

**Benchmarking Ab Initio Computational Methods for the  
Quantitative Prediction of Sunlight-Driven Pollutant  
Degradation in Aquatic Environments**

*An Honors Thesis for the Department of Chemistry*

By Kasidet (Jing) Trerayapiwat

*Bowdoin College, 2016*

©2016 Kasidet Trerayapiwat

## Abstract

Understanding the changes in molecular electronic structure following the absorption of light is a fundamental challenge for the goal of predicting photochemical rates and mechanisms. Proposed here is a systematic benchmarking method to evaluate accuracy of a model to quantitatively predict photo-degradation of small organic molecules in aquatic environments. An overview of underlying computational theories relevant to understanding sunlight-driven electronic processes in organic pollutants is presented. To evaluate the optimum size of solvent sphere, molecular Dynamics and Time Dependent Density Functional Theory (MD-TD-DFT) calculations of an aniline molecule in different numbers of water molecules using CAM-B3LYP functional yielded excited state energy and oscillator strength values, which were compared with data from experimental absorption spectra. For the first time, a statistical method of comparing experimental and theoretical data is proposed. Underlying Gaussian functions of absorption spectra were deconvoluted and integrated to calculate experimental oscillator strengths. A Matlab code written by Soren Eustis was utilized to decluster MD-TD-DFT results. The model with 256 water molecules was decided to give the most accurate results with optimized computational cost and accuracy. MD-TD-DFT calculations were then performed on aniline, 3-F-aniline, 4-F-aniline, 3-Cl-aniline, 4-MeOacetophenone, and (1,3)-dimethoxybenzophenone with CAM-B3LYP, PBE0, M06-2X, LCBLYP, and BP86 functionals. BP86 functional was determined to be the best functional.

# Contents

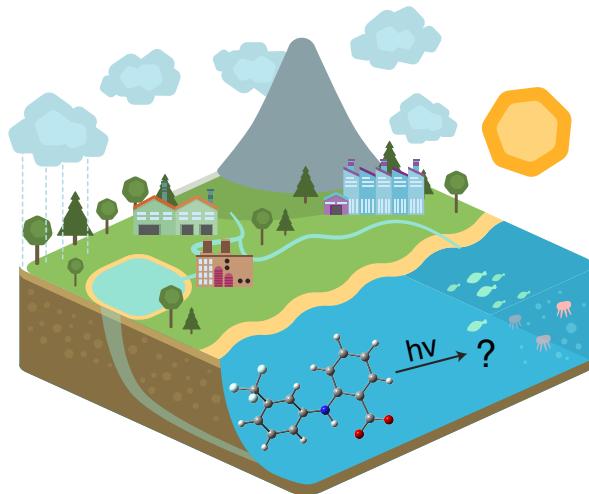
<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Environmental Photochemistry And Micropollutants . . . . .	1
1.2	Photo-excitation and UV-VIS Spectrum . . . . .	2
1.3	Computational Approach . . . . .	5
<b>2</b>	<b>Background</b>	<b>8</b>
2.1	Hartree-Fock Theory . . . . .	8
2.2	Basis Sets . . . . .	9
2.3	Density Functional Theory . . . . .	12
2.4	Solvent Models . . . . .	14
2.5	Molecular Dynamics . . . . .	16
2.6	Comparing Experimental And Theoretical Calculation Results . . . . .	16
<b>3</b>	<b>Methods</b>	<b>17</b>
3.1	Overview . . . . .	17
3.2	Computational . . . . .	17
3.2.1	Facilitating Each Steps Using Python Scripts . . . . .	17
3.2.2	Preparing Starting Geometry Using Packmol . . . . .	19
3.2.3	Computational Models . . . . .	19
3.3	Experimental . . . . .	20
<b>4</b>	<b>Results And Discussion</b>	<b>22</b>
4.1	Results from MD-TD-DFT . . . . .	22
4.1.1	Solvent Boundary Potential . . . . .	24
4.1.2	Aniline With 32 Water Molecules . . . . .	26
4.2	Determining The Optimum Solvent Environment . . . . .	27
4.3	Comparing Different Functionals . . . . .	29
<b>5</b>	<b>Conclusion</b>	<b>33</b>
<b>Acknowledgment</b>		<b>34</b>
<b>References</b>		<b>34</b>

<b>A Appendix</b>	<b>42</b>
A.1 MD-TD-DFT Result Tables . . . . .	42
A.2 Python Scripts . . . . .	44
A.2.1 Preparing MD Input Files . . . . .	44
A.2.2 MD Geometries Extraction . . . . .	47
A.2.3 Plot Potential Energy Of MD Run . . . . .	52
A.2.4 Find The Most Equilibrated Period . . . . .	56
A.2.5 Prepare TD-DFT Input . . . . .	58
A.2.6 Pull Excited State Energies From TD-DFT Log Files . . . . .	62
A.3 GAMESS Inputs . . . . .	64
A.3.1 MD Input File . . . . .	64
A.3.2 MD Restart . . . . .	66
A.3.3 TD-DFT Input File . . . . .	67

# 1 Introduction

## 1.1 Environmental Photochemistry And Micropollutants

During the past century, more and more synthetic chemicals have been created for commercial pharmaceuticals and personal care products. Each year, about 300 million tons of organic chemicals are being added into water systems from household septic tanks, waste water facilities, and many other paths.<sup>1</sup> These chemicals have been detected in ng/L to g/L quantity in natural water bodies: rivers, lakes, and the oceans.<sup>2</sup> The amount of chemicals being introduced into the water raises the need to understand the consequences of these compounds. If the water is polluted, eventually animals, humans, and other living beings in this world are affected. For example, in 1956, Minamata disease was identified to be caused by consumption of fish bio-accumulated by methylmercury from a nearby chemical factory.<sup>3</sup> While some major toxic chemicals, such as mercury,  $\text{CCl}_4$ , DDT, etc., have been constantly regulated or banned by governments for their effects on agriculture and the environment, many more chemical wastes have not been studied and regulated to the same degree. Even though some chemicals are approved as safe for daily usage in household products, their impacts after being disposed into water remain largely unknown. These low-concentration chemicals are collectively called micropollutants.<sup>1</sup> The small concentration, large number of different species in water, and difficulty of separating or concentrating them for analysis all retard the effort to understand their effects on in water systems. These are “The Challenge of Micropollutants in Aquatic Systems”.<sup>1</sup> Since each individual species of micropollutant is not classified as harmful, major concerns of introducing them to water bodies are their reactions in water. In water with exposure to sunlight, photoreaction can occur as photons from the sun energize organic pollutants from ground state to excited state. In rivers close to discharge areas, photolysis has been confirmed to have direct connection to the observed attenuation in concentration of micropollutants from upstream to downstream.<sup>5-8</sup> Observable decrease in concentrations means a significant portion of the pollutants undergo photo-reactions in exposure to sunlight. Eustis et al. carried out computational studies of triclosan in excited states.<sup>9</sup> Triclosan goes through photolysis under sunlight to form Dioxins and

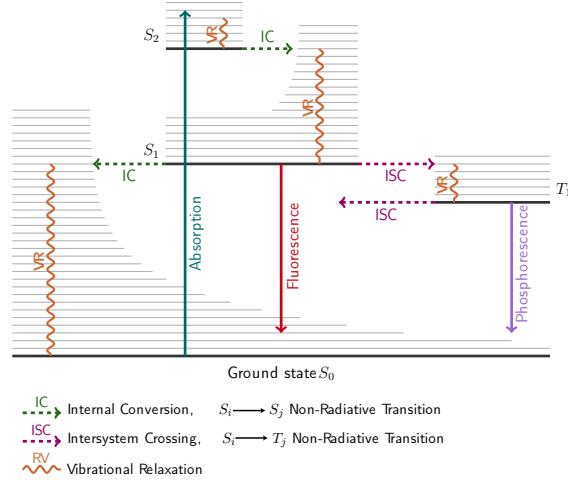


**Figure 1:** Environmental fate of micropollutants in water. Under sunlight, molecules can undergo unknown photo-reactions to form unknown products.<sup>4</sup>

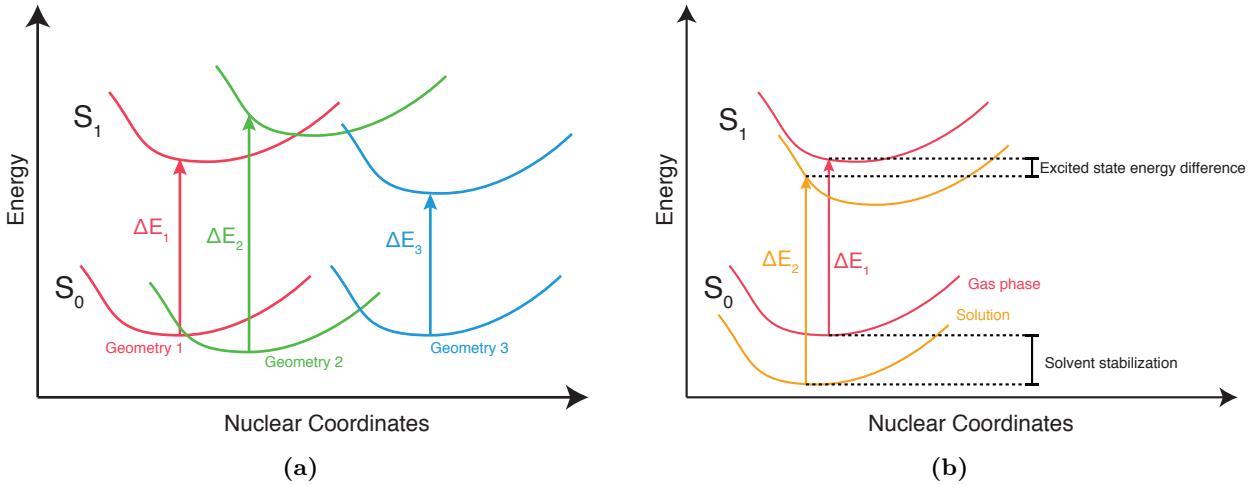
PCBs, well-known carcinogens.<sup>10</sup> This is but one of many examples of large-scale, unintended consequences as a direct result of persistent, low-level concentrations of pollutants in aquatic ecosystems<sup>11</sup>. One important step in triclosan's degradation process to its 4 major products is detachment of ortho-Cl from triclosan ring under 300 nm light. Confirmed by combining both experimental and computational results in the paper, ortho-chlorine photo-dissociates from triclosan through reductive cleavage pathway through biradical anion. Excited state bond energy of the breaking C-Cl reveals significantly smaller activation energy in reductive cleavage pathway compared to that of through homolytic mechanism. In this case, experimental method provides evidence of products from photo-reaction. Computational method provides reasons and insights into the microscopic world of atoms and electrons. Combing and analyzing results the two methods give a better understanding of photo-reactions in aquatic solutions.

## 1.2 Photo-excitation and UV-VIS Spectrum

Understanding each step in the photochemically driven reaction is critical to properly evaluate the impact of a pollutant's environmental fate. The first, and most important step in this process is light absorption of ground state molecules to form chemically activated excited state molecules. A molecule in its ground state ( $S_0$ ) can absorb a photon to its first singlet excited state ( $S_1$ ), and then relax to its first triplet excited state ( $T_1$ ) as illustrated in Figure 2. Experimentally, studying absorption of photons in chemical compounds is done through UV-VIS absorption spectroscopy, followed by analysis of molar absorptivity and wavelength. In contrast to common misconception, a value of molar absorptivity is not a characteristic property of a compound. Molar absorptivity as a function of wavelength is highly dependent on the resolution used in the measurement.<sup>13</sup> Solvent-solute interactions also play a role in changing the excited state energies.<sup>14-21</sup> One particular study done by Eilmes suggests that the effects solvent molecules have on solute can be separated into two main parts: solvent absorption peak shift and solvent broadening.<sup>22</sup> Two sets of Molecular Dynamics and Time Dependent Density Functional Theory (MD-TD-DFT) calculations were done controlling either solute or solvent molecules. The results were analyzed to study how either solvent or solute geometry variation af-



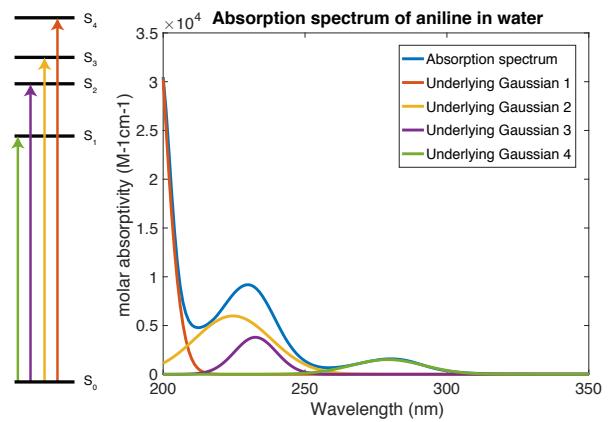
**Figure 2:** Jablonski diagram.<sup>12</sup> A possible pathway of photo-excitation starts at ground state to singlet excited state, then triplet excited states.



**Figure 3:** Illustration of solvent effects on potential energy surfaces of a solute. a) different ground state geometries give different potential energy surfaces for ground and excited states. Excited states energy calculation then gives different results due to this solute geometric fluctuation. The result is solvent peak broadening around the mean. b) Solvent stabilization of the ground state. Ground state geometry has lower potential energy surface as a result from dipole interaction from solvent molecules. Generally, the excited state geometries are not as well equilibrated by the ground state solute and solvent geometries.

fected the first excited state energy. First, the solute molecule were not static in the solution and geometry fluctuation creates a broad range of ground state Potential Energy Surfaces (PES) (see Figure 3a). Different starting ground states yield different predictions for the solute's excited states. This contributes to absorption peak broadening compared with an absorption in gas phase. In addition to solvent broadening, temperature broadening due to Boltzmann distribution of molecular energies around the global minimum energy as described in Molecular Dynamics calculation also contribute to absorption peak broadening.

Secondly, dipole moments from the solvent stabilize the ground state geometry lowering the potential energy surface. Since electrons move significantly faster than the nucleus, excited state geometries are less likely to be stabilized. Unequal lowering of the two potential energy surfaces causes blueshift in excited state energy prediction (see Figure 3b). In gas phase, simulation of UV-VIS spectrum can be done more easily due to a lack of solvent shift and solvent broadening. In solution, sharp and high-intensity gas-phase peaks become broadened and blueshifted. Each individual absorption



**Figure 4:** Absorption spectrum of aniline in water showing underlying electronic transitions. Underlying Gaussian components, fitted manually using Fityk, refer to 4 electronic transitions as shown in electronic transition diagram on the left. Each peak can be used to calculate oscillator strength using equation 16

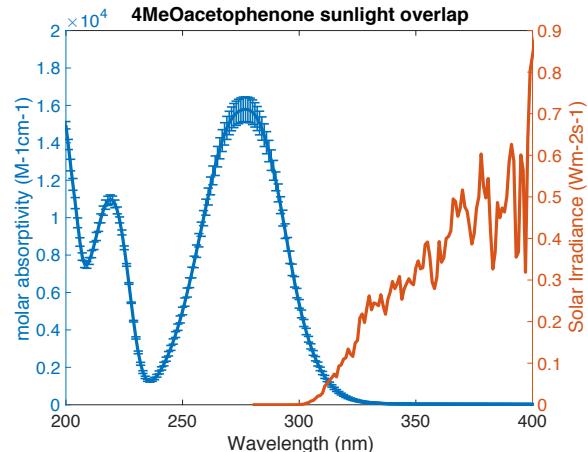
peak of an electronic transition over a period of wavelength may overlap. A molar absorptivity vs excitation energy plot is a linear combination of normal distribution plots from different electronic transition in a compound. Each individual peak combines together into a continuous experimental spectrum as seen in Figure 4. Notice different peak width at different excitation modes. Broadening effects from ground state geometry fluctuation affects their excited states differently. Once molecules are excited, solute molecules exist in the high energy state for a period of time. During that time a molecule can undergo different changes: undergoing a reaction by itself, emitting fluorescence light, or relaxing non-radiatively and de-energize back to the ground state. The rate constant of a reaction by the excited state ( $k_{rxn}$ ) and the rate of photon absorption ( $k_{abs}$ ) can determine the likelihood of a certain reaction to occur with respect to the amount of photons absorbed. A quantum yield ( $\Phi$ ) for a photo-reaction can then be calculated using Equation 1:

$$\Phi = \frac{\# \text{ of molecules undergoing a reaction}}{\# \text{ of photons absorbed}} = \frac{k_{rxn}}{k_{abs}} \quad (1)$$

Two main components needs to be determined to evaluate the quantum yield: the reaction rate constant and the light absorption rate constant. Reaction rate can be determined computationally from calculating two potential energy surfaces of the reactant and the product and solving for the conical intersection, where the two surfaces meet.

In Eustis lab, Nathan D. Ricke '14 worked on applying computational model to estimate quantum yields of photo-excited small organic molecules in the singlet excited states undergoing Intersystem Crossing (IC) from to the first triplet excited state using Landau-Zener model and outlining a possible method using Fermi's golden rule.<sup>23</sup> Landau-Zener model calculates the probability of IC using spin-orbit coupling constant at minimum crossing geometry, energy gradient, and crossing velocity. The rate is then calculated from probability from Arrhenius equation with activation energy being the

energy difference between the singlet excited state energy and the conical intersection. Quantum yields of aniline and benzene calculated using this theory were compared with the experimental values. Qualitative agreement with the experiments were observed and discussed. An alternative method is Fermi's golden rule, which calculates rates of transition from one energy state to another from electronic couplings,



**Figure 5:** Overlap between absorption spectrum of 4-MeOacetophenone and solar radiation (ASTM G173) showing at which wavelength photons from sunlight can be absorbed to energize the molecule to the excited states.

Franck-Condon factors, and solvent broadening factor. However, the calculation could not be done due to memory shortage caused by the number of terms generated at higher vibrational excitations. Nevertheless, Landau-Zener model can be used to calculate  $k_{rxn}$  of IC, the numerator term in Equation 1.

Since reactions of micropollutants in aquatic environment are the focus of this study, the main photon source for the photo-reactions is the sun. Sunlight at different wavelengths goes through the ozone layer, but most light in the UV range is absorbed by the ozone layer.<sup>24</sup> Some of the remaining light photons pass through the atmosphere before hitting the earth surface at a 90 degree. Some disperse in the atmosphere (commonly referred to as sky radiation<sup>25</sup>) before hitting the earth surface at different angles. Light intensity in terms of Solar irradiance (W) can be measured experimentally or looked up on a reference table: ASTM G173-03 Reference Spectra based on conditions considered average in most states in the US and tilt angle representative of the contiguous US states.<sup>26</sup> In water body, photon intensity decreases as a function of water depth and its attenuation characterized by the beam attenuation coefficient ( $\alpha$ ) and depth (z), both measured experimentally. These physical quantities of sunlight can be combined with the molar absorptivity ( $\epsilon$ ) of a species in water to give rate of light absorption ( $k_{abs}$ ) defined by Equation 2:

$$k_{abs}(\lambda) = \frac{W(\lambda) \cdot \epsilon(\lambda)}{z \cdot \alpha(\lambda)} \quad (2)$$

The focus of this study is to calculate molar absorptivity or equivalent in order to calculate the photon absorption rate in the aquatic environment. Study of UV-VIS absorption processes combined with the two methods outlined by Nathan Ricke illustrate what further calculations can be done to predict the quantum yield of sunlight-driven pollutant degradation in water.

### 1.3 Computational Approach

In recent years, high performance computers have become more and more accessible to researchers.<sup>27</sup> The power of High Performance Computing (HPC) extends the frontier of rigorous computational calculations to solve problems which a decade before were prohibited by the high computational cost (time and resources). To put “high cost” in perspective, wall clock time required to optimize geometry and vibration states of aniline in gas phase using different theories are shown in table 1. Computational time in the table was generated to provide guidance to understanding computation cost barrier faced when using higher-level theories. Theories are sorted in order of increasing computational time. Due to its accuracy, double excitation Coupled Cluster (CCSD)<sup>28,29</sup> is considered the hallmark of computational chemistry. Unfortunately, computational time required for CCSD limits its application to only 10-20 molecules; optimization of 14-atom aniline requires 15 computational days. Compared with normal Hartree Fock, CCSD takes 2000 times longer. For larger

systems, CCSD becomes impossible to be implemented. With larger number of CPU and more efficient CPU, wall clock time can be faster than computational time. If basis sets size doubles, the time required can be

**Table 1:** Comparing time required to optimize aniline’s geometry and vibrational states and ground state energy using Hartree Fock (HF), DFT, MP2, and CCSD level theory in Gaussian ’09<sup>30</sup> to demonstrate how computational time and accuracy varies with different theories. All calculation were done using 6-31G(d) basis sets on a 32-processor node and 10 GB connection.

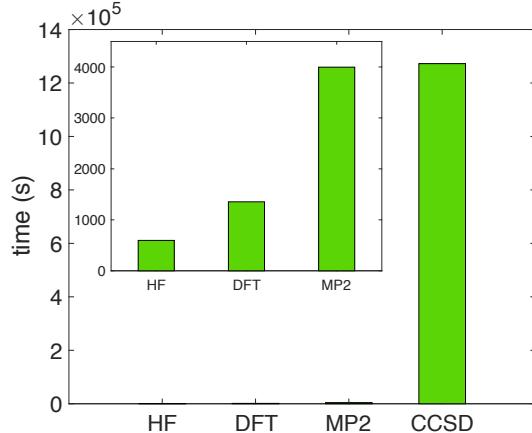
Theory	Scaling factor	Time	Time if basis sets double in size
HF	$N^4$	10 mins	2.7 hours
DFT	$N^3$	20 mins	2.7 hours
MP2	$N^5$	1 hours	32 hours
CCSD	$N^6$	15 days	2.6 years

calculated from the scaling factor (the last column in Table 1). It becomes clear that CCSD cannot be used for any calculation in a large system. None of the theories here scale linearly, which should reflect certain difficulties when the system becomes larger. MP2 also starts to take longer compared with HF or DFT.

Computational chemists can now implement more rigorous ab initio (from the first principles) methods.

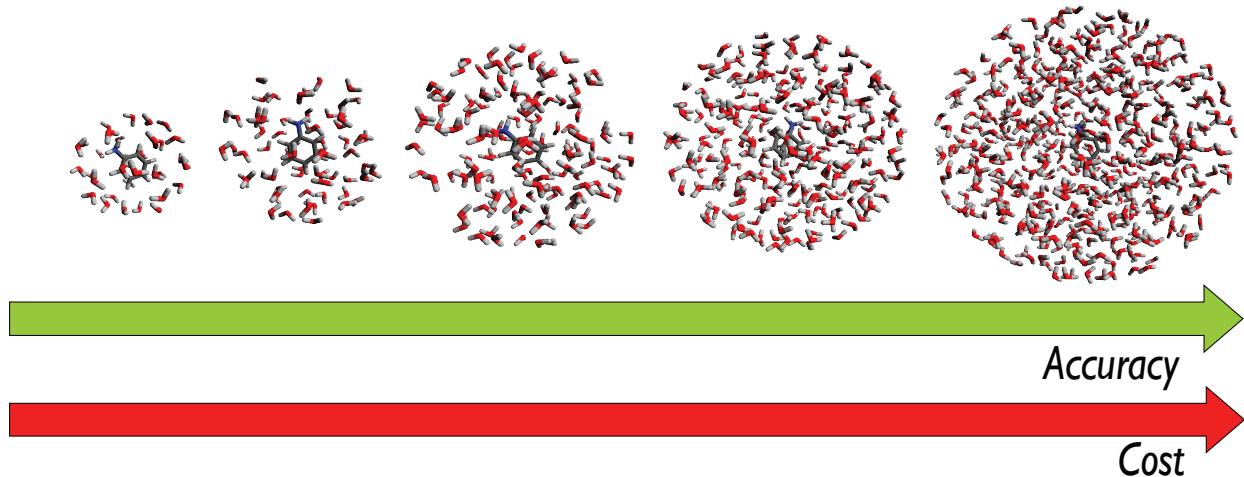
As more difficult problems are explored, there is a trend to incorporate both experimental and theoretical methods together for analysis because they have different advantages and disadvantages. For example when studying the photo-degradation of micropollutants, computational approach allows for a priori prediction of photo-products and their impact on the environment. Experimental approach is more limited to a posteriori study of negative impact micropollutants have on water system. Experimental

chemists face challenges to study individual species of micropollutants due to complex experiment planning. In natural waters, many species of micropollutants interfere with each other setting limits on what can be measured experimentally in a complex natural system. Experimental methods excel when used to study an observable natural phenomena and its consequences. While theoretical methods are great for studying simple system with good accuracy, they have gained attentions from scientists just in the past decades thanks to burgeoning of growth in computing speed. The computational theories are still in a stage of development as more and more researchers have more access to faster computers. They have become very accurate at calculating small and simple systems. However, ab initio theories remain incomplete and inaccurate for larger



**Figure 6:** Visualization of table 2. Due to their low computational costs compared with CCSD, time values of the other three theories are plotted in a box on the left hand corner.

and more complex systems. There is still a significant trade-off between cost and accuracy; the question is how much accuracy is needed? There are many factors that increase accuracy and computational cost: Theories, basis sets, solvent models...



**Figure 7:** Explicit  $2^n$  water model ( $n = 5-9$ ). The more water, the more complete solute-solvent model but the slower the calculation.

The objective of this project is to find a computational model with the best trade off between accuracy and computational cost to simulate excited states energy and probability of excitation of small organic molecules in water compared with experimental results. A work procedure will be developed to systematically evaluate the accuracy of different models. If successful, this procedure will be applicable to other larger molecules of different chemical properties. Criteria of being the best model include good accuracy when compared with the experimental data, reasonable computational cost, and applicability to different models. The best model to calculate the excited states will lead to rate constants of photon absorption, quantum yields, and general prediction of photo-transformation pathways of small organic pollutants.

## 2 Background

### 2.1 Hartree-Fock Theory

Schrödinger equation is an equation describing the relationship between the energy of system ( $E$ ) and the wave function ( $\Psi$ ) of a wave-particle. It is written in time-independent form:

$$H\Psi = E\Psi \quad (3)$$

$$H = -\frac{\hbar^2}{2m}\nabla^2 + V \quad (4)$$

with  $H$  is the Hamiltonian operator, which describe what the energy of the system is controlled, written here in italics instead of with a circumflex mark over the operator (in physics and computational chemistry,  $H$  is the convention more commonly used). The left operator is the kinetic energy operator with Laplace operator written using the del symbol. For simple external potential energy conditions like particle in a box or harmonic oscillator, the wave function can be deduced from the Schrödinger equation and its energy calculated. More information about Schrödinger equation can be found in any elementary quantum mechanics books.<sup>31,32</sup> For a system with more than two or more electrons, there is no direct way to solve the Schrödinger equation for an analytic wave function. Variational principle can be used to find the most “correct” numerical wave function that satisfies an eigenfunction in lieu of the Hamiltonian operator in a certain theories. With a trial wave function ( $\phi_{trial}$ ), energy calculated from the Hamiltonian operator is always higher than energy calculated using the true wave function ( $\psi_{real}$ ):

$$\langle\psi_{real}^*|H|\psi_{real}\rangle \leq \langle\phi_{trial}^*|H|\phi_{trial}\rangle \quad (5)$$

According to this principle, one can keep modifying parameters in the trial wave function until the calculated energy converges to the lowest value before concluding that the trial wave function is the best representation of the real wave function. The lowest energy calculated is by no means the “correct” value. Theories make assumptions and those assumptions limit how accurate the lowest energy of a theoretical calculation can achieve. Another useful assumption widely used in quantum chemistry is Linear Combination of Atomic Orbitals (LCAO) which allows *atomic* orbital wave functions ( $\chi_n$ ) to be summed to create a *molecular* wave function ( $\phi_i$ ):

$$\phi_i = a_{i1}\chi_1 + a_{i2}\chi_2 + a_{i3}\chi_3 + a_{i4}\chi_4 + \dots \quad (6)$$

Wave function theories mostly use this assumption. One fundamental theory, on which many other theories are base, in computational chemistry is Hartree-Fock Theory (HF). HF theory make assumptions about wave functions and the nature of quantum mechanical interactions in the atomic systems, including LCAO, single-electron wave functions decoupling, and non-relativistic electrons, to arrive at Fock operator ( $F$ ) for a single-electron wave function:

$$F_i = h_i + \sum_j^N (J_i - K_j) \quad (7)$$

with  $h$  being the kinetic and electron-nucleus Coulomb potential energy of an electron,  $J$  electron-electron Coulomb interactions,  $K$  electron-electron exchange (a term for Pauli Exclusion Principle). Eigenvalues of  $F$  are matrices with Lagrange coefficients. These coefficients can be used to calculate energies of each one-electron wave functions before being added back together to calculate total energy of the system. Iterations of guessing a trial wave function through LCAO coefficient adjustment, operating  $F$  on each single-electron wave functions, and energy calculations are done until the system with the lowest energy is achieved. Limitations of HF are in its assumptions. First, electrons are coupled so single-electron wave functions cannot truly represent one electron without taking into account electron-electron correlations. One of many post-HF theories, most notably Møller—Plesset perturbation theory (MP)<sup>33</sup>, attempts to add electron-correlation to correct errors introduced by approximations in HF.

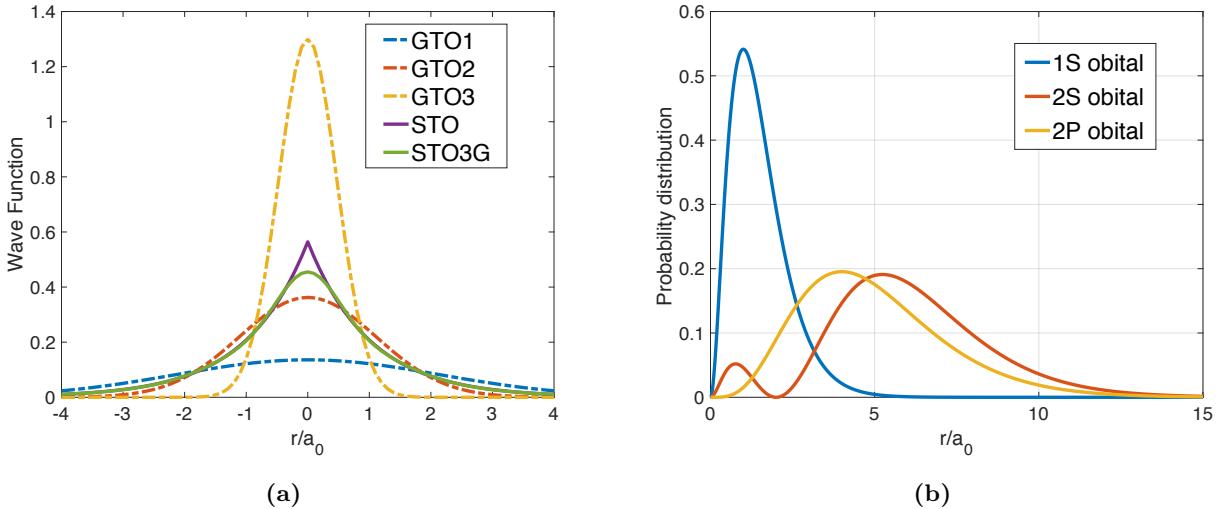
## 2.2 Basis Sets

To describe electron orbitals in a molecule, wave functions for quantum calculations are modeled using mathematical functions called basis sets. Simple Gaussian-type basis sets are a combination of Gaussian distributions modeling different orbitals in Cartesian coordinates. They were designed by John Pople to be optimized for speed.<sup>34</sup> Smaller basis sets have fewer basis functions per atom. In addition to having higher functions per atom, larger basis sets may also include polarized and diffuse basis sets. Equation 6 can be rewritten to include two 1S orbitals of two hydrogen atoms in hydrogen gas:

$$\phi_{H_2} = a_1 \chi_{H_A}^{1S} + a_2 \chi_{H_B}^{1S} \quad (8)$$

While this wave function is adequate to describe ground state hydrogen gas, it will certainly fail to describe the excited state which requires orbitals in the second electron shell for the electrons in the excited states. Therefore, more orbitals can be added to better describe the molecular orbitals:

$$\phi_{H_2} = a_1 \chi_{H_A}^{1S} + a_2 \chi_{H_A}^{2S} + a_3 \chi_{H_A}^{2P} + a_4 \chi_{H_B}^{1S} + a_5 \chi_{H_B}^{2S} + a_6 \chi_{H_B}^{2P} \quad (9)$$



**Figure 8:** Visualization of basis sets. a)plots of STO and STO-3G with its Gaussian components. All three Gaussian plots with 0.44, 0.54, and 0.15 mixing coefficients are added to STO describe atomic orbitals better. First, there is no more discontinuity in the first derivative of the wave function as in STO. Gaussian functions are also easier to integrate compare with exponential functions. Data replotted with permission from Magalhães<sup>35</sup>. Copyright 2014 American Chemical Society. b) probability density of different atomic orbitals. Diffuse functions can be added to light and heavy atoms to better describe more diffuse orbitals when electrons move further away from the nucleus. In this figure, 2S and 2P can increase the distance an electron wave function describe with expanded probability distribution.

The atomic orbitals can also be described by analytical mathematical expressions controlled by numerical parameters like Gaussian distributions. For example, in describing the molecular orbitals in aniline molecule, there are different level Gaussian basis set to choose from: STO-3G, 3-21G, 6-31+G(d), ..., 6-311++G(2d,p). STO<sup>36</sup> is a Slater-type minimal basis set designed for fast calculation. It is the solution to 1S Hydrogen atom using only an exponential function. In contrast, STO-3G mixes STO with 3 more Gaussian functions to describe the wave function more accurately (see figure 8a). 3-21G is also a small model (cheap and simple). It contains 3 Gaussian functions for inner shell orbitals, 2 and 1 different level of Gaussian functions describing valence electrons. 6-311++(2d,p) is more complete but also more expensive. It contains 6 Gaussian functions for inner shell orbitals, 3, 1, and 1 different level of Gaussian functions describing valence electrons. In addition, ++ means it adds two large s orbitals and one large p orbitals *diffuse* functions. This is important in calculating reactions with electron being further away from the nucleus such as electron addition, electronic transition, or charge transfer. Finally, (2d,p) means heavier atoms also contain two d orbitals polarized to the P orbitals and hydrogen atoms p functions to s orbital. Polarized basis sets becomes important when describing physical properties involving electron excitations.

Another one of the most well-known basis sets is the Dunning's correlation-consistent polarized Core and Valence (Double/Triple...) Zeta (cc-pVXZ) type basis sets. They are basis sets in spherical coordinates

that contain variationally optimized coefficients for electron correlations (in contrast to Pople's which was optimized for Hartree Fock). Dunning's cc-pVXZ basis sets may include a large number of polarized basis sets based on what level of Zeta is used. For example, cc-pV6Z describes a second-row atom using 140 basis sets.<sup>37</sup> Diffuse basis sets can be added to cc-pVXZ with aug- keyword in front of the basis sets. In comparing among theoretical calculation, aug-cc-pVQZ (ACCQ) can be used as theoretical benchmark.<sup>38,39</sup> ACCQ's computational cost is however prohibitive for solute-solvent model with Bowdoin HPC.

**Table 2: Demonstrating how time required to optimize hydrogen gas geometry and vibrational states and ground state energy using different basis sets in Gaussian '09. All calculation were done with HF level of theory on a 24-processor node and 64 GB of RAM. This table underscores trade-off between accuracy and cost using aug-cc-pvQz as the theoretical benchmark**

Basis Set	Time (seconds)	Energy (Eh)	R (A)
3-21G	134	-1.1230	0.735
6-31G	133	-1.1268	0.730
6-311G	137	-1.1280	0.732
6-31G(d,p)	154	-1.1313	0.732
6-311G(d,p)	153	-1.1325	0.735
6-311G(d,2p)	163	-1.1330	0.734
6-311G(d,3p)	161	-1.1331	0.734
6-311G(d,3pd)	194	-1.1331	0.734
aug-cc-pvQz	3250	-1.1335	0.734

To put computational time required by different basis sets in perspective, hydrogen gas geometry can be optimized using different Gaussian-type basis sets and aug-cc-pv6z. This computational demonstration illustrates the main points addressed below about how size of different basis sets affect time and accuracy. As basis sets size increased, computational time increased and energy converged closer and closer to that of the theoretical benchmark. 6-311G(d,2p) is the best basis sets for geometry optimization because it used 20% less time to calculate compared with 6-311G(2d,p), but with the same degree of accuracy. It is worth noting that depending on the type of calculations and the size of the systems, different basis sets should be chosen to best balance accuracy and cost. Adding (d,p) polarized basis sets helped get more accurate H-H radius both in 6-31G to 6-31G(d,p) and 6-311G to 6-311G(d,p). However, within Gaussian-type basis sets group, changes in time as more basis sets added was not as drastic as moving from 6-311G(d,3pd) to ACCQ which takes 15 times longer time to run. When the system becomes larger, the calculation time scaled would become longer. Increasing the size did not always increased accuracy. Note that the radius calculated using 6-311G(d,p) converged to 0.734 Angstroms, close to theoretical benchmark at 0.735 Angstroms. Increasing the size of the basis sets after did very little to improve the value. Different types of calculation then would be more suitable with different basis sets: 6-311G(d,p) for geometry optimization and 6-311G(d,2p) for energy calculation. In calculating excited state energies of organic molecules<sup>39,40</sup>, while having cheaper cost, an average-sized

basis set 6-311++G(2d,p)<sup>41</sup> performs better than aug-cc-pVDZ (ACCD)<sup>41</sup>. For example, transition energies calculated of CN molecule as calculated by ACCD deviates 1117-1669 cm<sup>-1</sup> from experimental value while those by 6-311++G(2d,p) only deviate 220-470 cm<sup>-1</sup>.<sup>39</sup>

### 2.3 Density Functional Theory

In theory, Density Functional Theory (DFT) does not require making assumptions about electron-electron decoupling because it does not encounter problems with degrees of freedom as wave-function-based methods do. It therefore has potential to be more exact compared to HF or post-HF wave function methods. Despite its low costs, DFT's accuracy in certain applications has been accepted and acclaimed.<sup>42,43</sup> Wave-function-based methods locate each particles in wave functions with 3 parameters on cartesian coordinates.

$$\Psi(x_1, y_1, z_1, x_2, y_2 \dots) \quad (10)$$

In multi-particle systems, computational cost rises exponentially with the number of particles. For a large solute-solvent system, the calculation cost becomes expensive. DFT uses electron density function with only 3 Cartesian parameters. Despite having fewer parameters, density functions also contain data about total number of electrons, position of nuclei, and atomic number; all required to calculate the Hamiltonian operator.

$$\rho(x, y, z) \quad (11)$$

Energy of a density function is not calculated based on HF or other wave function methods. Historically, DFT had large errors and limited usage in computational chemistry because there is no analytical or numerical way to solve for electron density functions. In wave function based methods, variational principles governs that energy calculated from the operating Hamiltonian on a candidate wave function gives higher energy than on the correct wave functions. By adjusting the wave function, until its energy is minimized, a wave function as a solution to that Hamiltonian operator can be solved for numerically. Two important theorems developed by Pierre Hohenberg and Walter Kohn allow DFT to be used more broadly in computational chemistry community by providing a way to calculate the energy via variational principle.<sup>44</sup> The first theorem states that only ground state electron density is required to construct the external potential function and hence the Hamiltonian. This is important because potential for excited state wave functions can also be calculated. The second theorem, is the variational principle for DFT allowing to change parameters in the density

function to get lowest energy through wave function and Hamiltonian calculation.

$$\left\langle \Psi_{cand(DFT)}^* | H | \Psi_{cand(DFT)} \right\rangle = E_{cand(DFT)} \geq E_0 \quad (12)$$

The problem with DFT as described is variational principle calculation on the wave functions because it requires solving the Schrödinger's equation to obtain the wave function and make assumptions/approximations as discussed above. To solve this, Kohn-Sham Self-consistent Field (Kohn-Shan SCF), developed in 1965, made use of the two theorems by breaking down electron density to non-interacting single-electron density functions.<sup>45</sup> Energy functionals can then operate on linear combination of these electron functionals to give the energy of the system. Functionals are functions which operate on another function. Instead of using the Hamiltonian operator and a wave function to calculate the energy (Equation 3), Kohn-Shan SCF calculates the energy of the system directly from electron density function using different functionals. For example, the Kohn-Sham one-electron operator as presented in Cramer's Essentials of Computational Chemistry<sup>37</sup>:

$$h_i^{KS} = -\frac{1}{2}\nabla_i^2 - \sum_k^{Nuclei} \frac{Z_k}{|r_i - r'|} + \int \frac{\rho(r')}{|r_i - r'|} dr' + V_{xc} \quad (13)$$

with the familiar looking first two terms (see Equation 4). The last term on the right hand side of the equation refers to exchange-correlation functional; it includes quantum mechanical exchange, electron correlation, and correction in kinetic and potential terms. Exchange interaction in fermions like electrons follows Pauli Exclusion Principle, constraining electrons from being described by the same quantum numbers. Correlation is an electron-electron interaction term ignored in HF. In modern DFT functionals, experimentally fitted parameters are used to describe these energy functional components.<sup>37</sup> Different functionals have different fitted parameters and their accuracy dependent on the system and physical properties being calculated. To find the best group of functionals for calculating excited state energies, there needs to be a benchmark specifically on excited state energy calculation on similar systems and electronic states. In modern functionals, a functional of electron density gradient function is added to the energy functionals. These functionals are gradient corrected and the method called generalized gradient approximation (GGA). Pure GGA functionals is known to underestimate system energy<sup>46-49</sup> or excited state energy<sup>39,50</sup>. In addition, a percent of contribution calculated from HF can be mixed in with the pure GGA results to create GGA-HF hybrid functionals. According to Hellmann-Feynman theorem (Equation 14), an actual exchange-correlation energy can be described as an energy plot over extent of correlation ( $\lambda$ ), as shown with green shaded area in Figure 9 recreated here based on figures in two literature sources.<sup>37,51</sup> Area of box A and B respectively refers to pure HF exchange energy and pure GGA exchange-correlation energy. Since the energy as a function of extent

of correlation is not known, experimentally fitted value ( $z$ ) is needed to find the integral using Equation 15. Adding HF exchange energy reduces overestimation of energy calculated from GGA functionals.

$$E_{xc} = E_{xc} = \int_0^1 \langle \Psi(\lambda) | V_{xc} | \Psi(\lambda) \rangle d\lambda \quad (14)$$

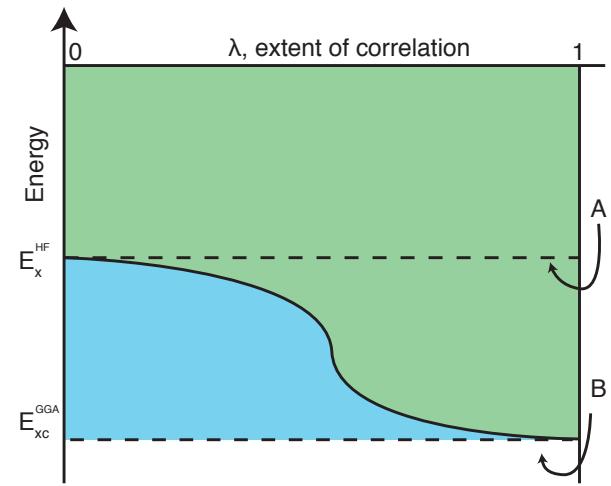
$$E_{xc} = E_x^{HF} + z(E_{xc}^{GGA} - E_x^{HF}) \quad (15)$$

Functionals with HF component are called hybrid functionals. Some of the more widely-used hybrid functionals include: B3LYP, CAM-B3LYP<sup>52</sup>, M06-2X<sup>53</sup>, and PBE0<sup>54</sup>. New functionals can also be created from mixing one functional for exchange and one for correlation. For example, BP86 is a combination of B88 for exchange and P86 for correlation. Unlike basis sets, which can become more complete as number of basis sets per atom increases, there is no 'full,' or complete, functional to use. Among all current theories to calculate the excited states, Time-Dependent Density Functional Theory (TD-DFT) is the most promising with its high accuracy when used with appropriate functionals and low computational cost<sup>55</sup>. TD-DFT calculate the electron density response to excitation photon modeled as an external perturbing field varying with time.<sup>56</sup>

Despite its under-performance in certain systems like organic dyes<sup>57</sup> and sulfur organic compounds<sup>58</sup>, studies suggest limitations can be overcome if appropriate functionals are used.<sup>49,59</sup> According to one study of 0-0 transition in radicals, the two best-performing DFT functionals are CAM-B3LYP and M06-2X. Pure GGA functional BP86<sup>60,61</sup> and long range corrected GGA LCBLYP<sup>60,62,63</sup> are also attractive to have a diverse group of functionals.

## 2.4 Solvent Models

Despite the recent advent of growth in computer speed and the burgeoning interest in incorporating computational models to further understand the natural world, large systems such as solvation models remain



**Figure 9:** Illustration of hybrid exchange-correlation term from HF and DFT created here based on figures in two literature sources.<sup>37,51</sup> Box A drawn by the first dotted line is an integral of exchange-only energy calculated from HF from  $\lambda$  value 0 to 1. Box B drawn by the second dotted line is the same integral but of exchange-correlation energy calculated from pure-GGA DFT. The value of actual exchange-correlation energy (shaded green) is likely to be described by the value between area of box A and B.

a big challenge.<sup>64</sup> In describing the solvents, one can either apply bulk solvent effect on the solute through different potentials without specifying the solvent molecules (implicit solvent model) or one can explicitly place the solvent molecules around the solute molecule (explicit solvent model). When calculating excited states energy of a solute, solute-solvent interactions will contribute to modeling solvent effects on absorption spectra as previously discussed (Figure 3a and 3b). In the past, implicit solvation models were frequently implemented because they gave acceptable results while maintains low computational cost. In contrast to explicit model, interactions between each individual solvent molecule and the solute molecule are not calculated. The most commonly used implicit model is the Polarizable Continuum Model (PCM).<sup>65–67</sup> Instead of explicitly handling each solvent molecule quantum mechanically, PCM expresses their bulk effects on solute molecule in means of dielectric continuum field surrounding molecule of interest. Dielectric field can polarize and affect the wave function of the solute system. Its disadvantage is its lack of dynamic solvent contribution on excited states properties.<sup>59</sup> Implicit solvent models also neglect hydrogen-bonding as they assume implicit implementation in dispersion forces and electrostatics.<sup>68</sup> In molecules with hydrogen bonding, implicit solvent model fail at quantum mechanical level to fully model solute-solvent interactions. Having more rigorous calculations, explicit solvent models, require more computational cost compared with implicit water models. Cost for treating a solute molecule in a large solvent sphere with full quantum mechanics can become prohibitive depending on basis sets, theory, the number of explicit solvent molecules (see Table 1 and Figure 7). One recent notable model, Effective Fragment Potentials (EFP), can be used to model explicit solvent molecules with non-bonded van der Waals interactions, hydrogen bonding. It implements Coulomb interactions, polarization, and exchange repulsion without high computational expense of explicit models.<sup>69,70</sup> Implementing EFP solvent model, TD-DFT can be used to accurately calculate excited state energy of acetone in water.<sup>70</sup> In incorporating the solvent model, there are also two different methods: static or dynamics. In static model, geometry optimization of solute molecule is carried out, then followed by calculation of excited state energies. This static ground state molecule however does not accurately represent solute in water.<sup>71</sup> As discussed earlier, the system molecules are dynamic. The difference in solute geometries can broaden the solute's excited state energies in solvent. Instead, Molecular Dynamics (MD) of solute and EFP solvent fragments can be used to obtain a range of equilibrated structures at a specific temperature (see Figure 12a). Mark Gordon et al. averaged the calculated energies of each excited state to arrive at a final excited states energy.<sup>71</sup>

## 2.5 Molecular Dynamics

Canonical Ensemble Molecular Dynamics (NVT)<sup>72</sup> provides a pool of equilibrated ground state geometries under constant number of particles, volume, and room temperature, 298K (see Appendix). Starting with a higher-energy initial geometry, a system relaxes to lower potential energy state until it reaches the dynamic equilibrium. However, there is no accepted measure of equilibrium. One study by Chodera suggests looking at how density of liquid argon changes over time and using variance-bias trade off to find equilibrium period.<sup>73</sup> By eliminating non-equilibrated geometries, the rest of the equilibrated geometries become “equilibrated” samples. As requirements for the equilibrium determination are relaxed, more “equilibrated” samples can be used for further calculations, the variance of the value calculated decreases. However, error from including non-equilibrated geometries (bias) increases as more samples, which has not really reached the equilibrium enter the pool. The best trade-off were determined to be a cut off-time when the ratio of sampling size to estimated statistical inefficiency is maximized. For a large system like organic solvent-solute model, this algorithm cuts off too much of the sampling size so it follows that its estimate of statistical inefficiency is not accurate. According to the study, the accuracy of cut-off should be questioned if the whole sample time is not at least an order of magnitude larger than cut-off time value.

## 2.6 Comparing Experimental And Theoretical Calculation Results

As discussed above, experimental molar absorptivity is not a characteristic of a chemical property; the value depends on solvent broadening and measurement resolution. Molar absorptivity at a single wavelength value( $\epsilon$ ,  $M^{-1}cm^{-1}$ ), with a value range of 0 to 50000  $M^{-1}cm^{-1}$ , may qualitatively imply probability of transition. The larger the molar absorptivity, the larger the absorptivity, the more photons absorbed. Molar absorptivity does not have an upper limit. Results obtained from theoretical calculations are composed of excited states energy (E, nm) and oscillator strength (no unit).<sup>74</sup> Oscillator strength, with range from 0 to 1, is a measure of how probable a specific electronic transition can occur. There is no solvent broadening involved in calculating static oscillator strength from a geometry, therefore theoretical oscillator strength cannot be compared with molar absorptivity directly. Oscillator strength tells the degree of allowness of transition between two electronic states. Mathematically, oscillator strength can be calculated from integrating molar absorptivity over a range of wavenumber within an electronic transition (see equation 16).<sup>75</sup>

$$f = 4.3 \times 10^{-9} \int_{\tilde{\nu}_a}^{\tilde{\nu}_b} \epsilon(\tilde{\nu}) d\tilde{\nu} \quad (16)$$

## 3 Methods

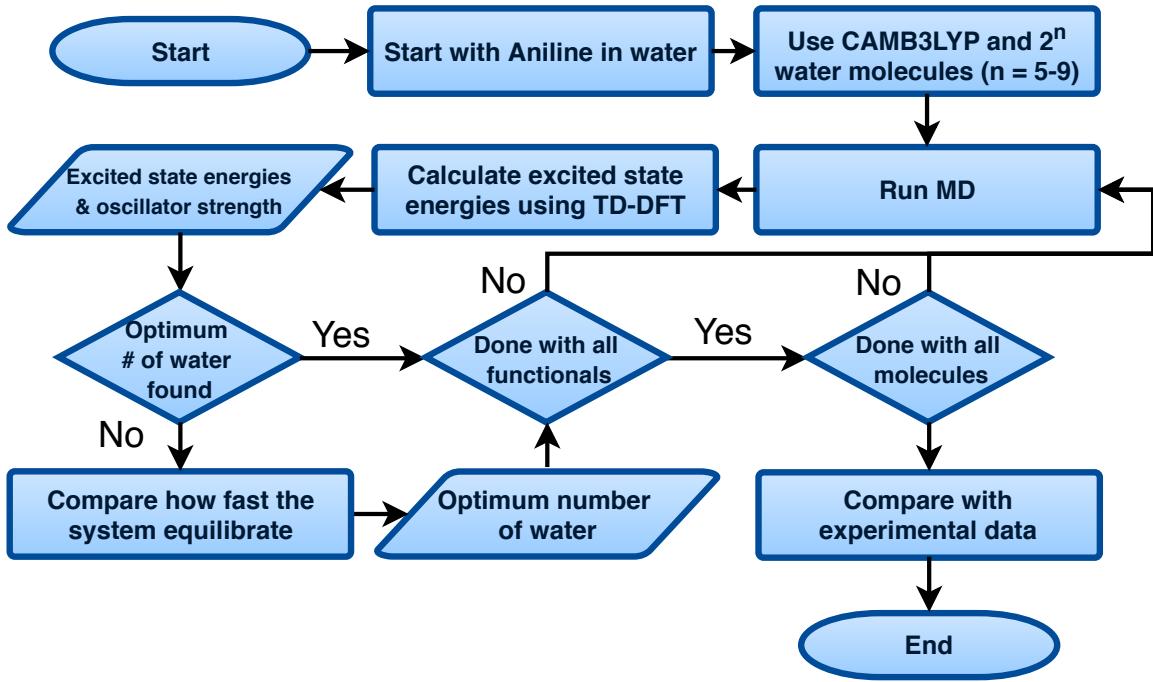
### 3.1 Overview

Systematic approach to developing computational model is outlined in this section. Python<sup>76</sup> scripts were used to facilitate some of the steps (see Appendix). Aniline was chosen as the initial test solute due to its hydrogen bonding capability, its small size, and extensive previous experimental and theoretical work on aniline within the Eustis lab. A successful model of aniline was expected to be applicable even to larger molecules of similar class, organic and aromatic. Starting molecule geometries were obtained from optimizing the molecular geometry of the solute and solvent molecules yielding two optimized geometries for solvent and solute. These two geometry files were then inputed into Packmol program<sup>77</sup> with empirical guess initial solvent radius and number of water molecules to create starting geometry of solvated solute. MD-TD-DFT using the geometry from Packmol was performed using GAMESS quantum computational package. All MD run used 6-311++G(2d,p) basis sets, calculating velocity vector every 1 fs while calculating other physical properties displaying every 10 fs at bath temp = 298 K. All TD-DFT calculations use 6-31+G(d,p) basis set. Water molecules were modeled explicitly using EFP1 and positioned initially using Packmol. Criterion for MD equilibrium is that the average change of systems potential energy over time crosses zero. After equilibrium has been reached, a 10 ps of equilibrated geometries are used in calculating excited state energies using TD-DFT in GAMESS. Matlab<sup>78</sup> was used to deconvolute experimental spectra and theoretical data cluster. The flowchart for finding the best model is outlined below in Figure 10. First, the optimal number of water molecules in the model was evaluated by using different numbers of water molecules surrounding aniline with CAM-B3LYP functional. The results were compared with experimental data to evaluate the accuracy of the model. The model which contained the lowest number of water, but still provided the most accurate results compared with experimental data is chosen, then applied to other molecules with different functionals before compared again to experimental data.

### 3.2 Computational

#### 3.2.1 Facilitating Each Steps Using Python Scripts

Python scripts play an important role in both data collection and smoothing the process between each computational steps. In order to automatically generate input files and cultivate output data from output files, many python scripts were written from scratch (see Appendix). Since each script is specific to each calculation task, there is a limited number of scripts available on the internet (virtually none for this project) Log files obtained from GAMESS contains both valuable theoretical data and useless text strings. For



**Figure 10:** Flowchart of work in this project. Python scripts were utilized to facilitate certain steps. Section and page number of important steps are listed here: reason for “Start with aniline in water,” Section 3.1 on page 17; procedure for “Use CAMB3LYP and  $2^n$  water molecules,” Section 3.2.2 on page 19; results and discussion of “Run MD” and “Calculate excited state energies using TDDFT” Section 4.1 on page 22; “Compare how fast the system equilibrate,” Section 4.2 on page 27; “Compare with experimental data,” Section 4.3 on page 29;

example, even though WebMO<sup>79</sup> automatically generates last geometry calculated from MD run, extracting a geometry from each MD step requires one to manually open the log file, copy-paste, reformat each geometry into an input file for TD-DFT calculation one geometry at a time. The python script postMDDataPull2.py is designed to pull thousands of geometries and generate GAMESS input files for TD-DFT energy calculation within seconds. In some cases, using Python scripts was required. When running MD of 512 water molecules, the size of the log file became so large (15 GB). The text files that size are unsupported by any regular graphical interface text editor. Extra care were taken when designing these scripts to deal with large data to make sure fast speed is maintained and results consistent when applied to different files. Generating these python scripts will also allow unified program to be developed in order to automate the whole project without any manual input.

### 3.2.2 Preparing Starting Geometry Using Packmol

As in many theoretical calculation, an initial guess is required. In MD, initial geometries of the solute and solvent were required to start MD calculations in GAMESS. This guess geometry was calculated using Packmol by inputting optimized geometries of solute and solvent calculated from using DFT with CAM-B3LYP functional and 6-311G(2d,p) basis sets. A sample of Packmol input file for solvating aniline in 16 water molecules is displayed below. The program does not return system geometry unless solvent molecules can fit in smallest, user-specified, radius without exceeding steric strain controlled by tolerance value (line 25). The value, 2.5 Angstrom, controlled the smallest distance atoms need to be far away from each other (line 6). If molecules could fit, solvent radius is increased slightly until a geometry was obtained. Notice this input command asks for spherical solvent sphere, which is convenient to use when implementing SSBP. If periodic boundary condition is to be used, line 25 should be changed to “inside cube 0. 0. 0. 5.” with a cube of dimension 5 x 5 x 5 Angstrom<sup>3</sup>.

```
1 #
2 # aniline solvated by 16 waters
3 #
4
5 # define minimum distance between atoms to 2.5 Angstrom
6 tolerance 2.5
7 # number of loops for fitting
8 nloop 10000
9 # number of trial per loop
10 maxiter 100
11
12 filetype xyz
13 # specify the solute
14 structure aniline.xyz
15   number 1
16   # set center of mass at (0,0,0) without any rotation
17   fixed 0. 0. 0. 0. 0. 0.
18   centerofmass
19 end structure
20
21 # specify the solute
22 structure h2o.xyz
23   number 16
24   # set center of solvent sphere at (0,0,0) with 5 Angstrom radius
25   inside sphere 0. 0. 0. 5.
26 end structure
27
28 output aniline16.xyz
```

### 3.2.3 Computational Models

Packmol was used to place water molecules around solute molecules to create starting geometry. GAMESS computational chemistry software was used in calculating MD and excited state energies. Fityk program<sup>80</sup> was used to fit Gaussian under absorption spectrum (Figure 4), and theoretical data (Figure 12b). In order to

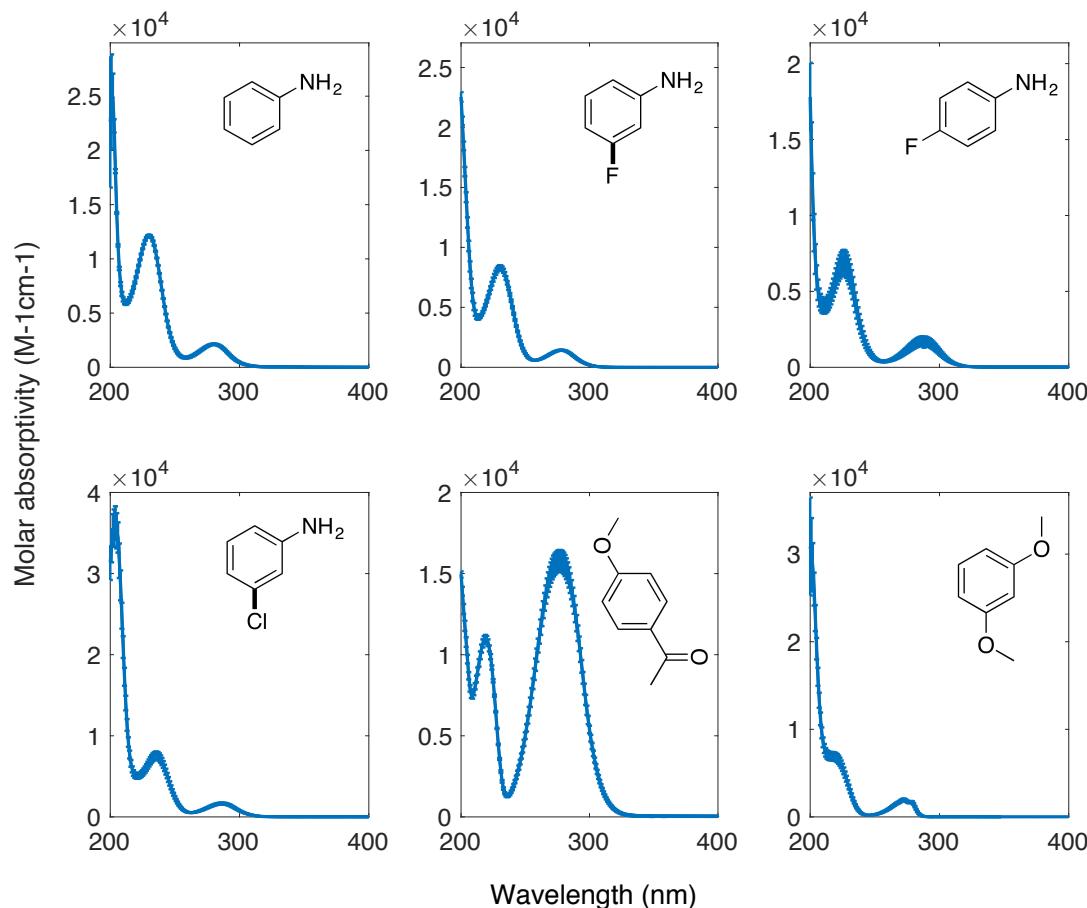
accurately calculate the excited energies, the 6-311++G(2d,p) basis sets was chosen to run vertical excitation energies on a set of equilibrated geometries. It was decided to reduce the basis set in running MD. A smaller basis set 6-31+G(d)<sup>41</sup> was used in order to cut computational cost. In MD calculation, diffuse and polarized basis sets are not as important since electrons do not move as much compared with TDDFT where charge transfer occurs when electrons move from HOMO to LUMO. It is an accepted knowledge that smaller basis sets produces as accurate geometry compared with larger basis sets (see also Table 2). Smaller basis sets also converge more easily giving better accuracy-cost trade off. EFP1 model of water molecules was chosen to implement an explicit solvent model and calculate excited states energies. The appropriate number of water molecules in the model had never been evaluated. Too many number of water molecules require more expensive computational cost. Too few water will not give a full model of solute-solvent interaction. A binary framework was used to test how many water molecules were needed to fully solvate the aniline: 2<sup>n</sup>, n=5 to 9 (see Figure 7).

After number of solvent molecules suitable to model solute-solvent for photo-excitations in aniline was decided, MD-TD-DFT excited state energies and oscillator strengths were calculated for 3-F-aniline, 4-F-aniline, 3-Cl-aniline, 4-MeOacetophenone, and (1,3)-dimethoxybenzophenone using that model. Once excited state energies for each system were calculated, the results were compared with experimental spectral data to find which functional performed best at calculating MD-TD-DFT excited states energy. Since CAM-B3LYP is a good functional for calculating ground state energies, equilibrated geometries will be calculated using CAM-B3LYP, and TD-DFT excited states using different functionals. Other functionals are PBE0, M06-2X, BP86, and LCBLYP. As previously discussed, M06-2X was expected to perform well when calculating excited state energies (see Introduction). PBE0 is a generic hybrid functional. BP86 is a pure GGA function (no HF). LCBLYP is a long range corrected functional.

### 3.3 Experimental

Experimental absorption spectra of 6 small organic molecules (aniline, 3-F-aniline, 4-F-aniline, 3-Cl-aniline, 4-MeOacetophenone, (1,3)-dimethoxybenzophenone) in pH 7 water was collected by Alex Poblete '17 and Holly Rudel '17 over the summer of 2015. Experimental methods as written in a currently submitted but unpublished article<sup>81</sup> is rewritten here in full. The organic analytes were purified via vacuum distillation under 99.998 % nitrogen atmosphere and/or as described in the literature<sup>82</sup>. Stock solutions (1 - 10 mM) were made in purified methanol and stored in amber vials at 4 degrees Celsius. Type I water (AquaSolutions,  $\geq 18 M\Omega^1 cm^{-1}$ ) was used for all aqueous solutions. Solutions were buffered at pH 7.0 with 10 mM phosphate buffer. Quartz glassware (Starna Cells, USA) was used for all spectroscopic measurements. UV-Vis

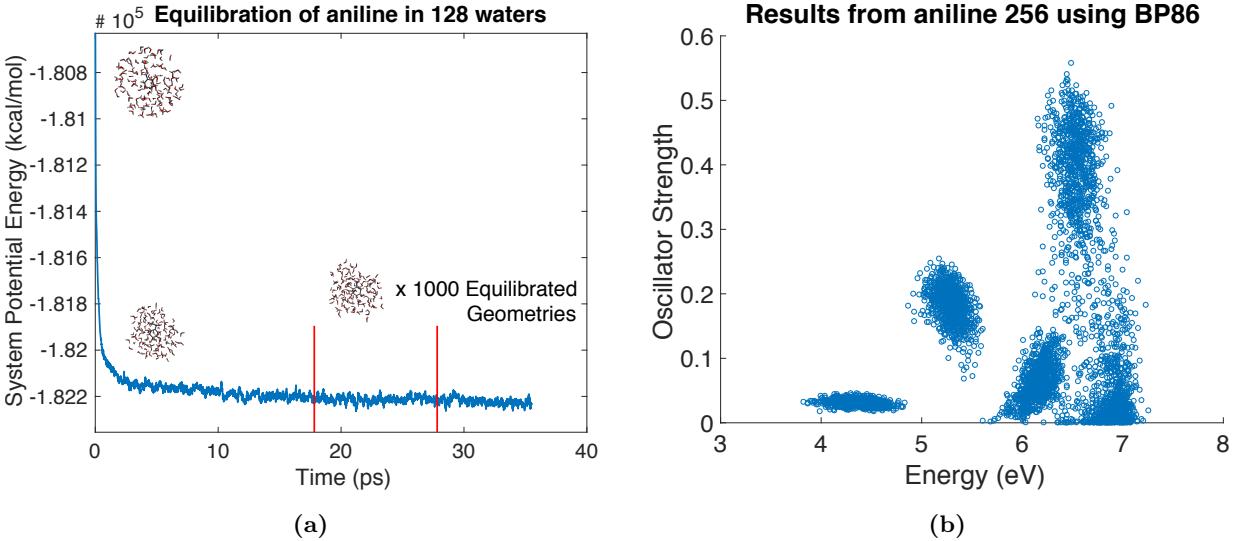
absorption measurements were performed using a dual-beam (sample/blank) Cary 400 spectrometer. The instrument is calibrated for energy accuracy (holmium oxide gel, Starna) and transmittance (Dichromate) on a regular basis. scanning from 190 - 400 nm at 1 nm increments at a rate of 50 nm per minute. The slit width was fixed at 1 nm resolution. Blank solutions were prepared from the same 10 mM phosphate buffered Type I water. Instrument correction was performed by inclusion of a blank in each concentration series, and via real-time reference subtraction against an additional blank during each measurement. Each analyte series consisted of a blank and five samples in the concentration range of 5 - 100  $\mu$ M. The concentration range chosen for each analyte was chosen to avoid concentrations approaching the solubility limit and to provide  $\epsilon_{\lambda}$  values  $\leq 0.1$  absorbance units for the most concentrated sample. Issues with insolubility at pH 7 / 25 degrees C were avoided by examining available solubility data<sup>83</sup>, assuring all stock and working solutions are fully stirred at room temperature before transfer.



**Figure 11:** Experimental UV-VIS spectra of 6 organic molecules measured in pH 7 water and 298K collected by Alex Poblete '17 and Holly Rudel '17 over the summer of 2015. Each spectrum represents the fitted slope of six samples with values of  $c$  from 0 to 100M, with the Y-error bar representing the standard error or 95% confidence interval.

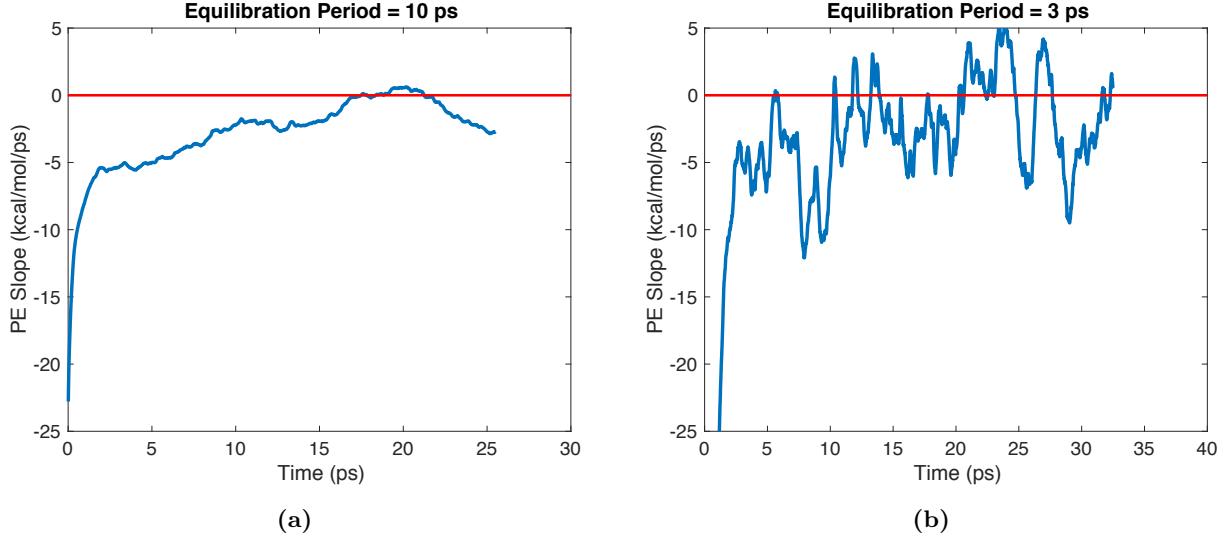
## 4 Results And Discussion

### 4.1 Results from MD-TD-DFT



**Figure 12:** MD run and excited state data calculated from TDDFT of equilibrated geometries from MD. a) plot of system's potential energy vs time in MD with geometries shown at  $t=0, 0.2, 1.78, 2.78$  ps. System energy and solvent radius decreased very rapidly in the first 0.2 ps, before converging to the equilibrium value. Fluctuation in energy is strong, around 0.1% but the equilibrium was determined to be from 1.78-2.78 ps (see methods). b) theoretical data from MD-TDDFT calculation using 1000 equilibrium geometries. Roughly five groups of data can be seen here. It would be erroneous however to assume each excited state energy calculated from one geometry comes from the same excitation mode.

To explain procedures and methods used to analyzed the results, MD-TD-DFT result using aniline in 128 water molecules using BP86 functionals was used to as an example case. MD's system potential energy as a function of time was plotted here in Figure 12a. Starting solute-solvent geometry at 0 picoseconds (ps) was obtained from Packmol program. Within 20 ps, system energy dropped 1%. Water sphere radius also appeared to have shrunk from its starting geometry. The decrease in solvent radius was expected. Hydrogen bonding and other intermolecular interactions stabilize the system causing system energy to lower to its equilibrium value. After the first 2 ps, system energy decreased, taking time to slowly equilibrate. In this case, 34.5 ps worth of MD was calculated took 9 days on 32-cpu and 10 Gb connection. Chodera's code<sup>73</sup> was used to calculate the cut-off time to be at 30 ps. The sample size should be at least one order of magnitude larger than the equilibrated one for the assumption made to estimate statistical efficiency used in the code to hold true. Assuming this cut-off point to be accurate, at least 300 ps of MD data would be obtained before deciding to choose use 30 ps as a cut-off point with confident. Projecting this forward, it would take 3 months to calculate MD run with time span larger enough to determine the equilibrium

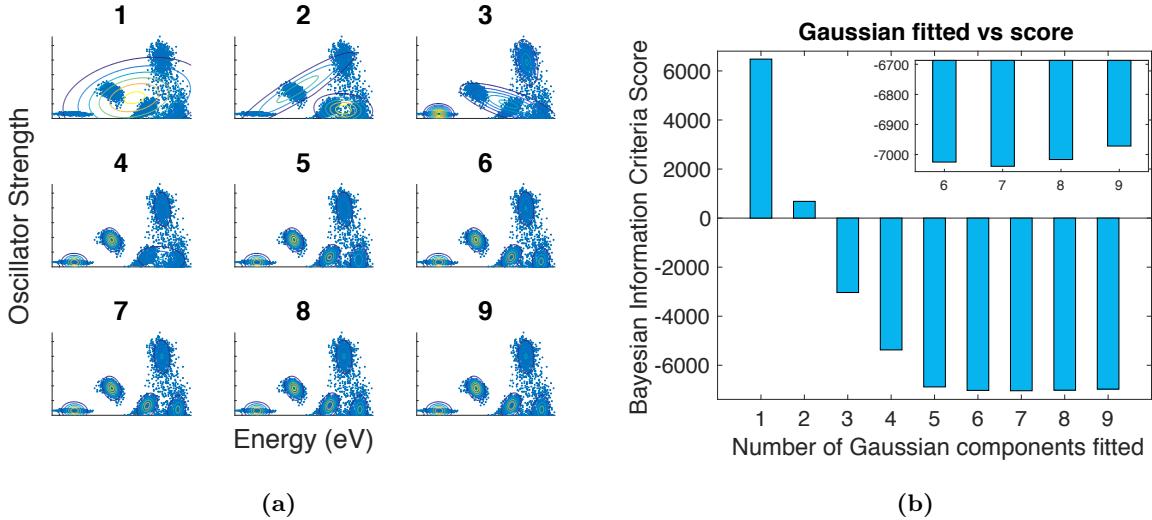


**Figure 13:** Slope plot of aniline with 128 water BP86 using different equilibrium period. a) 10 ps. b) 3 ps. Notice the higher in slope fluctuation in (b) compared to (a). This is due to smaller sampling size

using Choreda's method. Instead of using Chodera's code, the equilibrium was determined by plotting linear regression slope of system's potential energy over time with an arbitrarily set interval of 10 ps (equilibration period). The consistently low fluctuation indicated the start of equilibrium at 18 ps. Once the system's energy fluctuation reached zero, that interval should be considered at equilibrium. Slope of linear fit for each 10 ps interval starting at time on x-axis was displayed on Figure 13a. When equilibration period was changed to 3 ps (Figure 13b), energy slope reached zero earlier compared to 10-ps. The new equilibrium period yields equilibrium to be established at 6 ps instead of 20 ps. However, reducing equilibration period created larger energy slope fluctuation due to smaller sample size. Systematic criteria of measuring the equilibrium should still be properly explored and studied.

After MD equilibrium was found in the 1.78 to 2.78 ps region, system geometries of that region were used to calculate excited states using TD-DFT. Five excited energy states and its oscillator strengths from the ground state were calculated and shown in Figure 12b. Since pairs of excited state energy and oscillator strength data are collected over a range of equilibrated starting geometries, statistical analysis of raw theoretical data is required. A Matlab code was developed by Soren Eustis using Bayesian probability to fit Gaussian distribution to raw energy-oscillator-strength pair from Figure 12b to Figure 14a. The number of Gaussian plots fitted ranges from 1-9 in each sub-figure. The accepted number of plots was determined and the excited state energies and oscillator strengths given by Bayesian Information Criteria score. The lower the score of a sub-figure, the better fit. According to Figure 14b, 7 Gaussian components fitted to theoretical data was the best. Results of this fit is shown in the next section.

Experimental absorption spectra needed to be deconvoluted into a mixture of Gaussian functions for



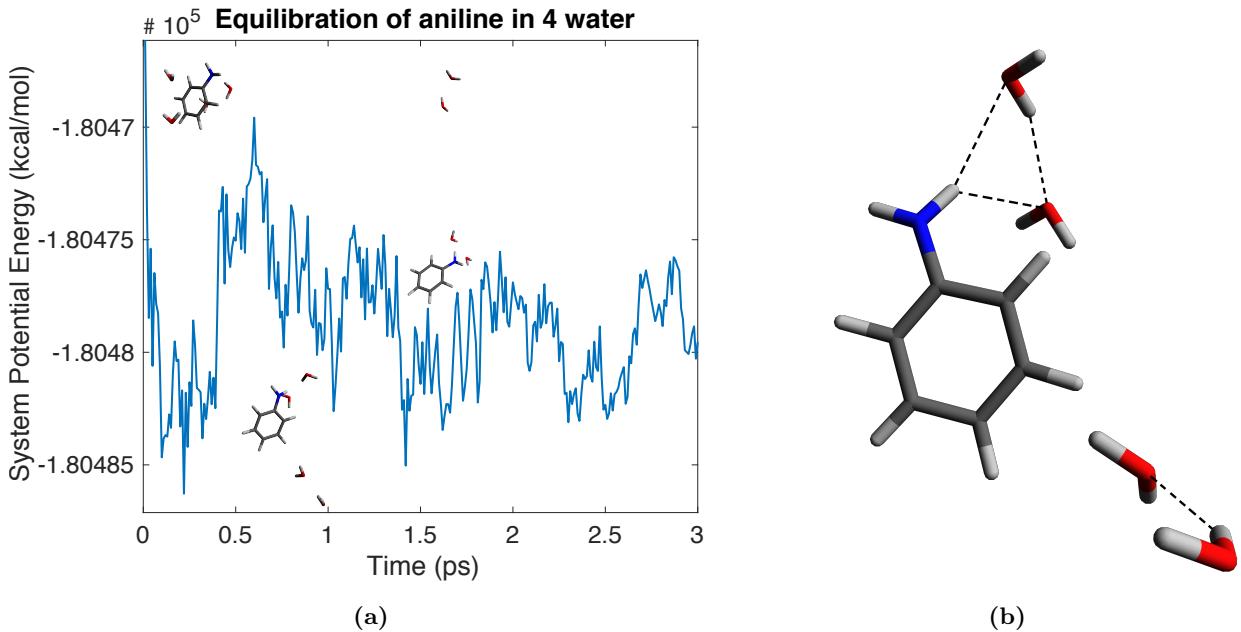
**Figure 14:** Results from Bayesian analysis using a Matlab code developed by Soren Eustis. a) deconvoluted theoretical excited state energies ( $E$ , eV) and oscillator strength ( $f$ , unitless) with varying number of Gaussian components from 1 to 9. Notice even after 5 apparent clusters were plotted, more Gaussian functions were can still be added. b) Bayesian information criteria score as a function of increasing number of fitted Gaussian functions. The lower the score, the better the fit. In this particular plot, the best number of Gaussian plots is 7.

oscillator strength calculation of each individual electronic transition (see Figure 4). Subsequently, each individual function was integrated numerically over the entire energy range of the function to give the total transition intensity using Equation 16. After calculation of excited state energies and the oscillator strength, computational results were compared with values calculated from experimental UV-VIS spectrum to evaluate accuracy of the models tested. There are currently no studies using full statistical method on quantitative calculation of excited state energies of organic molecule in water. Currently, Matlab code utilizing Bayesian probability as a statistical technique to properly find the Gaussian plots are being developed by professor Soren Eustis from previous work in R<sup>84</sup> with Peter Cohen '18. This allows room for a systematic approach to develop an appropriate computational model that would allow for further understanding of aquatic pollutants. In special cases, more information can be extracted from MD energy plot and system geometries. Two example of this will be discussed in the following sections.

#### 4.1.1 Solvent Boundary Potential

At room temperature, volatile solvent molecules like water can evaporate off the solvent shell during MD calculation. In preliminary runs, MD of aniline with 4 water molecules were calculated using CAM-B3LYP functional. System potential energy shown in Figure 15a as a function of time did not follow “well-behaved” trend seen previously (Figure 12a). Instead, the energy increased to almost initial energy before eventually

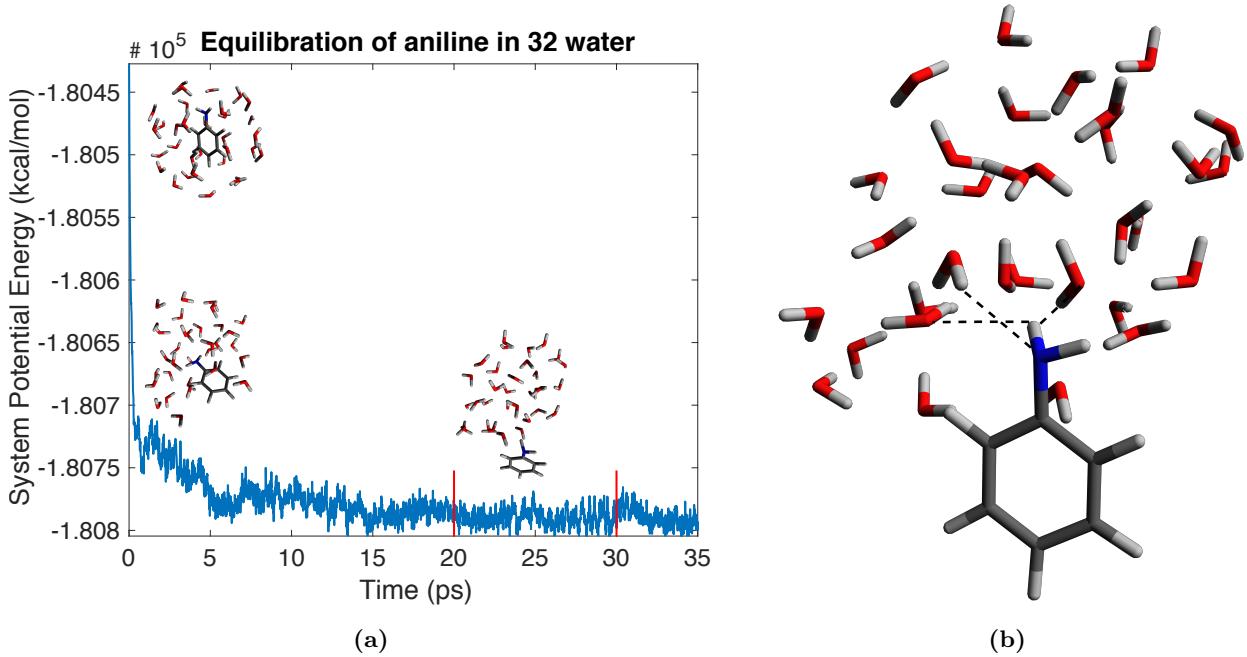
settling into equilibrium. System geometries were inspected to find problems with the MD run. After 1.5 ps into MD simulation, water molecules evaporated because the size of the solvent shell was not large enough to form a fully Hydrogen-bonded network around the solute. During the evaporation, the potential energy of the system increased momentarily, before slowly decreased as the system equilibrated. This emphasized a problem in running simulation at room temperature which provided enough potential energy for water molecules to evaporate from the water sphere. A method of using a potential boundary was proposed to



**Figure 15:** Molecular dynamics of aniline with 4 water without solvent boundary condition. a) A plot of system energy as a function of time with system geometries at 0, 0.43, and 1.5 picoseconds. Notice by 0.4 ps, the two water already started evaporating out from solute b) Hydrogen bonding between hydrogen and oxygen atoms in aniline with 4 water 0.4 ps into the MD simulation. The H-bonding ranged from 2.2 to 2.8 Angstrom

simulate the bulk solvent force field around the solute.<sup>85</sup> When a solvent fragment crosses a empirically-assigned radius, a potential “wall” will exert slight force to push the molecule back toward the solvent sphere. The Spherical Solvent Boundary Potential (SSBP) radius is set empirically from trial and error (see Appendix for `prepareMD2.py` for estimating SSBP radius). In the ab initio quantum chemistry package GAMESS<sup>72</sup>, the potential is static and does not move as the solute molecule translates. The potential boundary should not influence the structure of the system, but rather should serve as a preventive measure against evaporation. Solvent boundary condition is implemented, albeit incompletely, in GAMESS with command `SSBP=T` (see Appendix). The function for boundary potential is

$$V = 0.5 \times SFORCE \times (r_{SSBP})^2 \quad (17)$$



**Figure 16:** Molecular dynamics of aniline with 32 water without solvent boundary condition. a) A plot of system energy as a function of time with system geometries at 0, 2.5, and 20 ps. Notice the by 2.5 ps, the solute already started to get move out of the sphere b) Hydrogen bonding between hydrogen and oxygen atoms in aniline at 30 ps into the MD simulation. This showed that even though aniline was not stable being surrounded by 32 water molecules, the molecule does not completely lose contact with the solvent. More water molecules would indeed more fully solvate aniline.

Alternatively, the possibility of implementing periodic boundary conditions<sup>86</sup> should also be explored in the future. Periodic boundary conditions model an object on the other side of the box as the object moves through a wall. By limiting the size of solvent box, periodic boundary condition may allow for a more logical model compared with solvent boundary potentials. Both of these techniques are implemented in GAMESS.

#### 4.1.2 Aniline With 32 Water Molecules

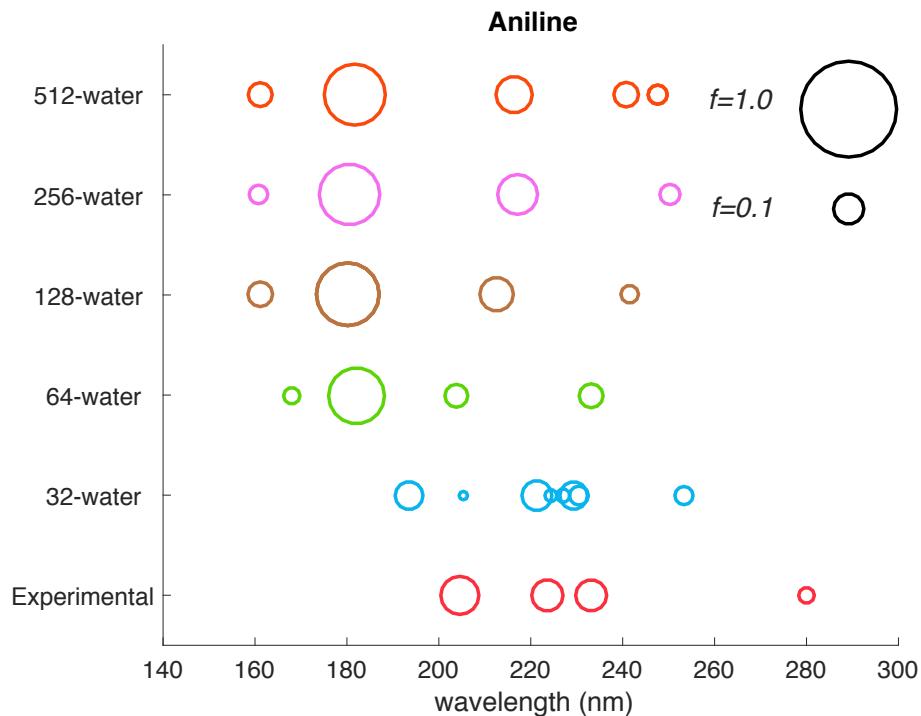
32-water model is the smallest solvent model in this study. When used with aniline, MD system energy is also not “well-behaved” (see Figure 16a). At 1ps and 7 ps, system energy increased slightly before decreased back down again. Even though these increases was not to the same degree as in 4-water case, a closer look at system geometry revealed a problem with this MD calculation. Aniline molecule appeared to be moving out of the water sphere at 2.5 ps. At equilibrium, aniline molecule could be seen pushed outside of the water cluster. This fact contradict the conclusion from Plugatyr that 32 water molecules should represent the first solvation shell for aniline.<sup>87</sup> If aniline molecule was not stable being surrounded by 32 water molecules, then the number of water molcules in the first solvation shell should be higher than 32. Thus, MD-TD-DFT of results from 32-water model be analyzed with caution due to incomplete solvation shell.

## 4.2 Determining The Optimum Solvent Environment

Experimental results and MD-TD-DFT excited states energy and oscillator strength pairs for aniline with 32, 64, 128, 256, and 512 surrounding water molecules were tabulated in Table 3. Clearer representation of the data is an area plot in Figure 17 with excitation energy (E, nm) on the x axis, and the oscillator strength (f, no unit) as an area of the data point.

**Table 3: Wavelength and oscillator strength from MD-TD-DFT calculation of aniline in 32-512 water molecules. Energy has unit of nm. Oscillator strength has no unit.**

Experimental		32-water		64-water		128-water		256-water		512-water	
E	f	E	f	E	f	E	f	E	f	E	f
204.5	0.1578	193.6	0.0842	167.9	0.0293	161.0	0.0668	160.9	0.0378	161.1	0.0627
223.6	0.1067	205.3	0.0079	182.3	0.3391	180.2	0.4280	180.5	0.3984	181.9	0.4046
233.3	0.1049	221.3	0.0975	203.7	0.0573	212.5	0.1226	217.1	0.1730	216.3	0.1437
280.0	0.0269	224.4	0.0157	233.1	0.0627	241.7	0.0343	250.1	0.0446	240.9	0.0666
		227.0	0.0195							247.4	0.0400
		229.3	0.0829								
		230.3	0.0392								
		253.2	0.0361								



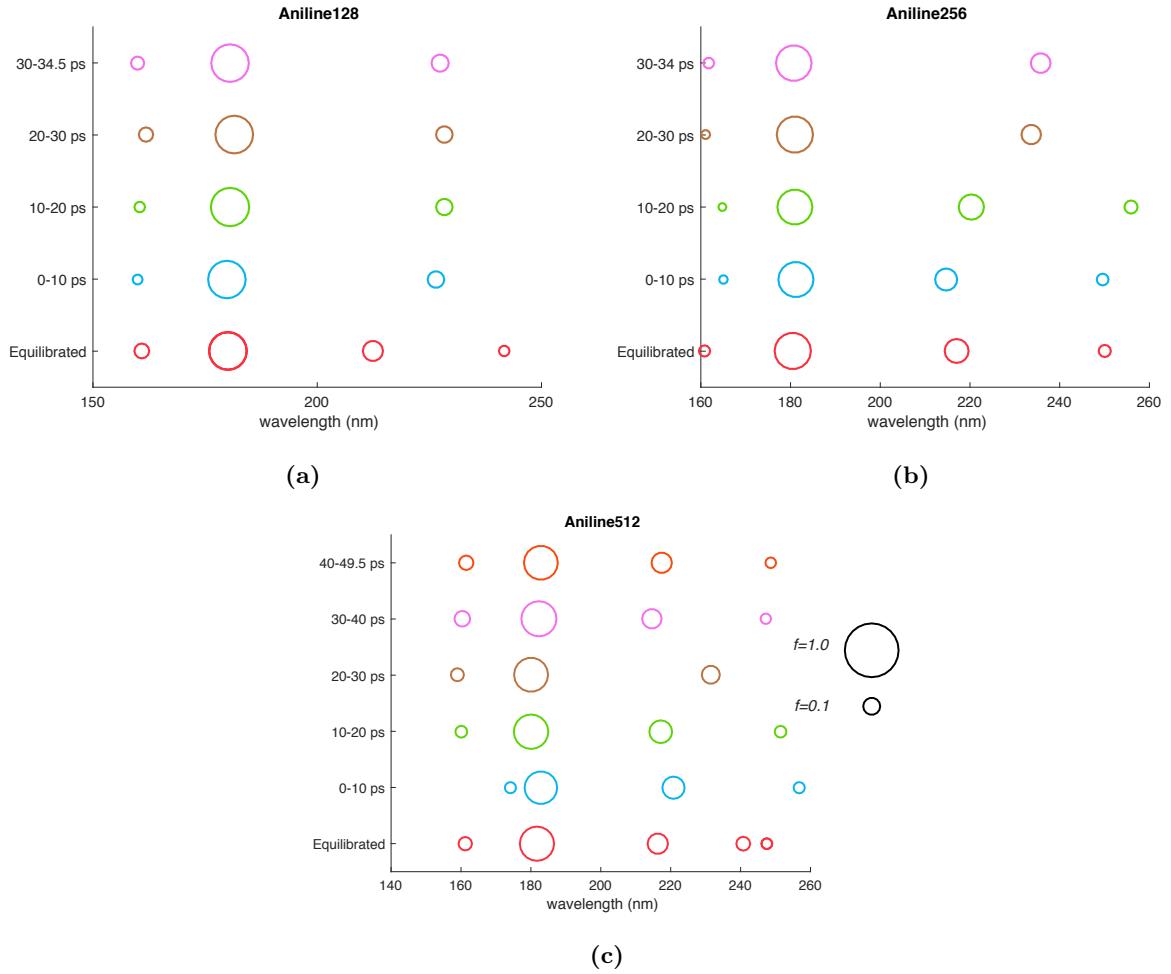
**Figure 17:** Predicted excitation energies (x-axis) and oscillator strength (area of marker) varying by the number of explicit water molecules using CAM-B3LYP functional. The two black dots on the top right shows the size of circle representing oscillator strength = 0.1 and 1. Different colors represent results from different water models.

Energy and oscillator strength values calculated from the 64-water model to 512-water model showed

convergence toward the values from 256-water model. (E, f) pairs of the absorption near 210 nm for 64 to 512-model were (203.7 nm, 0.0573), (212.5 nm, 0.1226), (217.1 nm, 0.1730), (216.3 nm, 0.1437). 256-model being the most redshifted and the most probable transition. As the number of solvent molecules increases, excited state energy from the models converged to 217 nm. Convergence in other electronic transitions in the circle plot could also be observed. Theoretical results started to converge beginning from 64-water to the values in 256 and 512-water model. 128, 256, and 512-water models gave similar results with 256-water being the most redshifted. Note that the two data points with similar values of E and f could be the same mode of excitation separated by faulty statistical analysis. For example, while other models only had one electronic transition at 240 nm, 512-model yielded two at (240.9 nm, 0.0666) and (247.4 nm, 0.0400). Matlab might put two Gaussian fits in 512-model but one in other models according to Bayesian statistics. This potential problem was acknowledged and will be explored in the next versions of Matlab code. The lowest energy transition in the 128-water model was still converging to the value of 256 and 512. It was noted that compared among theoretical data, 32-water model failed to correlate with the trend observed in models with more water molecules. 32-water model's downfall could be explained by a quick inspection at the system equilibrated geometries as discussed above. When compared with results from the experiment, theoretical results from CAM-B3LYP with any number of water molecules were mostly inaccurate. Experimental data ranged from 200 nm to 300 nm. All models missed to predict one particular transition of aniline at 280 nm. This could be because of CAM-B3LYP's failure at predicting a low-energy transition state. This problem will be discussed in the next Section.

More results in addition to the (E, f) calculated from the equilibrated region needed to be examined to determine the whether 128, 256, or 512 water molecules best model solute-solvent effects without too much cost. More excited state energies and oscillator strengths were calculated using a range of 128 to 512-model geometries obtained from MD. The results are then grouped by 10 ps time interval starting at 0 ps. This would show how quick the excited state energies and oscillator strengths value converge to the values calculated from geometries at the equilibrium. The values were calculated from geometries with interval of 0.5 ps in range of 10 ps. Tabulated data can be found in Table 4 in the Appendix. In Figure 18b, starting form 20-30 ps, the two transitions at 220 and 250 nm disappeared. This was speculated to be from Matlab's statistical method. When two transitions are close, depending on the sample, Matlab might fit one or two Gaussian components in. Excited state energy and oscillator strength values seemed to converge long before the currently determined equilibrium. In Figure 18a, the results from 128-water-model aniline were stable through out the time span and did not seem to converge as time progressed. This indicated that the equilibrium was established very early. On the other hand, 256-water-model aniline results did show

convergence in low-energy transitions. However, the apparent convergence could be a defect from Bayesian statistical analysis. In contrast, results from 512-model converged more slowly compared to 128 and 256-model. Note the fluctuation in data from 0-10 ps to 40-50 ps. Unlike, the previous two cases, 512-model had more inconsistency. Furthermore, MD calculation of 512-model took significantly longer time compared to 256-model. Combining this to the results from Figure 17, 256-water model was chosen to apply to other molecules.

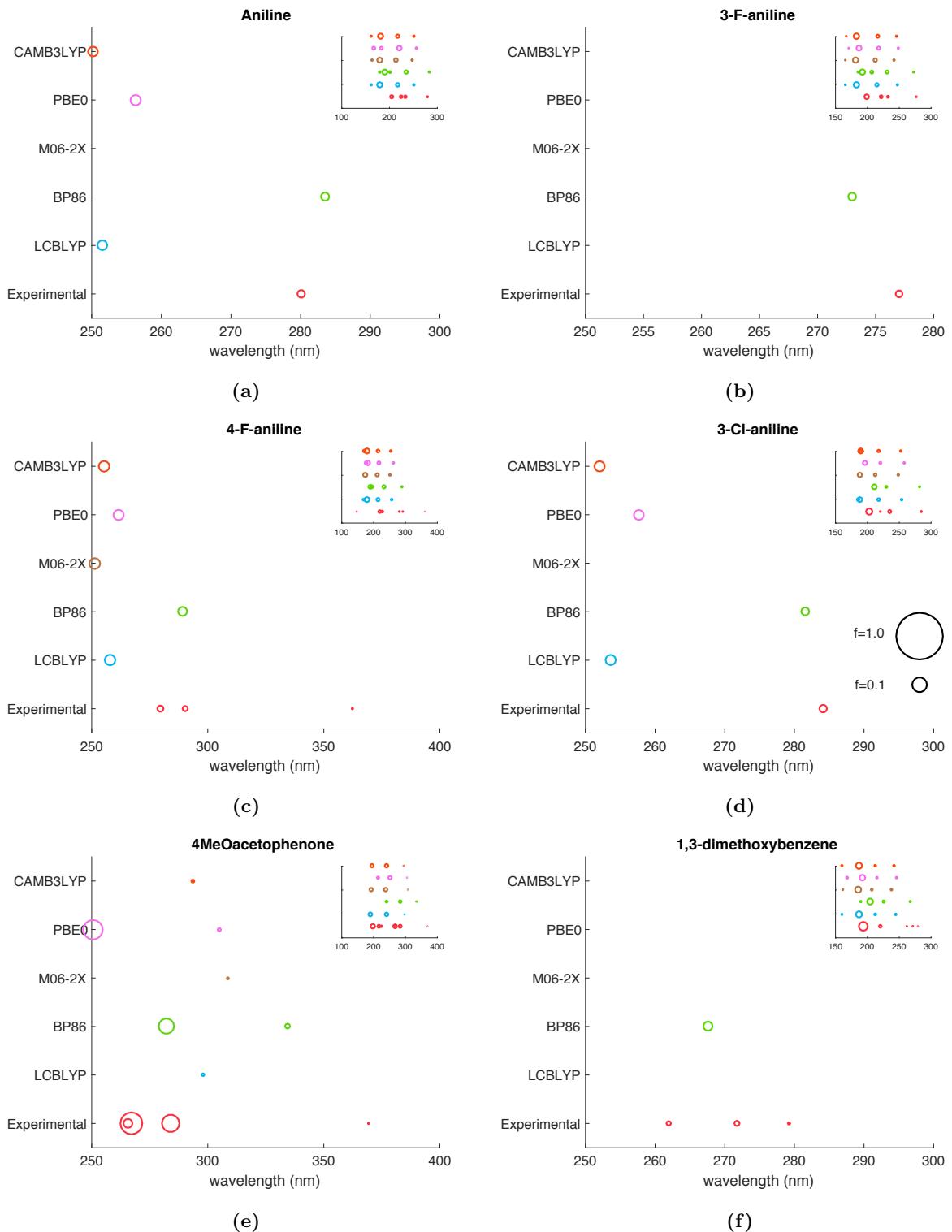


**Figure 18:** Predicted excitation energies and oscillator strengths varying by time (y-axis) for aniline in 128, 256, 512-water model

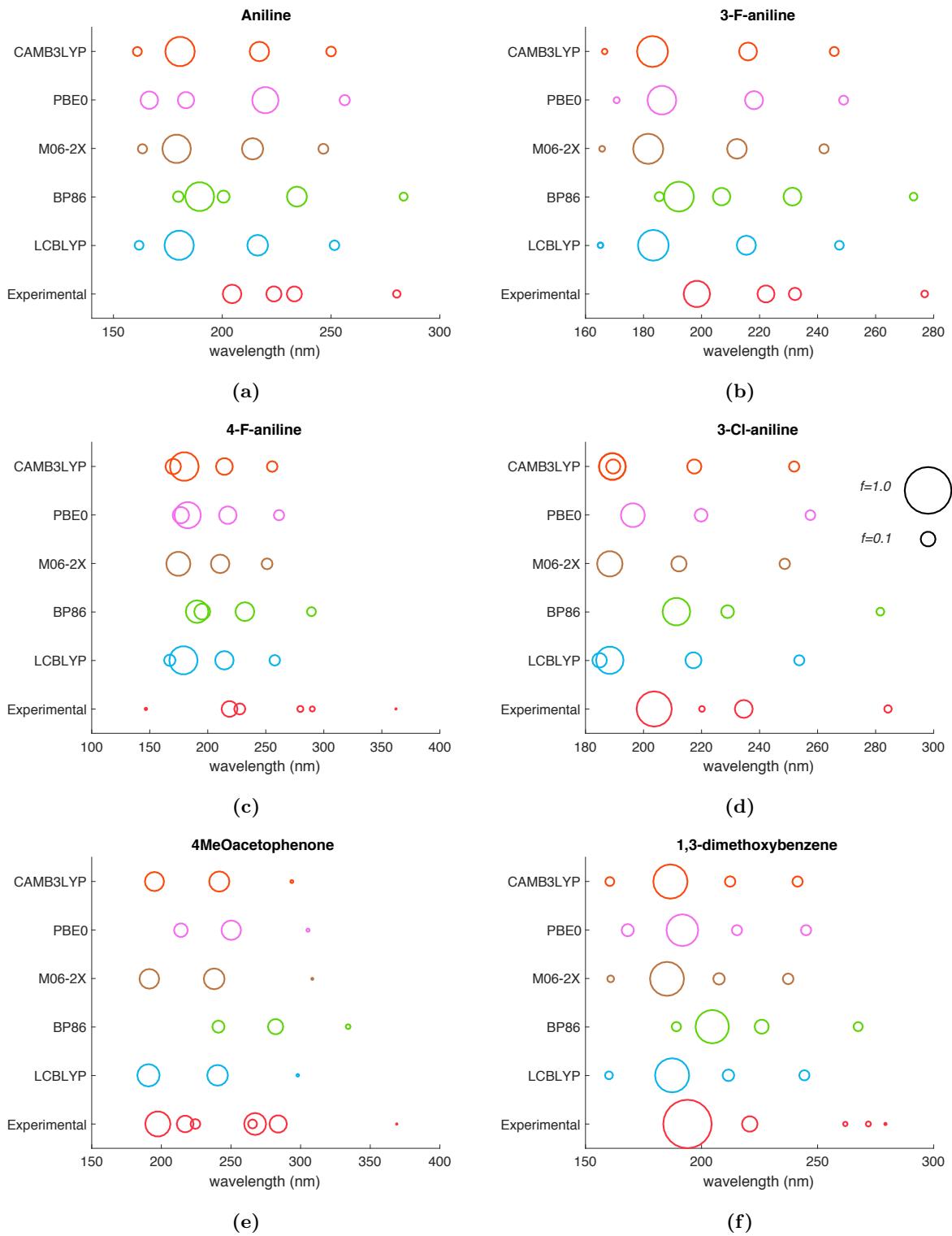
### 4.3 Comparing Different Functionals

MD-TD-DFT excited states of aniline, 3-F-aniline, 4-F-aniline, 3-Cl-aniline, 4-MeOacetophenone, and (1,3)-dimethoxybenzophenone using LCBLYP, BP86, M06-2X, PBE0, and CAM-B3LYP functionals were calculated and compared with the experimental results. In the environment, the main source of photon is the

sun. The absorption spectrum overlapping with the solar irradiance plot showed that low-energy electronic transitions above 250 nm are more relevant to modeling the light absorption of pollutants in aquatic environment (Figure 5). The circle plots in Figure 19 then had a domain starting from 250 nm to focus more on important electronic transition in nature. Table of the same data can be found in 5 in the Appendix. Most surprisingly, all functionals except BP86 missed low-energy transitions. Despite some small errors, BP86 was the best-performing functional. In Figure 19, BP86 in green calculated low energy transitions that other functionals failed to calculate. For aniline, for example, CAM-B3LYP, PBE0, and LCBLYP all predicted their lowest-energy transition at 250 nm, in contrast to BP86 and experimental data which predicted a transition at 280 nm. BP86’s lowest-energy excited state had the error in energy of 3.5 nm and error in oscillator strength of 0.005. In 3-F-aniline case, there was no other functional other than BP86 that predicted low energy transition at 275 nm with the error of 4 nm and 0.007. For 4-MeOacetophenone, BP86’s the second lowest excited state had the error of 0.5 nm and 0.03. BP86 is unique in this group for being the only pure GGA-type in all functionals considered. There is no other specific reason why BP86 might perform better than other functionals. BP86 s data in full range area plot, however, showed a consistent redshift in excited state energy calculation relative to other functionals. There were two possible explanation to this. It could mean BP86 performed badly compared to other functionals chosen here because it underestimated the excited energies and that all MD-TD-DFT excited states are blueshifted significantly. It could also mean all other hybrid functionals are inaccurate compared to pure GGA BP86 and the experimental results. If the former, improvement on current model such as EFP2<sup>88</sup> and dispersion correction<sup>89</sup> can reduce blueshift error by more correctly treat solvent-solute interaction. If the latter, more pure GGA functionals such as BPBE<sup>60,90</sup> and BLYP<sup>60,63</sup> should be explored.



**Figure 19:** Predicted excitation state energies and oscillator strength varying by DFT functionals starting from 250 nm with full Figure in a small box on the top right. A full Figure and full data can be seen in appendix.



**Figure 20:** Predicted excitation state energies and oscillator strength varying by DFT functionals.

## 5 Conclusion

An ideal model for calculating UV-VIS absorption spectrum has not yet been found, but a protocol and method of study to find it was proposed (see flowchart in Figure 10). In pioneering a project, development of method of work and system is crucial to the future study in the lab. For the years to come, Python scripts, Matlab codes and protocol will be used to continue searching for an ideal model to model photo-absorption part of the sunlight-dirven pollutant degradation in aquatic environments. For the first time, a systematic statistical approach was used to compare experimental molar absorptivity with theoretical MD-TD-DFT oscillator strength. Despite some issues with Matlab Bayesian criteria, a possibility has been demonstrated that a sensible set of data can be calculated. An approach to find the optimum number of explicit solvent molecules was also outlined. For small molecules like aniline, it was shown that 256 water molecules serve as a good trade-off between computational cost and accuracy. MD equilibrium determination was still a problem in the current protocol. MD-TD-DFT results demonstrated that the proposed method of equilibrium determination using `findEquilibrium.py` with 10 picoseconds as equilibration period is incorrect. There is then a need for statistical method to determine the equilibrium.

Even though no DFT functional was shown to have give results matching exactly with the experimental results, BP86 was identified as the best-performing functional for MD-TD-DFT calculation on aniline, 3-F-aniline, 4-F- aniline, 3-Cl-aniline, 4-MeOacetophenone, and (1,3)-dimethoxybenzophenone. More pure GGA and hybrid functionals should be explored and benchmarked with experimental results to observe more trends with different group of functionals. Larger molecules perhaps some well-studied micropollutants should be used for MD-TD-DFT calculation using the protocol proposed.

Proceeding forward, these initial results can be used to further move closer to the overall goal of developing fully computational model to predict behavior of pollutant in aquatic environment. When excited states energy and probability of excitation are accurately predicted, rate constant of photon absorption can be calculated, and combined with other data, used to predict the rate of photo-reactions.

## Acknowledgment

I would like to express my deepest appreciation professor Soren Eustis for being such a best mentor and helping me on almost all aspects of working within and without the lab. I wish to acknowledge the contribution provided by Alex Poblete, Holly Rudel, and Peter Cohen who have worked in Eustis lab and provide experimental data and mathematical insight to the project. I would like to offer my special thanks to Dj Merrill for supervising Bowdoin HPC and helping with troubleshooting throughout the years. I would like to thank James Stacy Coles / Littlefield Summer Research Fellowships, Grua/OConnell Fund, XSEDE project funded by NSF award number ACI 10-53575 for all the funding that comes to this project.

## References

- (1) Schwarzenbach, R. . P.; Escher, B. I.; Fenner, K.; Hofstetter, T. B.; Johnson, C. A.; von Gunten, U.; Wehrli, B. The Challenge of Micropollutants in Aquatic Systems. *Science* **2006**, *313*, 1072–1077.
- (2) Monteiro, S.; Boxall, A. *Reviews of Environmental Contamination and Toxicology SE - 2*; Reviews of Environmental Contamination and Toxicology; Springer New York, 2010; Vol. 202; pp 53–154.
- (3) Harada, M. Minamata disease: methylmercury poisoning in Japan caused by environmental pollution. *Critical reviews in toxicology* **1995**, *25*, 1–24.
- (4) Original image copyright belongs to noodledoodle and zhaolifang from vecteezy.com: non-exclusive right to use Eezy files to create new, derivative, unique works.
- (5) Conley, J. M.; Symes, S. J.; Schorr, M. S.; Richards, S. M. Spatial and temporal analysis of pharmaceutical concentrations in the upper Tennessee River basin. *Chemosphere* **2008**, *73*, 1178–1187.
- (6) Daneshvar, A.; Svanfelt, J.; Kronberg, L.; Weyhenmeyer, G. A. Winter accumulation of acidic pharmaceuticals in a Swedish river. *Environmental Science and Pollution Research* **2010**, *17*, 908–916.
- (7) Bonvin, F.; Rutler, R.; Chavre, N.; Halder, J.; Kohn, T. Spatial and temporal presence of a wastewater-derived micropollutant plume in Lake Geneva. *Environmental Science and Technology* **2011**, *45*, 4702–4709.
- (8) Carlson, J. C.; Stefan, M. I.; Parnis, J. M.; Metcalfe, C. D. Direct UV photolysis of selected pharmaceuticals, personal care products and endocrine disruptors in aqueous solution. *Water Research* **2015**, *84*, 350–361.

- (9) Kriegman, S.; Eustis, S. N.; Arnold, W. a.; McNeill, K. Experimental and theoretical insights into the involvement of radicals in triclosan phototransformation. *Environmental science & technology* **2013**, *47*, 6756–63.
- (10) Bedoux, G.; Roig, B.; Thomas, O.; Dupont, V.; Le Bot, B. Occurrence and toxicity of antimicrobial triclosan and by-products in the environment. *Environmental Science and Pollution Research* **2012**, *19*, 1044–1065.
- (11) Morss, P. C. The Photodegradation of Endocrine Disrupting Compounds in Aqueous Solutions. Ph.D. thesis, Bowdoin College, 2014.
- (12) Salvato-Vallverdu, G. The Perrin - Jablonski diagram. 2009.
- (13) Hirayama, S.; Steer, R. P. Understanding the Theory and Practice of Molecular Spectroscopy: The Effects of Spectral Bandwith. *Journal of Chemical Education* **2010**, *87*, 1344–1347.
- (14) Reichardt, C.; Welton, T. *Solvents and Solvent Effects in Organic Chemistry*; Wiley, 2011.
- (15) Turro, N. J.; Ramamurthy, V.; Scaiano, J. C. *Modern Molecular Photochemistry of Organic Molecules*; University Science Books, 2010.
- (16) Clemens, O.; Basters, M.; Wild, M.; Wilbrand, S.; Reichert, C.; Bauer, M.; Springborg, M.; Jung, G. Solvent effects on the absorption/emission spectra of an organic chromophore: A theoretical study. *Journal of Molecular Structure: THEOCHEM* **2008**, *866*, 15–20.
- (17) Suendo, V.; Viridi, S. Ab initio calculation of UV-Vis absorption spectra of a single chlorophyll a molecule: comparison study between RHF/CIS, TDDFT, and semi-empirical methods. *ITB Journal of Science* **2012**, *44A*, 93–112.
- (18) Kakitani, T.; Mataga, N. Comprehensive study on the role of coordinated solvent mode played in electron-transfer reactions in polar solutions. *The Journal of Physical Chemistry* **1987**, *91*, 6277–6285.
- (19) Marcus, R. A. Relation between Charge Transfer Absorption and Fluorescence Spectra and the Inverted Region. *The Journal of Physical Chemistry* **1989**, *93*, 3078–3086.
- (20) Matyushov, D. V. Solvent reorganization energy of electron-transfer reactions in polar solvents. *Journal of Chemical Physics* **2004**, *120*, 7532–7556.

- (21) Barthel, E. R.; Martini, I. B.; Schwartz, B. J. How does the solvent control electron transfer? Experimental and theoretical studies of the simplest charge transfer reaction. *Journal of Physical Chemistry B* **2001**, *105*, 12230–12241.
- (22) Eilmes, A. Solvatochromic probe in molecular solvents: implicit versus explicit solvent model. *Theoretical Chemistry Accounts* **2014**, *133*, 1–13.
- (23) Ricke, N. D. Application of the Landau-Zener Model and Fermi's Golden Rule to Estimate Triplet Quantum Yield for Organic Molecules. Ph.D. thesis, Bowdoin College, 2014.
- (24) Leifer, A. *The Kinetics of Environmental Aquatic Photochemistry: Theory and Practice*; ACS professional reference book; American Chemical Society, 1988.
- (25) Schwarzenbach, R. P.; Gschwend, P. M.; Imboden, D. M. *Environmental Organic Chemistry*; Wiley, 2005.
- (26) ASTM G173-03(2012), Standard Tables for Reference Solar Spectral Irradiances: Direct Normal and Hemispherical on 37 Tilted Surface, ASTM International, West Conshohocken, PA, 2012, [www.astm.org](http://www.astm.org).
- (27) Dykstra, C. E.; Frenking, G.; Kim, K. S.; Scuseria, G. E. *Theory and Applications of Computational Chemistry: The First Forty Years*; Elsevier, 2011.
- (28) Bartlett, R. J.; Purvis, G. D. Many-body perturbation theory, coupled-pair many-electron theory, and the importance of quadruple excitations for the correlation problem. *International Journal of Quantum Chemistry* **1978**, *14*, 561–581.
- (29) Pople, J.; Krishnan, R.; Schlegel, H. B.; Binkley, J. S. Electron correlation theories and their application to the study of simple reaction potential surfaces. *International Journal of Quantum Chemistry* **1978**, *XIV*, 545.
- (30) Frisch, M. J., et al. (2009). Gaussian 09. Wallingford, CT, USA, Gaussian, Inc.
- (31) Eisberg, R.; Resnick, R. *Quantum physics of atoms, molecules, solids, nuclei, and particles*, 2nd ed.; John Wiley & Sons: New York, 1985.
- (32) Griffiths, D. J. *Introduction to Quantum Mechanics*; Pearson Education India, 2005.

- (33) Møller, C.; Plesset, M. S. Note on an approximation treatment for many-electron systems. *Physical Review* **1934**, *46*, 618–622.
- (34) Hehre, W. J.; Radom, L.; Schleyer, P. v. R.; Pople, J. *AB INITIO Molecular Orbital Theory*; A Wiley-Interscience publication; Wiley, 1986.
- (35) Magalhães, A. L. Gaussian-type orbitals versus slater-type orbitals: A comparison. *Journal of Chemical Education* **2014**, *91*, 2124–2127.
- (36) Hehre, W. J.; Stewart, R. F.; Pople, J. A. SelfConsistent MolecularOrbital Methods. I. Use of Gaussian Expansions of SlaterType Atomic Orbitals. *The Journal of Chemical Physics* **1970**, *52*, 2769.
- (37) Cramer, C. J. *Essentials of Computational Chemistry: Theories and Models*, 2nd ed.; Wiley, 2005.
- (38) Rohrdanz, M. A.; Herbert, J. M. Simultaneous benchmarking of ground- and excited-state properties with long-range-corrected density functional theory. *Journal of Chemical Physics* **2008**, *129*, 1–9.
- (39) Barnes, L.; Abdul-Al, S.; Allouche, A.-R. TDDFT Assessment of Functionals for Optical 00 Transitions in Small Radicals. *The Journal of Physical Chemistry A* **2014**, *118*, 11033–11046.
- (40) Wiberg, K. B. Basis set effects on calculated geometries: 6-311++G\*\* vs. aug-cc-pVDZ. *Journal of Computational Chemistry* **2004**, *25*.
- (41) Feller, D. The Role of Databases in Support of Computational Chemistry Calculations. *Journal of Computational Chemistry* **1996**, *17*, 1571–1586.
- (42) Baerends, E. J.; Gritsenko, O. V. A quantum chemical view of density functional theory. *Journal of Physical Chemistry A* **1997**, *101*, 5383–5403.
- (43) Sousa, S. F.; Fernandes, P. A.; Ramos, M. J. General performance of density functionals. *Journal of Physical Chemistry A* **2007**, *111*, 10439–10452.
- (44) Hohenberg, P.; Kohn, W. Inhomogeneous Electron Gas. *Physical Review* **1964**, *136*, B864–B871.
- (45) Kohn, W.; Sham, L. J. Self-Consistent Equations Including Exchange and Correlation Effects. *Physical Review* **1965**, *140*, A1133–A1138.
- (46) Heßelmann, A. Molecular excitation energies from time-dependent density functional theory employing random-phase approximation Hessians with exact exchange. *Journal of Chemical Theory and Computation* **2015**, *11*, 1607–1620.

- (47) Hirao, H. Which DFT functional performs well in the calculation of methylcobalamin? Comparison of the B3LYP and BP86 functionals and evaluation of the impact of empirical dispersion correction. *Journal of Physical Chemistry A* **2011**, *115*, 9308–9313.
- (48) Parac, M.; Grimme, S. Comparison of multireference Møller-Plesset theory and time-dependent methods for the calculation of vertical excitation energies of molecules. *Journal of Physical Chemistry A* **2002**, *106*, 6844–6850.
- (49) Jacquemin, D.; Wathelet, V.; Perpète, E. a.; Adamo, C. Extensive TD-DFT benchmark: Singlet-excited states of organic molecules. *Journal of Chemical Theory and Computation* **2009**, *5*, 2420–2435.
- (50) Caricato, M.; Trucks, G. W.; Frisch, M. J.; Wiberg, K. B. Electronic transition energies: A study of the performance of a large range of single reference density functional and wave function methods on valence and rydberg states compared to experiment. *Journal of Chemical Theory and Computation* **2010**, *6*, 370–383.
- (51) Bersier, C.; Proetto, C.; Sanna, A.; Cangi, A.; Mirhosseini, S. H. Basic rules for functional construction at finite temperature. Ph.D. thesis, Max Planck Institute of Microstructure Physics.
- (52) Yanai, T.; Tew, D. P.; Handy, N. C. A new hybrid exchange correlation functional using the Coulomb-attenuating method (CAM-B3LYP). *Chemical Physics Letters* **2004**, *393*, 51–57.
- (53) Zhao, Y.; Truhlar, D. G. The M06 suite of density functionals for main group thermochemistry , thermochemical kinetics , noncovalent interactions , excited states , and transition elements : two new functionals and systematic testing of four M06-class functionals and 12 other fun. *Theoretical Chemistry Accounts* **2007**, *120*, 215–241.
- (54) Adamo, C.; Barone, V. Toward reliable density functional methods without adjustable parameters: The PBE0 model. *The Journal of Chemical Physics* **1999**, *110*, 6158.
- (55) Magyar, R. J.; Tretiak, S. Dependence of spurious charge-transfer excited states on orbital exchange in TDDFT: Large molecules and clusters. *Journal of Chemical Theory and Computation* **2007**, *3*, 976–987.
- (56) Runge, E.; Gross, E. K. U. Density-functional theory for time-dependent systems. *Physical Review Letters* **1984**, *52*, 997–1000.

- (57) Caricato, M.; Mennucci, B.; Tomasi, J. Solvent Effects on the Electronic Spectra: An Extension of the Polarizable Continuum Model to the ZINDO Method. *The Journal of Physical Chemistry A* **2004**, *108*, 6248–6256.
- (58) Fabian, J. Electronic excitation of sulfur-organic compounds - Performance of time-dependent density functional theory. *Theoretical Chemistry Accounts* **2001**, *106*, 199–217.
- (59) Barone, V.; Polimeno, A. Integrated computational strategies for UV/vis spectra of large molecules in solution. *Chemical Society Reviews* **2007**, *36*, 1724.
- (60) Becke, A. D. Density-functional exchange-energy approximation with correct asymptotic behavior. *Physical Review A* **1988**, *38*, 3098–3100.
- (61) Perdew, J. P. Density-functional approximation for the correlation energy of the inhomogeneous electron gas. *Physical Review B* **1986**, *33*, 8822–8824.
- (62) Iikura, H.; Tsuneda, T.; Yanai, T.; Hirao, K. A long-range correction scheme for generalized-gradient-approximation exchange functionals. *Journal of Chemical Physics* **2001**, *115*, 3540–3544.
- (63) Lee, C.; Yang, W.; Parr, R. G. Development of the Colle-Salvetti correlation-energy formula into a functional of the electron density. *Physical Review B* **1988**, *37*, 785–789.
- (64) Lin, H.; Truhlar, D. G. QM/MM: what have we learned, where are we, and where do we go from here? *Theoretical Chemistry Accounts* **2007**, *117*, 185–199.
- (65) Cossi, M.; Barone, V. Solvent effect on vertical electronic transitions by the polarizable continuum model. *The Journal of Chemical Physics* **2000**, *112*, 2427.
- (66) Mennucci, B. Continuum solvation models: What else can we learn from them? *Journal of Physical Chemistry Letters* **2010**, *1*, 1666–1674.
- (67) Tomasi, J.; Mennucci, B.; Cammi, R. Quantum mechanical continuum solvation models. *Chemical Reviews* **2005**, *105*, 2999–3093.
- (68) Li, J.; Cramer, C. J.; Truhlar, D. G. Two-Response-Time Model Based on CM2/INDO/S2 Electrostatic Potentials for the Dielectric Polarization Component of Solvatochromic Shifts on Vertical Excitation Energies. *International Journal of Quantum Chemistry* **1999**, *77*, 264–280.

- (69) Day, P. N.; Jensen, J. H.; Gordon, M. S.; Webb, S. P.; Stevens, W. J.; Krauss, M.; Garmer, D.; Basch, H.; Cohen, D. An effective fragment method for modeling solvent effects in quantum mechanical calculations. *The Journal of Chemical Physics* **1996**, *105*, 1968–1986.
- (70) Yoo, S.; Zahariev, F.; Sok, S.; Gordon, M. S. Solvent effects on optical properties of molecules: A combined time-dependent density functional theory/effective fragment potential approach. *Journal of Chemical Physics* **2008**, *129*, 1–8.
- (71) Defusco, A.; Minezawa, N.; Slipchenko, L. V.; Zahariev, F.; Gordon, M. S. Modeling Solvent Effects on Electronic Excited States. *Journal of Physical Chemistry Letters* **2011**, *2*, 2184–2192.
- (72) Schmidt, M. W.; Baldridge, K. K.; Boatz, J. A.; Elbert, S. T.; Gordon, M. S.; Jensen, J. H.; Koseki, S.; Matsunaga, N.; Nguyen, K. A.; Shyjun, S. U. et al. General Atomic and Molecular Electronic Structure System. *Journal of Computational Chemistry* **1993**, *14*, 1347–1363.
- (73) Chodera, J. D. A simple method for automated equilibration detection in molecular simulations. *Journal of Chemical Theory and Computation* **2015**, 021659.
- (74) Hilborn, R. C. Einstein coefficients, cross sections, f values, dipole moments, and all that. *American Journal of Physics* **1982**, *50*, 982–986.
- (75) Klán, P.; Wirz, J. *Photochemistry of organic compounds: From concepts to practice*; John Wiley & Sons, 2009.
- (76) Python Software Foundation. Python Language Reference, version 2.7.
- (77) Martínez, L.; Andrade, R.; Birgin, E. G.; Martínez, J. M. Software News and Update Packmol : A Package for Building Initial Configurations. *Journal of Computational Chemistry* **2009**, *30*, 2157–2164.
- (78) Matlab version 7.14.0.739. Natick, Massachusetts, The MathWorks Inc.
- (79) Schmidt, J.R.; Polik, W.F. WebMO Enterprise, version 13.0; WebMO LLC: Holland, MI, USA, 2013; <http://www.webmo.net> (accessed March, 2013).
- (80) Wojdyr, M. Fityk: A general-purpose peak fitting program. *Journal of Applied Crystallography* **2010**, *43*, 1126–1128.

- (81) Eustis, S. N.; Trerayapiwat, K.; Ricke, N.; Cohen, P.; Poblete, A.; Rodel, H. Computational Environmental Photochemistry: A Roadmap for Applied Quantitative Calculations. *Royal Chemical Society* **2016**,
- (82) Armarego, W. L. F.; Chai, C. L. L. *Purification of laboratory chemicals*, 7th ed.; Butterworth-Heinemann, 2013.
- (83) Yalkowsky, S. H.; He, Y.; Jain, P. *Handbook of Aqueous Solubility Data, Second Edition*; CRC Press, 2016.
- (84) R Development Core Team, R: A Language and Environment for Statistical Computing, 2008.
- (85) Beglov, D.; Roux, B. Finite representation of an infinite bulk system: Solvent boundary potential for computer simulations. *The Journal of Chemical Physics* **1994**, *100*.
- (86) Makov, G.; Payne, M. Periodic boundary conditions in ab initio calculations. *Physical Review B* **1995**, *51*, 4014–4022.
- (87) Plugatyr, A.; Svishchev, I. M. The hydration of aniline: Analysis of spatial distribution functions. *The Journal of chemical physics* **2009**, *130*, 114509.
- (88) Gordon, M. S.; Freitag, M. a.; Bandyopadhyay, P.; Jensen, J. H.; Kairys, V.; Stevens, W. J. The Effective Fragment Potential Method : A QM-Based MM Approach to Modeling. *Journal of Physical Chemistry A* **2001**, *105*, 293–307.
- (89) Guidez, E. B.; Gordon, M. S. Dispersion Correction Derived from First Principles for Density Functional Theory and HartreeFock Theory. *The Journal of chemical physics* **2015**, *119*, 21612168.
- (90) Perdew, J. P.; Burke, K.; Ernzerhof, M. Generalized Gradient Approximation Made Simple. *Physical Review Letters* **1996**, *77*, 3865–3868.

## A Appendix

### A.1 MD-TD-DFT Result Tables

**Table 4: Wavelength and Oscillator Strength from MD-TD-DFT calculation of aniline in 128-512 water molecules as a function of time.**

<i>128-water</i>		0-10 ps		10-20 ps		20-30 ps		30-34.5 ps	
Experimental		E	f	E	f	E	f	E	f
160.98	0.0668	159.88	0.0302	160.40	0.0345	161.93	0.0624	160.03	0.0522
180.20	0.4280	179.77	0.4255	180.53	0.4452	181.42	0.4297	180.57	0.4269
212.48	0.1226	226.39	0.0827	228.26	0.0815	228.35	0.0827	227.48	0.0894
241.70	0.0343								
<i>256-water</i>		0-10 ps		10-20 ps		20-30 ps		30-34 ps	
Experimental		E	f	E	f	E	f	E	f
160.90	0.0378	165.05	0.0224	164.81	0.0190	161.04	0.0246	161.83	0.0353
180.51	0.3984	181.14	0.3726	180.85	0.3690	180.94	0.3967	180.67	0.3826
217.07	0.1730	214.77	0.1465	220.35	0.1925	233.71	0.1150	235.77	0.1177
250.13	0.0446	249.57	0.0421	255.88	0.0522				
<i>512-water</i>		0-10 ps		10-20 ps		20-30 ps		30-40 ps	
Experimental		E	f	E	f	E	f	E	f
161.09	0.0627	174.24	0.0435	160.16	0.0489	158.87	0.0592	160.27	0.0847
181.86	0.4046	182.99	0.3670	179.95	0.4100	179.93	0.3994	182.21	0.4254
216.35	0.1437	220.86	0.1708	217.26	0.1790	231.48	0.1126	214.56	0.1313
240.86	0.0666	256.85	0.0438	251.30	0.0481			247.13	0.0380
247.44	0.0400							248.74	0.0394

**Table 5: Wavelength and Oscillator Strength from MD-TD-DFT calculation of aniline, 3-F-aniline, 4-F-aniline, 3-Cl-aniline, 4-MeOacetophenone, and (1,3)-dimethoxybenzophenone using LCBLYP, BP86, M06-2X, PBE0, and CAM-B3LYP functionals.**

Experiment		LCBLYP		BP86		M06-2X		PBE0		CAM-B3LYP	
<i>Aniline</i>		E	f	E	f	E	f	E	f	E	f
204.45	0.1578	161.76	0.0370	179.76	0.0531	163.29	0.0398	166.29	0.0523	160.90	0.0378
223.58	0.1067	180.22	0.3946	189.71	0.3844	179.08	0.3725	183.53	0.3305	180.51	0.3984
233.29	0.1049	216.38	0.1974	200.74	0.0694	213.88	0.2086	219.99	0.1818	217.07	0.1730
280.05	0.0269	251.51	0.0433	234.25	0.1829	246.59	0.0461	256.27	0.0419	250.13	0.0446
				283.49	0.0314						
<i>3-F-aniline</i>		E	f	E	f	E	f	E	f	E	f
198.42	0.3185	165.18	0.0144	185.47	0.0393	165.85	0.0162	170.68	0.0178	166.57	0.0146
222.31	0.1339	165.18	0.0144	192.23	0.4090	181.63	0.4154	186.18	0.3801	183.11	0.4379
232.29	0.0707	183.41	0.4332	206.97	0.1419	212.34	0.1780	218.13	0.1506	215.98	0.1467
277.03	0.0221	215.52	0.1672	231.44	0.1467	242.06	0.0406	248.86	0.0382	245.60	0.0385
		247.57	0.0365	272.99	0.0288						
<i>4-F-aniline</i>		E	f	E	f	E	f	E	f	E	f
146.36	0.0031	167.47	0.0600	190.53	0.2346	174.70	0.2661	176.95	0.1239	170.59	0.1098
218.88	0.1156	178.78	0.3663	195.05	0.1173	210.34	0.1570	182.48	0.3175	179.51	0.3779
227.74	0.0567	214.03	0.1560	231.99	0.1613	251.19	0.0545	217.30	0.1427	214.64	0.1298
279.66	0.0177	257.96	0.0516	289.10	0.0372			261.43	0.0495	255.29	0.0523
290.31	0.0129										
362.29	0.0011										
<i>3-Cl-aniline</i>		E	f	E	f	E	f	E	f	E	f
203.56	0.5666	184.97	0.0972	211.31	0.3493	188.35	0.2962	196.40	0.2627	189.36	0.3286
220.06	0.0156	188.32	0.3447	229.08	0.0765	212.30	0.1087	219.84	0.0771	189.36	0.3286
234.55	0.1461	217.16	0.1188	281.54	0.0293	248.60	0.0512	257.60	0.0450	189.60	0.0936
284.16	0.0267	253.58	0.0482							217.53	0.0932
<i>4MeOacetophenone</i>		E	f	E	f	E	f	E	f	E	f
197.60	0.2896	190.87	0.2276	241.21	0.0681	191.48	0.1769	214.30	0.0856	194.80	0.1665
216.87	0.1283	240.65	0.1964	282.30	0.1071	237.78	0.1983	250.42	0.1739	241.78	0.1865
224.69	0.0437	298.01	0.0034	334.28	0.0112	308.53	0.0023	305.07	0.0050	293.62	0.0044
265.81	0.0368										
267.15	0.2202										
283.86	0.1390										
369.32	0.0011										
<i>1,3-dimethoxybenzene</i>		E	f	E	f	E	f	E	f	E	f
193.84	1.0740	159.93	0.0293	189.16	0.0403	160.78	0.0217	168.15	0.0695	160.56	0.0371
220.80	0.1107	187.38	0.5333	204.48	0.5099	184.98	0.5289	191.91	0.4615	186.65	0.5361
261.93	0.0096	211.64	0.0624	225.89	0.0915	207.53	0.0620	215.21	0.0482	212.29	0.0506
271.73	0.0127	244.33	0.0484	267.60	0.0384	237.43	0.0527	245.18	0.0486	241.47	0.0492
279.19	0.0029					611.96	0.0081				

## A.2 Python Scripts

### A.2.1 Preparing MD Input Files

This script does two things. First (line 35-84), it calculates appropriate radius for solvent boundary potential without empirically fitting it. A simple model is proposed: at most solute will rotate around its outmost solute atom. In the code, the distance between the outmost solute atom to the solute's CG is called solute radius. The distance between the outmost solvent atom to the solute's CG is solvent radius. These two radius plus an extra 2-3 Angstrom gives SSBP radius for MD input file. Second (line 87-155), the script parses xyz file's geometry data into MD input file. Slight format change from xyz file type is required for GAMESS input files, but this python code automate that change. The output file is MD file which can be run on GAMESS. Output of this script can be seen below in MD Input File section.

```
1 ##### Create inp for MD run from xyz file from packmol #####
2 ### ssbp radius is also calculated roughly from #####
3 ### starting geometry #####
4 #####
5 #####
6
7 import sys
8 import csv
9 import os
10 import string
11
12 #for asking what the input in terminal should be
13 try:
14     if str(sys.argv[1])=='?':
15         print '\nCall function as: prepareMD.py input.xyz number_ofSoluteAtoms
16             ↪   number_ofSolventAtoms number_ofSolventMolecules \n'
17         sys.exit()
18 except IndexError:
19     print '\n!!!Input command Error. Call function as: prepareMD.py input.xyz number_ofSoluteAtoms
20             ↪   number_ofSolventAtoms number_ofSolventMolecules \n'
21     sys.exit()
22 #for assigning received input from terminal
23 try:
24     input=str(sys.argv[1])
25     number_ofSoluteAtoms=int(sys.argv[2])
26     number_ofSolventAtoms=int(sys.argv[3])
27     number_ofSolventMolecules=int(sys.argv[4])
28 except IndexError:
29     print '\n!!!Input command Error. Call function as: prepareMD.py input.xyz number_ofSoluteAtoms
30             ↪   number_ofSolventAtoms number_ofSolventMolecules \n'
31     sys.exit()
32 #generate output name
33 if input.endswith('.xyz'):
34     output = input[:-4]+'.inp'
35 #for safety - at worst the output will not overwrite the input
36 else:
37     output=input+'.inp'

#part one
```

```

38 #This part is for finding ssbp radius for inout file
39 #enumerate gets data in line - line and line index - n
40 radiusInSolute=0.0
41 radiusInSolvent=0.0
42 avgX=0.0
43 avgY=0.0
44 avgZ=0.0
45 X=[]
46 Y=[]
47 Z=[]
48
49 lineNumber=0
50 #open input
51 f2=open(input)
52 for line in f2:
53     lineNumber+=1
54     #first two line does not contain useful info - x y z start on the third line
55     if lineNumber>2:
56         #x y z
57         lineSplit=line.split()
58         X.append(float(lineSplit[1]))
59         Y.append(float(lineSplit[2]))
60         Z.append(float(lineSplit[3]))
61     #for looping through array below
62 size=len(X)
63 #find a CG for solute atoms
64 avgX=sum(X[:numberofSoluteAtoms-1])/numberofSoluteAtoms
65 avgY=sum(Y[:numberofSoluteAtoms-1])/numberofSoluteAtoms
66 avgZ=sum(Z[:numberofSoluteAtoms-1])/numberofSoluteAtoms
67 #looping to find radius of each atoms in relative to solute's CG
68 #also find the maximum value of them
69 for i in range(0,size):
70     d=((X[i]-avgX)**2+(Y[i]-avgY)**2+(Z[i]-avgZ)**2)**0.5
71     if i<numberofSoluteAtoms:
72         if radiusInSolute<d:
73             radiusInSolute=d
74         else:
75             if radiusInSolvent<d:
76                 radiusInSolvent=d
77
78 #radius should be a little bit larger than the two combined - 3 Angstrom larger - this does not
    ↵ need to be super accurate
79 radiusInSolute=radiusInSolute
80 radiusInSolvent=radiusInSolvent
81 ssbpRadius=radiusInSolute+radiusInSolvent+3
82 print '\n'
83 print 'Radius in solute is:\t'+str(radiusInSolute)
84 print 'Radius in solvent is:\t'+str(radiusInSolvent)
85 print 'ssbp Radius should be:\t'+str(ssbpRadius)
86 print '\n'
87 ##########
88
89 #Part two - this is where geometry data is taken from xyz, change into GAMESS input's format +
    ↵ other input
90 numberofAllSolventsAtoms=numberofSoventAtoms*numberofSolventMolecules
91 fragmentNumber=1;
92 atomLabel=1
93
94 #this dict is for generating atomic number from Acronym

```

```

95 atomicNumber={'LV': 116.0, 'BE': 4.0, 'FR': 87.0, 'BA': 56.0, 'BH': 107.0, 'BI': 83.0, 'BK': 97.0,
←   'EU': 63.0, 'FE': 26.0, 'BR': 35.0, 'ES': 99.0, 'FL': 114.0, 'FM': 100.0, 'RG': 111.0, 'RU':
←   44.0, 'NO': 102.0, 'NA': 11.0, 'NB': 41.0, 'ND': 60.0, 'NE': 10.0, 'RE': 75.0, 'RF': 104.0,
←   'LU': 71.0, 'RA': 88.0, 'RB': 37.0, 'NP': 93.0, 'RN': 86.0, 'RH': 45.0, 'B': 5.0, 'CO': 27.0,
←   'TH': 90.0, 'CM': 96.0, 'CL': 17.0, 'H': 1.0, 'CA': 20.0, 'CF': 98.0, 'CE': 58.0, 'N': 7.0,
←   'CN': 112.0, 'P': 15.0, 'GE': 32.0, 'GD': 64.0, 'GA': 31.0, 'V': 23.0, 'CS': 55.0, 'CR': 24.0,
←   'DS': 110.0, 'CU': 29.0, 'SR': 38.0, 'UUP': 115.0, 'UUS': 117.0, 'TC': 43.0, 'KR': 36.0, 'SI':
←   14.0, 'SN': 50.0, 'SM': 62.0, 'UUT': 113.0, 'SC': 21.0, 'SB': 51.0, 'TA': 73.0, 'OS': 76.0,
←   'PU': 94.0, 'SE': 34.0, 'AC': 89.0, 'HS': 108.0, 'YB': 70.0, 'DB': 105.0, 'C': 6.0, 'HO':
←   67.0, 'DY': 66.0, 'HF': 72.0, 'HG': 80.0, 'HE': 2.0, 'PR': 59.0, 'PT': 78.0, 'LA': 57.0, 'F':
←   9.0, 'UUO': 118.0, 'LI': 3.0, 'PB': 82.0, 'TL': 81.0, 'TM': 69.0, 'LR': 103.0, 'PD': 46.0,
←   'TI': 22.0, 'TE': 52.0, 'TB': 65.0, 'PO': 84.0, 'PM': 61.0, 'ZN': 30.0, 'AG': 47.0, 'NI':
←   28.0, 'I': 53.0, 'K': 19.0, 'IR': 77.0, 'AM': 95.0, 'AL': 13.0, 'O': 8.0, 'S': 16.0, 'AR':
←   18.0, 'AU': 79.0, 'AT': 85.0, 'W': 74.0, 'IN': 49.0, 'Y': 39.0, 'CD': 48.0, 'ZR': 40.0, 'ER':
←   68.0, 'MD': 101.0, 'MG': 12.0, 'PA': 91.0, 'SG': 106.0, 'MO': 42.0, 'MN': 25.0, 'AS': 33.0,
←   'MT': 109.0, 'U': 92.0, 'XE': 54.0}
96
97 #write out put the headers - all the commands for GAMESS + ssbp
98 #functional = M06-2X - DFTTYP=M06-2X
99 f = open(output, 'w');
100 f.write(''' $CONTRL SCFTYP=RHF RUNTYP=MD COORD=UNIQUE
101     DFTTYP=CAMB3LYP MAXIT=200 ICHARG=0 MULT=1 $END
102     $MD KEVERY=10 PROD=.T. NVTNH=2 MBT=.T. MBR=.T.
103         BATHT=298 RSTEMP=.T. DTEMP=25 NSTEPS=50000
104         SSBP=.T. SFORCE=1.0 DROFF='''+str(ssbpRadius)+''' $END
105     $DFT DC=.F. $END
106     $SYSTEM MWORDS=1000 MEMDDI=1000 $END
107     $SCF DIRSCF=.T. $END
108     $BASIS GBASIS=N31 NGAUSS=6 NDFUNC=2 NPFUNC=1
109     DIFFS=.TRUE. POLAR=POPN31 $END
110     $DATA\n'''+'MD INPUT for '+input+'\nC1 1\n''')
111
112 #geometry
113 with open(input) as f1:
114     #read by line
115     #readlines if okay to use bc xyz is not too big
116     lines = f1.readlines()
117     #enumerate gets data in line - line and line index - n
118     for n, line in enumerate(lines):
119         #take all solute molecules (in range of 2 (line 3 where packmol starts) to num+2)
120         #it's num+2 bc the range will go to num+1
121         if n == 2:
122             print 'Now Writing Solute:\n'
123         if n in range(2,numberofSoluteAtoms+2):
124             lineSplit=line.split();
125             lineSplit.insert(1,str(atomicNumber[lineSplit[0]]))
126             #convert coordinates to 10 decimals (add zeros if need be)
127             for index in [2,3,4]:
128                 lineSplit[index]=float(lineSplit[index])
129                 lineSplit[index]=format(lineSplit[index], '.10f')
130                 grandString=lineSplit[0]+'\t'+lineSplit[1]+'\t'+lineSplit[2] +
131                     '\t'+lineSplit[3]+'\t'+lineSplit[4]+'\n';
132             f.write(grandString)
133             print grandString
134         if n == numberofSoluteAtoms+2:
135             f.write(' $END\n\n $EFRAG\nCOORD=CART POSITION=OPTIMIZE\n')
136             print 'Now Writing Solvent:\n'
137             #now start doing solvent - (need to add fragment number and atom labels)
startPointOfSolvent=numberofSoluteAtoms+2

```

```

138     if n in range(startPointOfSolvent, startPointOfSolvent+number0fAllSolventsAtoms+1):
139         #atomlabel = O1, H2, H3 from O, H, H
140         if atomLabel%numberofSoventAtoms==1:
141             grandString='FRAGNAME=H2ODFT ! '+str(fragmentNumber)+'\n'
142             f.write(grandString)
143             print grandString
144             fragmentNumber+=1;
145             atomLabel%=numberofSoventAtoms
146             lineSplit=line.split();
147             lineSplit.insert(1,str(atomLabel))
148             atomLabel+=1
149             #convert coordinates to 10 decimals (add zeros if need be)
150             for index in [2,3,4]:
151                 lineSplit[index]=float(lineSplit[index])
152                 lineSplit[index]=format(lineSplit[index],'.10f')
153                 grandString=' '+lineSplit[0]+lineSplit[1]+'\t'+lineSplit[2] +
154                 ' '\t+lineSplit[3]+'\t'+lineSplit[4]+\n';
155                 f.write(grandString)
156                 print grandString
157             #close the inp with $END
158             f.write(' $END\n')
159 #####

```

### A.2.2 MD Geometries Extraction

One of the reasons, an MD run might fail is if solute molecule is pushed out of the water sphere. 3dExtract4.py allows geometries of the system at different time to be extracted from a large size log file into a xyz-movie file. xyz files, capable of containing more than one frame of geometries, allows one to follow MD through a combination of screenshot (each frame is 10 femtoseconds - in the current MD input file - see MD Input File section).

```

1 #####
2 ### 3dExtract pulls out geometries from MD run and make ###
3 ### an xyz-movie file for inspection MD progress #####
4 #####
5
6 import os as os
7 import sys
8
9 #for asking what the input in terminal should be
10 try:
11     if str(sys.argv[1])=='?':
12         print '\nCall function as: 3dExtract.py input.log number0fSoluteAtoms
13             → numberofSolventAtoms number0fSolventMolecules \n'
14         sys.exit()
15     except IndexError:
16         print '\n!!!Input command Error. Call function as: 3dExtract.py input.log number0fSoluteAtoms
17             → numberofSolventAtoms number0fSolventMolecules \n'
18         sys.exit()
19
20 #Call as 3dExtract.py inputFile #ofsoluteAtom #ofsolventAtom #ofsoluteMolecules
21 try:
22     input=str(sys.argv[1])
23     numberofSoluteAtoms=int(sys.argv[2])

```

```

22     numberofSoventAtoms=int(sys.argv[3])
23     numberofSolventMolecules=int(sys.argv[4])
24 except IndexError:
25     print '\n!!!Input command Error. Call function as: 3dExtract.py input.log numberofSoluteAtoms
26     →   numberofSolventAtoms numberofSolventMolecules \n'
27     sys.exit()
28 if input.endswith('.log'):
29     output = str(input[:-4]) + '.xyz'
30 else:
31     output=str(input)
32
33 numberOfAllSolventsAtoms=numberofSoventAtoms*numberofSolventMolecules
34 #This is for comparing files to be written
35 previousGrandString=''
36 collectionStarted=False
37 FoundFirstGeometry=False
38 foundTime=False
39 foundPE=False
40 time=' '
41 energy='0'
42 grandString=''
43 lineBwTimeAndEnergy=0
44 lineBwTimeAndGeometry=0
45 #1 is for cartesian line (useless), then 1 in 3(n+1) is for fragment H2O line (also useless)
46 numberOfLinesToBeCollected=numberofSoluteAtoms+1+numberOfSolventMolecules*(numberofSoventAtoms+1)
47
48 #number of molecules so far
49 timeCount=0
50 #total number of atoms (solute + solvent) - used later in checking if file is complete
51 atomCount=0
52 #define functions here
53 lineSinceTimeIsFound=0;
54 #do an input of solvent, solute atoms
55 molList=[]
56 #for finding PE
57 lineCountFromTime=0
58 #for printing time
59 def printAndReturnTime (thisLine):
60     lineComponents=thisLine.split();
61     timeString=str(lineComponents[3]);
62     #get rid off .00
63     timeString=timeString[:-3]
64     print "Analyzing t = "+timeString+" fsec\n"
65     return timeString
66
67 #to determine if line should be collected -
68 def shouldCollect():
69     #only check if collection is in progress - if it is, then continue to finish collecting the
70     →   lines
71     #collectionStarted is determined when 'QM ATOM COORDINATES (ANG)' is found
72     if collectionStarted:
73         #from first solute atom to the last fragment atom
74         if (atomCount>=0 and atomCount<numberOfLinesToBeCollected):
75             return True;
76         else:
77             return False;
78     else:
79         return False;

```

```

79 # only write when atomCount==numberOfLinesToBeCollected
80
81 def shouldWrite():
82     #only check if collection is in progress
83     if collectionStarted:
84         #solute
85         if (atomCount==numberOfLinesToBeCollected):
86             return True;
87         else:
88             return False;
89     else:
90         return False;
91
92 def shouldCollectTime(line):
93     #only check if collection is in progress
94     if ' *** AT T=' in line:
95         return True;
96     else:
97         return False;
98
99 def shouldCollectPE(line,reference,currentLine):
100    #check if line bw time and energy is known - this is written as reference
101    if reference>0:
102        if currentLine==reference:
103            return True
104        else:
105            return False
106    #if reference is not known, then it needs to be found by finding string POT EN...
107    elif reference==0:
108        if '    POT ENERGY' in line:
109            reference=int(currentLine)
110            return True;
111        else:
112            return False;
113
114 def shouldStartCollectGeometry(line,reference,currentLine):
115    #check if line bw time and geometry is known - reference
116    if reference>0:
117        if currentLine==reference:
118            return True
119        else:
120            return False
121    #if not then it needs to be found by searching for the string TEMPER...
122    elif reference==0:
123        if ' QM ATOM ' in line:
124            reference=int(currentLine)
125            return True;
126        else:
127            return False;
128
129 #clear output.xyz
130 f = open(output, 'w');
131 f.write('')
132 #open input
133 f1=open(input)
134 #enumerate gets data in line - line and line index - n
135 #readlines() is eliminated because it creates a huge array and python cannot handle it when log
136    → file get very large
#using for line in... alleviate the burden on memory and actually speed up the process

```

```

137  for line in f1:
138      #this keyword is usually before coordinate
139      #find out if we have found the first geometry
140      if not FoundFirstGeometry:
141          #if not then find it
142          if shouldCollect():
143              #split line
144              lineSplit=line.split()
145              atomCount+=1;
146              #append to molList
147              molList.append(lineSplit)
148              time=str(0)
149              if shouldWrite():
150                  atomCount=atomCount-(numberOfSolventMolecules+1);
151                  #for loop through a ***COPY*** of molList and delete some element from molList!
152                  #if you don't realize six asterisk then you should go back up - we do this so we
153                  #→ can remove element along the way without messing up the index
154                  for line in list(molList):
155                      #if line has 4 elements then it's a coordinate from solvent fragment - we have
156                      #→ to drop number behind atom - O1 to O
157                      if len (line) == 4:
158                          #store string
159                          oldString = line[0]
160                          #replacement string
161                          newString=','
162                          #loop to check if it's a alphabet or not
163                          for character in range(len(oldString)):
164                              #do substring of 1 character
165                              subString = oldString[character:character+1]
166                              #check if it's an alphabet - yes? then add to newString
167                              if subString.isalpha():
168                                  newString = newString + subString
169                                  #replace 'O1' with 'O'
170                                  line[0]=newString
171                                  #if it's 5 then it's solute coordinate - we have to get rid of atomic number
172                                  #→ behind atomic representation
173                                  elif len(line) == 5:
174                                      # 'N 7.0 ...' will become 'N ...'
175                                      del line[1]
176                                      #the rest are crap - just remove it out of the line
177                                      else:
178                                          #there's a reason why this is remove - not del - since we are iterating if
179                                          #→ we delete using index we are gonna be screwed
180                                          molList.remove(line)
181                                          #this is for if we have an incomplete file or inconsistant number of atoms we
182                                          #→ should only use the one before and break for loop without appending to
183                                          #→ grandString
184                                          if len(molList) != atomCount:
185                                              print 'error'
186                                              #xyz file has a format that we need atomCount at the top followed by snapshot
187                                              #→ number(timeCount) on the next line before adding any coordinates
188                                              grandString=grandString+str(atomCount)+'\n'+str(timeCount)+'\tfs\t'+str(energy)+'
189                                              #→ KCAL/MOL\n'
190                                              f.write(grandString)
191                                              #loop tho molList to add data - molList = [['N','1','1','1'], ['C',...], ... ] And
192                                              #→ element = ['N', '1', '1', '1']
193                                              for element in molList:
194                                                  #loop through element in molList data = 'N', '1', '1', '1'
195                                                  for data in element:

```

```

187             #add to grandString and don't forget tab, return
188             f.write(data+'\t')
189             #end one element with a return
190             f.write('\n')
191             #reset all values after writing
192             atomCount=0;
193             molList=[]
194             collectionStarted=False;
195             FoundFirstGeometry=True;
196             timeCount+=1
197             if ' QM ATOM COORDINATES (ANG)' in line:
198                 collectionStarted=True
199
200             if FoundFirstGeometry:
201                 if not foundTime:
202                     if shouldCollectTime(line):
203                         time=printAndReturnTime(line);
204                         foundTime=True;
205                         if time=='0.00':
206                             foundTime=False
207                             grandString=''
208
209             elif not foundPE:
210                 if shouldCollectPE(line,lineBwTimeAndEnergy,lineCountFromTime):
211                     #split line using space - sample( POT ENERGY      =      -1.804578585E+05
212                     #                           ← KCAL/MOL)
213                     #this will be split to [..., '=', '-1.804578585E+05'...] -time = element 4
214                     lineComponents=line.split();
215                     energy=str(lineComponents[3])
216                     foundPE=True
217
218             elif shouldCollect():
219                 #split line
220                 lineSplit=line.split()
221                 atomCount+=1;
222                 #append to molList
223                 molList.append(lineSplit)
224                 if shouldWrite():
225                     atomCount=atomCount-(number_of_SolventMolecules+1);
226                     #for loop through a ***COPY*** of molList and delete some element from molList!
227                     #if you don't realize six asterisk then you should go back up - we do this so we
228                     #   can remove element along the way without messing up the index
229                     for line in list(molList):
230                         #if line has 4 elements then it's a coordinate from solvent fragment - we have
231                         #   to drop number behind atom - O1 to O
232                         if len(line) == 4:
233                             #store string
234                             oldString = line[0]
235                             #replacement string
236                             newString=''
237                             #loop to check if it's a alphabet or not
238                             for character in range(len(oldString)):
239                                 #do substring of 1 character
240                                 subString = oldString[character:character+1]
241                                 #check if it's an alphabet - yes? then add to newString
242                                 if subString.isalpha():
243                                     newString = newString + subString
244                                     #replace 'O1' with 'O'
245                                     line[0]=newString
246                                     #if it's 5 then it's solute coordinate - we have to get rid of atomic number
247                                     #   behind atomic representation

```

```

242     elif len(line) == 5:
243         # 'N 7.0 ...' will become 'N ...'
244         del line[1]
245         #the rest are crap - just remove it out of the line
246     else:
247         #there's a reason why this is remove - not del - since we are iterating if
248         #→ we delete using index we are gonna be screwed
249         molList.remove(line)
250         #this is for if we have an incomplete file or inconsistant number of atoms we
251         #→ should only use the one before and break for loop without appending to
252         #→ grandString
253         if len(molList) != atomCount:
254             print 'error'
255             #xyz file has a format that we need atomCount at the top followed by snapshot
256             #→ number(timeCount) on the next line before adding any coordinates
257             grandString=grandString+str(atomCount)+'\n'+str(time)+'\tfs\t'+str(energy)+'
258             #→ KCAL/MOL\n'
259             f.write(grandString)
260             #loop tho molList to add data - molList = [['N', '1', '1', '1'], ['C', ...], ... ] And
261             #→ element = ['N', '1', '1', '1']
262             for element in molList:
263                 #loop through element in molList data = 'N', '1', '1', '1'
264                 for data in element:
265                     #add to grandString and don't forget tab, return
266                     f.write(data+'\t')
267                     #end one element with a return
268                     f.write('\n')
269                     #reset all values after writing
270                     atomCount=0;
271                     molList=[]
272                     collectionStarted=False;
273                     foundTime=False
274                     foundPE=False
275                     lineCountFromTime=0;
276                     grandString=''
277                     timeCount+=1
278                     elif shouldStartCollectGeometry(line,lineBwTimeAndGeometry,lineCountFromTime):
279                         collectionStarted=True
280                         lineCountFromTime+=1
281                     f.close()
282                     #sanity check
283                     print 'Done. Extract ' + str(timeCount) + ' snapshots total.'

```

### A.2.3 Plot Potential Energy Of MD Run

plotEnergyMD6.py script extracts potential energy and temperature of each MD frame to determine the if the system is close to equilibration. This script and 3dExtract are very essential to the first stage of the project: they determine whether MD has failed or reached equilibrium based on the geometry and potential energy of the system. Like 3dExtract, many versions of this code has been developed and modified and they are the most refined pieces of code for their purpose. In the future when the codes are unified, improvement can be made on plotting the plot on Matlab instead of matplotlib.

```

1 ##########
2 ### Use this to plot energy vs time to see if MD has #####
3 ### run its course. Generates: csv of PE and #####
4 ### temperature vs time, pdf of the plot #####
5 ##########
6
7 import matplotlib
8 matplotlib.use('Agg')
9 import matplotlib.pyplot as plt
10 import csv
11 import sys
12 import string
13
14 #call as plotEnergyMD3.py inputfile
15 #for asking what the input in terminal should be
16 try:
17     if str(sys.argv[1])=='?':
18         print '\nCall function as: plotEnergyMD.py input.log \n'
19         sys.exit()
20     else:
21         input=str(sys.argv[1])
22 except IndexError:
23     print '\n!!!Input command Error. Call function as: plotEnergyMD.py input.log \n'
24     sys.exit()
25 output=str(input) + '_energies.csv'
26
27 #initiate variables
28 lineBwTimeAndEnergy=0;
29 lineBwTimeAndTemp=0;
30 lineCountFromTime=0;
31 collectionStarted=False
32 foundTime=False
33 foundPE=False
34 foundTemp=False
35 firstTime=True
36 grandString=' '
37 time='0'
38
39 #####
40 #functions
41
42 #check if time is in line
43 def shouldCollectTime(line):
44     #only check if collection is in progress
45     if ' *** AT T=' in line:
46         return True;
47     else:
48         return False;
49
50 # for printing time so one can keep track of the progress
51 def printTime (thisLine):
52     lineComponents=thisLine.split();
53     timeString=str(lineComponents[3]);
54     print "Analyzing t = "+timeString+" fsec\n"
55
56 #Are we currently looking potential energy?
57 def shouldCollectPE(line,reference,currentLine):
58     #check if line bw time and energy is known - this is written as reference
59     if reference>0:

```

```

60     if currentLine==reference:
61         return True
62     else:
63         return False
64     #if reference is not known, then it needs to be found by finding string POT EN...
65     elif reference==0:
66         if '    POT ENERGY' in line:
67             reference=int(currentLine)
68             return True;
69         else:
70             return False;
71     #Are we currently looking Temp?
72     def shouldCollectTemp(line,reference,currentLine):
73         #check if line bw time and energy is known - reference
74         if reference>0:
75             if currentLine==reference:
76                 return True
77             else:
78                 return False
79         #if not then it needs to be found by searching for the string TEMPER...
80         elif reference==0:
81             if '    TEMPER' in line:
82                 reference=int(currentLine)
83                 return True;
84             else:
85                 return False;
86
87     #once everything is found, we should write down before moving on to the next snapshot
88     def shouldWrite():
89         #only check if collection is in progress
90         if foundTime:
91             if foundPE:
92                 if foundTemp:
93                     return True;
94                 else:
95                     return False
96             else:
97                 return False;
98         else:
99             return False;
100
101 #####
102 #open csv and prepare for writing
103 f = open(output, 'w');
104 f.write('')
105
106 #open input
107 f1=open(input)
108 print 'finding patterns...'
109 #avoid using readlines() so there'll be no problem with large files
110 for line in f1:
111     #time keyword is before PE, PE is before Temp so the search should be in this order in order to
112     #→ be most efficient
113     #find out if collection is needed
114     if not foundTime:
115         if shouldCollectTime(line):
116             printTime(line)
117             #split line
118             lineComponents=line.split();

```

```

118     #append time (split) to string
119     #split line using space - sample(** AT T=          10.00 FSEC, THIS RUN'S STEP NO.=
120     →   10)
121     #this will be split to ['***', 'AT', 'T=', '10.00...] -time = element 4
122     newTime=str(lineComponents[3])[:-3]
123     if int(time)>=int(newTime):
124         print 'there is a duplication at '+time+' and '+newTime +' ignoring new values'
125     else:
126         time=str(newTime)
127         grandString=grandString+time;
128         foundTime=True;
129
130     elif not foundPE:
131         if shouldCollectPE(line,lineBwTimeAndEnergy,lineCountFromTime):
132             #append time (split) to string
133             #split line using space - sample(      POT ENERGY      =      -1.804578585E+05 KCAL/MOL)
134             #this will be split to [..., '=', '-1.804578585E+05...] -time = element 4
135             lineComponents=line.split();
136             grandString=grandString+', '+str(lineComponents[3])
137             foundPE=True
138
139     elif not foundTemp:
140         if not firstTime:
141             if shouldCollectTemp(line,lineBwTimeAndTemp,lineCountFromTime):
142                 #append time (split) to string
143                 #split line using space - sample(      TEMPERATURE      =      349.98666547 K)
144                 #this will be split to [..., '=', '-349.98666547...] -time = element 3
145                 lineComponents=line.split();
146                 grandString=grandString+', '+str(lineComponents[2])+'\n';
147                 foundTemp=True
148
149             #problem with this is - the first snapshot's temperature is not given in the log file
150             #set Temp to 0 to indicate the beginning
151
152         if firstTime:
153             grandString=grandString+', 0\n';
154             #once append, turn off the boolean
155             firstTime=False
156             foundTemp=True
157
158         #write after all data is collected for one snapshot
159         if shouldWrite():
160             with open(output, 'a') as f:
161                 f.write(grandString)
162
163             #reset the variables
164             lineCountFromTime=0;
165             collectionStarted=False
166             foundTime=False
167             foundPE=False
168             foundTemp=False
169             grandString=''
170             lineCountFromTime+=1;
171
172         #finish writing csv
173         f.close()
174
175     #plot
176     #pull out CSV
177     #use csv.reader bc csv has ',' and this automate the formatting
178     f = csv.reader(open(output))
179     #convert column to array using zip (a built in function)
180     Time, Energy, Temp = zip(*f)
181     #convert string to float
182     Time = map(float, Time)
183     Energy = map(float, Energy)

```

```

176 Temp = map(float, Temp)
177 minEnergy=min(Energy)
178 relativeEnergy=[]
180 for eachEnergy in Energy:
181     relativeEnergy.append(eachEnergy-minEnergy)
182
183 #plot
184 x = Time
185 y1 = Energy
186 y2 = Temp
187
188 plt.subplot(2, 1, 1)
189 plt.plot(x, y1, 'g-')
190 plt.title('MD PE and T plot of: ' + str(input))
191 plt.ylabel('Potential Energy (KCal/mol)')
192
193 plt.subplot(2, 1, 2)
194 plt.plot(x, y2, 'b-')
195 plt.xlabel('time (fs)')
196 plt.ylabel('System Temperature (K)')
197 plt.savefig(str(input) + '_EnergyPlot.pdf', format='pdf')
198
199 plt.gcf().clear()
200 y1 = relativeEnergy
201
202 plt.subplot(2, 1, 1)
203 plt.plot(x, y1, 'g-')
204 plt.title('MD relative PE and T plot of: ' + str(input))
205 plt.ylabel('Relative Potential Energy (KCal/mol)')
206
207 plt.subplot(2, 1, 2)
208 plt.plot(x, y2, 'b-')
209 plt.xlabel('time (fs)')
210 plt.ylabel('System Temperature (K)')
211 plt.savefig(str(input) + '_EnergyPlot_2.pdf', format='pdf')

```

#### A.2.4 Find The Most Equilibrated Period

There are currently no consensus as to how to determine if a system has reached the equilibrium in molecular dynamics. In previous works, plotEnergyMD (previous script) was used to indicate whether the potential energy of the system (solute and solvent) has stabilized. Arbitrariness in determining the equilibrium falls in the hands of users. findEquilibrium.py is designed to solve this subjectivity. With a list of potential energies at different time from plotEnergyMD, linear fit can be done in a fix interval to evaluate the rise or fall in energy. Currently, the limit value is taken, still empirically, from 15000 to 25000 fs interval in CAM-B3LYP aniline32.log. The lowest slope of the list is used to identify good range for excited state energy calculation. Further improvement can be done to find the bottom slope limit as a variable with molecule input.

```

1 #####
2 ### This is used to determined the equilibrium using #####
3 ### linear regression and an upper limit for the slope #####
4 ### The set range for this script is 10 ps #####

```

```

5 ##########
6
7 import matplotlib
8 matplotlib.use('Agg')
9 import matplotlib.pyplot as plt
10 import csv
11 import sys
12 import numpy as np
13
14 #for asking what the input in terminal should be
15 try:
16     if str(sys.argv[1])=='?':
17         print '\nCall function as: findEquilibrium.py input.log \n'
18         sys.exit()
19     else:
20         input=str(sys.argv[1])
21 except IndexError:
22     print '\n!!!Input command Error. Call function as: findEquilibrium.py input.log \n'
23     sys.exit()
24 output=str(input) + '_energies.csv'
25 #This portion is the same as in plotEnergyMD6
26 #pull out CSV
27 f = csv.reader(open(output))
28 #convert column to array using zip (a built in function)
29 Time, Energy, Temp = zip(*f)
30 #convert string to float
31 Time = map(float, Time)
32 Energy = map(float, Energy)
33 Temp = map(float, Temp)
34
35 #find Equilibrium using linear regression
36 #this number control the range of time to be used in energy fluctuation calculation
37 minNumberOfStep=1000
38 #This is the limit above which the script will report no equilibrium is found
39 #this is from aniline32.log - 15000 to 25000
40 maxSlope=1e-4
41 #for plotting
42 slope=[]
43 print 'Finding equilibrium using minimum number of steps = '+str(minNumberOfStep) +' and top limit
   ↪   of acceptable slope = '+str(maxSlope)
44 try:
45     #for looping
46     size=len(Time)
47     x=Time[0:minNumberOfStep]
48     y=Energy[0:minNumberOfStep]
49     #for using in loop
50     #set thio a high value - it can be any number bc we will replace it with the lowest slope value
      ↪   found in loop
51     lowestSlopeValue=1e5
52     #Same - this will be replaced
53     indexOfLowestSlope=-1
54     for i in range(0,size-1-minNumberOfStep):
55         print 'Finding equilibrium from t= '+str(Time[i])+', to '+str(Time[i+minNumberOfStep-1])
56         #poly fit is basically a linear fit - m=slope, b=y_intersect
57         m,b = np.polyfit(x, y, 1)
58         #for plotting - append to array of existing slope values
59         slope.append(m)
60         print 'slope = ' +str(m)
61         #take absolute value and see which interval does not fluctuate the least

```

```

62     if abs(m)<=abs(lowestSlopeValue):
63         lowestSlopeValue=float(m)
64         indexOfLowestSlope=int(i)
65         #this is similar to queue structure - room for improvement is to make x and y arrays into
66         #→ actual queues
67         #else remove the head and add next tail
68         del x[0]
69         x.append(Time[minNumberOfStep+i])
70         del y[0]
71         y.append(Energy[minNumberOfStep+i])
72         #if x or y does not have enough element (minNumberOfStep) then report error
73         #room for improvement - move this up top instead of having a long try
74     except IndexError:
75         print 'There is not enough data to determine the equilibrium'
76         #if lowestSlopeValue pass the top limit then report
77     if abs(lowestSlopeValue)<=abs(maxSlope):
78         report=str('Found best equilibrium starting from '+str(Time[indexOfLowestSlope])+', to
79         #→ '+str(Time[indexOfLowestSlope+minNumberOfStep-1])+' with slope =
80         #→ '+str(lowestSlopeValue))
81         print report
82     #if not then say so
83 else:
84     report = 'Equilibrium is not yet reach'
85     print report
86     print str('The current limit is at '+str(maxSlope) +' kcal/mol/fs and the lowest value of
87     #→ slope = '+str(lowestSlopeValue))
88
89 #plot slope vs time
90 x=slope
91 #align time with slope
92 y=Time[0:size-1-minNumberOfStep]
93 plt.plot(y,x)
94 plt.gca().set_position((.125, .2, .8, .7))
95 plt.xlabel('Starting Time (fs)')
96 plt.ylabel('Slope (KCal/mol/fs)')
97 plt.ticklabel_format(axis='y', style='sci', scilimits=(-2,2), useOffset=False)
98 plt.title('Slope vs Time for the MD File: ' + str(input))
99 plt.figtext(.01, .05,report)
100 plt.savefig(str(input) + '_SlopePlot.pdf', format='pdf')

```

### A.2.5 Prepare TD-DFT Input

After equilibrium is determined, fincut2.py can be used to create TD-DFT input files from xyz-movie file and a text file containing gmssub commands especially for Bowdoin HPC grid. The script was originally created by Ricke '14 for this work, but many improvement has been made. The updated script works faster and more efficient, even though it still has outdated syntax and methods.

```

1 #####
2 ### This script create a folder of inp for TDDFT with #####
3 ### user-defined time range and TDDFT functional #####
4 #####
5
6 import os
7 import sys
8

```

```

9  #call as fincut.py input.log startTime stopTime timePerFrame
10 #this will run from starting startTime+timePerFrame to stopTime
11 #for example 15010-25000 if input is 15000, 25000
12
13 #this dict is for generating atomic number from Acronym
14 atomicNumber={ 'LV': 116.0, 'BE': 4.0, 'FR': 87.0, 'BA': 56.0, 'BH': 107.0, 'BI': 83.0, 'BK': 97.0,
   ↪ 'EU': 63.0, 'FE': 26.0, 'BR': 35.0, 'ES': 99.0, 'FL': 114.0, 'FM': 100.0, 'RG': 111.0, 'RU':
   ↪ 44.0, 'NO': 102.0, 'NA': 11.0, 'NB': 41.0, 'ND': 60.0, 'NE': 10.0, 'RE': 75.0, 'RF': 104.0,
   ↪ 'LU': 71.0, 'RA': 88.0, 'RB': 37.0, 'NP': 93.0, 'RN': 86.0, 'RH': 45.0, 'B': 5.0, 'CO': 27.0,
   ↪ 'TH': 90.0, 'CM': 96.0, 'CL': 17.0, 'H': 1.0, 'CA': 20.0, 'CF': 98.0, 'CE': 58.0, 'N': 7.0,
   ↪ 'CN': 112.0, 'P': 15.0, 'GE': 32.0, 'GD': 64.0, 'GA': 31.0, 'V': 23.0, 'CS': 55.0, 'CR': 24.0,
   ↪ 'DS': 110.0, 'CU': 29.0, 'SR': 38.0, 'UUP': 115.0, 'UUS': 117.0, 'TC': 43.0, 'KR': 36.0, 'SI':
   ↪ 14.0, 'SN': 50.0, 'SM': 62.0, 'UUT': 113.0, 'SC': 21.0, 'SB': 51.0, 'TA': 73.0, 'OS': 76.0,
   ↪ 'PU': 94.0, 'SE': 34.0, 'AC': 89.0, 'HS': 108.0, 'YB': 70.0, 'DB': 105.0, 'C': 6.0, 'HO':
   ↪ 67.0, 'DY': 66.0, 'HF': 72.0, 'HG': 80.0, 'HE': 2.0, 'PR': 59.0, 'PT': 78.0, 'LA': 57.0, 'F':
   ↪ 9.0, 'UUO': 118.0, 'LI': 3.0, 'PB': 82.0, 'TL': 81.0, 'TM': 69.0, 'LR': 103.0, 'PD': 46.0,
   ↪ 'TI': 22.0, 'TE': 52.0, 'TB': 65.0, 'PO': 84.0, 'PM': 61.0, 'ZN': 30.0, 'AG': 47.0, 'NI':
   ↪ 28.0, 'I': 53.0, 'K': 19.0, 'IR': 77.0, 'AM': 95.0, 'AL': 13.0, 'O': 8.0, 'S': 16.0, 'AR':
   ↪ 18.0, 'AU': 79.0, 'AT': 85.0, 'W': 74.0, 'IN': 49.0, 'Y': 39.0, 'CD': 48.0, 'ZR': 40.0, 'ER':
   ↪ 68.0, 'MD': 101.0, 'MG': 12.0, 'PA': 91.0, 'SG': 106.0, 'MO': 42.0, 'MN': 25.0, 'AS': 33.0,
   ↪ 'MT': 109.0, 'U': 92.0, 'XE': 54.0}
15
16 #if not sure use ? to ask
17 try:
18     if str(sys.argv[1])=='?':
19         print '\nCall function as: fincut.py input.log numberofSoluteAtoms
   ↪ numberofSolventAtoms numberofSolventMolecules startTime stopTime timePerFrame
   ↪ Functional\n'
20         sys.exit()
21 except IndexError:
22     print '\n!!!Input command Error. Call function as: fincut.py input.log numberofSoluteAtoms
   ↪ numberofSolventAtoms numberofSolventMolecules startTime stopTime timePerFrame
   ↪ Functional\n'
23     sys.exit()
24
25 #Call as fincut.py input.log numberofSoluteAtoms numberofSolventAtoms numberofSolventMolecules
   ↪ startTime stopTime timePerFrame
26 try:
27     input=str(sys.argv[1])
28     numberofSoluteAtoms=int(sys.argv[2])
29     numberofSolventAtoms=int(sys.argv[3])
30     numberofSolventMolecules=int(sys.argv[4])
31     startTime =int(sys.argv[5])
32     stopTime =int(sys.argv[6])
33     timePerFrame=int(sys.argv[7])
34     Functional=str(sys.argv[8])
35 except IndexError:
36     print '\n!!!Input command Error. Call function as: fincut.py input.log numberofSoluteAtoms
   ↪ numberofSolventAtoms numberofSolventMolecules startTime stopTime timePerFrame
   ↪ Functional\n'
37     sys.exit()
38
39 #LCBLYP has special file format
40 if Functional=='LCBLYP':
41     FunctionalPrint='BLYP'
42     LC=true
43 else:
44     FunctionalPrint=Functional
45     LC=false

```

```

46
47
48 #should collect data when 1) time within interval specified
49 #2) they are x y z line - not line 1 and 2
50 def shouldCollect(numberOfLine):
51     frame=int(numberOfLine/totalNumberOfLineInOneSet)
52     time=frame*timePerFrame
53     #if within range
54     if time>startTime and time<=stopTime:
55         #0 and 1 (line 1 and 2) in each frame in xyz are not useful
56         if numberOfLine%totalNumberOfLineInOneSet!=0 and numberOfLine%totalNumberOfLineInOneSet!=1:
57             return True
58         else:
59             return False
60     else:
61         return False
62
63 #condition for writing to each inp
64 def shouldWrite(numberOfLine):
65     #should write when last line of a frame is read
66     if numberOfLine%totalNumberOfLineInOneSet==totalNumberOfLineInOneSet-1 :
67         return True
68     else:
69         return False
70
71 #write FRAGNAME=H2ODFT ! 1
72 def shouldWriteFragmentHeader(numberOfLineInSolvent):
73     #should write before writing geometry of solvent molecules
74     if numberOfLineInSolvent%numberOfSolventAtoms==0 :
75         return True
76     else:
77         return False
78
79 #create format of GAMESS inp
80 def insertAtomicNumberInto(line):
81     lineSplit=line.split()
82     lineSplit.insert(1,str(atomicNumber[lineSplit[0]]))
83     return
84     ↪ lineSplit[0]+\t+lineSplit[1]+\t+lineSplit[2]+\t+lineSplit[3]+\t+lineSplit[4]+\n
85 #create format of GAMESS inp O -> O1, H-> H2...
86 def insertCountsInto(line, count):
87     lineSplit=line.split()
88     lineSplit[0]=lineSplit[0]+str(count)
89     return lineSplit[0]+\t+lineSplit[1]+\t+lineSplit[2]+\t+lineSplit[3]+\n
90
91 #even tho input is received written in .log - it will ultimately use .xyz created by 3dExtract for
92 #efficiency
93 #this can be confusing - room for improvement
94 if input.endswith('.log'):
95     input = input[:-4]
96
97 #Path for storing input files
98 path = os.getcwd()
99 outputPath=path+'/'+input+'InputFiles_'+Functional
100 #create folders if not already done
101 if not os.path.exists(outputPath): os.makedirs(outputPath)
102
103 #initiate variables
104 grandString=[]

```

```

103  lineNumber=0
104  numberAllSolventsAtoms=numberSolventAtoms*numberSolventMolecules
105  totalNumberOfAtoms=numberSoluteAtoms+numberAllSolventsAtoms
106  totalNumberOfLineInOneSet=totalNumberOfAtoms+2
107
108  #if no xyz in folder - ask for it
109      #if no 3dExtract.py - print...
110      #if there is, call it from here + print s'th
111  #else create one?
112  #if there is, then proceed
113  #this is very confusing - if there is xyz but it's not updated when it's of no use - room for
    ↪ improvement
114
115  #open input.xyz
116  try:
117      f=open(input+'.xyz')
118  except IOError:
119      print '\nThere is currently no xyz file named '+input
120      print 'trying to call 3dExtract'
121      try:
122          os.system('python 3dExtract4.py'+ ' '+input+'.log'+ ' '+str(numberSoluteAtoms)+'
    ↪ '+str(numberSolventAtoms)+ ' '+str(numberSolventMolecules))
123  except IOError:
124      print '\n Error. There is no 3dExtract to call. Try copying 3dExtract here. \n'
125      print 'Process terminated abnormally'
126      sys.exit()
127  f=open(input+'.xyz')
128
129  #since we are reading from xyz - using readlines() is okay
130  numberAllAtoms=int(f.readline().strip())
131  #enumerate gets data in line - line and line index - n
132  for line in f:
133      lineNumber+=1
134      #check if the line should be written
135      if shouldCollect(lineNumber):
136          grandString.append(line)
137          #if about to write -create file with this name
138          if shouldWrite(lineNumber):
139              frame=int(lineNumber/totalNumberOfLineInOneSet)
140              time=frame*timePerFrame
141              f1 = open(inputPath+'/'+input+'_'+str(time)+'.inp','w')
142              #write header
143              headerString=""" $CONTRL SCFTYP=RHF TDDFT=EXCITE DFTTYP=""Functional+"""
144              RUNTYP=ENERGY ICHARG=0 MULT=1 COORD=UNIQUE
145              MAXIT=200 NPRINT=-5 ISKPRP=1 $END
146
147
148
149  !TDDFT requires lots of memory space
150  $SYSTEM MWORDS=200 MEMDDI=250 $END
151  $SCF DIRSCF=.T. $END
152  $TDDFT NSTATE=5 TPA=.f. $END\n"""
153  if LC==False:
154      headerString=headerString+' $DFT LC=.t. $END'
155  else:
156      headerString=headerString+' $DFT LC=.f. $END'
157  headerString=headerString+"""\\n $BASIS GBASIS=N311 NGAUSS=6 NDFUNC=2 NPFUNC=1
158  DIFFSP=.TRUE. DIFFS=.TRUE. POLAR=POPN311 $END
159  $DATA\\n"""+input+' at t= '+str(time)+ '\\nC1 1\\n'

```

```

160     print 'Making input file for '+input+' at t= '+str(time) +' with '+FunctionalPrint
161     #write the header
162     f1.write(headerString)
163     #write one line by one - before writing we need to add atomic number in using the
164     #    function defined above
165     for eachLine in grandString[:number0fSoluteAtoms]:
166         eachLine=insertAtomicNumberInto(eachLine)
167         f1.write (eachLine)
168     #end solute and go to solvent
169     stringBwSoluteAndSolvent=""" $END\n\n $EFRAG\nCOORD=CART POSITION=OPTIMIZE \n"""
170     f1.write(stringBwSoluteAndSolvent)
171     #write fragments
172     #for iteration
173     i=0
174     for number0fLineWithInSolvents, eachLine in
175         #    enumerate(grandString[number0fSoluteAtoms:]):
176             #is header needed - if so write
177             if shouldWriteFragmentHeader(number0fLineWithInSolvents):
178                 fragmentsNumber=number0fLineWithInSolvents/number0fSolventAtoms+1
179                 f1.write("FRAGNAME=H2ODFT ! "+str(fragmentsNumber)+'\n')
180             #write each line in fragments
181             atomNumberInFragments=i%number0fSolventAtoms+1
182             eachLine=insertCountsInto(eachLine, atomNumberInFragments)
183             f1.write (eachLine)
184             i+=1
185             #close and reset
186             f1.write (" $END\n")
187             grandString=[]
188             f1.close()
189
190     #create megaio for gmssub input
191     allFiles = os.listdir(inputPath)
192     f = open('megaio_'+input+'_hpc.txt', 'w')
193     program = 'gmssub'
194     processors = '32'
195
196     #write item
197     for item in allFiles:
198         outputName=str(item)[-4:]+'.log'
199         f.write(program + ' %s '%item + processors + ' '+outputName+' -l virtual_free=4g\n')
200
201     #for moving megaio to current folder and set permission to read
202     cmdstring = 'chmod 777 megaio_'+input+'_hpc.txt\n'+ 'mv megaio_'+input+'_hpc.txt '+inputPath
203     os.system(cmdstring)
204
205     #call for xsedeGenerator.py
206     #os.system('python xsedeShellGenerator.py '+input+'InputFiles_'+Functional)

```

### A.2.6 Pull Excited State Energies From TD-DFT Log Files

postMDDataPull2.py pulls out excited state energies and dipole moments. Energy output is in the format of time, S1, S2... Dipole output is in the format of time, X1, Y1, Z1, X2...

```

1 ##########
2 ### Use this pull out excited state energies and dipole #####
3 ### moments from all the files in specified folder #####
4 ### Output is time, S1, S2... or time, X1, Y1, Z1, X2... #####

```

```

5 ##########
6
7 import os as os
8 import numpy as numpy
9 import matplotlib.pyplot as plt
10 import sys
11
12 #if not sure use ? to ask
13 try:
14     if str(sys.argv[1])=='?':
15         print '\nCall function as: postMDDataPull2.py inputDirectory
16             ↪  NumberOfExcitedStateEnergies\n'
17         sys.exit()
18 except IndexError:
19     print '\n!!!Input command Error. Call function as: postMDDataPull2.py inputDirectory
20             ↪  NumberOfExcitedStateEnergies\n'
21     sys.exit()
22
23 #Call as fincut.py input.log numberOfSoluteAtoms numberOfSolventAtoms numberOfSolventMolecules
24 #           → startTime stopTime timePerFrame
25 try:
26     input=str(sys.argv[1])
27     numberofExcitedStates=int(sys.argv[2])
28 except IndexError:
29     print '\n!!!Input command Error. Call function as: postMDDataPull2.py inputDirectory
30             ↪  NumberOfExcitedStateEnergies\n'
31     sys.exit()
32 # to prevent / at the end of the input file if used terminal autofill
33 if input.endswith('/'):
34     input = input[:-1]
35
36 #make input path - folder containing the log files or out files
37 path = os.getcwd()
38 outputPath=path+'/' +input
39
40 #arrays for storing and manipulating the data extracted
41 outputList=[]
42 fileList=os.listdir(inputPath)
43 listAllEAndf=[]
44 energies=[]
45 oscillatorStrengths=[]
46 timeList=[]
47 energyFinal=numpy.array(['Energy(eV)'])
48 oscillatorStrengthFinal=numpy.array(['Oscillator Strength'])
49
50 #get all the file names
51 for fileName in fileList:
52     if ("out" in fileName or ".log" in fileName):
53         outputList.append(fileName)
54
55 #open 2 output files
56 f1=open(input+'MD_data.csv','w')
57 f2=open(input+'MD_dipole.csv','w')
58 #now read each file and collect data
59 for fileName in outputList:
60     eachFile=inputPath+'/'+fileName
61     f=open(eachFile,'r')
62     lines=f.readlines()
63     startingLine=0

```

```

60     fileIsComplete=False
61     #find if the file is complete
62     # use enumerate and reverse
63     for i, line in reversed(list(enumerate(lines))):
64         if ' STATE # 1 ENERGY =' in line:
65             startingLine=i
66             fileIsComplete=True
67             break
68     # now start collecting data
69     if fileIsComplete:
70         print 'Harvesting data from file named: '+fileName
71         counter=0
72         #loop
73         for n, line in enumerate(lines):
74             #time
75             if line.startswith(' RUN TITLE'):
76                 nextLine=lines[n+2]
77                 time=str(nextLine.split()[3])
78                 f1.write(time)
79                 f2.write(time)
80             #energy
81             if line.startswith(' STATE # '):
82                 E=line.split()[5]
83                 nextLine=lines[n+1]
84                 f=nextLine.split()[3]
85                 f1.write(','+E+',',+f)
86                 counter+=1
87             #stop collecting
88             if (counter==numberOfExcitedStates):
89                 f1.write('\n')
90                 counter=0
91             #dipole
92             if line.startswith(' SUMMARY OF TDDFT RESULTS'):
93                 #5 lines are from 'SUMMARY...' to ' 1 A ...' which starts to
94                 # contain dipole moments
95                 for linesAhead in range(5,5+numberOfExcitedStates):
96                     lineContainingDipoles=lines[n+linesAhead]
97                     lineSplit=lineContainingDipoles.split()
98                     [X,Y,Z]=[lineSplit[4], lineSplit[5], lineSplit[6]]
99                     f2.write(','+X+',,'+Y+',,'+Z)
100                    f2.write('\n')

```

## A.3 GAMESS Inputs

### A.3.1 MD Input File

MD run is core to modeling explicit solvent. The MD run is simulated every femtosecond but only record every 10 femtoseconds. The bath temperature is  $25 \pm 25$  degree Celsius. Solvent boundary potential is also activated using default Sforce value ( $1.0\text{kcal} \cdot \text{mol}^{-1} \cdot \text{\AA}^{-2}$ ), but with estimate ssbp radius. #####n##### are for restarting MD in case the calculation abruptly ends (see next section). In this version, dispersion correction is not turned on. Basis set = 6-31+G(2d,p)

```

1 # run type = MD, with functional = CAMB3LYP, COORD = UNIQUE is important
2 $CONTRL SCFTYP=RHF RUNTYP=MD COORD=UNIQUE
3 DFTTYP=CAMB3LYP MAXIT=200 ICHARG=0 MULT=1 $END
4 # MD is recording every 10 frames with default 1 frame =10 fs
5 # 25 degree celcius, RSTEMP is on for keeping the temp ~ +/-25
6 # ssbp is on with default SForce value and radius estimated from prepareMD2.py
7 $MD KEVERY=10 PROD=.T. NVTNH=2 MBT=.T. MBR=.T.
8 BATHT=298 RSTEMP=.T. DTEMP=25 NSTEPS=50000
9 SSBP=.T. SFORCE=1.0 DROFF=12.0632116659 $END
10 #####1#####
11 #####
12 # dispersion correction is off
13 $DFT DC=.F. $END
14 # memory requested at each node =1000 million words
15 # memory reserved for communication = 1000 million words
16 $SYSTEM MWORDS=1000 MEMDDI=1000 $END
17 $SCF DIRSCF=.T. $END
18 # Basis set = 6-31+(2d,p)
19 $BASIS GBASIS=N31 NGAUSS=6 NDFUNC=2 NPFUNC=1
20 DIFFS=.TRUE. POLAR=POPN31 $END
21 # solute geometry - C1 1 = symmetry data
22 $DATA
23 MD INPUT for aniline32
24 C1 1
25 #####2#####
26 N      7.0      -2.3128100000      -0.0046000000      -0.0894530000
27 C      6.0      -0.9197160000      -0.0031280000      -0.0360090000
28 C      6.0      -0.2076150000      1.2004070000      -0.0355880000
29 .
30 .
31 H      1.0      -2.7544730000      0.8437050000      0.2369240000
32 H      1.0      -2.7570450000      -0.8248640000      0.2993670000
33 #####
34 $END
35
36 # solvent geometry in EFP1 (EFP2 is still not available)
37 #####3#####
38 $EFrag
39 COORD=CART POSITION=OPTIMIZE
40 FRAGNAME=H2ODFT ! 1
41 O1      1.7760990000      4.8390610000      -2.1049530000
42 H2      0.9224740000      4.4173920000      -2.2628490000
43 H3      2.4128200000      4.1148590000      -2.1441810000
44 FRAGNAME=H2ODFT ! 2
45 O1      3.6783070000      3.8351060000      0.8717800000
46 H2      3.6020030000      4.1116030000      1.7932670000
47 H3      3.4138960000      4.6126330000      0.3648680000
48 .
49 .
50 FRAGNAME=H2ODFT ! 32
51 O1      3.7691430000      -1.4091250000      -4.0319510000
52 H2      2.8756120000      -1.5562470000      -4.3656710000
53 H3      3.6686840000      -1.3995500000      -3.0721400000
54 $END
55 #####

```

### A.3.2 MD Restart

#####n##### are for restarting MD. For example, if MD stops from errors at t= 39000 fs, a restart geometry and \$MD should be obtained from t= 38960 in the run's trj file. #####1##### from trj goes to #####1##### in MD input file and so on with 2 and 3.

```

1 .
2 .
3 #time = 38960 fs
4 ===== MD DATA PACKET =====
5 NAT= 14 NFRG= 32 NQMMM= 0
6 TTOTAL= 38960.00 FS TOT. E= -180710.543716 KCAL/MOL
7 POT. E= -180783.917894 KCAL/MOL BATHT= 298.000000
8 KIN. E= 73.374179 TRANS KE= 44.012966 ROT KE= 29.361213 KCAL/MOL
9 ----- QM PARTICLE COORDINATES FOR $DATA GROUP -----
10 ######2#####
11 N 7.0 3.8554852781 -1.5814882196 -3.7440432660
12 C 6.0 2.9445799802 -1.6872198850 -4.7513152920
13 .
14 .
15 H 1.0 3.8827088558 -0.7337182877 -3.1704671239
16 ######3#####
17 ----- EFP PARTICLE COORDINATES FOR $EFRAG GROUP -----
18 ######3#####
19 $EFRAG
20 COORD=CART POSITION=OPTIMIZE
21 FRAGNAME=H2ODFT ! 1
22 O1 0.4899097683 5.6815060052 1.3332597175
23 H2 -0.1774491436 5.6953737679 2.0005813955
24 H3
25 .
26 .
27 FRAGNAME=H2ODFT ! 32
28 O1 1.8397452656 -1.0633683830 1.4432099580
29 H2 2.6999380304 -1.1499324935 1.8219625372
30 H3 1.3989608369 -0.3812464183 1.9241435987
31 $END
32 #####
33 GRADIENT DATA (NOT USED BY RESTARTS)...
34 FRAGMENT # 1 H2ODFT
35 .
36 .
37 ----- RESTART VELOCITIES FOR $MD GROUP -----
38 ######1#####
39 $MD READ=.TRUE. MBT=.FALSE. MBR=.FALSE. TTOTAL= 3.90E-11
40 MDINT= VVERLET DT= 0.10E-14 NVTNH= 2 NSTEPS= 11040
41 RSTEMP=.T. DTEMP= 25.00 LEVERY= 50000
42 RSRAND=.F. NRAND= 1000 NVTOFF= 0 JEVERY= 10
43 PROD=.T. KEVERY= 10 DELR= 0.020
44 Batht(1)=298.00
45 SSBP=.T. SFORCE= 1.0 DROFF= 12.1****
46 TVELQM(1)= ! QM ATOM TRANS. VELOCITIES (BOHR/PS) !
47 -1.436664655E+00 -1.014912632E+00 1.731488079E+01
48 .
49 .
50 .
51 -4.714191860E+00 -3.065697154E+00 1.306355299E+01

```

```

52 TVEL(1)=      ! EFP TRANSLATIONAL VELOCITIES (BOHR/PS) !
53   -9.879998843E+00  -1.286351783E+01  -1.625361337E-01
54 .
55 .
56   4.573922683E+00  9.602851076E+00  6.875495095E+00
57 QUAT(1)=      ! EFP QUATERNIONS !
58   6.658493597E-01  -7.408965451E-01  1.545516153E-02  8.647587874E-02
59 .
60 .
61   8.096376474E-01  4.753557554E-01  1.423717372E-01  3.134550593E-01
62 RVEL(1)=      ! EFP ANGULAR VELOCITY (RAD/PS) !
63   1.392338986E+01  -1.145718732E+01  -1.039863492E+01
64 .
65 .
66   -1.217796312E+01  2.229922184E+00  2.270460047E+01
67 QUAT1D(1)=      ! EFP QUATERNION 1ST DERIV. !
68   4.237107071E+12  4.149735474E+12  8.424472510E+12  1.422917746E+12
69 .
70 .
71   5.909828315E+12  -1.125106858E+13  -8.408659179E+11  2.179439222E+12
72 QUAT2D(1)=      ! EFP QUATERNION 2ND DERIV. !
73   -6.427981808E+26  -4.203604186E+26  -7.365308723E+25  1.102236012E+26
74 .
75 .
76   8.379527677E+26  4.204660564E+26  -1.448608909E+27  -2.676740098E+27
77 $END
78 #####
79
80 #time = 38970 fs
81 ===== MD DATA PACKET =====
82 NAT=      14 NFRG=      32 NQMMM=      0
83 TTOTAL=    38970.00 FS      TOT. E=    -180715.284973 KCAL/MOL
84 POT. E=    -180783.048242 KCAL/MOL  BATHT=      298.000000
85 KIN. E=     67.763269 TRANS KE=    41.754223  ROT KE=    26.009046 KCAL/MOL
86 ----- QM PARTICLE COORDINATES FOR $DATA GROUP -----
87 N          7.0          3.8518322659      -1.6019379232      -3.6477010146
88 .
89 .

```

### A.3.3 TD-DFT Input File

Excited state energies are calculated using TD-DFT. Direct SCF calculation is turned on. Basis set = 6-311++G(2d,p)

```

1  # run type = [excitation] energy, with functional = CAMB3LYP, and TDDFT
2  $CONTRL SCFTYP=RHF TDDFT=EXCITE DFTTYP=CAMB3LYP RUNTYP=ENERGY
3  ICHARG=0 MULT=1 COORD=UNIQUE MAXIT=200 $END
4  #TDDFT requires lots of memory space
5  # memory requested at each node =1000 million words
6  # memory reserved for communication = 1000 million words
7  $SYSTEM MWORDS=200 MEMDDI=250 $END
8  #activate direct SCF calculation
9  $SCF DIRSCF=.T. $END
10 # find 5 excited states - the current setting is purely driven by its lower cost
11 # Previous experience shows that 10 states gives only a few strong peak.
12   $TDDFT NSTATE=5 TPA=.f. $END

```

```

13  # Basis set = 6-311++(2d,p)
14  $BASIS GBASIS=N311 NGAUSS=6 NDFUNC=2 NPFUNC=1
15      DIFFSP=.TRUE. DIFFS=.TRUE. POLAR=POPN311 $END
16  # solute geometry - C1 1 = symmetry data
17  $DATA
18  aniline32 at t= 15010
19  C1 1
20  N      7.0      2.4008547653      5.9114221893      -1.1412310058
21  C      6.0      1.9371475177      5.9223533811      -2.4157851823
22  C      6.0      0.6209366009      6.2361805033      -2.7812041720
23  .
24  .
25  H      1.0      1.7323956480      5.6459040114      -0.4329519914
26  H      1.0      3.3564486514      5.6102990678      -1.0242941608
27  $END
28
29  # solvent geometry in EFP1 (EFP2 is still not available)
30  $EFRAG
31  COORD=CART POSITION=OPTIMIZE
32  FRAGNAME=H2ODFT ! 1
33  O1      2.335993939511      3.751856628604      1.427418842826
34  H2      1.439322965739      3.833168541986      1.710699607266
35  H3      2.854351938959      3.613582975500      2.203990690907
36  FRAGNAME=H2ODFT ! 2
37  O1      3.266753260182      2.613267395174      3.808546039622
38  H2      3.514552920456      1.768271678250      3.468757839439
39  H3      3.822697636667      2.780022240614      4.552855871622
40  .
41  .
42  FRAGNAME=H2ODFT ! 32
43  O1      -0.041331901955      -3.906598195206      1.282021099515
44  H2      -0.790424236756      -4.406259423929      1.565001283174
45  H3      -0.374908444181      -3.111491773717      0.898079798177
46  $END

```