to the environment polarizing the two monomers or chemical fragments involved in the interaction; (2) statistical effects due to averaging over many accessible arrangements of the system’s atoms; (3) *indirect, “effective” modification* of an interaction due to competition between the original interaction and interactions with the chemical environment. The first and second mechanisms are concerned with the affect of the chemical environment on modulating the strength of interaction once complexation has occurred; the third mechanism, on the other hand, is concerned with the favorability of complexation itself. This may be illustrated by consid-

ering the formation of a complex  $\mathcal{AB}$  in solution, from solvated monomers  $\mathcal{A}$ ,  $\mathcal{B}$ :



The free energy of binding for this process,  $\Delta G_{solv}^{\text{bind}}(\mathcal{AB})$ , can be written in terms of the gas-phase binding free energy,  $\Delta G_{gas}^{\text{bind}}(\mathcal{AB})$ , and the solvation energies of each species, as

$$\Delta G_{solv}^{\text{bind}}(\mathcal{AB}) = \Delta G_{gas}^{\text{bind}}(\mathcal{AB}) + \Delta G_{solv}(\mathcal{AB}) - [\Delta G_{solv}(\mathcal{A}) + \Delta G_{solv}(\mathcal{B})]$$

In this work, we will not address the thermodynamic cycle encompassed by this third mechanism. Instead, we will primarily explore the first mechanism and, to a lesser extent, the second, by investigating the manner in which chemical environment modulates  $\pi - \pi$  interactions by examining T-shaped configurations of eight mono-functionalized aromatic molecules ( $\text{ArX}$ ;  $\text{Ar}$  = benzene, pyridine,  $\text{X}$  = H,  $\text{NH}_2$ ,  $\text{NO}_2$ ,  $\text{OCH}_3$ ,  $\text{CH}_3$ ) interacting with benzene (Bz), solvated by statistically significant configurations of one or two hydration shells (with 28 and 50 water molecules, respectively). For each solvent configuration, sampled from molecular dynamics trajectories where the solute molecules were kept rigid, the tuning of the  $\text{PhX-Bz}$  interaction will be assessed by computing this interaction directly within the solvent environment by leveraging two recent extensions of SAPT, namely (i) its functional-group partition (F-SAPT)<sup>66</sup> and (ii) its intramolecular formulation (ISAPT).<sup>213</sup>

By accumulating contributions to the SAPT0 interaction energy and components from pairs of atoms on opposite monomers, F-SAPT provides a decomposition of the SAPT0 interaction energy (and its components) into contributions from the interaction of each monomer with chemical fragments (i.e., functional groups) on the other monomer (an order-1 partition), and furthermore into contributions from contacts between functional groups on opposite monomers (an order-2 partition). ISAPT, on the other hand, provides for the computation of the SAPT0 interaction energy between functional groups of the same molecule, rather than the traditionally rigid two-body formulation of SAPT. ISAPT does



this by first partitioning a single molecule  $\mathcal{X}$  into interacting fragments  $\mathcal{A}$  and  $\mathcal{B}$ , separated (but linked to one another) by a third fragment,  $\mathcal{C}$ . The zeroth-order wavefunctions for  $\mathcal{A}$  and  $\mathcal{B}$  are then prepared via a HF-in-HF embedding approach inspired by the procedure of Manby and coworkers, in which the orbitals of  $\mathcal{A}$  and  $\mathcal{B}$  are electronically deformed by the presence of  $\mathcal{C}$  before a standard SAPT0 computation is performed. The effect of the linker  $\mathcal{C}$  is therefore *effectively* captured, since the resulting ISAPT0 interaction energy and components are computed between the pre-polarized electron densities of fragments  $\mathcal{A}$  and  $\mathcal{B}$ . Each of these approaches can be made extensible to the computation of interactions embedded in a chemical environment, thereby offering complementary perspectives against which they may be mutually validated.

For all systems, we have obtained interaction energies using F-SAPT, ISAPT, and MP2, allowing us to quantify both the extent of environment tuning of ArX–Bz interactions and the pairwise additivity of the ArX–solvent, Bz–solvent, and ArX–Bz two-body interactions for recovering the total interaction energy of the full, mutually interacting system. We have found that for all but one solvated dimer (the hydrated pyridinium–benzene complex), the many-body interaction is well approximated by the sum of these two-body interactions. Where this pairwise additivity is exhibited, solvation by a single hydration shell of 28 explicit water molecules does not significantly modulate the ArX–Bz interaction, with deviations from the gas phase of only  $\pm 0.5$  kcal mol<sup>-1</sup> on average for both total ISAPT0 IEs and components. Furthermore, similarly small deviations are observed between different solvent configurations of the same ArX–Bz dimer. For the hydrated pyridinium–benzene complex, however, the ArX–Bz IE is reduced by up to  $\sim 2.5$  kcal mol<sup>-1</sup> on average relative to the gas phase, due almost exclusively to screening of the electrostatic component. We also observe large variations in the total ISAPT0 IE between different solvent configurations of this system of up to  $\sim 3$  kcal mol<sup>-1</sup>. From these findings, we conclude that solvent environment does not significantly tune ArX–Bz interactions and that a gas-phase treatment provides a good first approximation to describing these interactions, so long as the system

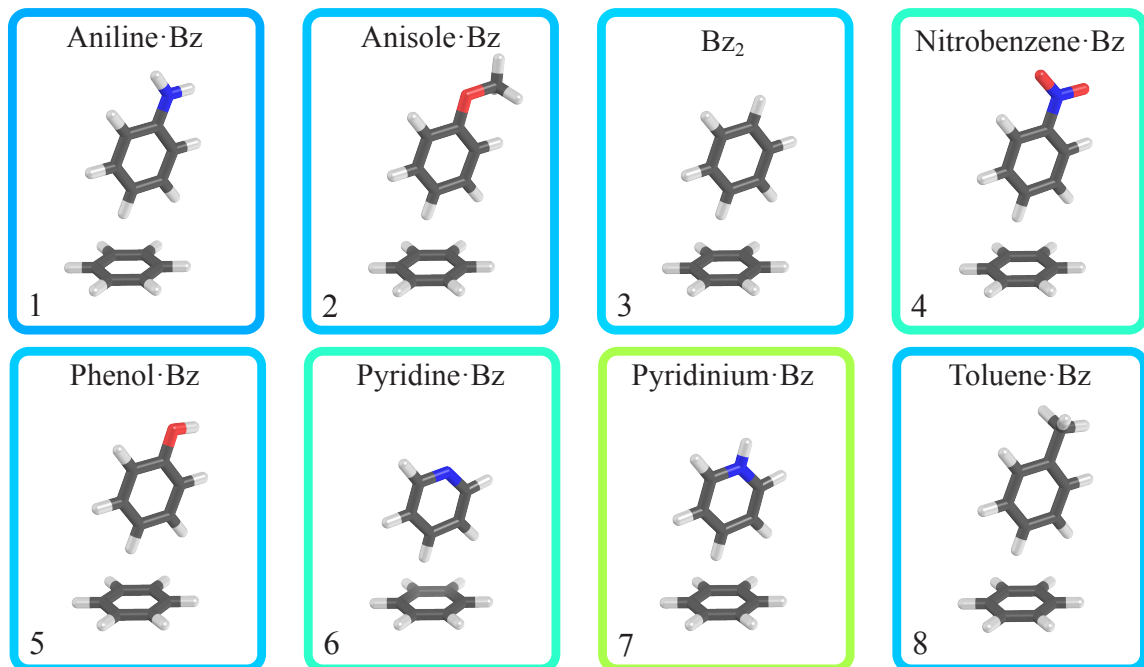


Figure 7.1: Bimolecular complexes from which the HYD8 test set is constructed. Structures prepared via functionalization of the T-shaped pyridine–benzene complex from the S66 test set,<sup>76,104</sup> before re-optimizing the structures at the B3LYP-D3M(BJ)/aug-cc-pVDZ level of theory within enforced  $C_s$  symmetry. Box coloring is based on SAPT0/jun-cc-pVDZ results computed in the gas phase, and indicates the interaction type: blue for dispersion-dominated interactions and yellow-green for mixed electrostatics and dispersion contributions.

does not experience marked non-additive many-body effects. In that event, however, we have also presented a robust framework for estimating the extent to which chemical environment tunes non-covalent interactions relative to the gas phase, given that complexation has already occurred.

### 7.3 Computational Methods

#### 7.3.1 Preparation of Geometries for Functionalized Complexes

Eight bimolecular complexes, each consisting of one benzene molecule and one benzene derivative, were prepared via functionalization of the tilted T-shaped pyridine–benzene complex from the S66 data set of Hobza and coworkers,<sup>76,104</sup> using Maestro v11.<sup>207</sup> Each

structure was then optimized within enforced  $C_s$  symmetry using a development version of the PSI4 electronic structure package,<sup>31</sup> using the dispersion-corrected B3LYP density functional and aug-cc-pVDZ basis set.<sup>22,147,148</sup> Optimizations were performed with default convergence thresholds, as recommended previously,<sup>128</sup> and employed the recently modified<sup>39</sup> parameters for the “-D3” dispersion correction of Grimme,<sup>35</sup> together with Becke-Johnsson damping.<sup>36,37</sup> For clarity, we will denote this combination of density functional, dispersion correction, and damping scheme here as “B3LYP-D3M(BJ)”. This test set of functionalized benzene dimer complexes was constructed to provide both structural diversity and differentiation with respect to (i) local substituent dipole (e.g., toluene vs. nitrobenzene), (ii) ability to form hydrogen bonds with the water solvent environment (e.g., phenol vs. benzene), (iii) molecular polarizability (e.g., anisole vs. benzene), and (iv) polarizing effect (e.g., pyridinium vs. benzene). In this way, we hope that conclusions drawn can be extended to the broader chemical space spanned by biologically relevant bimolecular complexes.

### 7.3.2 Hydration of Functionalized Complexes

Hydration of the eight complexes prepared above was carried out in the NAMD software package<sup>214</sup> within a TIP3P<sup>215</sup> water box measuring  $27 \times 27 \times 27$  Å under periodic boundary conditions (PBC). Restrained electrostatic potential (RESP) charges based on AM1-BCC<sup>216</sup> were fit using ligand geometries from the optimized dimers generated above; Amber topology and parameters based on the General Amber Force Field (GAFF) version 9 and the Amber94 force-field parameters were generated using the Antechamber program from Amber Tools 9.<sup>217–219</sup> An initial minimization was performed over 480 conjugate gradient steps with all solute atoms fixed and all O–H bonds held constant using the SETTLE algorithm, using (i) a vdW and electrostatic switching function from 8 to 10 Å, (ii) scaling factor of  $0.8\overline{33}$  for 1-4 interactions, and (iii) pairwise electrostatic and vdW interactions cutoff of 12.0 Å. Minimization was followed by 125 ps of NPT equilibration via Langevin

dynamics at 1 atm and 298 K using the Velocity–Verlet integrator with 2 fs time steps. Finally, NPT molecular dynamics trajectories were propagated for 10 ns under the same conditions, with coordinates saved every 5 ps. Five such trajectories were generated for each complex prepared above using different starting velocities.

For each complex, 20 snapshots were extracted from each of the five molecular dynamics trajectories. To eliminate any energetically spurious solvent configurations from among the gathered snapshots, energy minimization of each solvated bimolecular complex was performed according to an identical procedure as the initial minimization described above. For each of the raw (unminimized) and relaxed (minimized) solvent configurations surrounding all eight solute complexes, the maximum number of water molecules within baseline distances of 3.0 Å from the nearest solute atom was determined to correspond to target first solvation shell consisting of 28 water molecules. The cutoff distances were calibrated for all snapshots with a search resolution of up to 0.01 Å and a maximum search radius of  $r = 4.5$  Å in order to maintain the same number of waters across all snapshots and all complexes. Snapshots which could not attain this target number of water molecules below the search radius were discarded.

To maximize the diversity in water geometry among remaining snapshots, the structures of extracted water molecules were clustered based on both shape and atom-type similarities using the Tanimoto Combo similarity score, evaluated using the ROCS program of OpenEye Software.<sup>220</sup> This clustering was performed with Tanimoto Combo cutoff of 0.75,<sup>221</sup> for all minimized solvent configurations surrounding each bimolecular complex as well as for unminimized solvent configurations surrounding the benzene dimer, aniline–benzene complex, pyridine–benzene complex, and pyridinium–benzene complex. Finally, centroids of the 10 most populated clusters for each dimer were isolated. These 120 structures, corresponding to the 10 most populated cluster centroids for snapshots of minimized solvent configurations surrounding each of the eight bimolecular complexes, together with the 10 most populated cluster centroids for snapshots of the unminimized solvent configurations

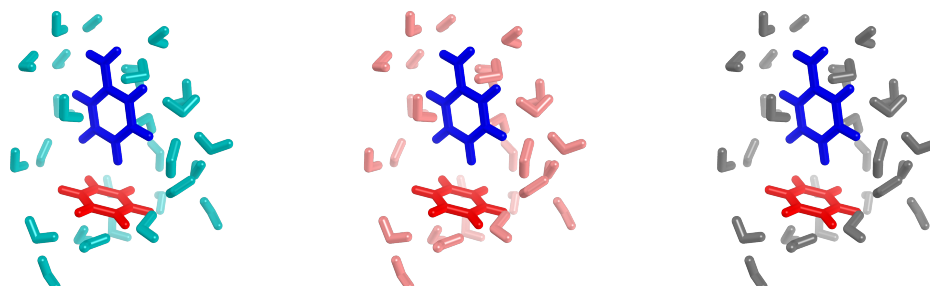


Figure 7.2: Environment binning schemes employed in this work, illustrated for the HYD8-1 (solvated aniline–benzene) complex. A–B interactions in “EnvC” scheme are computed directly using ISAPT, while A–B interactions in “EnvA” and “EnvB” are computed via F-SAPT post analysis,<sup>222</sup> via accumulation of functional group interactions.

surrounding those complexes mentioned previously, comprise the HYD8 test set. Furthermore, to investigate the convergence of the effect of chemical environment towards the “bulk” solvent, we also include the same 120 structures where we have extended the cut-off to include water molecules within 7 Å of the solute to represent a second solvation shell of 50 water molecules. These 240 structures are labeled according to the the format HYD8- $\mathbb{I}\times\mathbb{J}$ -wN, where  $\mathbb{I} = 1-8$  identifies the complex (see numbering in Fig. 7.1),  $\times = m, u$  indicates the solvent molecules were relaxed (*m*) or unrelaxed (*u*),  $\mathbb{J} = 1-10$  identifies the particular snapshot, and  $N = 28, 50$  indicates whether the dimer complexes are either singly (*w28*) or doubly (*w50*) solvated. For example, the label HYD8-1m4-w50 decodes to the relaxed, doubly-solvated nitrobenzene–benzene complex with solvent molecules in configuration #4. Geometries for each of these complexes are provided as a zipped archive in the Supplementary Information.

### 7.3.3 Quantifying Electronic Effect of Solvent on $\pi$ - $\pi$ Interactions and Energy Components via F-/ISAPT

To quantify the effect of solvation on the noncovalent ArX–Bz interactions within HYD8 complexes, both functional-group partitioned symmetry-adapted perturbation theory (F-SAPT)<sup>138</sup> and intramolecular symmetry-adapted perturbation theory (ISAPT)<sup>223</sup> were applied using the PSI4 electronic structure package,<sup>31</sup> all F-/ISAPT computations employ the

“zeroth-order” SAPT truncation (F-/ISAPT0), and utilize the truncated jun-cc-pVDZ basis set,<sup>22,224</sup> which has been recommended previously<sup>58</sup> for pairing with SAPT0. Unlike in the traditional formulation of SAPT, which describes the interactions between two distinct monomers (commonly denoted “A” and “B”), ISAPT can describe these A–B interactions in the presence of another chemical fragment, which we will denote monomer “C.” We can therefore incorporate the effects of solvation by including all solvent molecules, i.e., the environment, into the ISAPT monomer C. On the other hand, we can also incorporate the solvent into a functional group within the F-SAPT procedure, within either monomer A or B. We denote these three schemes for inclusion of solvent molecules within an F-/ISAPT computation as “EnvX” (X = A, B, C). For clarity, we have illustrated these three schemes in Fig. 7.2 for the solvated aniline–benzene complex (HYD8-1), and provide explicit description below:

- (a) “EnvA”: All solvent molecules placed within effective F-SAPT monomer A (aniline); particular aniline–benzene interaction energy and SAPT components computed via F-SAPT, with both aniline and solvent molecules treated as “functional groups”,<sup>222</sup>
- (b) “EnvB”: All solvent molecules placed within effective F-SAPT monomer B (benzene); particular aniline–benzene interaction energy and SAPT components computed via F-SAPT, with both benzene and solvent molecules treated as “functional groups”,<sup>222</sup> and
- (c) “EnvC”: All solvent molecules placed with ISAPT monomer C; particular aniline–benzene interactions and components computed directly with ISAPT.

Throughout this work, we will adopt the convention that for all HYD8 complexes, the substituted aromatic monomer (ArX) will be denoted monomer A, while the unsubstituted benzene (Bz) will be denoted as monomer B. In the case of HYD8-3 (the solvated benzene dimer), the benzene with C–H bond pointing towards the  $\pi$  cloud of its partner will be monomer A, for consistency with the other complexes. In addition to describing substituted

benzene–benzene interactions in solution, we performed a conventional two-body F-SAPT analysis of the interactions between monomers in the gas phase. In this way, we may investigate directly the effect of solvation on the interactions of interest.

#### 7.3.4 MP2 Results for Solute-Solute, Solute-Solvent, and Three-Body Interactions

To understand the basic physics of the interactions in the solvated dimers, we have supplemented our novel approach for quantifying the “tuning” of interactions embedded in a chemical environment with a traditional many-body approach, wherein we computed interaction energies within the solvated clusters according to a “three body” decomposition. In this picture, “monomer  $\mathcal{A}$ ” is defined as the substituted benzene donating a hydrogen to the C–H/ $\pi$  interactions in Figure 7.1, “monomer  $\mathcal{B}$ ” is the unsubstituted benzene at the base of the T-shaped dimer, and group “ $\mathcal{C}$ ” is the collection of the solvating water molecules. The overall interaction energy between monomers/groups  $\mathcal{A}$ ,  $\mathcal{B}$ , and  $\mathcal{C}$  may be defined as

$$\Delta E_{ABC}^{\text{IE}} = E_{ABC}(\mathcal{ABC}) - E_{\mathcal{A}}(\mathcal{ABC}) - E_{\mathcal{B}}(\mathcal{ABC}) - E_{\mathcal{C}}(\mathcal{ABC}), \quad (7.1)$$

where the subscripts denote the identity of the species and the parenthetical ( $\mathcal{ABC}$ ) denotes that each of the total energies in the expression have been computed in the trimer basis set according to the counterpoise correction scheme of Boys and Bernardi<sup>54</sup> to mitigate basis set superposition error. The counterpoise correction entails computing all required energies using the union of all basis functions in the entire cluster (all three monomers/groups), even when some of the atoms are not required in the computation. Note again that we have grouped all the water solvent molecules together as a single group “ $\mathcal{C}$ ” in this study, meaning that  $\Delta E_{ABC}^{\text{IE}}$  computed as above will be smaller in magnitude than if we computed the interaction relative to the limit in which all molecules (including H<sub>2</sub>O molecules) are infinitely separated; this leads to a simpler analysis because our primary concern here is the interaction between the two solute molecules,  $\Delta E_{AB}^{(2)}$ , and how it is

affected by the environment.

The three-body interaction energy,  $\Delta E_{ABC}^{\text{IE}}$ , can also be computed according to the many-body expansion as

$$\Delta E_{ABC}^{\text{IE}}(\mathcal{ABC}) = \sum_{\mathcal{I} < \mathcal{J}} \Delta E_{\mathcal{IJ}}^{(2)}(\mathcal{ABC}) + \Delta E_{ABC}^{(3)}(\mathcal{ABC}), \quad (7.2)$$

where each of the  $\Delta E_{\mathcal{IJ}}^{(2)}(\mathcal{ABC})$  are the standard two-body interaction energies between monomers  $\mathcal{I}$  and  $\mathcal{J}$ , and  $\Delta E_{ABC}^{(3)}(\mathcal{ABC})$  is the non-additive three-body contribution to the interaction energy. This non-additive contribution can be written as the difference between the overall interaction energy and the sum of the interactions between all pairs:

$$\Delta E_{ABC}^{(3)} = \Delta E_{ABC}^{\text{IE}} - \Delta E_{AB}^{(2)}(\mathcal{ABC}) - \Delta E_{BC}^{(2)}(\mathcal{ABC}) - \Delta E_{AC}^{(2)}(\mathcal{ABC}). \quad (7.3)$$

By computing this quantity, we will investigate the extent to which mutual three-body interaction (which is not directly included in the F-/ISAPT partitioning schemes discussed above) is present for each system. If this quantity is nontrivial, then we hypothesize that the choice of environment binning scheme (visualized in Fig. 7.2) will matter, i.e., the choice of in which ‘‘monomer’’ to include the solvent molecules will not be equivalent.

To construct this non-additive three-body energy correction via Eqn. 7.3 above for each of the 240 complexes in the HYD8 test set, we must compute seven individual computations (dimer energies for AB, BC, AC; monomer energies for A, B, C; trimer energy ABC all in the trimer basis set) for a total of 1,680 individual single-point computations. Considering that these complexes are comprised of up to 64 heavy atoms (for 178 atoms total) and each computation must be performed in the trimer basis set, choosing a level of theory (combination of method and basis set) which can be afforded is of critical concern. Since interaction energies are surprisingly sensitive to the choice of theoretical method,<sup>55,58,64,74,210,225,226</sup> we choose second-order Møller-Plesset perturbation theory (MP2) computations using the juncc-pVDZ basis set, which is Dunning’s correlation-consistent polarized double- $\zeta$  basis set,



augmented with diffuse  $s$  and  $p$  functions for heavy (non-hydrogen) atoms.<sup>224,227</sup> This level of theory represents a compromise between computational accuracy and speed, and it is also expected to yield the most similar interaction energies to the SAPT/jun-cc-pVDZ results discussed above that are the primary focus of this work. The largest of the interaction energy computations comprised 178 atoms (HYD8-2), with 1,786 orbital basis functions and 8,674 auxiliary basis functions.

## 7.4 Results and Discussion

### 7.4.1 Gas-Phase Interactions

Before considering the tuning of ArX–Bz interactions by solvent environment, it is important to first understand the interaction motif of the HYD8 dimers in the gas phase; these are provided in Table ???. In general, these eight ArX–Bz complexes are electrostatically attractive, but with an even larger dispersion term (which is sometimes up to twice as large), and a small stabilizing induction interaction. This is expected for T-shaped  $\pi - \pi$  interactions, as in an idealized T-shaped benzene dimer, SAPT2/jun-cc-pDVZ computes electrostatics, exchange, induction, and dispersion components to be -2.2, 4.9, -0.7, and -4.4 kcal mol<sup>-1</sup>, respectively, yielding a total IE of -2.4 kcal mol<sup>-1</sup>.<sup>228</sup> Of course, differences in geometry and substituents adjust these values somewhat, but they remain similar for each neutral HYD8 complex. For the cationic HYD8-6 (pyridinium–benzene), however, both the total SAPT0 IE and components are enhanced relative its neutral counterparts, with the electrostatic component overshadowing dispersion as the dominant contributor to this increase in total attraction.

Among the HYD8 complexes involving a functionalized benzene (PhX) — i.e., discounting HYD8-6 & HYD8-7 — the magnitude of the electrostatic attraction is smallest for the aniline–benzene complex, increasing in strength as the substituent becomes progressively more electron withdrawing. This is consistent with the fact that, since the *para*-Hydrogen of monomer  $\mathcal{A}$  is the atom closest to the  $\pi$  face of monomer  $\mathcal{B}$ , it stands to reason

that the attraction felt between this increasingly electron deficient site and the electron-rich  $\pi$  face would also increase. This is most pronounced for the nitrobenzene–benzene complex (HYD8-4), where the electrostatic interaction is a full kcal mol<sup>-1</sup> more attractive than for the next dimer in the series, the benzene dimer (HYD8-3). Interestingly, the complex for which the dispersion energy is most attractive is also the nitrobenzene–benzene complex, despite the fact that nitrobenzene is the least polarizable of the benzene derivatives represented in HYD8; this is likely due instead to the simple fact that nitrobenzene also has the largest number of electrons of any of the PhX molecules included in HYD8, since the dispersion energy is known to scale with the number of correlated electrons. Additionally, nitrobenzene–benzene also has the largest exchange-repulsion, most likely due to the polarization of electron density from the  $\pi$  cloud of the benzene ring towards the nitrobenzene *para*-Hydrogen, which is supported by the slightly increased magnitude in the induction energy for this complex relative to other PhX–Bz complexes. Finally, both the pyridine–benzene and pyridinium–benzene complexes exhibit larger magnitude total interaction energies and components relative to complexes involving PhX.

#### 7.4.2 Quantifying ArX–Bz Interactions in Solution via F-/ISAPT

Among the mechanisms of interest by which chemical environment may modulate NCI are (i) the electronic deformation of interacting species by the presence of the interaction, and (ii) statistical averaging due to multiple configurations of the environment. To ensure that conclusions drawn using the F-/ISAPT approach to quantify these effects are not artifactual due to the hydration procedure for our test systems, we must first examine the effect of relaxing the geometries of explicit solvent molecules on the ArX–Bz interactions.

##### *Solvent Molecule Relaxation*

In addition to the relaxed snapshots of solvent configurations considered for all dimers in the HYD8 test set, we have also computed the F-/ISAPT0 interaction energies for HYD8-1,

3, 6, and 7 (aniline–benzene, benzene dimer, pyridine–benzene, and pyridinium–benzene) with *unrelaxed* solvent molecules. Presented in Table ?? are differences between average F-/ISAPT0 IEs and components for each binning scheme for these complexes. Solvent molecule relaxation does not significantly affect the computed ArX–Bz interactions or its components, as these differences are nearly always below 0.1 kcal mol<sup>-1</sup> with only a few exceptions; namely, differences of up to about 0.3 kcal mol<sup>-1</sup> are observed for the EnvA grouping of the hydrated pyridinium–benzene complex. This finding is not surprising, as in the EnvA grouping the water molecules are expressly included with the pyridinium cation in a conventional two-body F-SAPT computation, whereby their electron densities are fully interacting in the preparation of the zeroth-order wavefunction for monomer A. Therefore, even small changes in the positions of the water molecules or their internal geometry can result in large effects in the F-SAPT interaction with benzene. When expanding our environment to include 50 water molecules, however, the effect of relaxing the positions of solvent molecules is damped, where differences no larger than 0.1 kcal mol<sup>-1</sup> are observed, even for the hydrated pyridinium–benzene complex. As a result of this finding, we are confident that the choice of relaxing or not relaxing solvent molecules will not affect the validity of our conclusions about the tuning of solute interactions.

#### *Environment Binning Scheme*

Visualized in Fig. 7.3 are total F-/ISAPT0 IEs and components for each environment grouping (EnvX; X = A, B, C) of the hydrated benzene dimer (HYD8-3; top panel) and the hydrated pyridinium–benzene complex (HYD8-7; bottom panel). Most strikingly, the hydrated benzene dimer (Fig. 7.3.a) exhibits very little variation between the average IE or components depending on environment binning scheme; furthermore, F-/ISAPT IEs and components for each binning scheme are quite similar to those computed for the gas-phase benzene dimer with conventional SAPT0. For the hydrated pyridinium–benzene complex (Fig. 7.3.b), on the other hand, notable variations between binning schemes exist for both

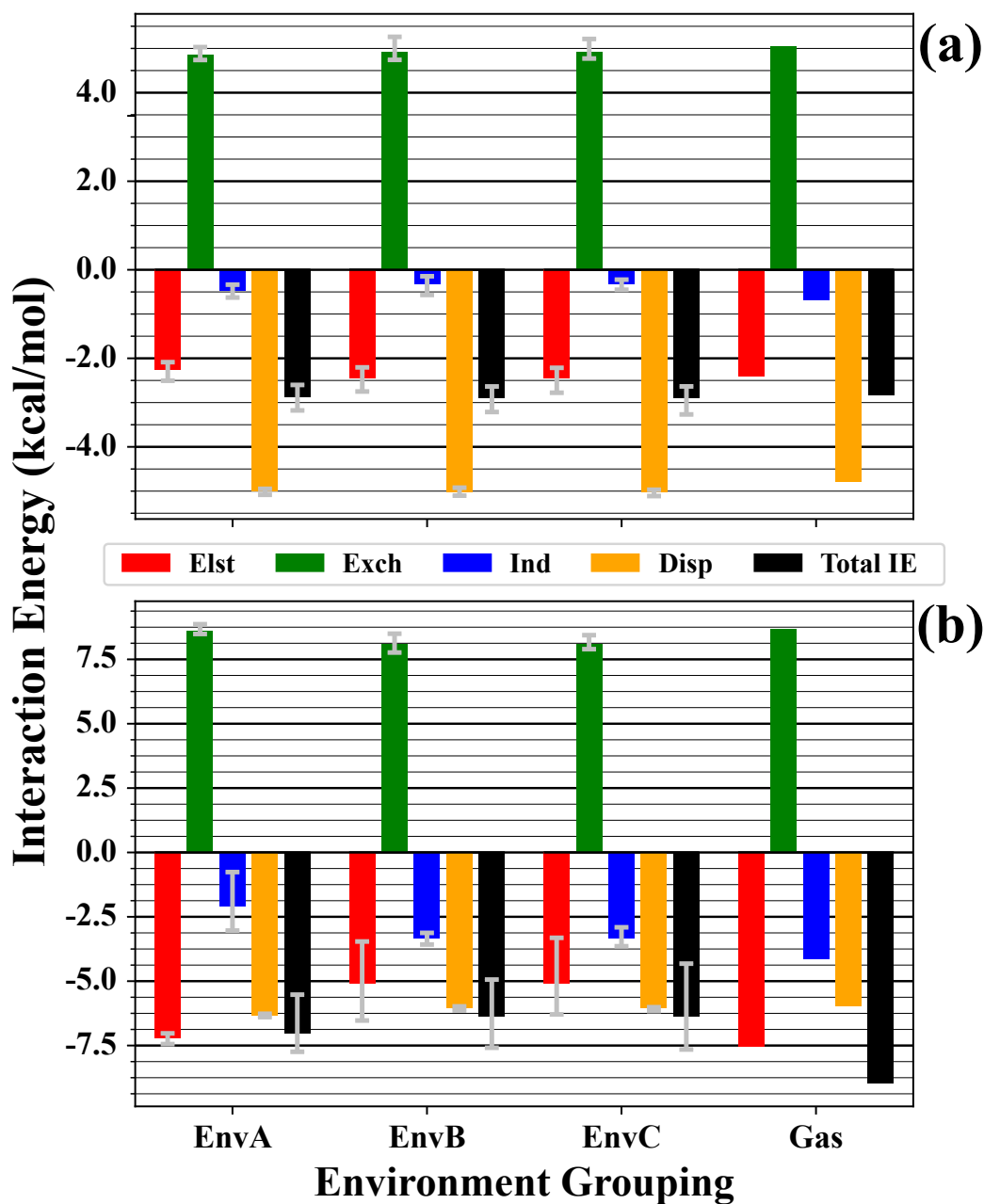


Figure 7.3: Total interaction energies and SAPT components for (a) HYD8-3m-w50 (benzene dimer) and (b) HYD8-7m-w50 (pyridinium–benzene) complexes solvated by 50 explicit solvent molecules, computed at the F-/ISAPT0/jun-cc-pVDZ level of theory and averaged over all ten relaxed solvent configurations. “EnvX” labels indicate that explicit solvent molecules are contained within monomer “X” during the SAPT computation (see text). Error bars encompassing the full range of values across all snapshots are also provided for SAPT terms and total IEs. Furthermore, we have provided a set of bars corresponding to the conventional two-body F-SAPT computation in the gas phase, i.e., in the absence of explicit solvent molecules. See Section II C for additional details regarding our nomenclature and details of the F-/ISAPT computations.

total IEs and components, and each solvated interaction exhibits some degree of screening relative to the interaction in the gas phase. For this complex, the largest difference from the gas phase is found for the electrostatic component in EnvB and EnvC of  $\sim 1$  and  $\sim 0.9$  kcal mol<sup>-1</sup>, respectively, and for the induction component of EnvA of  $\sim 0.9$  kcal mol<sup>-1</sup>. The other six solvated dimers examined here behave like the benzene dimer, in that the average interaction energies and components hardly differ between the various environment binning schemes, with deviations between EnvA/B/C typically of just a few tenths of one kcal mol<sup>-1</sup>. Furthermore, deviations from the gas phase interactions are similarly small, leading us to conclude that other than for the cationic pyridinium–benzene system, each binning scheme is essentially equivalent, and furthermore that the gas phase SAPT0 computation provides an adequate description even of the solvated interactions.

#### *Ranges due to Solvent Configuration*

Consulting Figures S-1–S12 and Tables S-4–S-15 in the SI, the range in the interaction energy, or its components, due to the different solvent configurations is fairly constant for most solvated dimers considered, regardless of the environment binning scheme (EnvA/B/C). For each F-/ISAPT component, differences between solvent configurations typically range from  $\sim 0.25$ – $0.5$  kcal/mol for exchange or induction,  $0.5$ – $1.0$  kcal/mol for electrostatics,  $0.1$  kcal/mol for dispersion, and  $0.5$ – $1.5$  kcal mol<sup>-1</sup> for total interaction energies. Somewhat larger ranges are seen in some cases, including the pyridinium–benzene system in panel (b) of Fig. 7.3, which demonstrates spreads of  $3$  kcal/mol among the electrostatic energies of the different solvent configurations for EnvB/C, and spreads of  $2.5$ – $3$  kcal/mol in the total interaction energies. The range in electrostatics is also somewhat larger than normal for some environmental binnings of HYD8-4 (nitrobenzene–benzene) and the range in induction values across snapshots increases to  $1$  kcal/mol for some environment binnings for HYD8-4 and HYD8-6 (pyridine–benzene), and can grow to more than  $2$  kcal/mol for HYD8-7 (pyridinium–benzene).

The increased susceptibility of solute interactions to solvent configuration for nitrobenzene–benzene, pyridine–benzene, and pyridinium–benzene seems counterintuitive; indeed, these complexes are comprised of the ArX molecules with largest dipole moments among all HYD8 complexes, with gas-phase dipole  $\mu_{\mathcal{A}}$  for monomer  $\mathcal{A}$  of  $\mu_{\mathcal{A}} = 4.22, 2.19,$  and XX Debye for HYD8-4, 6, 7, respectively, as compared to all other complexes in our test set for which  $\mu_{\mathcal{A}} \leq 1.5$  Debye. Since these molecules have the largest permanent dipole moments, should they not be the *least* polarizable by the solvent, and therefore their interactions with benzene less susceptible to changes in the solvent configuration? As it turns out, the opposite is actually true: regardless of gas-phase dipole moment, solvent actually *enhances* the molecular dipole moment by between 30–40%.<sup>229</sup> Furthermore, we believe that the variation of solute interactions between different solvent configurations is due not to the polarization of, e.g., nitrobenzene by the solvent environment, but rather the reverse. Since the polarizability of the solvent environment is highly dependent on the configuration of individual solvent molecules, the electronic deformation of the solvent by solute molecules (and therefore, the solvent’s tuning of the solute interactions) is also highly dependent on the solvent configuration. This higher-order effect seems only to be present when a solute monomer has a permanent electronic dipole moment of  $\mu \geq 2$  Debye, as the other HYD8 complexes do not exhibit the same variability of F-/ISAPT interactions and components with respect to solvent configuration.

### 7.4.3 Many-Body Analysis of Solvated Interactions

To investigate the possibility that variation in F-/ISAPT IEs and components between different solvent configurations could be due to some higher-order interactions between the “monomers,” as well as to validate our three-body picture (wherein we group all solvent molecules together into a single SAPT monomer), we have computed the non-additive three-body component of the total trimer energy according to Eqn. 7.3 at the HF/jun-cc-pVDZ and MP2/jun-cc-pVDZ levels of theory, as these combinations of methods and basis

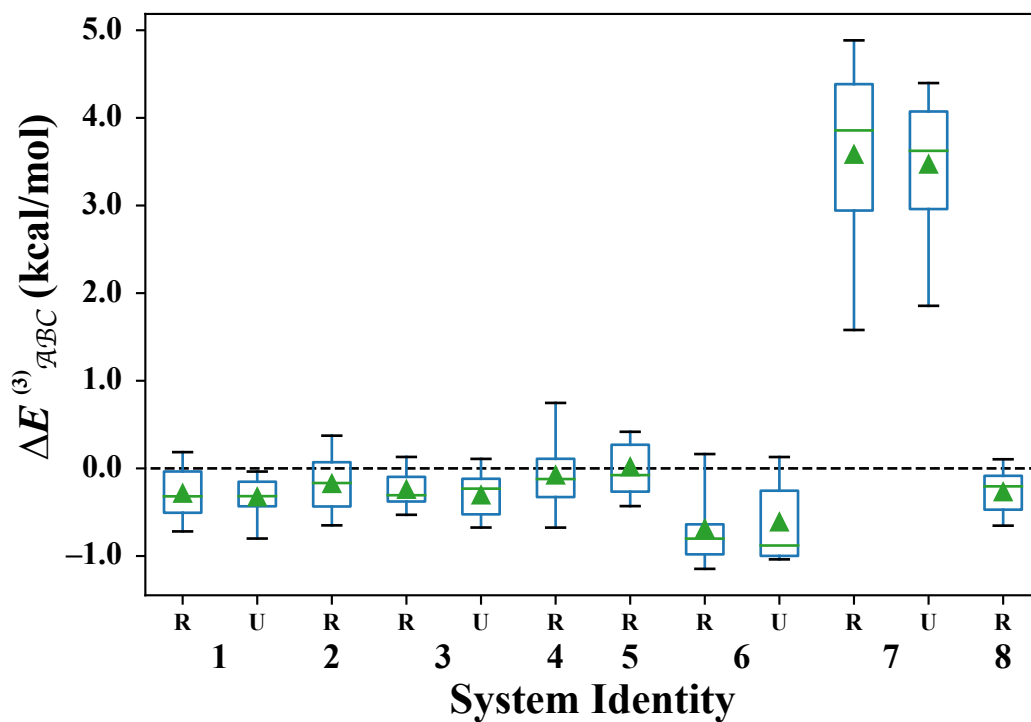


Figure 7.4: Box-and-whisker plots representing the non-additive three-body correction to total “trimer” energy ( $\Delta E_{ABC}$ ; kcal mol<sup>-1</sup>), for both relaxed (R) and unrelaxed (U) solvent configurations of each doubly-solvated HYD8 complex, computed at the HF/jun-cc-pVDZ level of theory (see text for details of solvent configuration selection and preparation). Boxes encompass the first (Q1) through third (Q3) quartiles of  $\Delta E_{ABC}$ , with values corresponding to the median (Q2) and mean  $\Delta E_{ABC}$  indicated as a solid green bar and green triangle, respectively. Additionally, whiskers encompass the full range of  $\Delta E_{ABC}$  values for all solvent configurations.

set provide the most direct comparison to the HF-in-HF embedding present in ISAPT0 and the SAPT0 computations themselves. Presented in Fig. 7.4 are box-and-whisker plots visualizing the distribution of non-additive three-body energy contribution,  $\Delta E_{ABC}^{(3)}$ , over different solvent configurations of each HYD8 complex. Immediately, it is apparent that the complex for which the largest non-additive behavior is present is for HYD8-7, the cationic pyridinium–benzene complex, which exhibits  $\Delta E_{ABC}^{(3)}$  between 3-4 kcal mol<sup>-1</sup> larger in magnitude on average versus all other HYD8 complexes.

This significant non-additive behavior is likely the cause of variation seen in pyridinium–benzene interactions both between different binning schemes (EnvA/B/C) and between different solvent configurations. Interestingly, neither HYD8-4 (nitrobenzene–benzene) nor HYD8-6 (pyridine–benzene) seem to exhibit a large non-additive interaction, as seemed to be indicated by the variation between solute interactions within different solvent configurations. Instead, however, the ranges of  $\Delta E_{ABC}^{(3)}$  between different solvent configurations are slightly larger for these systems ( $\sim 1-1.5$  kcal mol<sup>-1</sup>) than for other HYD8 complexes ( $\sim 0.5-0.8$  kcal mol<sup>-1</sup>). Instead of non-additive behavior on average, this variation in non-additivity between different solvent configurations of HYD8-4 and HYD8-6 may be the cause of the variations present in F-/ISAPT0 IEs and components between solvent configurations. This is supported by the fact that for HYD8-7, even larger ranges in  $\Delta E_{ABC}^{(3)}$  of  $\sim 2.3-3.5$  kcal mol<sup>-1</sup> are present, which matches the behavior for the magnitude of variation in F-/ISAPT IEs and components observed for this complex relative to other HYD8 members. It appears, therefore, that differences between environment binning scheme is due largely to the *permanent* non-additivity present in a given hydrated dimer, while variations in IE and components between different solvent configurations are due to changes in the non-additivity with respect to the solvent configuration.

While the preceding analysis has been performed for three-body interactions at the HF/jun-cc-pVDZ level of theory, it is worth noting that the two-body ArX–Bz interaction energy at this level is in some cases repulsive (see, e.g., Table S-15–S-22 in the Sup-



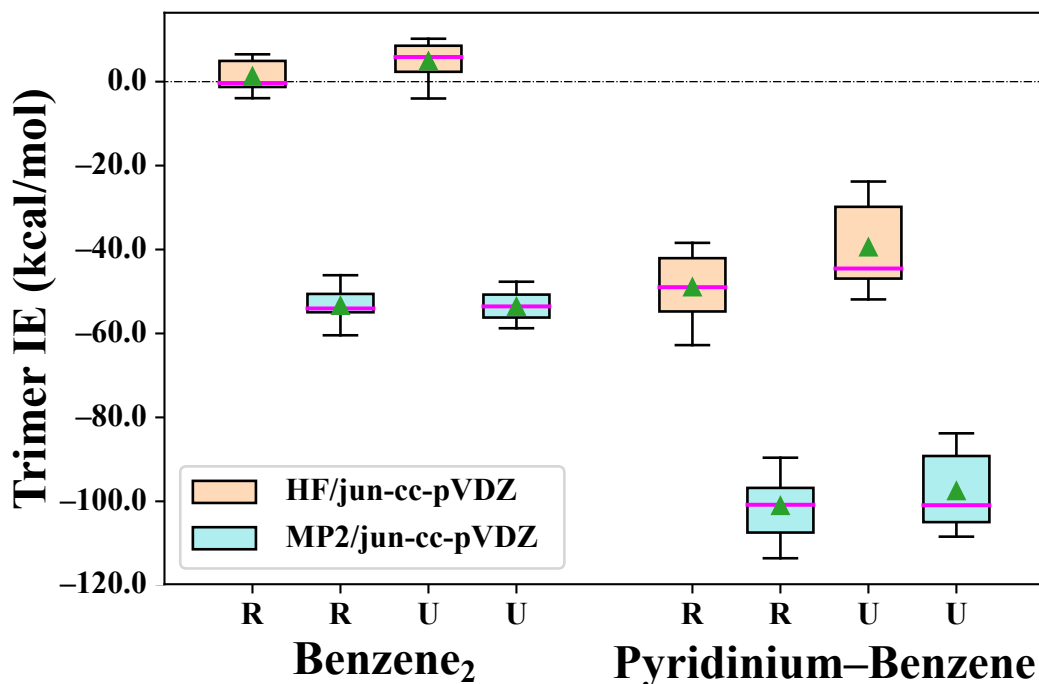


Figure 7.5: Box-and-whisker plots representing the total “trimer” interaction energy ( $\Delta E_{ABC}^{\text{IE}}$ ; kcal mol<sup>-1</sup>) for both relaxed (R) and unrelaxed (U) solvent configurations of benzene dimer (HYD8-3) and pyridinium–benzene (HYD8-7) complexes hydrated by 50 explicit water molecules, computed at the HF/jun-cc-pVDZ (orange boxes) and MP2/jun-cc-pVDZ (blue boxes) levels of theory. Boxes encompass the first (Q1) through third (Q3) quartiles of  $\Delta E_{ABC}^{\text{IE}}$ , with values corresponding to the median (Q2) and mean  $\Delta E_{ABC}^{\text{IE}}$  indicated as a solid green bar and green triangle, respectively. Additionally, whiskers encompass the full range of  $\Delta E_{ABC}^{\text{IE}}$  values for all solvent configurations.

plementary Materials). This is unphysical, as we know the T-shaped structure of these complexes are optimal at the B3LYP-D3M(BJ)/aug-cc-pVDZ level of theory. The repulsiveness of the ArX–Bz interactions computed with HF/jun-cc-pVDZ is due to the fact that at the Hartree–Fock level, no description of dispersion is present, as only electrostatics, exchange-repulsion, and lower-order induction are accounted for. Therefore, to ensure our conclusions for the non-additive three-body interactions based on energies computed at the Hartree–Fock level are relevant, we have performed an identical analysis at the MP2/jun-cc-pVDZ level of theory, which is the closest supermolecular wavefunction method to SAPT0. Provided in Fig. 7.5 is a comparison of the full trimer IE computed with HF/jun-cc-pVDZ and MP2/jun-cc-pVDZ for the hydrated benzene dimer (HYD8-3) and hydrated pyridinium–benzene complex (HYD8-7) (all components of the three-body MBE are given in Tables S-23–S-30 in the Supplementary Information). When using MP2/jun-cc-pVDZ, all two-body ArX–Bz IEs become attractive; furthermore, all trimer IEs are *also* more attractive, by a constant shift of approximately 50-60 kcal mol<sup>-1</sup> on average for all systems. Despite the shift in value for the full two- and three-body interaction energies, neither the ranges of three-body MBE components over different solvent configurations nor the non-additive three-body interaction is shifted between HF and MP2 descriptions. Therefore, the analysis of non-additive three-body interactions at the Hartree–Fock level above is consistent with MP2, and the conclusions thereof are retained.

#### 7.4.4 Effect of Multiple Hydration Shells

In EnvA and EnvB, the solvent molecules are grouped into a traditional SAPT monomer; therefore, they contribute to the  $\mathcal{O}(N^5)$  computational scaling of F-/ISAPT0, where  $N$  is proportional to the system size, whereas for EnvC, the solvent molecules are only treated at the  $\mathcal{O}(N^4)$ -scaling Hartree–Fock level. It is of critical interest, therefore, to determine exactly “how large” an environment is necessary to include in the F-/ISAPT computation to ensure the tuning of solute interactions is captured properly. We have therefore examined

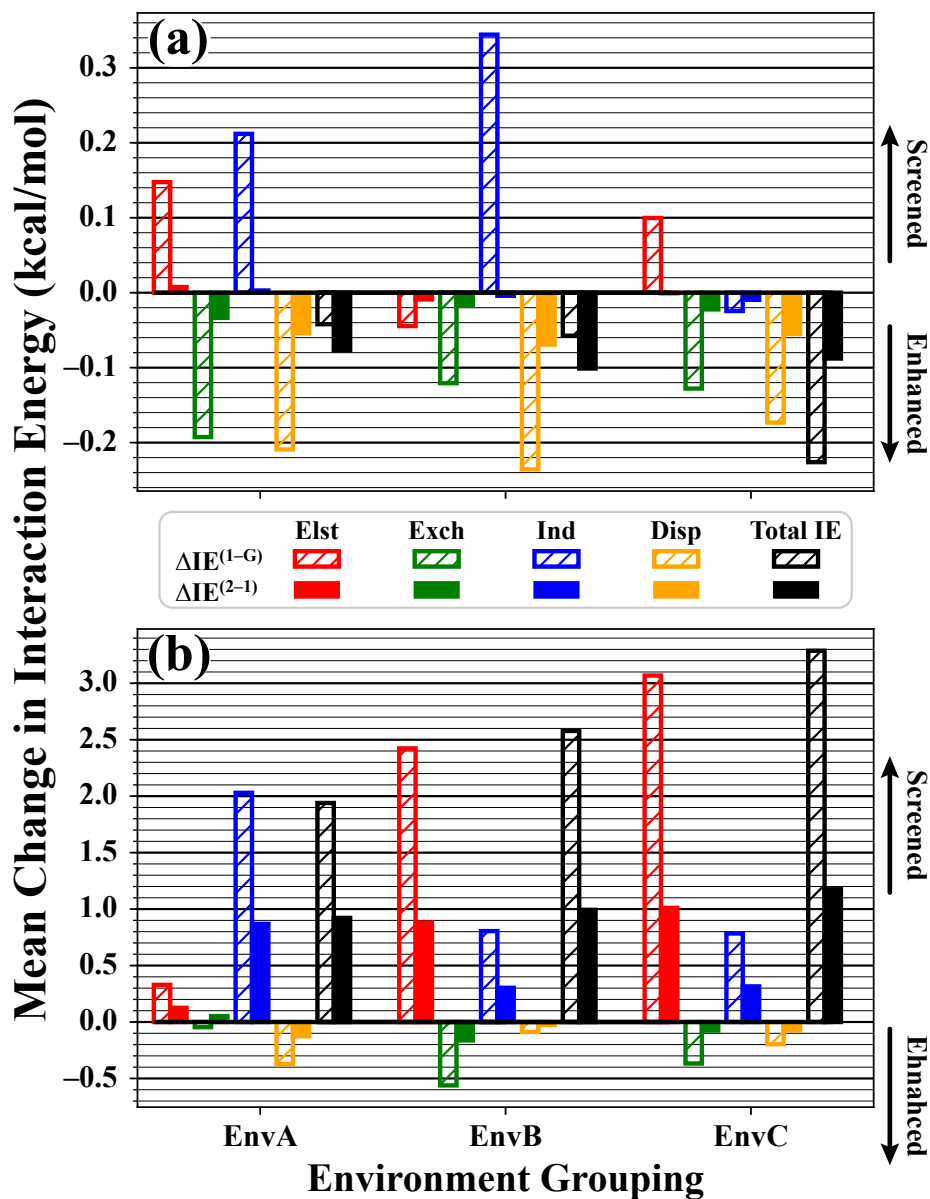


Figure 7.6: Mean change in total interaction energies and SAPT components upon first-shell solvation ( $\Delta IE^{1-G}$ ; striped bars) and second-shell solvation ( $\Delta IE^{2-1}$ ; solid bars) for the (a) HYD8-3mX (benzene dimer) and (b) HYD8-7mX (pyridinium-benzene) complexes, averaged over values computed at the F-/ISAPT0/jun-cc-pVDZ level of theory for all ten relaxed solvent configurations. “EnvX” labels indicate that explicit solvent molecules are contained within monomer “X” during the SAPT computation (see text). See Section II C for additional details regarding our nomenclature and details of the F-/ISAPT computations.

the difference between total ArX–Bz IEs and components upon the addition of the first solvation shell ( $\Delta E_{\text{int}}^{1-\text{G}}$ ) and upon addition of the second solvation shell ( $\Delta E_{\text{int}}^{2-1}$ ), visualized in Fig. 7.6 for each environment binning of the hydrated benzene dimer (HYD8-3; top panel, a) and hydrated pyridinium–benzene complex (HYD8-7; bottom panel, b). For systems where the non-additive trimer interaction is small (e.g., benzene dimer, Fig. 7.6.a), we find that the differences upon addition of further solvation shells is quite small. Additional levels of solvation (i.e., gas→1, 1→2) slightly enhances the attractiveness of IEs, thanks to slightly smaller or negligible changes in exchange repulsion and progressively more favorable dispersion interactions. Both electrostatics and induction are screened slightly by the addition of the first solvation shell, but no additional screening is caused by the addition of the second solvation shell. For these “additive” systems, the  $\Delta \text{IE}^{1-\text{G}}$  is less than 0.3 kcal mol<sup>-1</sup> (~10% of gas-phase IE), with  $\Delta \text{IE}^{2-1}$  only slightly larger.

For significantly non-additive systems (e.g., pyridinium–benzene, Fig. 7.6.b), the effects of adding both the first and second solvation shells are more notable. The electrostatics and induction are again screened upon addition of the first solvation shell, but to a larger magnitude; this leads to an overall screening of total IE of between ~1.9-3.3 kcal mol<sup>-1</sup>. Exchange-repulsion and dispersion exhibit identical behavior to the benzene dimer (Fig. 7.6.a), however since these effects are on the order of only tenths of one kcal mol<sup>-1</sup>, their effect is negligible compared to electrostatics and induction. Upon the addition of the second solvation shell, similar behavior is observed for each component, however the difference is only about half as large as was observed for the addition of the first solvation shell. Overall, the changes in IEs and components upon adding solvation shells indicates that while the solute interactions were roughly “converged” for additive systems after only a single solvation shell was included, this is not the case for non-additive systems, where in fact not even two solvation shells is sufficient to converge the solute IEs or components. This indicates that, especially for non-additive systems, the presence of significantly longer-range effects likely necessitates that larger environments be considered.

## 7.5 Summary and Conclusions

The extent to which chemical environment “tunes” non-covalent interactions (NCI), as well as the best manner in which to account for this effect, are open questions in the computational molecular sciences. To address this, we have presented an approach based on the functional-group partition and intramolecular formulation of symmetry-adapted perturbation theory (F-/ISAPT) which can (i) compute non-covalent interactions embedded in a chemical environment, and (ii) quantify the tuning of these interactions due to the environment relative to the interactions in the gas phase. We have applied our approach to quantify the extent to which explicit water solvent modulates  $\pi - \pi$  interactions in several functionalized, T-shaped arene–benzene complexes, hydrated by a statistically diverse set of solvent configurations. We have found that, for systems wherein no significant non-additive three-body interaction between the monomers and the collective solvent environment are present, the solvent environment does not significantly tune  $\pi - \pi$  interactions, either due to the choice of system partitioning or solvent configuration. For systems where the non-additive three-body interaction is significant, however — i.e., where it either is greater than  $\sim 2$  kcal mol<sup>-1</sup> or where it deviates between solvent configurations by greater than  $\sim 1$  kcal mol<sup>-1</sup> — the solvent environment *does* tune the interaction, sometimes by up to several kcal mol<sup>-1</sup> for both total interaction energies and F-/ISAPT components. Finally, we have shown that for these non-additive systems, even *two* hydration shells of 50 explicit water molecules within 7 Å of the solute complex may not be sufficient to ensure convergence of the solute–solute interactions towards the continuum limit, whereas for additive systems, only a single shell of 28 water molecules within 3 Å is necessary for convergence.

## 7.6 Acknowledgements

The authors gratefully acknowledge financial support from Bristol-Myers Squibb, and from the U.S. National Science Foundation through grant CHE-1566192.

## **PART V**

### **CONCLUSIONS**

## CHAPTER 8

### CONCLUSIONS AND OUTLOOK

#### 8.1 Conclusions

In this Thesis, we have taken a three-pronged approach towards the quantum chemical investigation of non-covalent interactions in extended chemical systems and diverse environments: first, we have developed protocols by which the structural and energetic properties of non-bonded complexes may be benchmarked, in order to establish best practices by which to obtain “the right answers for the right reasons;” second, we have developed semi-empirical perturbative approaches based on symmetry-adapted perturbation theory (SAPT), simultaneously reducing the computational expense of SAPT as well as to increase its accuracy and applicability for diverse non-covalent interaction motifs in extended chemical systems; and finally, we have applied these approaches, together with the functional group partitioned and intramolecular formulations of SAPT, to quantify the tuning of non-covalent interactions by their chemical environment. Taken together, these advances lay the foundation for future efforts leveraging computation to quantitatively study the effects of non-covalent interactions in new arenas, e.g., the long-range effects on enzyme activity due to allosteric inhibition or missense mutation.

In a broader sense, however, this Thesis feels very much like the end of an era in the field of non-covalent interactions. Gone are the days where the choice of computational molecular scientists is between quantum mechanics or classical mechanics, and so to are the days where simply assessing the quality of some method or another is of interest. Much of this space has already been explored, and the corners of the map have been largely filled in. Much of the current state of the field — and indeed, even future of quantum chemistry itself — seems dominated by the development of data-driven methodologies leveraging

machine learning (ML) and artificial intelligence (AI) to predict NCI and a host of other properties, trained to reproduce QM computations at a fraction of the price. Many of the concerns for ML methods going forward are ones of data curation, feature engineering, and defining appropriate, transferable models. Fortunately, the construction of massive datasets of very high-quality energetics (or properties) against which ML models may be fit is facilitated by the field’s collective expertise in benchmarking, to which this Thesis contributes. ML models are only as good as their training, however, so when completely novel chemical phenomena are being studied, neither their accuracy, nor even their appropriateness, is guaranteed. While ML may be unreliable in these scenarios, quantum mechanics can still provide definitive predictions of chemical behavior, so long as it is not prohibitively expensive. As such, the approaches developed in this Thesis for studying NCI in extended chemical systems and diverse chemical environments are hoped to provide insight where significantly more empirical approaches are likely to break down, thereby keeping the door open to new chemical discovery.

## **8.2 Outlook**

Based on the advances made in this Thesis, there exist a number of possible future avenues, including the two very briefly proposed here.

### 8.2.1 Towards a Multi-Level Embedded SAPT

In Chapter 7, we have observed that for systems where there exists a nontrivial mutual, many-body attraction or repulsion, the environment can significantly tune both the total interaction strength as well as interaction components. Also, when including a progressively larger environment in the computation, the values for IE components (and even the IEs themselves) begin to converge towards some “bulk” limit, which can be interpreted as the real interaction strength between two molecules embedded within an infinitely large environment. In order to investigate, e.g., the effects of allosteric control on enzyme function,



the size of the environment which must be included in the ISAPT computation may prove prohibitively expensive even for this approach. In order to both accelerate convergence towards the bulk limit and reduce the computational expense incurred by including long-range contacts in the chemical environment, the ISAPT methodology could be extended to be compatible with a multi-level embedding scheme amenable to very large systems, thereby enabling the study of embedded interactions natively in systems which were previously inaccessible. For example, the HF-in-HF embedding scheme for generating the ISAPT zeroth-order wavefunctions could be extended by further mechanically embedding the supersystem wavefunction into effective fragment potentials or classical force fields, or a field of distributed multipoles, point charges, or even a polarizable continuum model (PCM). Additionally, these schemes can be further combined to create a hierarchical series of embeddings, ordered by the relative exactness of the method.

### 8.2.2 The Influence of Long-Range Contacts in Drug-Protein Binding Specificity

By leveraging the functional group partition of SAPT (F-SAPT), Parrish *et al.* showed that the local contact model, where direct interactions between residue side chains in a protein's binding pocket and ligand functional groups are hypothesized to be the driving force for protein-ligand binding, was insufficient to justify the relative binding affinity of chloro versus methyl aryl substituted factor Xa (fXa) inhibitor drugs.<sup>69</sup> Instead, it was shown that mid- to long-range contacts, in particular involving peptide bonds in the protein backbone itself, were the cause of binding specificity. The F-SAPT0-D method developed in Chapter 6 can be applied in order to examine these long-range contacts more directly by including a larger subsystem of fXa in each computation. Furthermore, it will be of interest to investigate a wider range of conformation space by sampling structures along various binding trajectories and submitting these structures to F-SAPT. While we hypothesize that the importance of peptide bond contacts will decay with distance from the binding pocket, it is possible that these contacts are cooperative and more important than previously expected.

## REFERENCES

- <sup>1</sup>A. Chalmers, in *The stanford encyclopedia of philosophy*, edited by E. N. Zalta (Metaphysics Research Lab, Stanford University, 2019).
- <sup>2</sup>D. Lindley, *the myth of a unified theory* (Basic Books, May 1993).
- <sup>3</sup>P. A. M. Dirac, *Proceedings of the Royal Society A* **123**, 714–733 (1929).
- <sup>4</sup>N. S. Ostlund and A Szabo, 1982.
- <sup>5</sup>M Born and R Oppenheimer, *Annalen der Physik* **389**, 457–484 (1927).
- <sup>6</sup>D. G. A. Smith, L. A. Burns, D. A. Sirianni, D. R. Nascimento, A. Kumar, A. M. James, J. B. Schriber, T. Zhang, B. Zhang, A. S. Abbott, E. J. Berquist, M. H. Lechner, L. A. Cunha, A. G. Heide, J. M. Waldrop, T. Y. Takeshita, A. Alenaizan, D. Neuhauser, R. A. King, A. C. Simmonett, J. M. Turney, H. F. Schaefer, F. A. Evangelista, A. E. DePrince III, T. D. Crawford, K. Patkowski, and C. D. Sherrill, *J. Chem. Theory Comput.*, 1–8 (2018).
- <sup>7</sup>K. D. Vogiatzis, D. Ma, J. Olsen, L. Gagliardi, and W. A. de Jong, *J. Chem. Phys.* **147**, 184111–14 (2017).
- <sup>8</sup>C. Miller and M. S. Plesset, *Phys. Rev.* **46**, 618–622 (1934).
- <sup>9</sup>G. D. Purvis and R. J. Bartlett, *J. Chem. Phys.* **76**, 1910–1918 (1982).
- <sup>10</sup>K. Raghavachari, G. W. Trucks, J. A. Pople, and M. Head-Gordon, *Chem. Phys. Lett.* **157**, 479–483 (1989).
- <sup>11</sup>M. L. Leininger, W. D. Allen, H. F. Schaefer, and C. D. Sherrill, *J. Chem. Phys.* **112**, 9213–9222 (2000).
- <sup>12</sup>D. Smith, P Jankowski, M. S. J. o. chemical, and 2014, ACS Publications 10.1021/ct500347q.
- <sup>13</sup>B. W. Hopkins and G. S. Tschumper, *J. Phys. Chem. A* **108**, 2941–2948 (2004).
- <sup>14</sup>J. Řezáč and P. Hobza, *J. Chem. Theory Comput.* **9**, 2151–2155 (2013).
- <sup>15</sup>J. Řezáč, L. Simova, and P. Hobza, *J. Chem. Theory Comput.* **9**, 364–369 (2013).

- <sup>16</sup>R. Podeszwa, K. Patkowski, and K. Szalewicz, *Phys. Chem. Chem. Phys.* **12**, 5974–5979 (2010).
- <sup>17</sup>G. S. Tschumper, M. L. Leininger, B. C. Hoffman, E. F. Valeev, H. F. Schaefer, and M. Quack, *J. Chem. Phys.* **116**, 690–701 (2002).
- <sup>18</sup>M. Schutz, S. Brdarski, P. O. Widmark, R. Lindh, and G. Karlstrom, *J. Chem. Phys.* **107**, 4597–4605 (1997).
- <sup>19</sup>T. Kato, *Commun. Pure Appl. Math.* **10**, 151–177 (1957).
- <sup>20</sup>A. Halkier, T. Helgaker, P. Jørgensen, W. Klopper, H. Koch, J. Olsen, and A. K. Wilson, *Chem. Phys. Lett.* **286**, 243–252 (1998).
- <sup>21</sup>A. Halkier, T. Helgaker, P. Jørgensen, W. Klopper, and J. Olsen, *Chem. Phys. Lett.* **302**, 437–446 (1999).
- <sup>22</sup>T. H. Dunning, *J. Chem. Phys.* **90**, 1007–1023 (1989).
- <sup>23</sup>A. L. L. East and W. D. Allen, *J. Chem. Phys.* **99**, 4638–4650 (1993).
- <sup>24</sup>A. G. Császár, W. D. Allen, and H. F. Schaefer, *J. Chem. Phys.* **108**, 9751–9764 (1998).
- <sup>25</sup>H. Koch, B. Fernández, and O. Christiansen, *J. Chem. Phys.* **108**, 2784–2790 (1998).
- <sup>26</sup>S. Tsuzuki, K. Honda, T. Uchimaru, M. Mikami, and K. Tanabe, *J. Am. Chem. Soc.* **124**, 104–112 (2002).
- <sup>27</sup>M. O. Sinnokrot, E. F. Valeev, and C. D. Sherrill, *J. Am. Chem. Soc.* **124**, 10887–10893 (2002).
- <sup>28</sup>P. Jurečka and P. Hobza, *Chem. Phys. Lett.* **365**, 89–94 (2002).
- <sup>29</sup>M. O. Sinnokrot and C. D. Sherrill, *J. Phys. Chem. A* **108**, 10200–10207 (2004).
- <sup>30</sup>M. S. Marshall, L. A. Burns, and C. D. Sherrill, *J. Chem. Phys.* **135**, 194102 (2011).
- <sup>31</sup>R. M. Parrish, L. A. Burns, D. G. A. Smith, A. C. Simmonett, A. E. DePrince, E. G. Hohenstein, U. Bozkaya, A. Y. Sokolov, R. Di Remigio, R. M. Richard, J. F. Gonthier, A. M. James, H. R. McAlexander, A. Kumar, M. Saitow, X. Wang, B. P. Pritchard, P. Verma, H. F. Schaefer, K. Patkowski, R. A. King, E. F. Valeev, F. A. Evangelista, J. M. Turney, T. D. Crawford, and C. D. Sherrill, *J. Chem. Theory Comput.* **13**, 3185–3197 (2017).
- <sup>32</sup>E. A. Hylleraas, *Z. Phys.* **54**, 347 (1929).

- <sup>33</sup>S. Ten-no, Chem. Phys. Lett. **398**, 56–61 (2004).
- <sup>34</sup>S. Grimme, A. Hansen, J. G. Brandenburg, and C. Bannwarth, Chem. Rev. **116**, 5105–5154 (2016).
- <sup>35</sup>S. Grimme, J. Antony, S. Ehrlich, and H. Krieg, J. Chem. Phys. **132**, 154104 (2010).
- <sup>36</sup>S. Grimme, S. Ehrlich, and L. Goerigk, J. Comput. Chem. **32**, 1456–1465 (2011).
- <sup>37</sup>E. R. Johnson and A. D. Becke, J. Chem. Phys. **123**, 024101 (2005).
- <sup>38</sup>J. Chai and M. Head-Gordon, Phys. Chem. Chem. Phys. **10**, 6615 (2008).
- <sup>39</sup>D. G. A. Smith, L. A. Burns, K. Patkowski, and C. D. Sherrill, J. Phys. Chem. Lett. **7**, 2197–2203 (2016).
- <sup>40</sup>J. Yang, W. Hu, D. Usvyat, D. Matthews, M. Schuetz, and G. K. Chan, Science **345**, 640–643 (2014).
- <sup>41</sup>K. D. Nanda and G. J. O. Beran, J. Chem. Phys. **137**, 174106 (2012).
- <sup>42</sup>P. J. Bygrave, N. L. Allan, and F. R. Manby, J. Chem. Phys. **137**, 164102 (2012).
- <sup>43</sup>A. Tkatchenko, D. Alfè, and K. S. Kim, J. Chem. Theory Comput. **8**, 4317–4322 (2012).
- <sup>44</sup>J. Friedrich, H. Yu, H. R. Leverentz, P. Bai, J. I. Siepmann, and D. G. Truhlar, J. Phys. Chem. Lett. **5**, 666–670 (2014).
- <sup>45</sup>M. J. Gillan, D. Alfe, P. J. Bygrave, C. R. Taylor, and F. R. Manby, J. Chem. Phys. **139**, 114101 (2013).
- <sup>46</sup>G. R. Medders, V. Babin, and F. Paesani, J. Chem. Theory Comput. **9**, 1103–1114 (2013).
- <sup>47</sup>A. Stone (Oxford University Press, 2013).
- <sup>48</sup>C Sherrill and K. Merz, in *Many-body effects and electrostatics in biomolecules* (Pan Stanford, Boca Raton, FL, Mar. 2016), pp. 65–120.
- <sup>49</sup>E. G. Hohenstein and C. D. Sherrill, WIREs Computational Molecular Science **2**, 304–326 (2012).
- <sup>50</sup>A. O. de la Roza and G. A. DiLabio (Elsevier Science, 2017).
- <sup>51</sup>F London, Trans. Faraday Soc. **33**, 8b–26 (1937).

- <sup>52</sup>F London, *Z. Physik. Chem. B* **14**, 222–251 (1930).
- <sup>53</sup>R. M. Richard, B. W. Bakr, and C. D. Sherrill, *J. Chem. Theory Comput.*, [acs.jctc.7b01232–15](#) (2018).
- <sup>54</sup>S. F. Boys and F. Bernardi, *Mol. Phys.* **19**, 553–566 (1970).
- <sup>55</sup>L. A. Burns, M. S. Marshall, and C. D. Sherrill, *J. Chem. Phys.* **141**, 234111 (2014).
- <sup>56</sup>D. M. Bates and G. S. Tschumper, *J. Phys. Chem. A* **113**, 3555–3559 (2009).
- <sup>57</sup>L. A. Burns, M. S. Marshall, and C. D. Sherrill, *J. Chem. Theory Comput.* **10**, 49–57 (2014).
- <sup>58</sup>T. M. Parker, L. A. Burns, R. M. Parrish, A. G. Ryno, and C. D. Sherrill, *J. Chem. Phys.* **140**, 094106 (2014).
- <sup>59</sup>C. D. Sherrill, in *Non-covalent interactions in quantum chemistry and physics*, edited by A. O. de la Roza and G. A. DiLabio (Elsevier, 2017), pp. 137–168.
- <sup>60</sup>E. R. Johnson, in *Non-covalent interactions in quantum chemistry and physics*, edited by A. O. de la Roza and G. A. DiLabio (Elsevier, 2017), pp. 169–194.
- <sup>61</sup>E. Schröder, V. R. Cooper, K. Berland, B. I. Lundqvist, P. Hyldgaard, and T. Thonhauser, in *Non-covalent interactions in quantum chemistry and physics*, edited by A. O. de la Roza and G. A. DiLabio (Elsevier, 2017), pp. 241–274.
- <sup>62</sup>L. Goerigk, in *Non-covalent interactions in quantum chemistry and physics*, edited by A. O. de la Roza and G. A. DiLabio (Elsevier, 2017), pp. 195–219.
- <sup>63</sup>K. S. Thanthiriwatte, E. G. Hohenstein, L. A. Burns, and C. D. Sherrill, *J. Chem. Theory Comput.* **7**, 88–96 (2011).
- <sup>64</sup>M. S. Marshall, L. A. Burns, and C. D. Sherrill, *J. Chem. Phys.* **135**, 194102 (2011).
- <sup>65</sup>K. Patkowski, *WIREs Comput. Mol. Sci.* **119**, 123401–47 (2019).
- <sup>66</sup>R. M. Parrish, T. M. Parker, and C. D. Sherrill, *J. Chem. Theory Comput.* **10**, 4417–4431 (2014).
- <sup>67</sup>B. W. Bakr and C. D. Sherrill, *Phys. Chem. Chem. Phys.* **20**, 18241–18251 (2018).
- <sup>68</sup>B. W. Bakr and C. D. Sherrill, *Phys. Chem. Chem. Phys.* **18**, 10297–10308 (2016).

- <sup>69</sup>R. M. Parrish, D. F. Sitkoff, D. L. Cheney, and C. D. Sherrill, *Chem. Eur. J.* **23**, 7887–7890 (2017).
- <sup>70</sup>R. M. Parrish and C. D. Sherrill, *J. Chem. Phys.* **141**, 044115–22 (2014).
- <sup>71</sup>R. M. Parrish, J. F. Gonthier, C. Corminboeuf, and C. D. Sherrill, *J. Chem. Phys.* **143**, 051103–6 (2015).
- <sup>72</sup>J. F. Gonthier and C. Corminboeuf, *J. Chem. Phys.* **140**, 154107 (2014).
- <sup>73</sup>C. Sutton, M. S. Marshall, C. D. Sherrill, C. Risko, and J. L. Brédas, *J. Am. Chem. Soc.* **137**, 8775–8782 (2015).
- <sup>74</sup>D. A. Sirianni, L. A. Burns, and C. D. Sherrill, *J. Chem. Theory Comput.* **13**, 86–99 (2017).
- <sup>75</sup>L. A. Burns, Á. Vázquez-Mayagoitia, B. G. Sumpter, and C. D. Sherrill, *J. Chem. Phys.* **134**, 084107 (2011).
- <sup>76</sup>J. Řezáč, K. E. Riley, and P. Hobza, *J. Chem. Theory Comput.* **7**, 2427–2438 (2011).
- <sup>77</sup>L. A. Burns, M. S. Marshall, and C. D. Sherrill, *J. Chem. Phys.* **141**, 234111 (2014).
- <sup>78</sup>P. Hobza and J. Šponer, *J. Am. Chem. Soc.* **124**, 11802–11808 (2002).
- <sup>79</sup>P. Jurečka, J. Šponer, J. Černý, and P. Hobza, *Phys. Chem. Chem. Phys.* **8**, 1985–1993 (2006).
- <sup>80</sup>M. O. Sinnokrot and C. D. Sherrill, *J. Phys. Chem. A* **110**, 10656–10668 (2006).
- <sup>81</sup>A. D. Boese, J. M. L. Martin, and W. Klopper, *J. Phys. Chem. A* **111**, 11122–11133 (2007).
- <sup>82</sup>T. Janowski and P. Pulay, *Chem. Phys. Lett.* **447**, 27–32 (2007).
- <sup>83</sup>M. Pitoňák, T. Janowski, P. Neogrády, P. Pulay, and P. Hobza, *J. Chem. Theory Comput.* **5**, 1761–1766 (2009).
- <sup>84</sup>T. Takatani, E. G. Hohenstein, M. Malagoli, M. S. Marshall, and C. D. Sherrill, *J. Chem. Phys.* **132**, 144104 (2010).
- <sup>85</sup>J. C. Faver, M. L. Benson, X. He, B. P. Roberts, B. Wang, M. S. Marshall, M. R. Kennedy, C. D. Sherrill, and K. M. Merz Jr., *J. Chem. Theory Comput.* **7**, 790–797 (2011).

- <sup>86</sup>E. J. Carrell, C. M. Thorne, and G. S. Tschumper, *J. Chem. Phys.* **136**, 014103 (2012).
- <sup>87</sup>L. Kong, F. A. Bischoff, and E. F. Valeev, *Chem. Rev.* **112**, 75–107 (2012).
- <sup>88</sup>W. Klopper, F. R. Manby, S. Ten-no, and E. F. Valeev, *Int. Rev. Phys. Chem.* **25**, 427–468 (2006).
- <sup>89</sup>J. Noga, W. Klopper, and W. Kutzelnigg, in *Recent advances in coupled-cluster methods*, Vol. 3, edited by R. J. Bartlett, *Recent Advances in Computational Chemistry* (World Scientific, Singapore, 1997), pp. 1–48.
- <sup>90</sup>J. Noga, S. Kedzuch, J. Simunek, and S. Ten-no, *J. Chem. Phys.* **128**, 174103 (2008).
- <sup>91</sup>T. B. Adler, G. Knizia, and H.-J. Werner, *J. Chem. Phys.* **127**, 221106 (2007).
- <sup>92</sup>G. Knizia, T. B. Adler, and H.-J. Werner, *J. Chem. Phys.* **130**, 054104 (2009).
- <sup>93</sup>C. Hättig, D. P. Tew, and A. Kohn, *J. Chem. Phys.* **132**, 231102 (2010).
- <sup>94</sup>H.-J. Werner, G. Knizia, and F. R. Manby, *Mol. Phys.* **109**, 407–417 (2011).
- <sup>95</sup>O. Marchetti and H.-J. Werner, *J. Phys. Chem. A* **113**, 11580–11585 (2009).
- <sup>96</sup>M. S. Marshall and C. D. Sherrill, *J. Chem. Theory Comput.* **7**, 3978–3982 (2011).
- <sup>97</sup>K. D. Vogiatzis, Klopper, and W, *Mol. Phys.* **111**, 2299–2305 (2013).
- <sup>98</sup>A. D. Boese, *Mol. Phys.* **113**, 1618 (2015).
- <sup>99</sup>A. D. Boese, G. Jansen, M. Torheyden, S. Hofener, and W. Klopper, *Phys. Chem. Chem. Phys.* **13**, 1230–1238 (2011).
- <sup>100</sup>K. Patkowski, *J. Chem. Phys.* **138**, 154101 (2013).
- <sup>101</sup>B. Brauer, M. K. Kesharwani, and J. M. L. Martin, *J. Chem. Theory Comput.* **10**, 3791–3799 (2014).
- <sup>102</sup>G. Schmitz, C. Hättig, and D. P. Tew, *Phys. Chem. Chem. Phys.* **16**, 22167 (2014).
- <sup>103</sup>M. Torheyden and E. F. Valeev, *Phys. Chem. Chem. Phys.* **10**, 3410 (2008).
- <sup>104</sup>J. Řezáč, K. E. Riley, and P. Hobza, *J. Chem. Theory Comput.* **10**, 1359–1360 (2014).
- <sup>105</sup>J. Řezáč, K. E. Riley, and P. Hobza, *J. Chem. Theory Comput.* **7**, 3466–3470 (2011).

- <sup>106</sup>J. Řezáč and P. Hobza, *J. Chem. Theory Comput.* **8**, 141–151 (2012).
- <sup>107</sup>L. Gráfová, M. Pitoňák, J. Řezáč, and P. Hobza, *J. Chem. Theory Comput.* **6**, 2365–2376 (2010).
- <sup>108</sup>K. E. Yousaf and K. A. Peterson, *J. Chem. Phys.* **129**, 184108 (2008).
- <sup>109</sup>K. A. Peterson, M. K. Kesharwani, and J. M. L. Martin, *Mol. Phys.* **113**, 1551–1558 (2015).
- <sup>110</sup>H. Fliegl, W. Klopper, and C. Hättig, *J. Chem. Phys.* **122**, 084107 (2005).
- <sup>111</sup>MOLPRO, version 2010.1, a package of ab initio programs, H.-J. Werner, P. J. Knowles, F. R. Manby, M. Schütz, P. Celani, G. Knizia, T. Korona, R. Lindh, A. Mitrushenkov, G. Rauhut, T. B. Adler, R. D. Amos, A. Bernhardsson, A. Berning, D. L. Cooper, M. J. O. Deegan, A. J. Dobbyn, F. Eckert, E. Goll, C. Hampel, A. Hesselmann, G. Hetzer, T. Hrenar, G. Jansen, C. Köppl, Y. Liu, A. W. Lloyd, R. A. Mata, A. J. May, R. Tarroni, T. Thorsteinsson, M. Wang, and A. Wolf, see <http://www.molpro.net>.
- <sup>112</sup>E. F. Valeev, *Chem. Phys. Lett.* **395**, 190–195 (2004).
- <sup>113</sup>B. Brauer, M. K. Kesharwani, and J. M. L. Martin, *J. Chem. Theory Comput.* **10**, 3791–3799 (2014).
- <sup>114</sup>J. A. Pople, in *Energy, structure and reactivity: proceedings of the 1972 boulder summer research conference on theoretical chemistry*, edited by D. W. Smith and W. B. McRae (Wiley, New York, 1973), p. 51.
- <sup>115</sup>J. Řezáč, M. Dubecký, P. Jurečka, and P. Hobza, *Phys. Chem. Chem. Phys.* **17**, 19268 (2015).
- <sup>116</sup>K. A. Peterson, T. B. Adler, and H.-J. Werner, *J. Chem. Phys.* **128**, 084102 (2008).
- <sup>117</sup>M. S. Marshall, J. S. Sears, L. A. Burns, J. L. Brédas, and C. D. Sherrill, *J. Chem. Theory Comput.* **6**, 3681–3687 (2010).
- <sup>118</sup>D. G. A. Smith and K. Patkowski, *J. Chem. Theory Comput.* **9**, 370–389 (2013).
- <sup>119</sup>C. Sosa, J. Geersten, G. W. Trucks, R. J. Barlett, and J. A. Franz, *Chem. Phys. Lett.* **159**, 148–154 (1989).
- <sup>120</sup>W. Klopper, J. Noga, H. Koch, and T. Helgaker, *Theor. Chem. Acc.* **97**, 164–176 (1997).
- <sup>121</sup>A. G. Taube and R. J. Bartlett, *Collect. Czech. Chem. Commun.* **70**, 837–850 (2005).



- <sup>122</sup>A. Landau, K. Khistyayev, S. Dolgikh, and A. I. Krylov, *J. Chem. Phys.* **132**, 014109 (2010).
- <sup>123</sup>A. E. DePrince and C. D. Sherrill, *J. Chem. Theory Comput.* **9**, 293–299 (2013).
- <sup>124</sup>A. E. DePrince and C. D. Sherrill, *J. Chem. Theory Comput.* **9**, 2687–2696 (2013).
- <sup>125</sup>E. Epifanovsky, D. Zuev, X. Feng, K. Khistyayev, Y. Shao, and A. I. Krylov, *J. Chem. Phys.* **139**, 134105 (2013).
- <sup>126</sup>A. E. DePrince, M. R. Kennedy, B. G. Sumpter, and C. D. Sherrill, *Mol. Phys.* **112**, 844–852 (2014).
- <sup>127</sup>J. M. Turney, A. C. Simmonett, R. M. Parrish, E. G. Hohenstein, F. A. Evangelista, J. T. Fermann, B. J. Mintz, L. A. Burns, J. J. Wilke, M. L. Abrams, N. J. Russ, M. L. Leininger, C. L. Janssen, E. T. Seidl, W. D. Allen, H. F. Schaefer, R. A. King, E. F. Valeev, C. D. Sherrill, and T. D. Crawford, *WIREs Comput. Mol. Sci.* **2**, 556–565 (2012).
- <sup>128</sup>D. A. Sirianni, A. Alenaizan, D. L. Cheney, and C. D. Sherrill, *J. Chem. Theory Comput.* **14**, 3004–3013 (2018).
- <sup>129</sup>R. Z. Khaliullin, M. Head-Gordon, and A. T. Bell, *J. Chem. Phys.* **124**, 204105 (2006).
- <sup>130</sup>R. Z. Khaliullin, E. A. Cobar, R. C. Lochan, A. T. Bell, and M. Head-Gordon, *J. Phys. Chem. A* **111**, 8753–8765 (2007).
- <sup>131</sup>R. Z. Khaliullin, A. T. Bell, and M. Head-Gordon, *J. Chem. Phys.* **128**, 184112 (2008).
- <sup>132</sup>P. R. Horn, E. J. Sundstrom, T. A. Baker, and M. Head-Gordon, *J. Chem. Phys.* **138**, 134119 (2013).
- <sup>133</sup>B. Jeziorski, R. Moszynski, and K. Szalewicz, *Chem. Rev.* **94**, 1887–1930 (1994).
- <sup>134</sup>K. Szalewicz, *WIREs Comput. Mol. Sci.* **2**, 254–272 (2012).
- <sup>135</sup>E. G. Hohenstein and C. D. Sherrill, *WIREs Comput. Mol. Sci.* **2**, 304–326 (2012).
- <sup>136</sup>C. D. Sherrill, *Acc. Chem. Res.* **46**, 1020–1028 (2013).
- <sup>137</sup>R. M. Parrish and C. D. Sherrill, *J. Chem. Phys.* **141**, 044115 (2014).
- <sup>138</sup>R. M. Parrish, T. M. Parker, and C. D. Sherrill, *J. Chem. Theory Comput.* **10**, 4417–4431 (2014).

- <sup>139</sup>R. M. Parrish, D. F. Sitkoff, D. L. Cheney, and C. D. Sherrill, *Chem. Eur. J.* **23**, 7887–7890 (2017).
- <sup>140</sup>C. D. Sherrill, T. Takatani, and E. G. Hohenstein, *J. Phys. Chem. A* **113**, 10146–10159 (2009).
- <sup>141</sup>P. Jurečka, J. Černý, P. Hobza, and D. R. Salahub, *J. Comput. Chem.* **28**, 555–569 (2007).
- <sup>142</sup>H. Valdés, V. Klusák, M. Pitoňák, O. Exner, I. Starý, P. Hobza, and L. Rulíšek, *J. Comput. Chem.* **29**, 861–870 (2008).
- <sup>143</sup>L. Demovicova, P. Hobza, and J. Řezáč, *Phys. Chem. Chem. Phys.* **16**, 19115–19121 (2014).
- <sup>144</sup>J. Řezáč, K. E. Riley, and P. Hobza, *J. Chem. Theory Comput.* **8**, 4285 (2012).
- <sup>145</sup>K. S. Thanthiriwatte, E. G. Hohenstein, L. A. Burns, and C. D. Sherrill, *J. Chem. Theory Comput.* **7**, 88–96 (2011).
- <sup>146</sup>S. Grimme, *J. Comput. Chem.* **27**, 1787–1799 (2006).
- <sup>147</sup>A. D. Becke, *J. Chem. Phys.* **98**, 5648–5652 (1993).
- <sup>148</sup>P. J. Stephens, F. J. Devlin, C. F. Chabalowski, and M. J. Frisch, *J. Phys. Chem.* **98**, 11623–11627 (1994).
- <sup>149</sup>Y. Zhao, N. E. Schultz, and D. G. Truhlar, *J. Chem. Theory Comput.* **2**, 364–382 (2006).
- <sup>150</sup>A. Ambrosetti, N. Ferri, R. A. DiStasio Jr., and A. Tkatchenko, *Science* **351**, 1171–1176 (2016).
- <sup>151</sup>A. Ambrosetti, A. M. Reilly, R. A. DiStasio Jr., and A. Tkatchenko, *J. Chem. Phys.* **140**, 18A508 (2014).
- <sup>152</sup>A. Ambrosetti, D. Alfè, R. A. DiStasio Jr., and A. Tkatchenko, *J. Phys. Chem. Lett.* **5**, 849–855 (2014).
- <sup>153</sup>A. J. Misquitta, J. Spencer, A. J. Stone, and A. Alavi, *Phys. Rev. B.* **82**, 075312 (2010).
- <sup>154</sup>J. F. Dobson, A. White, and A. Rubio, *Phys. Rev. Lett.* **96**, 073201 (2006).
- <sup>155</sup>A. D. Becke and E. R. Johnson, *J. Chem. Phys.* **122**, 154104 (2005).
- <sup>156</sup>E. R. Johnson and A. D. Becke, *J. Chem. Phys.* **124**, 174104 (2006).

- <sup>157</sup>J. Kong, Z. T. Gan, E. Proynov, M. Freindorf, and T. R. Furlani, *Phys. Rev. A* **79**, 042510 (2009).
- <sup>158</sup>M. Dion, H. Rydberg, E. Schröder, D. C. Langreth, and B. I. Lundqvist, *Phys. Rev. Lett.* **92**, 246401 (2004).
- <sup>159</sup>D. C. Langreth, M. Dion, H. Rydberg, E. Schroder, P. Hyldgaard, and B. I. Lundqvist, *Int. J. Quantum Chem.* **101**, 599–610 (2005).
- <sup>160</sup>K. Lee, É. Murray, L. Kong, B. I. Lundqvist, and D. C. Langreth, *Phys. Rev. B* **82**, 081101(R) (2010).
- <sup>161</sup>O. A. Vydrov and T. V. Voorhis, *J. Chem. Phys.* **133**, 244103 (2010).
- <sup>162</sup>Y. Zhao and D. G. Truhlar, *J. Chem. Theory Comput.* **3**, 289–300 (2007).
- <sup>163</sup>E. G. Hohenstein, S. T. Chill, and C. D. Sherrill, *J. Chem. Theory Comput.* **4**, 1996–2000 (2008).
- <sup>164</sup>R. Sedlak, T. Janowski, M. Pitoňák, J. Řezáč, P. Pulay, and P. Hobza, *J. Chem. Theory Comput.* **9**, 3364–3374 (2013).
- <sup>165</sup>J. P. Perdew and K Schmidt, *AIP Conference Proceedings* **577**, 1 (2001).
- <sup>166</sup>S. Grimme, A. Hansen, J. G. Brandenburg, and C. Bannwarth, *Chem. Rev.* **116**, 5105–5154 (2016).
- <sup>167</sup>Y. Shao, Z. Gan, E. Epifanovsky, A. T. B. Gilbert, M. Wormit, J. Kussmann, A. W. Lange, A. Behn, J. Deng, X. Feng, D. Ghosh, M. Goldey, P. R. Horn, L. D. Jacobson, I. Kaliman, R. Z. Khaliullin, T. Kus, A. Landau, J. Liu, E. I. Proynov, Y. M. Rhee, R. M. Richard, M. A. Rohrdanz, R. P. Steele, E. J. Sundstrom, H. L. Woodcock, P. M. Zimmerman, D. Zuev, B. Albrecht, E. Alguire, B. Austin, G. J. O. Beran, Y. A. Bernard, E. Berquist, K. Brandhorst, K. B. Bravaya, S. T. Brown, D. Casanova, C. Chang, Y. Chen, S. H. Chien, K. D. Closser, D. L. Crittenden, M. Diedenhofen, R. A. DiStasio, H. Do, A. D. Dutoi, R. G. Edgar, S. Fatehi, L. Fusti-Molnar, A. Ghysels, A. Golubeva-Zadorozhnaya, J. Gomes, M. W. D. Hanson-Heine, P. H. P. Harbach, A. W. Hauser, E. G. Hohenstein, Z. C. Holden, T. Jagau, H. Ji, B. Kaduk, K. Khistyayev, J. Kim, J. Kim, R. A. King, P. Klunzinger, D. Kosenkov, T. Kowalczyk, C. M. Krauter, K. U. Lao, A. D. Laurent, K. V. Lawler, S. V. Levchenko, C. Y. Lin, F. Liu, E. Livshits, R. C. Lochan, A. Luenser, P. Manohar, S. F. Manzer, S. Mao, N. Mardirossian, A. V. Marenich, S. A. Maurer, N. J. Mayhall, E. Neuscamman, C. M. Oana, R. Olivares-Amaya, D. P. O’Neill, J. A. Parkhill, T. M. Perrine, R. Peverati, A. Prociuk, D. R. Rehn, E. Rosta, N. J. Russ, S. M. Sharada, S. Sharma, D. W. Small, A. Sodt, T. Stein, D. Stueck, Y. Su, A. J. W. Thom, T. Tsuchimochi, V. Vanovschi, L. Vogt, O. Vydrov, T. Wang, M. A. Watson, J. Wenzel, A. White, C. F. Williams, J. Yang, S. Yeganeh, S. R. Yost, Z. You, I. Y. Zhang, X. Zhang,

- Y. Zhao, B. R. Brooks, G. K. L. Chan, D. M. Chipman, C. J. Cramer, W. A. Goddard, M. S. Gordon, W. J. Hehre, A. Klamt, H. F. Schaefer, M. W. Schmidt, C. D. Sherrill, D. G. Truhlar, A. Warshel, X. Xu, A. Aspuru-Guzik, R. Baer, A. T. Bell, N. A. Besley, J. Chai, A. Dreuw, B. D. Dunietz, T. R. Furlani, S. R. Gwaltney, C. Hsu, Y. Jung, J. Kong, D. S. Lambrecht, W. Liang, C. Ochsenfeld, V. A. Rassolov, L. V. Slipchenko, J. E. Subotnik, T. V. Voorhis, J. M. Herbert, A. I. Krylov, P. M. W. Gill, and M. Head-Gordon, *Mol. Phys.* **113**, 184–215 (2015).
- <sup>168</sup>E. R. Johnson, A. D. Becke, C. D. Sherrill, and G. A. DiLabio, *J. Chem. Phys.* **131**, 034111 (2009).
- <sup>169</sup>J. L. Whitten, *J. Chem. Phys.* **58**, 4496–4501 (1973).
- <sup>170</sup>B. I. Dunlap, J. W. D. Connolly, and J. R. Sabin, *Int. J. Quantum Chem. Symp.* **11**, 81 (1977).
- <sup>171</sup>B. I. Dunlap, J. W. D. Connolly, and J. R. Sabin, *J. Chem. Phys.* **71**, 3396–3402 (1979).
- <sup>172</sup>M. Feyereisen, G. Fitzgerald, and A. Komornicki, *Chem. Phys. Lett.* **208**, 359–363 (1993).
- <sup>173</sup>O. Vahtras, J. Almlöf, and M. W. Feyereisen, *Chem. Phys. Lett.* **213**, 514–518 (1993).
- <sup>174</sup>A. P. Rendell and T. J. Lee, *J. Chem. Phys.* **101**, 400–408 (1994).
- <sup>175</sup>R. A. Kendall and H. A. Fruchtl, *Theor. Chem. Acc.* **97**, 158–163 (1997).
- <sup>176</sup>F. Weigend, *Phys. Chem. Chem. Phys.* **4**, 4285–4291 (2002).
- <sup>177</sup>J. Řezáč and P. Hobza, *J. Chem. Theory Comput.* **9**, 2151–2155 (2013).
- <sup>178</sup>N. J. Singh, S. K. Min, D. Y. Kim, and K. S. Kim, *J. Chem. Theory Comput.* **5**, 515–529 (2009).
- <sup>179</sup>S. van der Walt, S. C. Colbert, and G. Varoquaux, *Computing in Science and Engineering* **13**, 22–30 (2011).
- <sup>180</sup>D. A. Sirianni, D. G. A. Smith, L. A. Burns, D. F. Sitkoff, D. L. Cheney, and C. D. Sherrill, “Optimized Damping Parameters for Empirical Dispersion Corrections to Symmetry-Adapted Perturbation Theory,” *in preparation*.
- <sup>181</sup>E. G. Hohenstein and C. D. Sherrill, *J. Chem. Phys.* **133**, 014101 (2010).
- <sup>182</sup>H. L. Williams, K. Szalewicz, R. Moszynski, and B. Jeziorski, *J. Chem. Phys.* **103**, 4586–4599 (1995).

- <sup>183</sup>T. Korona, *Mol. Phys.* **111**, 3705–3715 (2013).
- <sup>184</sup>A. Heßelmann, G. Jansen, and M. Schütz, *J. Chem. Phys.* **122**, 014103 (2005).
- <sup>185</sup>A. J. Misquitta, R. Podeszwa, B. Jeziorski, and K. Szalewicz, *J. Chem. Phys.* **123**, 214103 (2005).
- <sup>186</sup>E. G. Hohenstein and C. D. Sherrill, *J. Chem. Phys.* **132**, 184111 (2010).
- <sup>187</sup>E. G. Hohenstein, R. M. Parrish, C. D. Sherrill, J. M. Turney, and H. F. Schaefer, *J. Chem. Phys.* **135**, 174107 (2011).
- <sup>188</sup>R. M. Parrish, K. C. Thompson, and T. J. Martínez, *J. Chem. Theory Comput.* **14**, 1737–1753 (2018).
- <sup>189</sup>R. Podeszwa, K. Pernal, K. Patkowski, and K. Szalewicz, *J. Phys. Chem. Lett.* **1**, 550–555 (2010).
- <sup>190</sup>A. Hesselmann, *J. Phys. Chem. A* **115**, 11321–11330 (2011).
- <sup>191</sup>K. U. Lao and J. M. Herbert, *J. Phys. Chem. Lett.* **3**, 3241–3248 (2012).
- <sup>192</sup>J. Hepburn, G. Scoles, and R. Penco, *Chem. Phys. Lett.* **36**, 451–456 (1975).
- <sup>193</sup>R. Ahlrichs, R. Penco, and G. Scoles, *Chem. Phys.* **19**, 119–130 (1977).
- <sup>194</sup>C. Douketis, G. Scholes, S. Marchetti, M. Zen, and A. J. Thakkar, *J. Chem. Phys.* **76**, 3057–3063 (1982).
- <sup>195</sup>*DFTD3*, A dispersion correction for density functionals, Hartree–Fock, and semi-empirical quantum chemical methods, version 3.2 Rev. 0; Grimme Research Group: Mulliken Center for Theoretical Chemistry, Universität Bonn, 2016. <https://www.chemie.uni-bonn.de/pctc/mulliken-center/software/dft-d3/> (accessed August 23, 2019).
- <sup>196</sup>F.-Y. Lin, C.-I. Liu, Y.-L. Liu, Y. Zhang, K. Wang, W.-Y. Jeng, T.-P. Ko, R. Cau, A. H. J. Wang, and E. Oldfield, *Proc. Natl. Acad. Sci.* **107**, 21337–21342 (2010).
- <sup>197</sup>T. Warne, P. C. Edwards, A. S. Doré, A. G. W. Leslie, and C. G. Tate, *Science* **364**, 775–778 (2019).
- <sup>198</sup>T. M. Parker and C. D. Sherrill, *J. Chem. Theory Comput.* **11**, 4197–4202 (2015).
- <sup>199</sup>L. A. Burns, J. C. Faver, Z. Zheng, M. S. Marshall, D. G. A. Smith, K. Vanommeslaeghe, A. D. MacKerell, K. M. Merz, and C. D. Sherrill, *J. Chem. Phys.* **147**, 161727 (2017).

- <sup>200</sup>R. Bukowski, K. Szalewicz, G. C. Groenenboom, and A. van der Avoird, *Science* **315**, 1249–1252 (2007).
- <sup>201</sup>R. Bukowski, K. Szalewicz, G. C. Groenenboom, and A. van der Avoird, *J. Chem. Phys.* **128**, 094313 (2008).
- <sup>202</sup>D. G. A. Smith and K. Patkowski, unpublished.
- <sup>203</sup>S. Li and K. Patkowski, unpublished.
- <sup>204</sup>D. G. A. Smith and K. Patkowski, *J. Phys. Chem. C* **118**, 544–550 (2014).
- <sup>205</sup>S. Li, D. G. A. Smith, and K. Patkowski, *Phys. Chem. Chem. Phys.* **17**, 16560–16574 (2015).
- <sup>206</sup>D. G. A. Smith and K. Patkowski, *J. Phys. Chem. C* **119**, 4934 (2015).
- <sup>207</sup>Maestro, Schrödinger, LLC, New York, NY, 2019.
- <sup>208</sup>G Madhavi Sastry, M. Adzhigirey, T. Day, R. Annabhimoju, and W. Sherman, *J Comput Aided Mol Des* **27**, 221–234 (2013).
- <sup>209</sup>T. Warne, R. Moukhametzianov, J. G. Baker, R. Nehmé, P. C. Edwards, A. G. W. Leslie, G. F. X. Schertler, and C. G. Tate, *Nature* **469**, 241–244 (2011).
- <sup>210</sup>D. G. A. Smith, L. A. Burns, K. Patkowski, and C. D. Sherrill, *J. Phys. Chem. Lett.* **7**, 2197–2203 (2016).
- <sup>211</sup>H. W. Qi and H. J. Kulik, *Journal of Chemical Information and Modeling* **59**, 2199–2211 (2019).
- <sup>212</sup>R. M. Parrish, K. C. Thompson, and T. J. Martínez, *J. Chem. Theory Comput.* **14**, 1737–1753 (2018).
- <sup>213</sup>R. M. Parrish, J. F. Gonthier, C. Corminboeuf, and C. D. Sherrill, *J. Chem. Phys.* **143**, 051103–6 (2015).
- <sup>214</sup>J. C. Phillips, R. Braun, W. Wang, J. Gumbart, E. Tajkhorshid, E. Villa, C. Chipot, R. D. Skeel, L. Kalé, and K. Schulten, *J. Comput. Chem.* **26**, 1781–1802 (2005).
- <sup>215</sup>W. L. Jorgensen, J. Chandrasekhar, J. D. Madura, R. W. Impey, and M. L. Klein, *J. Chem. Phys.* **79**, 926–935 (1983).
- <sup>216</sup>A. Jakalian, D. B. Jack, and C. I. Bayly, *J. Comput. Chem.* **23**, 1623–1641 (2002).

- <sup>217</sup>D. A. Case, T. E. Cheatham, T. Darden, H. Gohlke, R. Luo, K. M. Merz, A. Onufriev, C. Simmerling, B. Wang, and R. J. Woods, *J. Comput. Chem.* **26**, 1668–1688 (2005).
- <sup>218</sup>J. Wang, W. Wang, P. A. Kollman, and D. A. Case, *Journal of Molecular Graphics and Modelling* **25**, 247–260 (2006).
- <sup>219</sup>J. Wang, R. M. Wolf, J. W. Caldwell, P. A. Kollman, and D. A. Case, *J. Comput. Chem.* **25**, 1157–1174 (2004).
- <sup>220</sup>P. C. D. Hawkins, A. G. Skillman, and A. Nicholls, *J. Med. Chem.* **50**, 74–82 (2007).
- <sup>221</sup>For further details regarding the Tanimoto combo scoring, see <https://docs.eyesopen.com/rocs/usage.htm>
- <sup>222</sup>We refer the reader to the documentation in the PSI4 manual for the FISAPT module at <http://psicode.org/psi4manual/1.2/fisapt.html> for specific details regarding the post-processing required for F-SAPT analysis.
- <sup>223</sup>R. M. Parrish, J. F. Gonthier, C. Corminboeuf, and C. D. Sherrill, *J. Chem. Phys.* **143**, 051103 (2015).
- <sup>224</sup>E. Papajak, J. Zheng, X. Xu, H. R. Leverentz, and D. G. Truhlar, *J. Chem. Theory Comput.* **7**, 3027–3034 (2011).
- <sup>225</sup>L. A. Burns, M. S. Marshall, and C. D. Sherrill, *J. Chem. Theory Comput.* **10**, 49–57 (2014).
- <sup>226</sup>L. A. Burns, Á. V. Mayagoitia, B. G. Sumpter, and C. D. Sherrill, *J. Chem. Phys.* **134**, 084107 (2011).
- <sup>227</sup>R. A. Kendall, T. H. Dunning, and R. J. Harrison, *J. Chem. Phys.* **96**, 6796–6806 (1992).
- <sup>228</sup>A. L. Ringer, M. O. Sinnokrot, R. P. Lively, and C. D. Sherrill, *Chem. Eur. J.* **12**, 3821–3828 (2006).
- <sup>229</sup>D. J. Tannor, B. Marten, R. Murphy, R. A. Friesner, D. Sitkoff, A. Nicholls, B. Honig, M. Ringnalda, and W. A. Goddard, *J. Am. Chem. Soc.* **116**, 11875–11882 (1994).

## VITA

Dominic A. (Dom) Sirianni was born May 6, 1994 to James (Jim) and Cathy Sirianni, who are both secondary school teachers, and grew up in the small rural town of Kane, Pennsylvania. Dominic attended first through fourth grades at the Holy Rosary School in Johnsonburg, PA, before moving to Kane public school district in fifth grade after the closure of Holy Rosary. Dominic graduated from Kane Area High School in 2011, having been taught by both his parents (drawing, ceramics, and painting by Cathy, and world cultures by Jim). After graduation, Dominic attended Edinboro University of Pennsylvania with the intent of following in his parents' footsteps, to become a high school chemistry teacher. After discovering that he loved learning as much as he did teaching, Dominic switched his major to pursue simultaneous Bachelors of Science degrees in Chemistry and Mathematics, while also serving as a student peer tutor and teaching assistant in the Departments of Chemistry and Mathematics & Computer Science. Dominic graduated summa cum laude from EUP in 2015, after which he attended the Georgia Institute of Technology in to pursue his PhD in Chemistry in the research group of Prof. C. David Sherrill.

### Publications

5. "PSI4 1.4: Open Source Software for High-Throughput Quantum Chemistry," D. G. A. Smith, L. A. Burns, A. C. Simmonett, R. M. Parrish, M. C. Schieber, R. Galvelis, P. Kraus, H. Kruse, R. Di Remigio, A. Alenaizan, A. M. James, S. Lehtola, J. P. Misiewicz, M. Scheurer, R. A. Shaw, J. B. Schriber, Y. Xie, Z. L. Glick, D. A. Sirianni, J. S. O'Brien, J. M. Waldrop, A. Kumar, E. G. Hohenstein, B. P. Pritchard, B. R. Brooks, H. F. Schaefer III, A. Yu. Sokolov, K. Patkowski, A. E. DePrince III, U. Bozkaya, R. A. King, F. A. Evangelista, J. M. Turney, T. D. Crawford, and C. D. Sherrill, *J. Chem. Phys.* (Submitted)



4. "Tipping the Balance between S- $\pi$  and O- $\pi$  Interactions," J. Whang, P. Li, M. D. Smith, C. E. Warden, D. A. Sirianni, E. C. Vik, J. M. Maier, C. J. Yehl, C. D. Sherrill, and K. D. Shimizu, *J. Am. Chem. Soc.* **140**, 13301-13307 (2018) (doi: 10.1021/jacs.8b07617)
3. "PSI4NUMPY: An Interactive Quantum Chemistry Programming Environment for Reference Implementations and Rapid Development," D. G. A. Smith, L. A. Burns, D. A. Sirianni, D. R. Nascimento, A. Kumar, A. M. James, J. B. Schriber, T. Zhang, B. Zhang, A. S. Abbott, E. Berquist, M. H. Lechner, L. dos A. Cunha, A. G. Heide, R. A. King, A. C. Simmonett, J. M. Turney, H. F. Schaefer, F. A. Evangelista, A. E. DePrince III, T. D. Crawford, K. Patkowski, and C. D. Sherrill, *J. Chem. Theory Comput.* **14**, 3504-3511 (2018) (doi: 10.1021/acs.jctc.8b00286)
2. "Assessment of Density Functionals for Optimization of Bimolecular van der Waals Complexes," D. A. Sirianni, A. Alenaizan, D. L. Cheney, and C. D. Sherrill, *J. Chem. Theory Comput.* **14**, 3004-3013 (2018) (doi: 10.1021/acs.jctc.8b00114)
1. "Comparison of Explicitly Correlated Methods for Computing High-Accuracy Benchmark Energies for Noncovalent Interactions," D. A. Sirianni, L. A. Burns, and C. D. Sherrill, *J. Chem. Theory Comput.* **13**, 86-99 (2017) (doi: 10.1021/acs.jctc.6b00797)