

# NH<sub>4</sub><sup>+</sup> Resides Inside the Water 20-mer Cage As Opposed to H<sub>3</sub>O<sup>+</sup>, Which Resides on the Surface: A First Principles Molecular Dynamics Simulation Study

Soohaeng Yoo Willow, N. Jiten Singh,\* and Kwang S. Kim\*

Center for Superfunctional Materials, Department of Chemistry, Pohang University of Science and Technology, San 31, Hyojadong, Namgu, Pohang 790-784, Korea

**S** Supporting Information

**ABSTRACT:** Experimental vibrational predissociation spectra of the magic NH<sub>4</sub><sup>+</sup>(H<sub>2</sub>O)<sub>20</sub> clusters are close to those of the magic H<sub>3</sub>O<sup>+</sup>(H<sub>2</sub>O)<sub>20</sub> clusters. It has been assumed that the geometric features of NH<sub>4</sub><sup>+</sup>(H<sub>2</sub>O)<sub>20</sub> clusters might be close to those of H<sub>3</sub>O<sup>+</sup>(H<sub>2</sub>O)<sub>20</sub> clusters, in which H<sub>3</sub>O<sup>+</sup> resides on the surface. Car–Parrinello molecular dynamics simulations in conjunction with density functional theory calculations are performed to generate the infrared spectra of the magic NH<sub>4</sub><sup>+</sup>(H<sub>2</sub>O)<sub>20</sub> clusters. In comparison with the experimental vibrational predissociation spectra of NH<sub>4</sub><sup>+</sup>(H<sub>2</sub>O)<sub>20</sub>, we find that NH<sub>4</sub><sup>+</sup> is *inside* the cage structure of NH<sub>4</sub><sup>+</sup>(H<sub>2</sub>O)<sub>20</sub> as opposed to *on* the surface structure. This shows a clear distinction between the structures of NH<sub>4</sub><sup>+</sup>(H<sub>2</sub>O)<sub>20</sub> and H<sub>3</sub>O<sup>+</sup>(H<sub>2</sub>O)<sub>20</sub> as well as between the hydration phenomena of NH<sub>4</sub><sup>+</sup> and H<sub>3</sub>O<sup>+</sup>.

## INTRODUCTION

Since solvated ions play a pivotal role in the chemical and physical properties of chemical and biological systems<sup>1,2</sup> and in the environment of the upper atmosphere,<sup>3</sup> diverse experimental and theoretical approaches have been carried out to understand the intriguing phenomena of interactions between ions and solvent molecules. Useful information has been gained from studies of the solvation of cations<sup>4–10</sup> and anions<sup>11–19</sup> in the gas phase as well as water clusters.<sup>20–28</sup>

Solvated ammonium (NH<sub>4</sub><sup>+</sup>) cations have been intensively studied experimentally<sup>29–34</sup> and theoretically<sup>35–43</sup> due to their similarity to the solvated hydronium (H<sub>3</sub>O<sup>+</sup>) in that both clusters have an excess proton. The similarity and difference between H<sub>3</sub>O<sup>+</sup>(H<sub>2</sub>O)<sub>*n*</sub> and NH<sub>4</sub><sup>+</sup>(H<sub>2</sub>O)<sub>*n*</sub> could shed light on the intriguing role of how protonated water systems are related to proton transfer in aqueous chemistry and biology.<sup>44–49</sup> One interesting experimental IR spectrum of NH<sub>4</sub><sup>+</sup>(H<sub>2</sub>O)<sub>4–6</sub> showed the peak of NH<sub>d</sub> stretching clearly in the range of 2900–3450 cm<sup>-1</sup>, where NH<sub>d</sub> indicates a dangling NH bond.<sup>34,35</sup> For theoretical interpretation of the experimental IR spectra, some low-lying structures were identified to show a dangling NH<sub>d</sub> appearance.<sup>39</sup> The global minimum structures of NH<sub>4</sub><sup>+</sup>(H<sub>2</sub>O)<sub>4–6</sub> for the complete basis set (CBS) limit at the CCSD(T) level of theory have no dangling NH<sub>d</sub> since NH<sub>4</sub><sup>+</sup> is fully solvated in the global minimum isomers. Hence, several low-lying structures would be required for interpretation of the experimental IR spectra.

The experimental IR spectra for larger clusters of NH<sub>4</sub><sup>+</sup>(H<sub>2</sub>O)<sub>*n*</sub> (*n* > 8) did not show the peak of NH<sub>d</sub> stretching near 2900–3450 cm<sup>-1</sup>.<sup>32,34,35</sup> Furthermore, vibrational predissociation spectra of the magic NH<sub>4</sub><sup>+</sup>(H<sub>2</sub>O)<sub>20</sub> clusters in 2500–3900 cm<sup>-1</sup> are close to those of the magic H<sub>3</sub>O<sup>+</sup>(H<sub>2</sub>O)<sub>20</sub> clusters. It has been assumed that the geometric features of NH<sub>4</sub><sup>+</sup>(H<sub>2</sub>O)<sub>20</sub> clusters might be the same as those of H<sub>3</sub>O<sup>+</sup>(H<sub>2</sub>O)<sub>20</sub> clusters, in which H<sub>3</sub>O<sup>+</sup> prefers to be on the surface of clusters. Diken et al.

suggested several “handmade” structures of NH<sub>4</sub><sup>+</sup>(H<sub>2</sub>O)<sub>20</sub>, in which the isomer with NH<sub>4</sub><sup>+</sup> on the surface of clusters was lower in internal energy than that with the fully solvated NH<sub>4</sub><sup>+</sup>.<sup>32</sup> On the basis of their “handmade” structures, they discussed that the experimental vibrational spectra failed to display the dangled NH<sub>d</sub> stretch near 3450 cm<sup>-1</sup> since its peak might overlap with OH<sub>d</sub> stretching transitions, where OH<sub>d</sub> indicates a dangling OH bond of water. Since Brutschy and co-workers<sup>8</sup> reported the intriguing mass spectra of the magic clusters of NH<sub>4</sub><sup>+</sup>(H<sub>2</sub>O)<sub>20</sub>, K<sup>+</sup>(H<sub>2</sub>O)<sub>20</sub>, and Cs<sup>+</sup>(H<sub>2</sub>O)<sub>20</sub>, we briefly addressed that, for NH<sub>4</sub><sup>+</sup>(H<sub>2</sub>O)<sub>20</sub>, the structure with the NH<sub>4</sub><sup>+</sup> ion inside the cage is more stable than that on the surface of the cage by about 2 kcal/mol, and K<sup>+</sup>(H<sub>2</sub>O)<sub>20</sub> and Cs<sup>+</sup>(H<sub>2</sub>O)<sub>20</sub> also show internal structures.<sup>49</sup> Then, Douady et al.<sup>40</sup> studied low-lying isomers of NH<sub>4</sub><sup>+</sup>(H<sub>2</sub>O)<sub>*n*</sub> (*n* ≤ 24). They confirmed that NH<sub>4</sub><sup>+</sup> is fully solvated in the global minimum structures for clusters NH<sub>4</sub><sup>+</sup>(H<sub>2</sub>O)<sub>*n*</sub> (*n* ≥ 6). In addition, the isomers NH<sub>4</sub><sup>+</sup>(H<sub>2</sub>O)<sub>20</sub> with the fully solvated NH<sub>4</sub><sup>+</sup> were lower in internal energy by 1–4 kcal/mol than the isomers with a dangling NH<sub>d</sub> of NH<sub>4</sub><sup>+</sup> residing on their surfaces. In this letter, we present a study of low-lying isomers of NH<sub>4</sub><sup>+</sup>(H<sub>2</sub>O)<sub>20</sub> and calculations of their infrared (IR) spectra to compare them with the experimental spectra.

## COMPUTATIONAL APPROACH

**Global Search of NH<sub>4</sub><sup>+</sup>(H<sub>2</sub>O)<sub>20</sub>.** We searched for low-lying energy structures using the density-functional tight binding theory (DFTB).<sup>50</sup> The basin-hopping global optimization method<sup>51,52</sup> was used to search for the geometries of low-lying isomers. A key idea of the basin-hopping method is to generate the transformed potential-energy surface (PES)  $\tilde{U}$  using the

Received: July 12, 2011

Published: October 11, 2011

following mapping equation

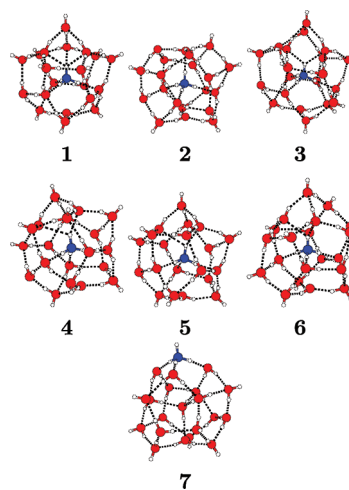
$$\tilde{U}(N, \mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_N) = \min\{U(N, \mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_N)\}$$

where min denotes the energy minimization with starting configuration of  $\{\mathbf{r}_1, \dots, \mathbf{r}_N\}$  and  $U$  is the PES. This transformed PES  $\tilde{U}$  removes the potential well existing in PES  $U$  and allows the system “hop” directly between different local minima at each step.

Since DFTB is much less reliable in determining the relative stability between isomers, the following three steps were employed to seek possible low-lying isomers: (1) The database of the top  $\sim 100$  low-lying isomers from the global search with DFTB was categorized into distinct structural families according to their different pentagonal dodecahedron  $(\text{H}_2\text{O})_{20}$  clusters. Note that the low-lying isomers with the fully solvated  $\text{NH}_4^+$  have the pentagonal dodecahedron  $(\text{H}_2\text{O})_{20}$  clusters. (2) The geometries within each different structural family furthermore were optimized using the Becke–3–Lee–Yang–Parr (B3LYP) hybrid functional<sup>53,54</sup> to screen the top 10 low-lying isomers. (3) Finally, the resolution of identity MP2 (RIMP2) calculations were performed for more-accurate energies and to determine the reliable low-lying isomers.<sup>55</sup> DFT calculations with the B3LYP hybrid functional were mainly used as a screening tool since it can yield more reliable energy rankings than DFTB. Recently, it has been reported that the dispersion interaction correction should be included into DFT to accurately predict thermochemical properties, electronic excitations, infrared vibrational spectra, and solvent effects.<sup>56–59</sup> Hence, we employed the M06-2X functional in order to include the corrected dispersion interactions into the potential energy calculations.<sup>59</sup> Thresholds of  $10^{-6}$  au (convergence of the potential energy) and 0.001 au/bohr (convergence of the gradient) were used during the geometry optimizations. The geometry optimization and vibrational frequency calculations were carried out at the level of B3LYP/aug-cc-pVDZ' (aVDZ')<sup>60,61</sup> (in which ' denotes that the diffuse basis function of the hydrogen atom was removed) and M06-2X/aug-cc-pVDZ (aVDZ) theories using the Gaussian 03 suite of programs.<sup>62</sup> TURBOMOLE<sup>63</sup> was used for geometry optimization at the level of RIMP2/aVDZ theory.

**Simulated Infrared spectrum.** We performed Car–Parrinello molecular dynamics (CP-MD) simulations<sup>64,65</sup> at a temperature of 125 K to generate the simulated infrared spectra. The CP-MD simulations were carried out at the level of the plane-wave-pseudopotential density functional theory with the Becke exchange<sup>53</sup> and Lee–Yang–Parr correlation<sup>54</sup> (BLYP) functionals. The core–valence interaction was described by a norm-conserving Troullier–Martins pseudopotential,<sup>66</sup> and the wave function energy cutoff value was 90 Ry. A fictitious electron mass of 600 au and an integration step of  $dt = 0.1$  fs were used. A Nose–Hoover thermostat was employed to generate the canonical ensemble (constant volume and constant temperature). During the simulations of 10 ps, we kept the molecules at the center of isolated cubic boxes of side lengths  $L = 15$  Å for  $\text{NH}_4^+(\text{H}_2\text{O})_{20}$ . From the last 6 ps trajectory of the CP-MD simulations, we evaluated the time correlation function to investigate the spectra of the clusters in the equilibrium state. During NVT simulations, we monitored whether other isomers would be sampled, confirming that the hydrogen bond network of isomers remained. The Fourier transform of dipole moment autocorrelation functions was carried out. The infrared absorption spectrum can be computed from FT-DACF as

$$I(\omega) = (\hbar\beta/2\pi)\omega^2 \int dt e^{-i\omega t} \langle \mu(0) \mu(t) \rangle$$



**Figure 1.** Possible low-lying isomers of  $\text{NH}_4^+(\text{H}_2\text{O})_{20}$  optimized at RIMP2/aVDZ. Atoms of O, N, and H are highlighted in red, blue, and white, respectively. Six isomers with the fully solvated  $\text{NH}_4^+$  are labeled as 1–6. The isomer labeled 7 has one dangling  $\text{NH}_4^+$  due to  $\text{NH}_4^+$  residing on the surface.<sup>32</sup> Isomer 3 was suggested by Douady et al.<sup>40</sup>

Here, the symbols are used to denote intensity ( $I$ ), frequency ( $\omega$ ), Planck constant ( $\hbar = h/2\pi$ ), inverse of the Boltzmann constant multiplied by temperature ( $\beta = 1/kT$ ), time ( $t$ ), and dipole moment ( $\mu$ ), which is the total dipole moment of the clusters rather than the dipole moment of each molecule. For computational and interpretative purposes, it is more convenient to compute the autocorrelation function of the time derivative of the dipole moment:

$$I(\omega) = (\hbar\beta/2\pi) \int dt e^{-i\omega t} \langle \dot{\mu}(0) \dot{\mu}(t) \rangle$$

as discussed by Schmitt and Voth.<sup>67</sup> Hence, this method was employed in our calculations.

## RESULTS AND DISCUSSION

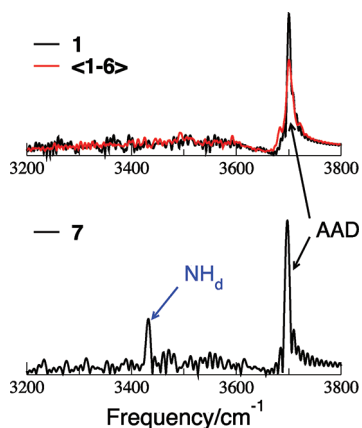
Even though the basin-hopping method is a very efficient global search method, there is a limitation in searching for the global minimum structures using the basin-hopping method coupled with DFTB since the number of possible isomers of  $\text{NH}_4^+(\text{H}_2\text{O})_n$  increases dramatically with the increase of  $n$  (the number of water molecules) and DFTB is not as fast as the empirical water potentials with regard to the computational speed of the energy calculation. Since DFTB is much less reliable in determining the relative stability between isomers, we do not claim that we sampled the true global minimum energy structure of  $\text{NH}_4^+(\text{H}_2\text{O})_{20}$  within the top  $\sim 100$  low-lying isomers from the global search with DFTB. Instead, we present low-lying energy structures at RIMP2/aVDZ in comparison with previously reported low-lying energy structures as the most likely candidate for the global minimum energy structure.

The selected low-lying isomers for the calculation of IR spectra are shown in Figure 1. Structures 1–6 are low-lying energy isomers with fully solvated  $\text{NH}_4^+$ . Isomers 1–6 have different pentagonal dodecahedron  $(\text{H}_2\text{O})_{20}$  clusters. Isomer 7 has  $\text{NH}_4^+$  residing on the surface, as suggested by Diken et al.<sup>32</sup> Relative energies ( $\Delta E_e$  and  $\Delta E_0$  in kcal/mol) of the B3LYP, M06-2X, and RIMP2 calculations are listed in Table 1. The effect of harmonic zero-point energies (ZPE) at the M06-2X/aVDZ level of theory

**Table 1. Relative Energies (in kcal/mol) without ( $\Delta E_e$ ) and with ( $\Delta E_0$ ) Harmonic Zero-Point Energy (ZPE) Corrections for the Low-Lying Isomers of  $\text{NH}_4^+(\text{H}_2\text{O})_{20}$  Shown in Figure 1<sup>a</sup>**

Isomer	B3LYP/aVDZ'		M06-2X/aVDZ		RIMP2/aVDZ	
	$\Delta E_e$	$\Delta E_0$	$\Delta E_e$	$\Delta E_0$	$\Delta E_e$	$\Delta E_0$
1	1.20	0.40	0.71	0.41	<b>0.00</b>	<b>0.00</b>
2	0.85	0.08	1.85	1.36	0.51	0.32
3	1.21	0.19	0.26	0.22	0.13	0.40
4	1.41	0.45	<b>0.00</b>	<b>0.00</b>	0.18	0.48
5	2.39	1.29	0.97	0.31	1.26	0.91
6	2.26	1.35	1.71	1.72	1.33	1.65
7	<b>0.00</b>	<b>0.00</b>	6.67	6.21	3.51	3.36

<sup>a</sup> Isomers 3 and 7 were suggested by Douady et al.<sup>40</sup> and Diken et al.,<sup>32</sup> respectively. The boldface energies denote the lowest-lying isomer. The ZPE-corrected relative energies ( $\Delta E_0$ ) at the RIMP2/aVDZ level were obtained using harmonic ZPE estimates at the M06-2X/aVDZ level of theory.



**Figure 2.** Simulated infrared (IR) spectra of  $\text{NH}_4^+(\text{H}_2\text{O})_{20}$  (frequencies scaled by 1.049). The AAD-type  $\text{OH}_d$  peak appears at  $\sim 3700 \text{ cm}^{-1}$ . The AAD-type  $\text{OH}_d$  peak is from the water molecules on the surface. The cluster 1 and <1–6> with fully solvated  $\text{NH}_4^+$  have no  $\text{NH}_d$  peak, while the cluster 7 with  $\text{NH}_4^+$  on the surface shows the  $\text{NH}_d$  peak at  $\sim 3430 \text{ cm}^{-1}$ . Here, <1–6> indicates the average spectra from six isomers of 1–6.

was added to the RIMP2/aVDZ relative energetics ( $\Delta E_0$ ). The DFT with the B3LYP hybrid functional gave the isomer 7 as the lowest-energy isomer, while M06-2X and RIMP2 gave the highest energy among the given 7 structures (Figure 1). Furthermore, the geometries optimized at B3LYP/aVDZ' were quite different from those optimized at M06-2X/aVDZ and RIMP2/aVDZ (these optimized geometries are provided in the Supporting Information). In detail, the former prefers to form a four-hydrogen-bond network of  $\text{NH}_4^+$  with the pentagonal dodecahedron ( $\text{H}_2\text{O}$ )<sub>20</sub> cages, while the latter prefers a five- or more hydrogen-bond network. This wide deviation in the relative energies of the B3LYP functional from the RIMP2 level of theory is mainly due to underestimation of the dispersion interaction in DFT.<sup>56–59</sup> Hence, the M06-2X functional with the corrected dispersion interaction gives a highly improved description for both potential energies and geometric features of  $\text{NH}_4^+(\text{H}_2\text{O})_{20}$  clusters. We confirmed that isomer 7 has much higher energy than the

solvated  $\text{NH}_4^+$  clusters (isomers 1–6) at the level of RIMP2/aVDZ. In order to make sure that isomer 1 is one candidate for the lowest isomer, we performed geometry optimization at the level of RIMP2/aVDZ for the low-lying isomer (isomer 3 in Table 1) suggested by Douady et al.<sup>40</sup> and computed its relative energy of  $\Delta E_e = 0.13 \text{ kcal/mol}$  and  $\Delta E_0 = 0.4 \text{ kcal/mol}$ . In summary, isomers 1, 2, 3, and 4 with fully solvated  $\text{NH}_4^+$  are nearly isoenergetic lowest-energy structures.

Since the simulated vibrational spectra can give insight into the origin of the “no peaks” of  $\text{NH}_d$  stretching of the vibrational spectra of  $\text{NH}_4^+(\text{H}_2\text{O})_{20}$ , we performed Car–Parrinello molecular dynamics (CP-MD) simulations<sup>64,65</sup> at a temperature of 125 K for the seven low-lying isomers shown in Figure 1. Note that the CP-MD IR spectra are more realistic, being closer to the experimental data since the CP-MD IR spectra reflect both the anharmonic potential surfaces and the contribution of temperature. Figure 2 shows the simulated vibrational spectrum in the broad range  $3200\text{--}3800 \text{ cm}^{-1}$  for two isomers 1 and 7. The dangling AAD-type  $\text{OH}_d$  peak appears at  $\sim 3700 \text{ cm}^{-1}$ , indicating that both isomers 1 and 7 have the AAD-type  $\text{OH}_d$  on the water cages of  $\text{NH}_4^+(\text{H}_2\text{O})_{20}$ . Note that the AAD-type  $\text{OH}_d$  peak from  $\text{NH}_4^+(\text{H}_2\text{O})_{20}$  is also shown in the IR spectra of  $\text{H}_3^+\text{O}(\text{H}_2\text{O})_{20}$ .<sup>49</sup> Both  $\text{NH}_4^+(\text{H}_2\text{O})_{20}$  and  $\text{H}_3^+\text{O}(\text{H}_2\text{O})_{20}$  have the water cage with the dangling  $\text{OH}_d$  on their cage surface. In comparison with isomer 1, the additional IR peak near  $\sim 3430 \text{ cm}^{-1}$  appears in isomer 7. This additional peak is mainly due to the dangling  $\text{NH}_d$  stretching of  $\text{NH}_4^+$ , which was confirmed in the vibrational spectra of smaller clusters of  $\text{NH}_4^+(\text{H}_2\text{O})_n$  ( $n < 8$ ).<sup>34</sup> Hence, the  $\text{NH}_d$  peak near  $\sim 3430 \text{ cm}^{-1}$  becomes the unique peak to identify whether  $\text{NH}_4^+$  is fully solvated in water clusters. The experimental vibrational spectrum<sup>32</sup> of  $\text{NH}_4^+(\text{H}_2\text{O})_{20}$  did not show the dangling  $\text{NH}_d$  stretching peak near  $3430 \text{ cm}^{-1}$ . We suggest that the absence of the experimental peak of  $\text{NH}_d$  near  $3430 \text{ cm}^{-1}$  in  $\text{NH}_4^+(\text{H}_2\text{O})_{20}$  indicates clearly that the ammonium cation  $\text{NH}_4^+$  is fully solvated in most structures of  $\text{NH}_4^+(\text{H}_2\text{O})_{20}$ .

Why is the ammonium cation  $\text{NH}_4^+$  fully solvated in  $\text{NH}_4^+(\text{H}_2\text{O})_n$  ( $n \geq 6$ )? And why can the hydronium cation  $\text{H}_3\text{O}^+$  not be fully solvated in  $\text{H}_3\text{O}^+(\text{H}_2\text{O})_{20}$ ? All hydrogen atoms of  $\text{NH}_4^+$  have a favorable hydrogen bonding interaction with the oxygen atom of a water molecule. Thus, its solubility is similar to that of alkali cations ( $\text{Na}^+$ ,  $\text{K}^+$ , and  $\text{Cs}^+$ ). In contrast, the oxygen atom of  $\text{H}_3\text{O}^+$  is hydrophobic, as in our previous work,<sup>68</sup> since it is no longer a good electron donor or proton acceptor for the formation of a hydrogen bonding interaction. Hence, the hydronium shows an amphiphilic behavior. The  $\text{H}_3\text{O}^+$  ion favors the surface to maximize the polarization-driven binding energy, as shown in amphiphilic species such as lipids. Though the H atoms in  $\text{H}_3\text{O}^+$  are involved in the H bonding, the O atom is not involved in the H bonding. Thus, the  $\text{H}_3\text{O}^+$  remains on the surface of the cluster with three H bonds by three H atoms, while there is no H bond by the O atom. On the other hand, the  $\text{NH}_4^+$  prefers the internal structure forming a five or more H-bond network with the pentagonal dodecahedron ( $\text{H}_2\text{O}$ )<sub>20</sub> cages.

In summary, we performed the CP-MD simulations at 125 K in conjunction with DFT calculations to generate the simulated IR spectra of  $\text{NH}_4^+(\text{H}_2\text{O})_{20}$ . When the experimental vibrational spectrum of  $\text{NH}_4^+(\text{H}_2\text{O})_{20}$  was compared with those of 1 and 7, the simulated spectrum of 1 (the isomer with fully solvated  $\text{NH}_4^+$ ) was consistent with the experimental vibrational spectra. This result indicates that in experimentally measured structures of  $\text{NH}_4^+(\text{H}_2\text{O})_{20}$  the ammonium cation  $\text{NH}_4^+$  is fully solvated

inside the cage structure of  $(\text{H}_2\text{O})_{20}$  against on the surface structure. Even though the experimental vibrational spectra of  $\text{NH}_4^+(\text{H}_2\text{O})_{20}$  are no different from those of  $\text{H}_3\text{O}^+(\text{H}_2\text{O})_{20}$  in the broad range  $3200\text{--}3800\text{ cm}^{-1}$ , the geometric features of  $\text{NH}_4^+(\text{H}_2\text{O})_{20}$  are very different from those of  $\text{H}_3\text{O}^+(\text{H}_2\text{O})_{20}$ .

## ■ ASSOCIATED CONTENT

**S Supporting Information.** Total electronic energies ( $E_e$ ) and zero-point corrected energies ( $E_0$ ) and the optimized geometries of low-lying isomers shown in Figure 1. This material is available free of charge via the Internet <http://pubs.acs.org>.

## ■ AUTHOR INFORMATION

### Corresponding Author

\*Tel.: +82-54-279-2110. Fax: +82-54-279-8137. E-mail: kim@postech.ac.kr (K.S.K.), jiten@postech.ac.kr (N.J.S.).

## ■ ACKNOWLEDGMENT

This work was supported by NRF (National Honor Scientist Program: 2010-0020414, WCU: R32-2008-000-10180-0) and KISTI (KSC-2011-G3-02).

## ■ REFERENCES

- Stace, A. *Science* **2001**, *294*, 1292–1293.
- Singh, N. J.; Olleta, A. C.; Kumar, A.; Park, M.; Yi, H.-B.; Bandyopadhyay, I.; Lee, H. M.; Tarakeshwar, P.; Kim, K. S. *Theor. Chem. Acc.* **2005**, *115*, 127–135.
- Knipping, E. M.; Lakin, M. J.; Foster, K. L.; Jungwirth, P.; Tobias, D. J.; Gerber, R. B.; Dabdub, D.; Finlayson-Pitts, B. J. *Science* **2000**, *288*, 301–306.
- Miller, D. J.; Lisy, J. M. *J. Chem. Phys.* **2006**, *124*, 184301.
- Vaden, T. D.; Lisy, J. M.; Carnegie, P. D.; Dinesh Pillai, E.; Duncan, M. A. *Phys. Chem. Chem. Phys.* **2006**, *8*, 3078.
- Kolaski, M.; Lee, H. M.; Choi, Y. C.; Kim, K. S.; Tarakeshwar, P.; Miller, D. J.; Lisy, J. M. *J. Chem. Phys.* **2007**, *126*, 074302.
- Reinhardt, B. M.; Niedner-Schatteburg, G. *Phys. Chem. Chem. Phys.* **2002**, *4*, 1471–1477.
- Sobott, F.; Wattenberg, A.; Barth, H.-D.; Brutschy, B. *Int. J. Mass Spectrom.* **1999**, *187*, 271–279.
- Lee, H. M.; Min, S. K.; Lee, E. C.; Min, J.-H.; Odde, S.; Kim, K. S. *J. Chem. Phys.* **2005**, *122*, 064314.
- Karthikeyan, S.; Park, M.; Shin, I.; Kim, K. S. *J. Phys. Chem. A* **2008**, *112*, 10120–10124.
- Lehr, L.; Zanni, M. T.; Frischkorn, C.; Weinkauff, R.; Neumark, D. M. *Science* **1999**, *284*, 635–638.
- Robertson, W. H.; Johnson, M. A. *Science* **2002**, *298*, 69–69.
- Robertson, W. H.; Diken, E. G.; Price, E. A.; Shin, J.-W.; Johnson, M. A. *Science* **2003**, *299*, 1367–1372.
- Hurley, S. M.; Dermota, T. E.; Hydutsky, D. P.; Castleman, A. W., Jr. *Science* **2002**, *298*, 202–204.
- Kim, J.; Lee, H.; Suh, S.; Majumdar, D.; Kim, K. S. *J. Chem. Phys.* **2000**, *113*, 5259–5272.
- Odde, S.; Mhin, B. J.; Lee, S.; Lee, H. M.; Kim, K. S. *J. Chem. Phys.* **2004**, *120*, 9524.
- Wang, X.-B.; Kowalski, K.; Wang, L.-S.; Xantheas, S. S. *J. Chem. Phys.* **2010**, *132*, 124306.
- Yates, B. F.; Schaefer, H. F., III; Lee, T. J.; Rice, J. E. *J. Am. Chem. Soc.* **1988**, *110*, 6327–6332.
- Kemp, D. D.; Gordon, M. S. *J. Phys. Chem. A* **2005**, *109*, 7688–7699.
- Gruenloh, C. J.; Carney, J. R.; Arrington, C. A.; Zwier, T. S.; Fredericks, S. Y.; Jordan, K. D. *Science* **1997**, *276*, 1678–1681.
- Losada, M.; Leutwyler, S. *J. Chem. Phys.* **2002**, *117*, 2003–2016.
- Day, P.; Pachter, R.; Gordon, M. S.; Merrill, G. N. *J. Chem. Phys.* **2000**, *112*, 2063–2073.
- Bulusu, S.; Yoo, S.; Aprà, E.; Xantheas, S.; Zeng, X. C. *J. Phys. Chem. A* **2006**, *110*, 11781–11784.
- Yoo, S.; Aprà, E.; Zeng, X. C.; Xantheas, S. S. *J. Phys. Chem. Lett.* **2010**, *1*, 3122–3127.
- Lagutschenkov, A.; Fanourgakis, G.; Niedner-Schatteburg, G.; Xantheas, S. S. *J. Chem. Phys.* **2005**, *122*, 194310.
- Lenz, A.; Ojamäe, L. *J. Phys. Chem. A* **2006**, *110*, 13388–13393.
- Lee, H. M.; Suh, S. B.; Lee, J. Y.; Tarakeshwar, P.; Kim, K. S. *J. Chem. Phys.* **2000**, *112*, 9759.
- Lee, H.; Suh, S.; Lee, J.; Tarakeshwar, P.; Kim, K. *J. Chem. Phys.* **2001**, *113*, 3343.
- Shinohara, H.; Nagashima, U.; Nishi, N. *Chem. Phys. Lett.* **1984**, *111*, 511–513.
- Perrin, C. L.; Gipe, R. K. *Science* **1987**, *238*, 1393–1394.
- Fox, B. S.; Beyer, M. K.; Bondybey, V. E. *J. Phys. Chem. A* **2001**, *105*, 6386–6392.
- Diken, E. G.; Hammer, N. I.; Johnson, M. A.; Christie, R. A.; Jordan, K. D. *J. Chem. Phys.* **2005**, *123*, 164309.
- Pankewitz, T.; Lagutschenkov, A.; Niedner-Schatteburg, G.; Xantheas, S. S.; Lee, Y.-T. *J. Chem. Phys.* **2007**, *126*, 074307.
- Wang, Y.-S.; Chang, H.-C.; Jiang, J.-C.; Lin, S. H.; Lee, Y. T.; Chang, Y.-T. *J. Am. Chem. Soc.* **1998**, *120*, 8777–8788.
- Jiang, J.-C.; Chang, H.-C.; Lee, Y. T.; Lin, S. H. *J. Phys. Chem. A* **1999**, *103*, 3123–3135.
- Brugé, F.; Bernasconi, M.; Parrinello, M. *J. Am. Chem. Soc.* **1999**, *121*, 10883–10888.
- Chang, T.; Dang, L. X. *J. Chem. Phys.* **2003**, *118*, 8813–8820.
- Lee, H. M.; Tarakeshwar, P.; Park, J.; Kolaski, M. R.; Yoon, Y. J.; Yi, H.-B.; Kim, W. Y.; Kim, K. S. *J. Phys. Chem. A* **2004**, *108*, 2949–2958.
- Karthikeyan, S.; Singh, J. N.; Park, M.; Kumar, R.; Kim, K. S. *J. Chem. Phys.* **2008**, *128*, 244304.
- Douady, J.; Calvo, F.; Spiegelman, F. *J. Chem. Phys.* **2008**, *129*, 154305.
- Zhao, Y.-L.; Meot-Ner Mautner, M.; Gonzalez, C. *J. Phys. Chem. A* **2009**, *113*, 2967–2974.
- Kim, H.; Lee, H. M. *J. Phys. Chem. A* **2009**, *113*, 6859–6864.
- Morrell, T. E.; Shields, G. C. *J. Phys. Chem. A* **2010**, *114*, 4266–4271.
- Miyazaki, M.; Fujii, A.; Ebata, T.; Mikami, N. *Science* **2004**, *304*, 1134–1137.
- Mizuse, K.; Fujii, A. *J. Phys. Chem. Lett.* **2011**, *2*, 2130–2134.
- Shin, J.-W.; Hammer, N. I.; Diken, E. G.; Johnson, M. A.; Walters, R. S.; Jaeger, T. D.; Duncan, M. A.; Christie, R. A.; Jordan, K. D. *Science* **2004**, *304*, 1137–1140.
- Wu, C.-C.; Lin, C.-K.; Chang, H.-C.; Jiang, J.-C.; Kuo, J.-L.; Klein, M. L. *J. Chem. Phys.* **2005**, *122*, 074315.
- Iyengar, S. S.; Petersen, M. K.; Day, T. J. F.; Burnham, C. J.; Teige, V. E.; Voth, G. A. *J. Chem. Phys.* **2005**, *123*, 084309.
- Singh, N. J.; Park, M.; Min, S. K.; Suh, S. B.; Kim, K. S. *Angew. Chem., Int. Ed.* **2006**, *118*, 3879–3884.
- Elstner, M.; Porezag, D.; Jungnickel, G.; Elsner, J.; Haugk, M.; Frauenheim, T.; Suhai, S.; Seifert, G. *Phys. Rev. B* **1998**, *58*, 7260–7268.
- Li, Z.; Scheraga, H. A. *Proc. Natl. Acad. Sci. U.S.A.* **1987**, *84*, 6611–6615.
- Wales, D.; Hodges, M. P. *Chem. Phys. Lett.* **1998**, *286*, 65–72.
- Becke, A. D. *Phys. Rev. A* **1988**, *38*, 3098–3100.
- Lee, C.; Yang, W.; Parr, R. *Phys. Rev. B* **1988**, *37*, 785–789.
- Weigend, F.; Haser, M.; Patzelt, H.; Ahlrichs, R. *Chem. Phys. Lett.* **1998**, *294*, 143–152.
- Grimme, S.; Antony, J.; Schwabe, T.; Mck-Lichtenfeld, C. *Org. Biomol. Chem.* **2007**, *5*, 741.
- Gräfenstein, J.; Cremer, D. *J. Chem. Phys.* **2009**, *130*, 124105.
- Johnson, E. R.; Mackie, I. D.; DiLabio, G. A. *J. Phys. Org. Chem.* **2009**, *22*, 1127–1135.
- Zhao, Y.; Truhlar, D. G. *Theor. Chem. Acc.* **2007**, *120*, 215–241.

- (60) Dunning, T., Jr. *J. Chem. Phys.* **1989**, *90*, 1007–1023.
- (61) Kendall, R.; Dunning, T., Jr.; Harrison, R. *J. Chem. Phys.* **1992**, *96*, 6796–6806.
- (62) Frisch, M. J.; Trucks, G. W.; Schlegel, H. B.; Scuseria, G. E.; Robb, M. A.; Cheeseman, J. R.; Montgomery, J. A., Jr.; Vreven, T.; Kudin, K. N.; Burant, J. C.; Millam, J. M.; Iyengar, S. S.; Tomasi, J.; Barone, V.; Mennucci, B.; Cossi, M.; Scalmani, G.; Rega, N.; Petersson, G. A.; Nakatsuji, H.; Hada, M.; Ehara, M.; Toyota, K.; Fukuda, R.; Hasegawa, J.; Ishida, M.; Nakajima, T.; Honda, Y.; Kitao, O.; Nakai, H.; Klene, M.; Li, X.; Knox, J. E.; Hratchian, H. P.; Cross, J. B.; Bakken, V.; Adamo, C.; Jaramillo, J.; Gomperts, R.; Stratmann, R. E.; Yazyev, O.; Austin, A. J.; Cammi, R.; Pomelli, C.; Ochterski, J. W.; Ayala, P. Y.; Morokuma, K.; Voth, G. A.; Salvador, P.; Dannenberg, J. J.; Zakrzewski, V. G.; Dapprich, S.; Daniels, A. D.; Strain, M. C.; Farkas, O.; Malick, D. K.; Rabuck, A. D.; Raghavachari, K.; Foresman, J. B.; Ortiz, J. V.; Cui, Q.; Baboul, A. G.; Clifford, S.; Cioslowski, J.; Stefanov, B. B.; Liu, G.; Liashenko, A.; Piskorz, P.; Komaromi, I.; Martin, R. L.; Fox, D. J.; Keith, T.; Al-Laham, M. A.; Peng, C. Y.; Nanayakkara, A.; Challacombe, M.; Gill, P. M. W.; Johnson, B.; Chen, W.; Wong, M. W.; Gonzalez, C.; Pople, J. A. *Gaussian 03*, Revision C.02; Gaussian, Inc.: Wallingford, CT, 2004.
- (63) Ahlrichs, R.; Bar, M.; Haser, M.; Horn, H.; Kolmel, C. *Chem. Phys. Lett.* **1989**, *162*, 165–169. {Available from TURBOMOLE-Program Package for ab initio Electronic Structure Calculations. <http://www.turbomole.com> (accessed Sep 27, 2011).}
- (64) Car, R.; Parrinello, M. *Phys. Rev. Lett.* **1985**, *55*, 2471–2474.
- (65) CPMD; IBM Corp.: Armonk, NY, 1990; MPI für Festkörperforschung Stuttgart: Stuttgart, Germany, 1997. <http://www.cpmc.org/> (accessed Sep 27, 2011).
- (66) Troullier, N.; Martins, J. *Phys. Rev. B* **1991**, *43*, 1993–2006.
- (67) Schmitt, U. W.; Voth, G. A. *J. Chem. Phys.* **1999**, *111*, 9361.
- (68) Park, M.; Shin, I.; Singh, N. J.; Kim, K. S. *J. Phys. Chem. A* **2007**, *111*, 10692–10702.

# Extensions of the S66 Data Set: More Accurate Interaction Energies and Angular-Displaced Nonequilibrium Geometries

Jan Řezáč,<sup>\*,†</sup> Kevin E. Riley,<sup>†</sup> and Pavel Hobza<sup>†,‡</sup>

<sup>†</sup>Institute of Organic Chemistry and Biochemistry, Academy of Sciences of the Czech Republic and Center for Biomolecules and Complex Molecular Systems, 166 10 Prague, Czech Republic

<sup>‡</sup>Regional Centre of Advanced Technologies and Materials, Department of Physical Chemistry, Palacky University, 771 46 Olomouc, Czech Republic

## S Supporting Information

**ABSTRACT:** We present two extensions of the recently published S66 data set [Řezáč, Riley, Hobza; DOI: 10.1021/ct2002946]. Interaction energies for the equilibrium geometry complexes have been recalculated using a triple- $\zeta$  basis set for the CCSD(T) term in the CCSD(T)/CBS scheme. This allows for the extrapolation of this term to the complete basis set limit, improving accuracy by almost 1 order of magnitude compared to the scheme previously used for the S66 set. Now, we estimate the largest error in the set to be about 1%. Validation of several methods against the new data indicates the exceptional robustness and accuracy of the SCS-MI-CCSD method. The second extension improves the coverage of nonequilibrium geometries. We introduce a new data set, S66a8, that samples intermolecular angular degrees of freedom in the S66 complexes. For each of the 66 complexes, eight displaced geometries have been constructed, systematically sampling possible rotations of the monomers. Interaction energies in this set are calculated at the CCSD(T)/CBS level consistently with the earlier introduced S66x8 data set that samples the intermolecular distance.

## INTRODUCTION

The importance of accurate *ab initio* calculations for evaluation of the performance of more approximate methods is now widely recognized. One of the fields where such calculations serve as very important benchmarks is the study of noncovalent interactions. Multiple databases of reference data covering this topic have been published in the past decade.<sup>1–5</sup> Recently, we introduced the S66 data set,<sup>6</sup> which was designed to overcome multiple limitations of the previously available sets. The most important improvements are the increased size of the set and more balanced coverage of different types of interactions, which was achieved by careful selection of the complexes. The interaction energies of these complexes were calculated consistently at the CCSD(T) level and extrapolated to the complete basis set limit (CBS). For a more detailed description of the S66 set, we refer the reader to the original paper.<sup>6</sup>

If these complexes are used as a model for interaction between atomic groups in large molecular systems, it is also necessary to consider nonequilibrium geometries. Although the interaction is strongest in equilibrium, there are much greater numbers of weaker interactions acting over longer distances in large condensed systems. To address this, we have published dissociation curves for all of the 66 complexes of the S66 data set, also calculated at the CCSD(T)/CBS level.<sup>6</sup>

In this work, we present two extensions of the S66 data set. First, we have recalculated the interaction energies in the S66 set using a larger basis set in order to improve the accuracy toward the complete basis set limit. The new benchmark interaction energies are based on extrapolation of the CCSD(T) term from double- and triple- $\zeta$  basis sets with diffuse functions on first row

atoms, improving the accuracy over the previously published results by almost 1 order of magnitude. Second, we introduce the S66a8 data set, which samples the intermolecular angular degrees of freedom in all of the 66 complexes. Together with the S66x8 set, sampling intermolecular distances, our results represent the largest body of accurate data available for nonequilibrium structures of molecular complexes, calculated consistently at the same level.

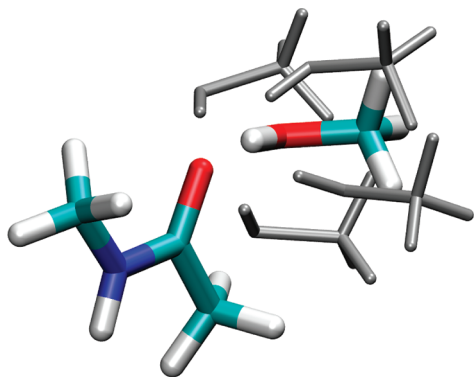
We use the new results as a reference for assessment of the accuracy of more approximate methods. We focus on two methods that have been shown in our previous paper to have very good performance to cost ratios, MP2.5<sup>7</sup> and SCS-MI-CCSD.<sup>8</sup> A comparison of results calculated with different basis sets allows us to discuss the transferability of the parameters used in these methods.

## METHODS

**Basis Sets.** Throughout this work, we use the correlation consistent basis sets of Dunning.<sup>9</sup> The full names of the basis sets, cc-pVXZ (X = D, T, Q) and aug-cc-pVXZ (augmented with diffuse functions<sup>10</sup>), are abbreviated as XZ and aXZ in this work.

**Benchmark Calculations.** The CCSD(T)/CBS interaction energy is obtained as a sum of the Hartree–Fock (HF) interaction energy, the MP2 correlation energy extrapolated to CBS from large basis sets, and a  $\Delta$ CCSD(T) correction ( $\Delta E^{\text{CCSD(T)}} - \Delta E^{\text{MP2}}$ ) calculated with a smaller basis set. For the large set of

**Published:** October 03, 2011



**Figure 1.** Angular-displaced geometries obtained by rotating one of the monomers (in gray), compared to the original geometry of the complex (in color).

angular displacement calculations, we use the same scheme as for the original S66 and S88x8 sets. The MP2 term is extrapolated from the aTZ and aQZ basis sets, and the  $\Delta\text{CCSD(T)}$  correction is calculated with the aDZ basis set. Extrapolations to the CBS limit are done using the formula of Helgaker et al.<sup>11</sup> All interaction energy calculations presented here are corrected for the basis set superposition error using the counterpoise scheme.

Here, we have recalculated the S66 binding energies using a more accurate scheme, in which the  $\Delta\text{CCSD(T)}$  correction is extrapolated to the CBS from two calculations. Because the calculation of the largest systems in the set using the aTZ basis set is not possible with our current resources, we use a smaller basis set that still retains most of the qualities of the aTZ basis set, combining the TZ basis for hydrogens and aTZ for all other atoms.<sup>12</sup> This basis set is often referred to as the heavy-augmented basis, and we use the abbreviation haTZ. Recently, it has also been published under the name jul-cc-pVTZ.<sup>13</sup> An analogously constructed haDZ basis set is used along with haTZ to extrapolate the  $\Delta\text{CCSD(T)}$  correction. This scheme achieves binding energy values close to the most accurate results available for the S22 set,<sup>14</sup> where aDZ and aTZ basis sets were used for the extrapolation of  $\Delta\text{CCSD(T)}$ .

To assess the improvement brought by this approach, we compare these two schemes, along with other possible combinations of basis sets, on a set of small model systems introduced earlier.<sup>15</sup> Description and geometries of these complexes are available in the original paper.

**S66a8 Geometries.** The angular-displaced geometries have been prepared from the S66 geometries in the following way: For each complex, the principal plane of the monomers was identified. In simple cases, these are defined by symmetry; in more complex molecules, the definition was only approximate but reflects the shape of the molecule. Eight geometries are built for each system: each of the monomers is rotated in both directions ( $\pm$ ) in the molecular plane (coordinate  $y$ ) and perpendicular to it (coordinate  $z$ ) by  $30^\circ$ . In order to sample the angular displacements but not the nonequilibrium intermolecular distances (already covered by the S66x8 set), the intermolecular distances in the S66a8 complexes are optimized at the RI-MP2/TZ level, with the counterpoise correction, while all other degrees of freedom are kept fixed. An example of the displaced geometries obtained by the rotation of one monomer is shown in Figure 1. We do not explore rotation around the intermolecular axis, the interaction energy is not sensitive to such a change of the

**Table 1.** Accuracy of the CCSD(T)/CBS Scheme with Different Basis Sets Used for the Calculation of the  $\Delta\text{CCSD(T)}$  Term, Measured As the Root Mean Square Error in a Set of 10 Small Complexes for Which Accurate Estimates of the CCSD(T)/CBS Interaction Energies Are Available

$\Delta\text{CCSD(T)}$ basis set(s)	RMSE (kcal/mol)
aDZ	0.080
aTZ	0.020
TZ	0.107
haTZ	0.033
haDZ $\rightarrow$ haTZ extrapolation	0.009

geometry in most of the studied complexes. No symmetry is assumed in the generation of the displaced geometries; therefore, the complete set contains several pairs of structures built from symmetrical minima that are equivalent in energy. For simplicity, we provide the complete set of geometries, but identical results can be achieved by eliminating the duplicates and setting the weight of the results to two in the statistical analysis.

**Methods Tested.** There are two methods that we investigate in this work in greater detail; these have been parametrized to describe noncovalent interactions. MP2.5<sup>7</sup> is a variant of MP3 where the third-order contribution is scaled by one-half. SCS-MI-CCSD<sup>8</sup> is a spin-component-scaled CCSD method optimized specifically for noncovalent interactions.

Several other methods have been tested on the S66a8 set. In addition to the ones described above, these are MP2 in multiple basis sets, spin-component-scaled MP2<sup>16</sup> (SCS-MP2) and its reparameterization for noncovalent interactions SCS-MI-MP2,<sup>17</sup> dispersion weighted MP2<sup>18</sup> (DW-MP2), MP3, CCSD, and the original version of spin-component-scaled CCSD<sup>19</sup> (SCS-CCSD). Details on these methods and values of the parameters used can be found in the original paper on the S66 data set.<sup>6</sup>

To extrapolate all of these methods to the CBS limit, we employ a scheme analogous to the CCSD(T)/CBS calculations; the result is built from accurate extrapolation of the MP2 energy and a higher order correction (e.g.,  $\Delta E^{\text{SCS-MI-CCSD}} - \Delta E^{\text{MP2}}$ ) calculated in a smaller basis set or extrapolated independently, using the same basis set(s) as the benchmark in the given data set.

**Computational Details.** Optimization of the intermolecular distance in the S66a8 set was carried out in Turbomole 6.2.<sup>20</sup> All interaction energy calculations were performed using the Molpro 2010 package,<sup>21</sup> using density fitting for the MP2 calculations.

## RESULTS AND DISCUSSION – S66 DATA SET

**CCSD(T)/CBS Extrapolation.** In the set of 10 model complexes,<sup>15</sup> we compare the effect of the basis sets and extrapolation schemes applicable to the calculation of the  $\Delta\text{CCSD(T)}$  correction in the S66 set (Table 1). The results are compared to our best estimates of the CCSD(T)/CBS interaction energies directly extrapolated from the aTZ and aQZ basis sets. First, to justify the use of the heavy-augmented basis sets, we discuss the differences between the TZ, aTZ, and haTZ basis sets, also including aDZ for comparison. It is clear that the heavy-augmented basis set yields results close to the fully augmented ones and that the improvement over the TZ basis set with no diffuse functions is substantial. These results are supported by previous studies of this basis set.<sup>12,13</sup>

**Table 2. New, More Accurate CCSD(T)/CBS Interaction Energies (in kcal/mol) for the Complexes in the S66 Data Set<sup>a</sup>**

hydrogen bonds		$\Delta E$
1	water...water	-5.01
2	water...MeOH	-5.70
3	water...MeNH <sub>2</sub>	-7.04
4	water...peptide	-8.22
5	MeOH...MeOH	-5.85
6	MeOH...MeNH <sub>2</sub>	-7.67
7	MeOH...peptide	-8.34
8	MeOH...water	-5.09
9	MeNH <sub>2</sub> ...MeOH	-3.11
10	MeNH <sub>2</sub> ...MeNH <sub>2</sub>	-4.22
11	MeNH <sub>2</sub> ...peptide	-5.48
12	MeNH <sub>2</sub> ...water	-7.40
13	peptide...MeOH	-6.28
14	Peptide...MeNH <sub>2</sub>	-7.56
15	peptide...peptide	-8.72
16	peptide...water	-5.20
17	uracil...uracil (BP)	-17.45
18	water...pyridine	-6.97
19	MeOH...pyridine	-7.51
20	AcOH...AcOH	-19.41
21	AcNH <sub>2</sub> ...AcNH <sub>2</sub>	-16.52
22	AcOH...uracil	-19.78
23	AcNH...uracil	-19.47
dispersion		$\Delta E$
24	benzene...benzene ( $\pi-\pi$ )	-2.72
25	pyridine...pyridine ( $\pi-\pi$ )	-3.80
26	uracil...uracil ( $\pi-\pi$ )	-9.75
27	benzene...pyridine ( $\pi-\pi$ )	-3.34
28	benzene...uracil ( $\pi-\pi$ )	-5.59
29	pyridine...uracil ( $\pi-\pi$ )	-6.70
30	benzene...ethene	-1.36
31	uracil...ethene	-3.33
32	uracil...ethyne	-3.69
33	pyridine...ethene	-1.80
34	pentane...pentane	-3.76
35	neopentane...pentane	-2.60
36	neopentane...neopentane	-1.76
37	cyclopentane...neopentane	-2.40
38	cyclopentane...cyclopentane	-2.99
39	benzene...cyclopentane	-3.51
40	benzene...neopentane	-2.85
41	uracil...pentane	-4.81
42	uracil...cyclopentane	-4.09
43	uracil...neopentane	-3.69
44	ethene...pentane	-1.99
45	ethyne...pentane	-1.72
46	peptide...pentane	-4.26
others		$\Delta E$
47	benzene...benzene (TS)	-2.83
48	pyridine...pyridine (TS)	-3.51
49	benzene...pyridine (TS)	-3.29

**Table 2. Continued**

others		$\Delta E$
50	benzene...ethyne (CH... $\pi$ )	-2.86
51	ethyne...ethyne (TS)	-1.54
52	benzene...AcOH (OH... $\pi$ )	-4.73
53	benzene...AcNH <sub>2</sub> (NH... $\pi$ )	-4.40
54	benzene...water (OH... $\pi$ )	-3.29
55	benzene...MeOH (OH... $\pi$ )	-4.17
56	benzene...MeNH <sub>2</sub> (NH... $\pi$ )	-3.20
57	benzene...peptide (NH... $\pi$ )	-5.26
58	pyridine...pyridine (CH...N)	-4.24
59	ethyne...water (CH...O)	-2.93
60	ethyne...AcOH (OH... $\pi$ )	-4.97
61	pentane...AcOH	-2.91
62	pentane...AcNH <sub>2</sub>	-3.53
63	benzene...AcOH	-3.75
64	peptide...ethene	-3.00
65	pyridine...ethyne	-4.10
66	MeNH <sub>2</sub> ...pyridine	-3.97

<sup>a</sup>The CCSD(T)/CBS is based on extrapolation of the MP2 correlation energy from the aTZ and aQZ basis sets with the  $\Delta$ CCSD(T) correction extrapolated from the haDZ and haTZ basis sets.

**Table 3. Errors of the MP2.5 and SCS-MI-CCSD Methods with Higher Order Terms Calculated in Different Basis Sets Compared to the CCSD(T)/CBS Reference with  $\Delta$ CCSD(T) Term Calculated in aDZ Basis Set and Extrapolated from haDZ and haTZ Basis Sets**

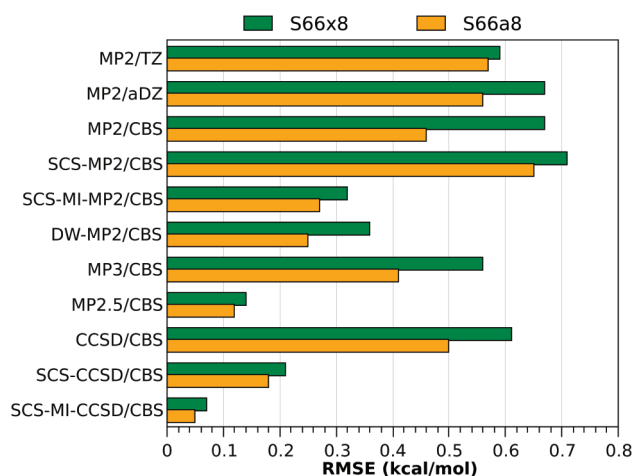
method	$\Delta$ CCSD(T) basis		
	set in reference basis set	aDZ RMSE (kcal/mol)	haDZ→haTZ RMSE (kcal/mol)
MP2.5	aDZ	0.16	0.21
MP2.5	aDZ→aTZ	0.16	0.22
SCS-MI-CCSD	aDZ	0.08	0.14
SCS-MI-CCSD	haDZ→haTZ	0.14	0.07
SCS-MI-CCSD	aDZ→aTZ	0.14	0.07

To improve the accuracy of the final CBS estimate further, it is now possible to extrapolate the  $\Delta$ CCSD(T) from the haDZ/haTZ series. This scheme yields a root-mean-square error (RMSE) of only 0.009 kcal/mol for the model complexes, lowering the error by almost 1 order of magnitude compared to the scheme used previously (using the aDZ basis set, RMSE 0.080 kcal/mol). The improvement obtained by the extrapolation of the  $\Delta$ CCSD(T) term, compared to the use of a single calculation in the haTZ basis set, is also substantial. We find this to be the most accurate approach practically applicable even to the largest complexes in the S66 set. Therefore, we consider these results to be the new benchmark for the S66 set.

The largest error in the set of model complexes is 0.7%, compared to 2.5% in the previously used scheme. Therefore, we estimate that the largest errors in the S66 set are lower than 1%.

**Benchmark Results.** Interaction energies in the S66 data set obtained using the new, more accurate CCSD(T)/CBS extrapolation are listed in Table 2. These results are also available online through the BEGDB database.<sup>22</sup>





**Figure 2.** Root mean square error of selected methods in the S66x8 (dissociation curves) and S66a8 (angular displacements) data sets.

Comparing the CCSD(T)/CBS scheme using the aDZ basis set for the  $\Delta$ CCSD(T) term used previously with the new results, we find an average unsigned error of 0.08 kcal/mol (1.5%) and a RMSE of 0.10 kcal/mol. The strength of hydrogen bonds had been systematically underestimated, while dispersion and mixed-type interactions had been overestimated. This reflects the sign of the  $\Delta$ CCSD(T) term, whose magnitude is smaller in the unsaturated basis set.

**Methods Tested.** The overall performance of the methods originally tested on the S66 set does not change significantly when the new, more accurate benchmark is used. Here, we will focus only on the effect of the basis set (used for the higher-order correction in the CBS scheme) on two methods found to be the most accurate in their categories, MP2.5 and SCS-MI-CCSD. Table 3 lists RMS errors in the S66 set for multiple combinations of basis sets, both in the studied method and in the reference.

The performance of MP2.5 is slightly worse when the more accurate reference is used, regardless of the basis set used for the MP2.5 calculation. This indicates that the source of the error comes from the approximations in the method itself and not the basis set. Some improvement can be achieved by optimization of the scaling factor, discussed in detail in a separate paper.<sup>23</sup>

The behavior of SCS-MI-CCSD is surprisingly consistent. This method is able to accurately reproduce the CCSD(T) results calculated using the same basis set or extrapolation scheme (RMSE 0.08 and 0.07 kcal/mol for aDZ and CBS extrapolation). When compared crosswise, the error of the method combined with the difference between the references, the RMSE is larger (0.14 kcal/mol). On the basis of these results, it is obvious that the SCS-MI-CCSD method is very robust and provides results very close to CCSD(T) in a given basis set.

## RESULTS AND DISCUSSION – S66A8 DATA SET

**Benchmark Results.** The complete set of S66a8 geometries, benchmark interaction energies, and results from the tested methods are available online in the BEGDB database,<sup>22</sup> where it is possible to browse, plot, and download the data.

The most important information that can be derived from these data is how the magnitude of the interaction in a given complex decreases when the geometry is displaced. We list the average interaction energy per complex as a percentage of the

interaction energy in equilibrium for the individual groups of complexes of the S66 set: hydrogen bonds, 77%; dispersion-dominated complexes, 67%; others, 77%. It seems counter-intuitive that hydrogen bonds, which are known to be sensitive to the mutual orientation of the interacting molecules, do not exhibit the largest decrease. On the other hand, the hydrogen bonding motif is conserved rather well upon rotation by 30°, while the rotation of larger molecules in the dispersion-dominated complexes leads to a large decrease of the contact surface that determines the strength of the interaction. Dispersion interactions also decay with interatomic distance faster than electrostatic interactions.

**Methods Tested.** We used the S66a8 data set to test the same set of methods studied on the S66x8 set. The results, plotted in Figure 2 and listed in Table S1 in the Supporting Information, are very similar. Therefore, we refer the reader to the discussion of the performance of individual methods in the original paper.<sup>6</sup> In general, the errors in the S66a8 set are slightly lower even when the relative error (with respect to average interaction energy in the set) is considered. The major source of this discrepancy stems from the shorter than equilibrium geometries in the S66x8 set, where the errors are larger than at (or above) equilibrium distances.

However, the similarity of the results presented here does not make the S66a8 set redundant. It is not surprising that high-level QM methods describe the entire potential energy surface similarly. The main reason for building the S66a8 set was to aid in the development and testing of more approximate methods, such as DFT-D, semiempirical QM methods, or force fields, where significant differences can be expected.

## CONCLUSIONS

We present new benchmark interaction energies for the S66 data set. Using the extrapolation of the  $\Delta$ CCSD(T) correction from haDZ and haTZ basis sets, we have improved the average accuracy by almost 1 order of magnitude. On the basis of small model calculations, we estimate the largest error in the S66 set to be approximately 1% of the interaction energy when compared to CCSD(T) complete basis set limit.

The SCS-MI-CCSD method was found to be very robust, as the scaling coefficients are transferable between different basis sets. For a given basis set or extrapolation scheme, the results reproduce reference CCSD(T) calculations using the same basis sets with high accuracy.

Extension of the S66 set to nonequilibrium geometries obtained by rotation of the monomers in the complex, the S66a8 set, is also presented. Here, the high-level QM methods tested yield errors similar to those in the S66x8 data set, but we expect the S66a8 set to be useful for development and validation of more approximate methods where larger differences can be found. When the S66x8 and S66a8 sets are combined, they constitute 1056 points covering the most important coordinates of the intermolecular potential energy surface of the complexes in the S66 set.

Geometries of the complexes, the benchmark CCSD(T)/CBS interaction energies, and results from the other methods discussed in the paper are freely available in the online database at [www.begdb.com](http://www.begdb.com).<sup>22</sup>

## ASSOCIATED CONTENT

**Supporting Information.** Table S1, listing errors of all of the methods tested on S66a8 and S66x8 sets and geometries and

benchmark interaction energies for the data sets featured in this paper are provided. This material is available free of charge via the Internet at <http://pubs.acs.org>.

## AUTHOR INFORMATION

### Corresponding Author

\*Fax: +420 220 410 320. E-mail: [rezac@uochb.cas.cz](mailto:rezac@uochb.cas.cz).

## ACKNOWLEDGMENT

This work was a part of Research Project No. Z40550506 of the Institute of Organic Chemistry and Biochemistry, Academy of Sciences of the Czech Republic, and was supported by Grants No. LC512 and MSM6198959216 from the Ministry of Education, Youth and Sports of the Czech Republic. It was also supported by the operational program Research and Development for Innovations of European Social Fund (CZ.1.05/2.1.00/03.0058). The support of Praemium Academiae, Academy of Sciences of the Czech Republic, awarded to P.H. in 2007 is also acknowledged.

## REFERENCES

- (1) Jurečka, P.; Šponer, J.; Černý, J.; Hobza, P. *Phys. Chem. Chem. Phys.* **2006**, *8*, 1985.
- (2) Gráfová, L.; Pitoňák, M.; Řezáč, J.; Hobza, P. *J. Chem. Theory Comput.* **2010**, *6*, 2365–2376.
- (3) Goerigk, L.; Grimme, S. *J. Chem. Theory Comput.* **2011**, *7*, 291–309.
- (4) Zhao, Y.; Truhlar, D. G. *J. Chem. Theory Comput.* **2005**, *1*, 415–432.
- (5) Burns, L. A.; Mayagoitia, A. V.; Sumpter, B. G.; Sherrill, C. D. *J. Chem. Phys.* **2011**, *134*, 084107.
- (6) Řezáč, J.; Riley, K. E.; Hobza, P. *J. Chem. Theory Comput.* **2011**, *7*, 2427–2438.
- (7) Pitoňák, M.; Neogrady, P.; Černý, J.; Grimme, S.; Hobza, P. *ChemPhysChem* **2009**, *10*, 282–289.
- (8) Pitoňák, M.; Řezáč, J.; Hobza, P. *Phys. Chem. Chem. Phys.* **2010**, *12*, 9611.
- (9) Dunning, T. H. *J. Chem. Phys.* **1989**, *90*, 1007.
- (10) Woon, D. E.; Dunning, T. H. *J. Chem. Phys.* **1994**, *100*, 2975.
- (11) Halkier, A.; Helgaker, T.; Jorgensen, P.; Klopper, W.; Koch, H.; Olsen, J.; Wilson, A. K. *Chem. Phys. Lett.* **1998**, *286*, 243–252.
- (12) ElSohly, A. M.; Tschumper, G. S. *Int. J. Quantum Chem.* **2009**, *109*, 91–96.
- (13) Papajak, E.; Zheng, J.; Xu, X.; Leverentz, H. R.; Truhlar, D. G. *J. Chem. Theory Comput.*, **2011**; DOI: 10.1021/ct200106a.
- (14) Takatani, T.; Hohenstein, E. G.; Malagoli, M.; Marshall, M. S.; Sherrill, C. D. *J. Chem. Phys.* **2010**, *132*, 144104.
- (15) Řezáč, J.; Hobza, P. *J. Chem. Theory Comput.* **2011**, *7*, 685–689.
- (16) Grimme, S. *J. Chem. Phys.* **2003**, *118*, 9095.
- (17) Distasio, R.; Head-Gordon, M. *Mol. Phys.* **2007**, *105*, 1073–1083.
- (18) Marchetti, O.; Werner, H.-J. *J. Phys. Chem. A* **2009**, *113*, 11580–11585.
- (19) Takatani, T.; Hohenstein, E. G.; Sherrill, C. D. *J. Chem. Phys.* **2008**, *128*, 124111.
- (20) TURBOMOLE, v6.2; University of Karlsruhe and Forschungszentrum Karlsruhe GmbH: Karlsruhe, Germany, 2010.
- (21) Werner, H.-J.; Knowles, P. J.; Manby, F. R.; Schütz, M.; et al. *MOLPRO*, version 2010.1; Cardiff University: Cardiff, U.K.; Universität Stuttgart: Stuttgart, Germany, 2010.
- (22) Řezáč, J.; Jurečka, P.; Riley, K. E.; Černý, J.; Valdes, H.; Pluháčková, K.; Berka, K.; Řezáč, T.; Pitoňák, M.; Vondrášek, J.; Hobza, P. *Collect. Czech. Chem. Commun.* **2008**, *73*, 1261–1270.
- (23) Riley, K. E.; Řezáč, J.; Hobza, P. Manuscript submitted.

# How Different Are Aromatic $\pi$ Interactions from Aliphatic $\pi$ Interactions and Non- $\pi$ Stacking Interactions?

Kwang S. Kim,\* S. Karthikeyan, and N. Jiten Singh

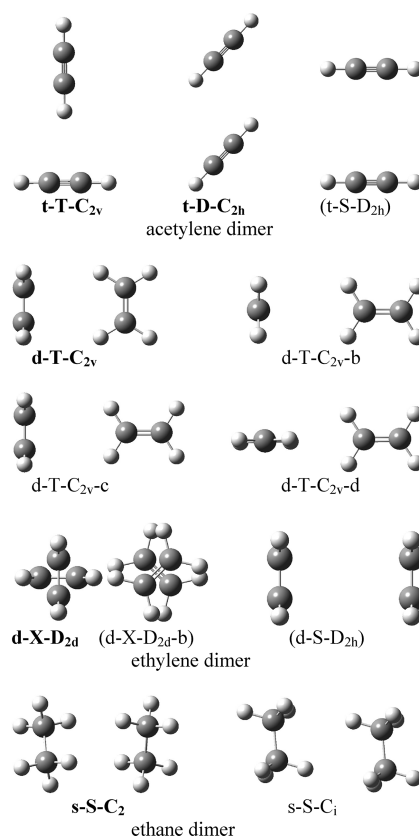
Center for Superfunctional Materials, Department of Chemistry, Pohang University of Science and Technology, San 31, Hyojadong, Namgu, Pohang 790-784, Korea

**S** Supporting Information

**ABSTRACT:** We compare aromatic  $\pi$  interactions with aliphatic  $\pi$  interactions of double- and triple-bonded  $\pi$  systems and non- $\pi$  stacking interactions of single-bonded  $\sigma$  systems. The model dimer systems of acetylene ( $C_2H_2$ )<sub>2</sub>, ethylene ( $C_2H_4$ )<sub>2</sub>, ethane ( $C_2H_6$ )<sub>2</sub>, benzene ( $C_6H_6$ )<sub>2</sub>, and cyclohexane ( $C_6H_{12}$ )<sub>2</sub> are investigated. The ethylene dimer has large dispersion energy, while the acetylene dimer has strong electrostatic energy. The aromatic  $\pi$  interactions are strong with particularly large dispersion and electrostatic energies, which would explain why aromatic compounds are frequently found in crystal packing and molecular self-engineering. It should be noted that the difference in binding energy between the benzene dimer (aromatic–aromatic interactions) and the cyclohexane dimer (aliphatic–aliphatic interactions) is not properly described in most density functionals.

Given that the interactions involved in  $\pi$  systems<sup>1–5</sup> are very important in molecular/biomolecular recognition,<sup>6–9</sup> assembly,<sup>10–12</sup> and engineering,<sup>13,14</sup> there have been numerous studies on  $\pi$  interactions.<sup>15–53</sup> One might speculate if there are significant differences in dispersion energies between aromatic  $\pi$  interactions and aliphatic  $\pi$  interactions<sup>54</sup> or non- $\pi$  stacking interactions and between single, double, and triple-bonded  $\pi$  systems. In this regard, it is necessary to compare the intermolecular interaction energies for dimers of acetylene HC≡CH (triple bond; t-), ethylene H<sub>2</sub>C=CH<sub>2</sub> (double bond; d-), ethane H<sub>3</sub>C–CH<sub>3</sub> (single bond; s-), benzene ( $\text{---CH---CH---}$ )<sub>3</sub> (aromatic bond; a-), and cyclohexane ( $\text{---CH}_2\text{---CH}_2\text{---}$ )<sub>3</sub> (cyclic single bond; h-). To this end, we need to focus our attention on the dimerization energy at high levels of theory and its energy decomposition. Since the dispersion energy is very important in  $\pi$  interactions, this study requires the complete basis set (CBS) limit binding energies at the level of coupled cluster theory with single, double, and perturbative triple excitations [CCSD(T)]. The strength of the  $\pi$ -interactions is determined by the combined effect of attractive forces (electrostatic, dispersive, and inductive) and repulsive forces (electrostatic, exchange repulsion). Each of these components shows distinctive differences in physical origin, magnitude, and directionality of the molecular interaction. This investigation is done by using symmetry adapted perturbation theory (SAPT) calculations. On the basis of the above calculations, we show the importance of aromatic  $\pi$  interactions in crystal packing and molecular self-engineering. In addition, since all of these results are very useful to test the reliability of various functionals for density functional theory (DFT), their validity for  $\pi$ -interactions along with their strength and weakness is assessed.

To search for the lowest energy structures of acetylene, ethylene, ethane, benzene, and cyclohexane dimers, we investigated diverse topologically different conformers using a few different types of DFT calculations. To confirm the minimum energy structures for the acetylene, ethylene, and ethane dimers,

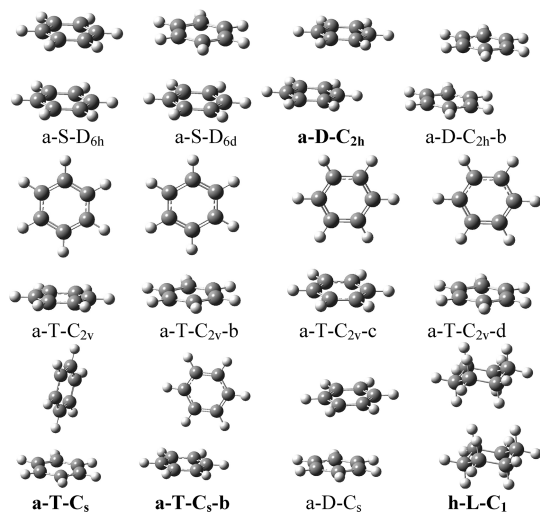


**Figure 1.** Low energy structures of the acetylene(t-) dimer, the ethylene(d-) dimer, and the ethane(s-) dimer; (t-, triple bonded; d-, double bonded; s-, single bonded).

**Received:** August 21, 2011

**Published:** September 22, 2011

frequency calculations were carried out at the DFT and Moller–Plesset second order perturbation (MP2) theory levels. Then, the low-lying energy structures were optimized with the basis set superposition error (BSSE) correction at the MP2 level using the aug-cc-pVDZ (aVDZ) and aug-cc-pVTZ (aVTZ) basis



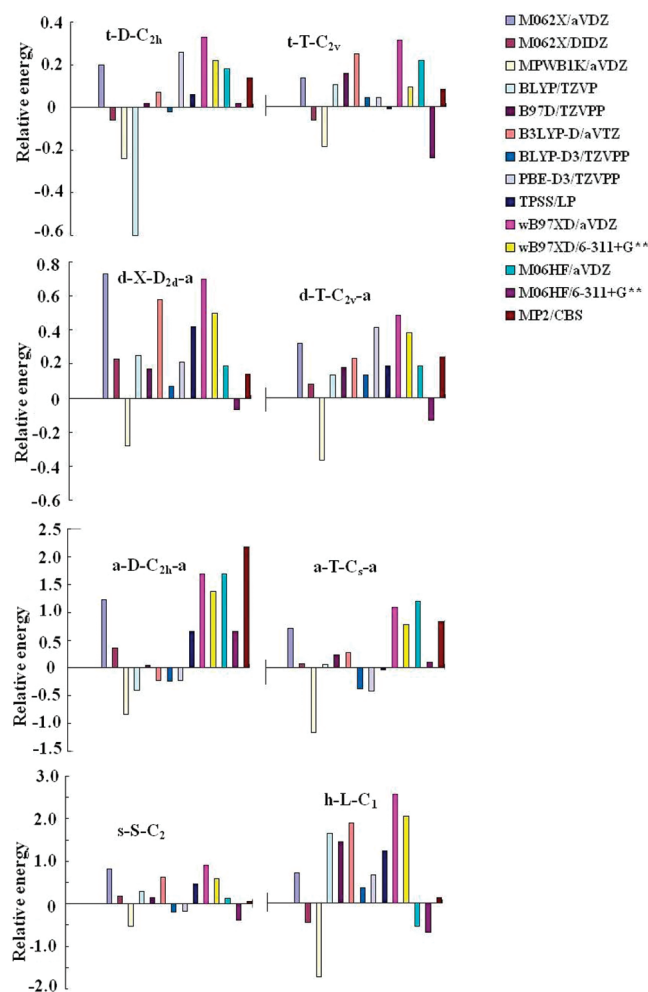
**Figure 2.** Low energy structures of the benzene(a-) dimer and the cyclohexane(h-) dimer (a-, aromatic bonded; h-, cyclohexane single bonded).

sets for the acetylene, ethylene, and ethane dimers and at the resolution of identity approximation (RI) of MP2 (RI-MP2)/aVDZ level for the benzene and cyclohexane dimers. Single-point MP2/aug-cc-pVQZ (aVQZ) calculations were performed to estimate the CBS limit binding energies. The estimated CBS limit values were evaluated on the basis of the extrapolation method exploiting the fact that the basis set error in the electron correlation energy is proportional to  $N^{-3}$  for the aug-cc-pVNZ (aVNZ) basis set ( $E_{\text{CBS}} = [E_N^* N^3 - E_{N-1}(N-1)^3] / [N^3 - (N-1)^3]$ ).<sup>55,56</sup> We also carried out CCSD(T)/aVDZ optimization and single point BSSE-corrected CCSD(T)/aVDZ and CCSD(T)/aVTZ calculations on the CCSD(T)/aVDZ geometries. However, in the case of cyclohexane dimer, we carried out the same calculations on the BSSE-corrected MP2/aVDZ (RIMP2) geometries. Given that the difference in binding energy between MP2/aVNZ and CCSD(T)/aVNZ does not change significantly with increasing basis set size, we obtained the estimated CCSD(T)/CBS energies by assuming that the difference in binding energies between MP2/aVDZ and MP2/CBS calculations is similar to that between CCSD(T)/aVDZ and CCSD(T)/CBS calculations ( $E_{\text{CCSD(T)/CBS}} = E_{\text{CCSD(T)/aVDZ}} + (E_{\text{MP2/CBS}} - E_{\text{MP2/aVDZ}})$ ).<sup>56–58</sup> However, the CCSD(T)/CBS values are also obtained from the extrapolation based on the CCSD(T)/aVDZ and CCSD(T)/aVTZ values. Since the MP2 dispersion energy corrections are not so reliable, the latter extrapolation method could be a better choice. The BLYP/TZVP, B97-D/TZVP,<sup>59</sup> B3LYP-D/aVTZ,<sup>60</sup> BLYP-D3/TZVPP,

**Table 1.** DFT, MP2/CBS, and CCSD(T)/CBS Binding Energies ( $-\Delta E_e$  in kcal/mol) for Low Energy Structures of the Triple-Bonded Acetylene (t-), Double-Bonded Ethylene (d-), Single-Bonded Ethane (s-), Aromatic-Bonded Benzene (a-) and Cyclic Single-Bonded Cyclohexane (h-) Dimers<sup>a</sup>

	M06-2X/aVDZ	M06-2X/DIDZ	MPWB1K/aVDZ	BLYP/TZVP	B97-D/TZV2P	B3LYP-D/aVTZ	MP2/CBS	CCSD(T)/CBS <sup>b</sup>
t-D-C <sub>2h</sub>	1.57(+0.20)	1.31(−0.06)	1.13(−0.24)	1.31(−0.6)	1.39(+0.02)	1.44(+0.07)	1.51(+0.14)	1.37[1.33]
t-T-C <sub>2v</sub>	1.68(+0.13)	1.49(−0.06)	1.37(−0.18)	1.65(+0.10)	1.70(+0.15)	1.79(+0.24)	1.63(+0.08)	1.55[1.45]
d-X-D <sub>2d</sub>	2.13(+0.73)	1.63(+0.23)	1.12(−0.28)	1.65(+0.25)	1.57(+0.17)	1.98(+0.58)	1.54(+0.14)	1.40[1.41]
d-T-C <sub>2v</sub>	1.34(+0.36)	1.07(+0.09)	0.57(−0.41)	1.13(+0.15)	1.18(+0.20)	1.24(+0.26)	1.25(+0.27)	0.98[1.00]
s-S-C <sub>2</sub>	2.25(+0.82)	1.61(+0.18)	0.90(−0.53)	1.71(+0.28)	1.57(+0.14)	2.09(+0.62)	1.37(+0.06)	1.43[1.33]
a-D-C <sub>2h</sub>	3.96(+1.23)	3.09(+0.36)	1.90(−0.83)	2.33(−0.40)	2.78(+0.05)	2.50(−0.23)	4.93(+2.18)	2.73[2.66]
a-T-C <sub>s</sub>	3.59(+0.75)	2.90(+0.06)	1.58(−1.26)	2.89(+0.05)	3.08(+0.24)	3.12(+0.28)	3.72(+0.88)	2.84[2.81]
h-L-C <sub>1</sub>	3.30(+0.68)	2.21(−0.41)	1.02(−1.60)	4.15(+1.53)	3.98(+1.36)	4.38(+1.76)	2.76(+0.14)	[2.62]
MAD <sup>c</sup>	0.61	0.18	0.67	0.42	0.29	0.50	0.48	
	BLYP-D3/TZVPP	PBE-D3/TZVPP	TPSS/LP	wB97XD/aVDZ	wB97XD/6-311+G**	M06HF/aVDZ	M06HF/MG3S	
t-D-C <sub>2h</sub>	1.35(−0.02)	1.63(+0.26)	1.43(+0.06)	1.70(+0.33)	1.59(+0.22)	1.55(+0.18)	1.39(+0.02)	
t-T-C <sub>2v</sub>	1.58(+0.03)	1.59(+0.04)	1.54(−0.01)	1.85(+0.30)	1.64(+0.09)	1.76(+0.21)	1.32(−0.23)	
d-X-D <sub>2d</sub>	1.47(+0.07)	1.61(+0.21)	1.82(+0.42)	2.12(+0.70)	1.90(+0.50)	1.59(+0.19)	1.33(−0.07)	
d-T-C <sub>2v</sub>	1.13(+0.15)	1.44(+0.46)	1.19(+0.21)	1.52(+0.54)	1.41(+0.43)	1.19(+0.21)	0.83(−0.15)	
s-S-C <sub>2</sub>	1.24(−0.19)	1.25(−0.18)	1.90(+0.47)	2.34(+0.91)	2.02(+0.59)	1.56(+0.13)	1.05(−0.38)	
a-D-C <sub>2h</sub>	2.49(−0.24)	2.51(−0.22)	3.38(+0.65)	4.42(+1.69)	4.10(+1.37)	4.42(+1.69)	3.39(+0.66)	
a-T-C <sub>s</sub>	2.42(−0.42)	2.38(−0.46)	2.79(−0.05)	3.99(+1.15)	3.66(+0.82)	4.11(+1.27)	2.94(+0.10)	
h-L-C <sub>1</sub>	2.95(+0.33)	3.22(+0.60)	3.75(+1.13)	4.99(+2.37)	4.51(+1.89)	2.10(−0.52)	1.57(−1.05)	
MAD <sup>c</sup>	0.18	0.30	0.38	1.00	0.74	0.55	0.33	

<sup>a</sup> While the geometries for MP2/CBS and CCSD(T)/CBS were optimized at the MP2/aVDZ and CCSD(T)/aVDZ levels of theory, respectively, those for all other methods were optimized at each given calculation method. For the CCSD(T)/CBS value of the cyclohexane dimer, the intermolecular distance of the MP2/aVDZ optimized geometry was optimized at the CCSD(T)/aVDZ level. The values in parentheses are the differences of the theoretical method-dependent binding energies from the CCSD(T)/CBS values. The bold characters indicate the largest upper and the smallest lower binding energy differences with respect to the CCSD(T)/CBS values. <sup>b</sup> CCSD(T)/CBS values are obtained from the extrapolation by using the CCSD(T)/aVDZ and CCSD(T)/aVTZ values. In the brackets, the values are obtained by assuming that the difference in binding energies between MP2/aVDZ and MP2/CBS calculations is similar to that between CCSD(T)/aVDZ and CCSD(T)/CBS calculations.<sup>56–58</sup> Since the MP2 dispersion energy corrections are not so reliable, the former extrapolation method could be a better choice. <sup>c</sup> MAD is mean absolute deviation.



**Figure 3.** Difference of the binding energies of DFT and MP2/CBS from the CCSD(T)/CBS values.

PBE-D3/TZVPP, and TPSS/LP calculations were performed with the Turbomol 5.10 suite of programs.<sup>61</sup> In the case of the TPSS functional, we used the 6-311++G\*\*(3df,3dp) basis set (which will be abbreviated as Large Pople's (LP) basis set) because the empirical dispersion is specially parametrized for this particular basis set. The suitable basis set for M06-2X method is the DIDZ basis set in that the results are basis set dependent. The M06-2X,<sup>62</sup> MPWB1K,<sup>63</sup> wB97XD, M06HF, and MP2 calculations were carried out with the Gaussian 09 suite of programs.<sup>64</sup> The CCSD(T) calculations were done with the MOLPRO suite.<sup>65</sup> The molecular structures were drawn with the POSMOL package.<sup>66</sup>

By using SAPT calculations,<sup>67–76</sup> the total interaction energy ( $E_{\text{tot}}$ ) is decomposed into electrostatic ( $E_{\text{es}}$ ), effective induction ( $E_{\text{in}}$ ), effective dispersion ( $E_{\text{dp}}$ ), and effective exchange repulsion ( $E_{\text{x}}$ ) energies, as in our earlier work<sup>77,78</sup> and others.<sup>67–76</sup> Here,  $E_{\text{in}}$  and  $E_{\text{dp}}$  include the exchange-induction term and exchange-dispersion term, respectively, while  $E_{\text{x}}$  excludes these terms from the exchange term. The coupled Hartree–Fock response term ( $\delta_{\text{int,resp}}^{\text{HF}}$ ) is added to  $E_{\text{in}}$ , since it tends to be more related to the induction than other terms. DFT-SAPT calculations were performed with the PBE0 functional<sup>79</sup> and aVDZ basis set.

Using the DFT level of theory, we investigated many low-energy structures for each dimer. Then, important low-energy structures (Figures 1 and 2) were further investigated by using

MP2 and CCSD(T) calculations. The dimers are named as “group-shape-sym-index”. Here, “group” denoted as “s/d/t/a/h” indicates “single/double/triple/aromatic/cyclohexane-single”-bonded; “type” as “S/D/X/T/L” indicates “stacked/displaced-stacked/cross-stacked/T-shaped/overlayered”; “sym” denotes the point group symmetry of molecular cluster; “index” as “/b/c...” distinguishes each isomer from the lowest energy structure for more than two isomers. The predicted binding energies for the important dimer structures at various DFT, MP2/CBS, and CCSD(T)/CBS levels are in Table 1. The binding energies in the literature are reported in Table S5, Supporting Information and the differences in binding energies of DFT and MP2/CBS from the CCSD(T)/CBS values are in Figure 3. Our discussion will be based on the CCSD(T)/CBS results unless otherwise specified, because these results are considered to be very reliable.

First, we briefly discuss the most stable structures of acetylene, ethylene, ethane, benzene, and cyclohexane and their competing stable structures. Although the geometrical search was carried out using various DFT and MP2 methods, these results are not quite consistent with the CCSD(T) results. Thus, we discuss the low energy conformers on the basis of the CCSD(T) results, and then the assessment of other methods will be given in comparison with the CCSD(T) results.

We begin with the discussion on the structures of the various types of dimers in terms of the energies ( $\Delta E_{\text{e}}$ ) on the Born–Oppenheimer potential surface at the level of CCSD(T)/CBS. The most stable stacked structures of acetylene, ethylene, ethane, benzene, and cyclohexane dimers are t-D-C<sub>2h</sub>, d-X-D<sub>2d</sub>, s-S-C<sub>2</sub>, a-D-C<sub>2h</sub>, and h-L-C<sub>1</sub>, respectively. Against these structures, there are competing stable T-shaped structures of acetylene, ethylene, and benzene, which are t-T-C<sub>2v</sub>, d-T-C<sub>2v</sub>, and a-T-C<sub>s</sub>, respectively. In most cases, the stacked or displaced-stacked structures are more stable, except for the cases of the acetylene dimer and the benzene dimer for which the T-shaped structures are slightly more stable. In the case of the acetylene dimer, the zero point energy (ZPE) correction makes both t-T-C<sub>2v</sub> and t-D-C<sub>2h</sub> nearly isoenergetic, resulting in the quantum probabilistic structure spanning both T-shaped and displaced-stacked structures<sup>80,81</sup> (see the Supporting Information).

Now, we compare various DFT results and MP2 results with the CCSD(T)/CBS results for a few important cases. For the acetylene dimer and the ethylene dimer, the BLYP-D3/TZVPP results are better than other DFT methods. Both BLYP-D3/TZVPP and MP2/CBS results are in reasonable agreement with the CCSD(T)/CBS results. The intermolecular distance between the centers of mass of two ethylene monomers is 3.82 Å at BLYP-D3/TZVPP and MP2/aVTZ, as compared with 3.78 Å at CCSD(T)/aVDZ. Among the DFT methods, BLYP-D3/TZVPP is in good agreement with CCSD(T)/CBS, though slightly overestimated. The MP2/CBS results are also in good agreement with the CCSD(T)/CBS results.

In the case of the benzene dimer, at the DFT level, the displaced-stacked structures are more stable at the M06-2X, MPWB1K, BLYP-D3, PBE-D3, TPSS, wB97XD, and M06HF levels, while the T-shaped isomers are more stable at the BLYP/TZVP, B97-D/TZVPP, and B3LYP-D/aVTZ levels. At the MP2/CBS level, the displaced-stacked structure is far more stable than the T-shaped structure, which is the weakest point of the MP2 level of theory on the  $\pi$ – $\pi$  interaction.

For the cyclohexane dimer, against the CCSD(T)/CBS binding energy of 2.62 kcal/mol for the isomer h-L-C<sub>1</sub>, we note

**Table 2. Center-to-Center Distance of the Low Energy Structures of the Acetylene(t-), Ethylene(d-), Ethane(s-), Benzene(a-), and Cyclohexane(h-) Dimers<sup>a</sup>**

	M06-2X/aVDZ	MPWB1K/aVDZ	BLYP/TZVP	B97-D/TZV2P	B3LYP-D/aVTZ	MP2/aVTZ	CCSD(T)/aVDZ	CCSD(T)/aVTZ
Center-to-Center Distance								
t-D-C <sub>2h</sub>	4.17(-0.15)	4.27(-0.05)	4.27(-0.05)	4.29(-0.03)	4.09(-0.23)	4.18(-0.14)	4.32	4.26
t-T-C <sub>2v</sub>	4.33(-0.12)	4.42(-0.03)	4.29(-0.16)	4.33(-0.12)	4.27(-0.18)	4.35(-0.10)	4.45	4.40
d-X-D <sub>2d</sub>	3.59(-0.19)	3.84(+0.06)	3.55(-0.23)	3.65(-0.13)	3.51(-0.27)	3.82(+0.04)	3.78	3.73
d-T-C <sub>2v</sub>	3.73(-0.23)	3.86(-0.10)	3.64(-0.32)	3.77(-0.19)	3.63(-0.33)	3.79(-0.17)	3.96	3.91
s-S-C <sub>2</sub>	3.52(-0.33)	3.73(-0.12)	3.46(-0.39)	3.60(-0.25)	3.40(-0.45)	3.70(-0.15)	3.85	3.80
a-D-C <sub>2h</sub>	3.80(-0.24)	3.92(-0.12)	3.87(-0.17)	3.93(-0.11)	3.84(-0.20)	3.69(-0.35)	4.04	3.97
a-T-C <sub>s</sub>	4.85(+0.18)	5.18(+0.15)	4.82(-0.21)	4.87(-0.16)	4.79(-0.24)	4.85(-0.18)	5.03	4.97
h-L-C <sub>1</sub>	4.44(-0.18)	4.64(+0.02)	4.20(-0.42)	4.26(-0.36)	4.22(-0.55)	4.62(-0.15)	4.77	
Vertical Distance								
t-D-C <sub>2h</sub>	2.84(-0.09)	2.90(-0.03)	2.90(-0.03)	2.90(-0.03)	2.79(-0.14)	2.77(-0.16)	2.93	2.83
t-T-C <sub>2v</sub>	2.66(-0.09)	2.75(0.00)	2.62(-0.13)	2.66(-0.09)	2.60(-0.15)	2.68(-0.07)	2.75	2.70
d-X-D <sub>2d</sub>	1.74(-0.17)	1.99(+0.08)	1.69(-0.22)	1.79(-0.12)	1.66(-0.25)	1.95(+0.04)	1.91	1.86
d-T-C <sub>2v</sub>	2.85(-0.22)	2.98(-0.09)	2.76(-0.31)	2.89(-0.18)	2.76(-0.31)	2.91(-0.16)	3.07	2.98
s-S-C <sub>2</sub>	1.76(-0.31)	2.00(-0.07)	1.69(-0.38)	1.84(-0.23)	1.64(-0.43)	1.94(-0.13)	2.07	2.02
a-D-C <sub>2h</sub>	3.33(-0.29)	3.57(-0.05)	3.39(-0.23)	3.47(-0.15)	3.37(-0.25)	3.34(-0.28)	3.62	3.56
a-T-C <sub>s</sub>	2.49(-0.08)	2.74(+0.17)	2.46(-0.11)	2.48(-0.09)	2.46(-0.11)	2.42(-0.15)	2.57	2.50
h-L-C <sub>1</sub>	1.78(-0.26)	1.99(-0.05)	1.54(-0.50)	1.60(-0.44)	1.57(-0.47)	1.94(-0.10)	2.04	
	BLYP-D3/TZVPP	PBE-D3/TZVPP	TPSS/LP	wB97XD/aVDZ	wB97XD/6-311+G**	M06HF/aVDZ	M06HF/MG3S	
Center-to-Center Distance								
t-D-C <sub>2h</sub>	4.27(-0.05)	4.27(-0.05)	4.32(0.00)	4.16(-0.16)	4.18(-0.14)	4.14(-0.18)	4.16(-0.16)	
t-T-C <sub>2v</sub>	4.42(-0.03)	4.42(-0.03)	4.52(0.07)	4.34(-0.11)	4.38(-0.07)	4.35(-0.10)	4.41(-0.04)	
d-X-D <sub>2d</sub>	3.82(0.04)	3.82(0.04)	3.73(-0.05)	3.65(-0.13)	3.65(-0.13)	3.71(-0.07)	3.65(-0.13)	
d-T-C <sub>2v</sub>	3.80(-0.16)	3.80(-0.16)	3.89(-0.07)	3.74(-0.22)	3.74(-0.22)	3.78(-0.18)	3.89(-0.07)	
s-S-C <sub>2</sub>	3.67(-0.18)	3.88(0.03)	3.67(-0.18)	3.56(-0.29)	3.59(-0.26)	3.60(-0.25)	3.57(-0.28)	
a-D-C <sub>2h</sub>	3.77(-0.27)	3.97(0.07)	3.85(-0.19)	3.82(-0.22)	3.81(-0.23)	3.75(-0.29)	3.74(-0.30)	
a-T-C <sub>s</sub>	5.13(.10)	5.18(.15)	5.14(.11)	4.87(-0.16)	4.89(-0.14)	4.75(-0.28)	4.82(-0.21)	
h-L-C <sub>1</sub>	4.51(-0.26)	4.59(-0.18)	4.45(-0.32)	4.30(-0.47)	4.31(-0.46)	4.47(-0.30)	4.51(-0.26)	
Vertical Distance								
t-D-C <sub>2h</sub>	2.89(-0.04)	2.90(-0.03)	2.90(-0.03)	2.85(-0.08)	2.85(-0.08)	2.77(-0.15)	2.88(-0.05)	
t-T-C <sub>2v</sub>	2.74(-0.01)	2.74(-0.01)	2.85(0.10)	2.66(-0.09)	2.71(-0.04)	2.68(-0.07)	2.75(.00)	
d-X-D <sub>2d</sub>	1.97(0.06)	1.96(0.05)	1.89(0.02)	1.80(-0.11)	2.73(-0.18)	1.79(-0.12)	1.81(-0.10)	
d-T-C <sub>2v</sub>	2.87(-0.20)	2.86(-0.21)	2.96(-0.11)	2.80(-0.27)	2.81(-0.26)	2.85(-0.22)	2.97(-0.10)	
s-S-C <sub>2</sub>	1.94(-0.13)	2.06(-0.01)	1.90(-0.17)	1.80(-0.27)	1.88(-0.19)	2.78(-0.29)	1.83(-0.24)	
a-D-C <sub>2h</sub>	3.43(-0.19)	3.64(0.02)	3.55(-0.07)	3.47(-0.15)	3.45(-0.17)	3.28(-0.44)	3.27(-0.35)	
a-T-C <sub>s</sub>	2.80(.23)	2.86(.29)	2.82(.25)	2.54(-0.3)	2.57(.00)	2.27(-0.30)	2.42(-0.15)	
h-L-C <sub>1</sub>	1.85(-0.19)	1.92(0.12)	1.80(-0.24)	1.64(-0.40)	1.67(-0.37)	1.80(-0.24)	1.86(-0.20)	

<sup>a</sup>The bold characters indicate the largest upper and the smallest lower distance differences with respect to the CCSD(T)/aVDZ values, which are close to the CCSD(T)/aVTZ values (t-, triple bonded; d-, double bonded; s-, single bonded; a-, aromatic bonded; h-, cyclohexane single bonded).

substantial energy differences between different methods. The binding energy is 1.02 kcal/mol at MPWB1K/TZVPP, 3.30 kcal/mol at M06-2X/aVDZ, 3.98 kcal/mol at B97-D/TZV2P, 4.15 kcal/mol at BLYP/TZVP, 2.95 kcal/mol at BLYP-D3/TZVPP, and 4.38 kcal/mol at B3LYP-D/aVTZ. Among the DFT methods, only the BLYP-D3/TZVPP method gives a reasonable value. The MP2/CBS value (2.76 kcal/mol) is in good agreement with the CCSD(T)/CBS. The intermolecular distance between two centers of mass of the two monomer units is 4.51 Å at BLYP-D3, 4.62 Å at RIMP2/aVTZ, and 4.77 Å at the BBSE-corrected CCSD(T)/aVDZ.

When various DFT and MP2/CBS binding energies are compared with the CCSD(T)/CBS binding energies for t-D-C<sub>2h</sub>, d-X-D<sub>2d</sub>, s-S-C<sub>2</sub>, a-D-C<sub>2h</sub>, a-T-C<sub>s</sub>, and h-L-C<sub>1</sub>, there are significant discrepancies in certain cases. For example, M06-2X/aVDZ, wB97XD/6-311+G\*\*, wB97XD/aVDZ, and MP2/CBS overestimate the binding energy of a-D-C<sub>2h</sub> (by 1.2, 1.4, 1.7, and 2.2 kcal/mol, respectively), while MPWB1K/aVDZ underestimates it (by 0.8 kcal/mol). B3LYP-D/aVTZ, BLYP/TZVP, and B97-D/TZV2P overestimate that of a-L-C<sub>1</sub> (by 1.7, 1.5, and 1.3 kcal/mol, respectively), while MPWB1K/aVDZ underestimates it (by 1.6 kcal/mol). The BLYP-D3/TZVPP binding energies

agree with the CCSD(T)/CBS energies within 0.4 kcal/mol, but the a-D-C<sub>2h</sub>/a-T-C<sub>s</sub> is underestimated by 0.2/0.4 kcal/mol, while the h-L-C<sub>1</sub> is overestimated by 0.3 kcal/mol. This results in a substantial binding energy difference (0.75 kcal/mol) between the two cases (i.e., benzene and cyclohexane) as compared with an insignificant binding energy difference between the two (0.1–0.2 kcal/mol) in CCSD(T)/CBS. Among various DFT methods, the M06-2X/DIDZ, B97-D/TZV2P, MP2/CBS, and M06HF/aVDZ methods are better for the ethane dimer. In the case of the T-shaped ethylene dimer, M06-2X/DIDZ, BLYP/TZVP, BLYP-D3/TZVPP, and M06HF/MG3S are better, while for the X-shaped ethylene dimer, MP2/CBS, BLYP-D3/TZVPP, and M06HF/MG3S are better. The M06-2X/DIDZ, BLYP/TZVP, B97-D/TZVPP, BLYP-D3/TZVPP, and TPSS/LP methods are better for the acetylene dimer. The M06-2X/DIDZ, BLYP/TZVP, and TPSS/LP methods are better for the T-shaped benzene dimer, while B97-D/TZVPP is better for the displaced stacked structure. The BLYP-D3/TZVPP and MP2/CBS methods are better for the cyclohexane dimer. Overall, BLYP-D3/TZVPP and M06-2X/DIDZ calculations properly reproduce the CCSD(T)/CBS binding energies as compared with other methods. However, we believe that the DFT functionals need to be further improved to describe the difference in binding energy between the benzene dimer and the cyclohexane dimer. M06-2X/DIDZ, BLYP/TZVP, B97-D/TZVPP, and BLYP-D3/TZVPP are working well, except for the cyclohexane dimer. Overall, M06-2X/DIDZ and BLYP-D3/TZVPP are working well. Its mean unsigned relative error or mean absolute deviation (MAD) is 0.18 kcal/mol. Such deviations in the center-to-center distances and vertical distances for the dimer structures (in particular, significantly shortened distances for the cyclohexane dimer) are also noted in the above DFT calculation methods (Table 2).

**Table 3. SAPT-DFT Interaction Energy Components (kcal/mol) for Important Conformers of the Acetylene, Ethylene, Ethane, Benzene, and Cyclohexane Dimers<sup>a</sup>**

	t-D-C <sub>2h</sub>	t-T-C <sub>2v</sub>	d-X-D <sub>2d</sub>	d-T-C <sub>2v</sub>	s-S-C <sub>2</sub>	a-D-C <sub>2h</sub>	a-T-C <sub>s</sub>	h-L-C <sub>1</sub>
$E_{\text{tot}}$	-1.40	-1.56	-1.29	-1.03	-1.08	-2.67	-2.51	-2.32
$E_{\text{es}}$	-1.82	-2.22	-1.18	-1.04	-0.84	-2.64	-1.99	-1.27
$E_{\text{id}}$	-0.36	-0.72	-0.23	-0.38	-0.19	-0.89	-0.60	-0.38
$E_{\text{dp}}$	-1.52	-1.65	-2.56	-2.41	-3.17	-8.68	-5.01	-5.54
$E_{\text{x}}$	2.30	3.04	2.68	2.81	3.13	9.54	5.09	4.88

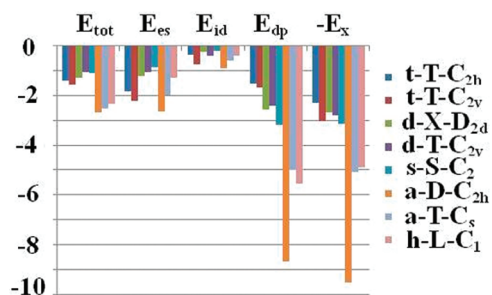
<sup>a</sup> PBE0 functional and aVDZ basis set are employed (t-, triple bonded; d-, double bonded; s-, single bonded; a-, aromatic bonded; h-, cyclohexane single bonded).

The energy components based on the SAPT/DFT are listed in Tables 3 and 4. The total interaction energies for the stacked structures of acetylene, ethylene, ethane, benzene (displaced-stacked/T-shaped), and cyclohexane dimers are -1.40, -1.29, -1.08, -2.67/-2.51, and -2.32 kcal/mol, respectively. As compared with the CCSD(T)/CBS energies, the SAPT/DFT energies of the T-shaped benzene dimer and the cyclohexane dimer are somewhat underestimated by 0.3 kcal/mol, while other cases are close to the CCSD(T) energies. The corresponding  $E_{\text{dp}}$  values of the above dimers are -1.52, -2.56, -3.17, -8.68/-5.01, and -5.54 kcal/mol, respectively, and the corresponding electrostatic energies ( $E_{\text{es}}$ ) are -1.82, -1.18, -0.84, -2.64/-1.99, and -1.27 kcal/mol, respectively. Although the total interaction energies for acetylene, ethylene, and ethane dimers are similar, the C≡C has much weaker dispersion energy than the C=C, which is again weaker than the C-C (for which the dispersion would arise from the interaction between C and H). On the other hand, the electrostatic energy is large for the acetylene dimer (because of significantly positive charge of H atoms in acetylene,  $q_{\text{H}}$ : 0.225 au) but small for the ethane dimer (Figure 4). Again, although the total interaction energies of the dimers for benzene and cyclohexane are similar, the (···C···C···)<sub>cyclic</sub> case has much stronger dispersion and electrostatic energies than the (-C-C-)<sub>cyclic</sub> case. Despite that the dispersion energy tends to be correlated with the number of valence electrons participating in the interaction between the molecules (acetylene, 20; ethylene, 22; ethane, 24; benzene, 60; cyclohexane, 66), one can note that the stacked ethane dimer and the displaced-stacked benzene dimer show particularly large dispersion energies. In the case of the benzene dimer, the electrostatic interaction energy (due to the quadrupole moments of benzene) is also large. This result would thus explain why aromatic compounds are easily found in crystals and self-assembled systems.

To obtain insight into these noncovalent interactions, the second order perturbation theory approximately gives the dispersion energies between two 1s electrons (1s-1s), between two 2s electrons (2s-2s), between two 2p<sub>z</sub> electrons (2p<sub>z</sub>-2p<sub>z</sub>), and between two 2p<sub>1</sub> electrons (2p<sub>1</sub>-2p<sub>1</sub>) for two hydrogenic atoms, which are separated by distance  $R$  along the  $z$  axis. These four values in units of  $(e^4/R^6)(a_0/Z_{\text{eff}})^4$  are estimated to be ~6 (1s-1s), ~1176 (2s-2s), ~1368 (2p<sub>z</sub>-2p<sub>z</sub>), and ~432 (2p<sub>1</sub>-2p<sub>1</sub>), where  $e$  is the electron charge,  $a_0$  is 1 Bohr, and  $Z_{\text{eff}}$  is the effective nuclear charge of the hydrogenic atom ( $Z_{\text{eff}}$  is 1 for H, 3.22 for C(2s), and 3.14 for C(2p)). The dispersion energy between 1s(H) and 2s(C) electrons is estimated to be ~84 ( $e^4/R^6)(a_0^4/Z_{\text{eff}}^2)$ ). In general, two closely contacted nonbonded carbon atoms are separated by  $R = \sim 3.5$  Å, while the closely contacted nonbonded distance between H and C is  $R = \sim 2.5$  Å. Then, the C···H dispersion energy is also strong. What is

**Table 4. SAPT-DFT Interaction Energy Components (kcal/mol) of the Stacked Conformers of Acetylene, Ethylene, and Benzene Dimers, As in Table 3 (t-, triple bonded; d-, double bonded; s-, single bonded; a-, aromatic bonded; h-, cyclohexane single bonded)**

DFT-SAPT	CCSD(T)/aVDZ <sub>opt</sub> interplane distance			interplane distance 3.41 Å			interplane distance 3.7 Å		
	t-S-C <sub>2h</sub>	d-S-D <sub>2h</sub>	a-S-D <sub>6h</sub>	t-S-C <sub>2h</sub>	d-S-D <sub>2h</sub>	a-S-D <sub>6h</sub>	t-S-C <sub>2h</sub>	d-S-D <sub>2h</sub>	a-S-D <sub>6h</sub>
$E_{\text{tot}}$	0.18	0.00	-1.18	1.67	1.31	0.49	0.71	0.36	-1.18
$E_{\text{es}}$	0.42	0.31	-0.46	-0.01	-0.88	-2.91	0.45	-0.03	-0.46
$E_{\text{id}}$	-0.02	-0.06	-0.30	-0.24	-0.33	-0.52	-0.12	-0.17	-0.30
$E_{\text{dp}}$	-0.33	-0.89	-6.70	-2.32	-3.23	-10.74	-1.38	-1.95	-6.70
$E_{\text{x}}$	0.12	0.63	6.28	4.23	5.75	14.66	1.76	2.51	6.28



**Figure 4.** SAPT-DFT/aVDZ energy contributions of the acetylene, ethylene, ethane, benzene, and cyclohexane dimers.

interesting is that the dispersion energy for the  $2p_z-2p_z$  electrons is stronger than that for the  $2s-2s$  electrons, which is much stronger than that for the  $2p_1-2p_1$  electrons. Similarly, in the accurate calculations of the dispersion energy between all electrons of two O atoms, one can note that the van der Waals coefficients of the  $3P_1-3P_1$  interaction and the  $3P_0-3P_0$  interaction are 18.0 and 16.7 au, respectively.<sup>82</sup> In this regard, we expect that the dispersion energy for  $2p_z-2p_z$  electrons would favor the maximally overlapped stacked conformation, while the electrostatic energy would disfavor this overlapping. In the ethene and benzene, the gain by the dispersion energy for  $2p_z-2p_z$  electrons is substantial due to their planar conformation, while in the ethane and cyclohexane, the dispersion energy between  $1s(H)$  and  $2s(C)$  electrons is significant. The present results including the competition and cooperation between dispersion and electrostatic energies are very important for investigating the reliability of density functionals in predicting diverse molecular interaction energies.

In summary, we have studied the structural isomers and interaction energy of triple-bonded acetylene, double-bonded ethylene, single-bonded ethane, aromatic benzene, and cyclic single-bonded cyclohexane dimers. In the case of the ethylene dimer, ethane dimer, benzene dimer, and cyclohexane dimer, the dispersion energy is dominant, while in the case of acetylene dimer, the electrostatic energy is dominant. The aromatic  $\pi$  interactions have particularly large dispersion and electrostatic energies among various types of  $\pi$  interactions. This phenomenon would be related to the fact that aromatic compounds are easily found in crystals, which is indeed very important for crystal packing and molecular self-engineering.

## ASSOCIATED CONTENT

**S Supporting Information.** Coordinates of the low energy structures at the CCSD(T)/aVDZ level; the various types of DFT, MP2, and CCSD(T) binding energies; previous results in the literature. This material is available free of charge via the Internet at <http://pubs.acs.org>.

## AUTHOR INFORMATION

### Corresponding Author

\*E-mail: kim@postech.ac.kr.

## ACKNOWLEDGMENT

This work was supported by NRF (National Honor Scientist Program: 2010-0020414, WCU: R32-2008-000-10180-0) and KISTI (KSC-2011-G3-02).

## REFERENCES

- (1) Meyer, E. A.; Castellano, R. K.; Diederich, F. *Angew. Chem., Int. Ed.* **2003**, *42*, 1210–1250.
- (2) Hoeben, F. J. M.; Jonkheijm, P.; Meijer, E. W.; Schenning, A. P. H. J. *Chem. Rev.* **2005**, *105*, 1491–1546.
- (3) Hunter, C. A.; Sanders, J. K. M. *J. Am. Chem. Soc.* **1990**, *112*, 5525–5534.
- (4) Kim, K. S.; Tarakeshwar, P.; Lee, J. Y. *Chem. Rev.* **2000**, *100*, 4145–4185.
- (5) Riley, K. E.; Pitonak, M.; Jurecka, P.; Hobza, P. *Chem. Rev.* **2010**, *110*, 5023–5063.
- (6) Burley, S. K.; Petsko, G. *Science* **1985**, *229*, 23–28.
- (7) Waters, M. L. *Cur. Opi. Chem. Biol.* **2002**, *6*, 736–741.
- (8) Hobza, P.; Selzle, H. L.; Schlag, E. W. *J. Am. Chem. Soc.* **1994**, *116*, 3500–3506.
- (9) Singh, N. J.; Lee, H. M.; Hwang, I.-C.; Kim, K. S. *Supramol. Chem.* **2007**, *19*, 321–332.
- (10) Hong, B. H.; Lee, J. Y.; Lee, C.-W.; Kim, J. C.; Bae, S. C.; Kim, K. S. *J. Am. Chem. Soc.* **2001**, *123*, 10748–10749.
- (11) Janssen, P. G. A.; Vandenberg, J.; van Dongen, J. L. J.; Meijer, E. W.; Schenning, A. P. H. J. *J. Am. Chem. Soc.* **2007**, *129*, 6078–6079.
- (12) Lee, E. C.; Kim, D.; Jurecka, P.; Tarakeshwar, P.; Hobza, P.; Kim, K. S. *J. Phys. Chem A* **2007**, *111*, 3446–3457.
- (13) Petitjean, A.; Khoury, R. G.; Kyritsakas, N.; Lehn, J. M. *J. Am. Chem. Soc.* **2004**, *126*, 6637–6647.
- (14) Lee, J. Y.; Hong, B. H.; Kim, W. Y.; Min, S. K.; Kim, Y.; Jouravlev, M. V.; Bose, R.; Kim, K. S.; Hwang, I.-C.; Kaufman, L. J.; Wong, C. W.; Kim, P.; Kim, K. S. *Nature* **2009**, *460*, 498–501.
- (15) Müller-Dethlefs, K.; Hobza, P. *Chem. Rev.* **2000**, *100*, 143–167.
- (16) Hobza, P.; Sponer, J. *J. Am. Chem. Soc.* **2002**, *124*, 11802–11808.
- (17) Sponer, J.; Jurecka, P.; Hobza, P. *J. Am. Chem. Soc.* **2004**, *126*, 10142–10151.
- (18) Rezac, J.; Fanfrlik, J.; Salahub, D.; Hobza, P. *J. Chem. Theory Comput.* **2009**, *5*, 1749–1760.
- (19) Sinnokrot, M. O.; Valeev, E. F.; Sherrill, C. D. *J. Am. Chem. Soc.* **2002**, *124*, 10887–10893.
- (20) Sinnokrot, M. O.; Sherrill, C. D. *J. Am. Chem. Soc.* **2004**, *126*, 7690–7697.
- (21) Sinnokrot, M. O.; Sherrill, C. D. *J. Phys. Chem.* **2006**, *110*, 10656–10668.
- (22) Hunter, C. A. *Angew. Chem., Int. Ed.* **1993**, *32*, 1584–1586.
- (23) Hunter, C. A. *Chem. Soc. Rev.* **1994**, *23*, 101–109.
- (24) Cockroft, S. L.; Hunter, C. A.; Lawson, K. R.; Perkins, J.; Urch, C. J. *J. Am. Chem. Soc.* **2005**, *127*, 8594–8595.
- (25) Hong, B. H.; Lee, J. Y.; Cho, S. J.; Yun, S.; Kim, K. S. *J. Org. Chem.* **1999**, *64*, 5661–5665.
- (26) Lee, E. C.; Hong, B. H.; Lee, J. Y.; Kim, J. C.; Kim, D.; Kim, Y.; Tarakeshwar, P.; Kim, K. S. *J. Am. Chem. Soc.* **2005**, *127*, 4530–4537.
- (27) Singh, N. J.; Min, S. K.; Kim, D. Y.; Kim, K. S. *J. Chem. Theor. Comput.* **2009**, *5*, 515–529.
- (28) Geronimo, I.; Lee, E. C.; Singh, N. J.; Kim, K. S. *J. Chem. Theor. Comput.* **2010**, *6*, 1931–1934.
- (29) Piacenza, M.; Grimme, S. *J. Am. Chem. Soc.* **2005**, *127*, 14841–14848.
- (30) Grimme, S. *J. Comput. Chem.* **2006**, *27*, 1787–1799.
- (31) Grimme, S. *Angew. Chem., Int. Ed.* **2008**, *47*, 3430–3434.
- (32) Tsuzuki, S.; Honda, K.; Uchimaru, T.; Mikami, M.; Tanabe, K. *J. Am. Chem. Soc.* **2002**, *124*, 104–112.
- (33) Tsuzuki, S.; Honda, K.; Fujii, A.; Uchimaru, T.; Mikami, M. *Phys. Chem. Chem. Phys.* **2008**, *10*, 2860–2865.
- (34) Podeszwa, R.; Szalewicz, K. *Chem. Phys. Lett.* **2005**, *412*, 488–493.
- (35) Podeszwa, R.; Bukowski, R.; Szalewicz, K. *J. Phys. Chem. A* **2006**, *110*, 10345–10354.
- (36) Sato, T.; Tsuneda, T.; Hirao, K. *J. Chem. Phys.* **2005**, *123*, 104307–10.
- (37) Adamovic, I.; Li, H.; Lamm, M. H.; Gordon, M. S. *J. Phys. Chem. A* **2006**, *110*, 519–525.



- (38) DiStasio, R. A., Jr.; Helden, G. V.; Steele, R. P.; Head-Gordon, M. *Chem. Phys. Lett.* **2007**, *437*, 277–283.
- (39) Dykstra, C. E.; Lisy, J. M. *J. Mol. Struct. (Theochem)* **2000**, *500*, 375–390.
- (40) Chenoweth, K.; Dykstra, C. E. *Theor. Chem. Acc.* **2003**, *110*, 100–104.
- (41) Janowski, T.; Pulay, P. *Chem. Phys. Lett.* **2007**, *447*, 27–32.
- (42) Tsuzuki, S.; Uchimaru, T.; Mikami, M.; Tanabe, K. *J. Phys. Chem. A* **1998**, *102*, 2091–2094.
- (43) Shuler, K.; Dykstra, C. E. *J. Phys. Chem. A* **2000**, *104*, 4562–4570.
- (44) Alberts, I. L.; Rowlands, T. W.; Handy, N. C. *J. Chem. Phys.* **1988**, *88*, 3811–3816.
- (45) Bone, R. G. A.; Handy, N. C. *Theor. Chim. Acta.* **1990**, *78*, 133–163.
- (46) Hobza, P.; Selzle, H. L.; Schlag, E. W. *Collect. Czech. Chem. Commun.* **1992**, *57*, 1186–1190.
- (47) Yu, J.; Shujun, Su.; Bloor, J. E. *J. Phys. Chem.* **1990**, *94*, 5589–5592.
- (48) Karpfen, A. *J. Phys. Chem. A* **1999**, *103*, 11431–11441.
- (49) Karpfen, A. *J. Phys. Chem. A* **1998**, *102*, 9286–9296.
- (50) Brenner, V.; Millie, Ph. *Z. Phys. D* **1994**, *30*, 327–340.
- (51) Hobza, P.; Selzle, H. L.; Schlag, E. W. *J. Chem. Phys.* **1996**, *100*, 18790–18796.
- (52) Sinnokrot, M. O.; Sherrill, C. D. *J. Phys. Chem. A* **2006**, *110*, 10656–10668.
- (53) Jaffe, R. I.; Smith, G. D. *J. Chem. Phys.* **1996**, *105*, 2780–2788.
- (54) Tarakeshwar, P.; Choi, H. S.; Kim, K. S. *J. Am. Chem. Soc.* **2001**, *123*, 3323–3331.
- (55) Helgaker, T.; Klopper, W.; Koch, H.; Noga, J. *J. Chem. Phys.* **1997**, *106*, 9639–9646.
- (56) Min, S. K.; Lee, E. C.; Lee, H. M.; Kim, D. Y.; Kim, D.; Kim, K. S. *J. Comput. Chem.* **2008**, *29*, 1208–1221.
- (57) Császár, A. G.; Allen, W. D.; Schaefer, H. F., III *J. Chem. Phys.* **1998**, *108*, 9751–9764.
- (58) Sinnokrot, M. O.; Sherrill, C. D. *J. Phys. Chem. A* **2004**, *108*, 10200–10207.
- (59) Grimme, S. *J. Comput. Chem.* **2006**, *27*, 1787–1799.
- (60) Pavone, M.; Rega, N.; Barone, V. *Chem. Phys. Lett.* **2008**, *452*, 333–339.
- (61) TURBOMOLE V6.02009; University of Karlsruhe and Forschungszentrum Karlsruhe GmbH: Karlsruhe, Germany, 2007. Available from <http://www.turbomole.com>.
- (62) Zhao, Y.; Truhlar, D. G. *J. Chem. Phys.* **2006**, *110*, 5121–5129.
- (63) Zhao, Y.; Truhlar, D. G. *J. Chem. Theory Comput.* **2005**, *1*, 415–432.
- (64) Gaussian 09, Revision A.1, Frisch, M. J.; Trucks, G. W.; Schlegel, H. B.; Scuseria, G. E.; Robb, M. A.; Cheeseman, J. R.; Scalmani, G.; Barone, V.; Mennucci, B.; Petersson, G. A.; Nakatsuji, H.; Caricato, M.; Li, X.; Hratchian, H. P.; Izmaylov, A. F.; Bloino, J.; Zheng, G.; Sonnenberg, J. L.; Hada, M.; Ehara, M.; Toyota, K.; Fukuda, R.; Hasegawa, J.; Ishida, M.; Nakajima, T.; Honda, Y.; Kitao, O.; Nakai, H.; Vreven, T.; Montgomery, Jr., J. A.; Peralta, J. E.; Ogliaro, F.; Bearpark, M.; Heyd, J. J.; Brothers, E.; Kudin, K. N.; Staroverov, V. N.; Kobayashi, R.; Normand, J.; Raghavachari, K.; Rendell, A.; Burant, J. C.; Iyengar, S. S.; Tomasi, J.; Cossi, M.; Rega, N.; Millam, J. M.; Klene, M.; Knox, J. E.; Cross, J. B.; Bakken, V.; Adamo, C.; Jaramillo, J.; Gomperts, R.; Stratmann, R. E.; Yazyev, O.; Austin, A. J.; Cammi, R.; Pomelli, C.; Ochterski, J. W.; Martin, R. L.; Morokuma, K.; Zakrzewski, V. G.; Voth, G. A.; Salvador, P.; Dannenberg, J. J.; Dapprich, S.; Daniels, A. D.; Farkas, Ö.; Foresman, J. B.; Ortiz, J. V.; Cioslowski, J.; Fox, D. J. Gaussian, Inc., Wallingford CT, 2009.
- (65) Werner, H.-J.; Knowles, P. J.; Lindh, R.; Manby, F. R.; Schütz, M.; Celani, P.; Korona, T.; Rauhut, G.; Amos, R. D.; Bernhardsson, A.; Berning, A.; Cooper, D. L.; Deegan, M. J. O.; Dobbyn, A. J.; Eckert, F.; Hampel, C.; Hetzer, G.; Lloyd, A. W.; McNicholas, S. J.; Meyer, W.; Mura, M. E.; Nicklass, A.; Palmieri, P.; Pitzer, R.; Schumann, U.; Stoll, H.; Stone, A. J.; Tarroni, R.; Thorsteinsson, T. MOLPRO, version 2006.1; University College Cardiff Consultants Limited: Cardiff, Wales, U.K., 2006. See <http://www.molpro.net>.
- (66) Lee, S. J.; Chung, H. Y.; Kim, K. S. *Bull. Korean Chem. Soc.* **2004**, *25*, 1061–1064.
- (67) Heßelmann, A.; Jansen, G. *Chem. Phys. Lett.* **2002**, *357*, 464–470.
- (68) Heßelmann, A.; Jansen, G. *Phys. Chem. Chem. Phys.* **2003**, *5*, 5010–5014.
- (69) Fiethen, A.; Jansen, G.; Heßelmann, A.; Schütz, M. *J. Am. Chem. Soc.* **2008**, *130*, 1802–1803.
- (70) Jansen, G.; Heßelmann, A. *J. Phys. Chem. A* **2001**, *105*, 11156–11157.
- (71) Heßelmann, A.; Jansen, G.; Schütz, M. *J. Chem. Phys.* **2005**, *122*, 014103–17.
- (72) Misquitta, A. J.; Bukowski, R.; Szalewicz, K. *J. Chem. Phys.* **2000**, *112*, 5308–5319.
- (73) Misquitta, A. J.; Szalewicz, K. *Chem. Phys. Lett.* **2000**, *357*, 301–306.
- (74) Misquitta, A. J.; Szalewicz, K. *J. Chem. Phys.* **2005**, *122*, 214109.
- (75) Misquitta, A. J.; Podeszwa, R.; Jezierski, B.; Szalewicz, K. *J. Chem. Phys.* **2005**, *123*, 214103.
- (76) Misquitta, A. J.; Stone, A. J. *J. Chem. Theory Comput.* **2008**, *4*, 7–18.
- (77) Lee, E. C.; Hong, B. H.; Lee, J. Y.; Kim, J. C.; Kim, D.; Kim, Y.; Tarakeshwar, P.; Kim, K. S. *J. Am. Chem. Soc.* **2005**, *127*, 4530–4537.
- (78) Kim, D. Y.; Singh, N. J.; Kim, K. S. *J. Chem. Theory Comput.* **2008**, *4*, 1401–1407.
- (79) Adamo, C.; Barone, V. *J. Chem. Phys.* **1999**, *110*, 6158–6170.
- (80) Karthikeyan, S.; Lee, H. M.; Kim, K. S. *J. Chem. Theory Comput.* **2010**, *6*, 3190–3197.
- (81) Prichard, D. G.; Nandi, R. N.; Muentzer, J. S. *J. Chem. Phys.* **1988**, *89*, 115–123.
- (82) Kelly, H. P. *Phys. Lett.* **1969**, *29A*, 30–31.

# Use of Direct Dynamics Simulations to Determine Unimolecular Reaction Paths and Arrhenius Parameters for Large Molecules

Li Yang, Rui Sun, and William L. Hase\*

Department of Chemistry and Biochemistry, Texas Tech University, Lubbock, Texas 79409, United States

**ABSTRACT:** In a previous study (*J. Chem. Phys.* **2008**, *129*, 094701) it was shown that for a large molecule, with a total energy much greater than its barrier for decomposition and whose vibrational modes are harmonic oscillators, the expressions for the classical Rice–Ramsperger–Kassel–Marcus (RRKM) (i.e., RRK) and classical transition-state theory (TST) rate constants become equivalent. Using this relationship, a molecule's unimolecular rate constants versus temperature may be determined from chemical dynamics simulations of microcanonical ensembles for the molecule at different total energies. The simulation identifies the molecule's unimolecular pathways and their Arrhenius parameters. In the work presented here, this approach is used to study the thermal decomposition of  $\text{CH}_3\text{—NH—CH=CH—CH}_3$ , an important constituent in the polymer of cross-linked epoxy resins. Direct dynamics simulations, at the MP2/6-31+G\* level of theory, were used to investigate the decomposition of microcanonical ensembles for this molecule. The Arrhenius  $A$  and  $E_a$  parameters determined from the direct dynamics simulation are in very good agreement with the TST Arrhenius parameters for the MP2/6-31+G\* potential energy surface. The simulation method applied here may be particularly useful for large molecules with a multitude of decomposition pathways and whose transition states may be difficult to determine and have structures that are not readily obvious.

## 1. INTRODUCTION

Computational chemistry is an important tool for studying unimolecular reactions.<sup>1,2</sup> To understand the dynamics and kinetics of a unimolecular reaction, it is necessary to know the atomic-level mechanism(s) by which a molecule dissociates.<sup>3</sup> Electronic structure calculations<sup>4</sup> are often used to identify the important unimolecular pathways and transition states (TSs). A classical trajectory chemical dynamics simulation<sup>5</sup> may be performed to investigate the molecule's atomistic intramolecular and unimolecular dynamics.<sup>1,2</sup> The potential energy surface (PES) for this simulation may be an analytic potential energy function,<sup>6,7</sup> or the simulation may be performed by direct dynamics,<sup>8,9</sup> in which the gradient and potential energy for calculating the trajectory is obtained directly from an electronic structure theory.

For large molecules and/or high energies (e.g., hyperthermal),<sup>10–12</sup> identifying reaction pathways and TS properties by electronic structure calculations becomes less practical and more challenging. This is because the important decomposition pathways may become less identifiable, and there is the possibility of a multitude of pathways.<sup>13,14</sup> Such effects are found when protonated peptide ions collide with hydrocarbon surfaces.<sup>15,16</sup> For collisions of protonated diglycine with the diamond {111} surface, at a collision energy of 100 eV, 88 different fragmentation pathways of the peptide ion are observed.<sup>15</sup> Similarly, protonated octaglycine dissociates via 304 pathways when it collides with the diamond {111} surface at 100 eV.<sup>16</sup> Identifying TSs for all of these pathways would be a formidable task and may also be impractical.

In this article a classical trajectory direct chemical dynamics simulation procedure is described and applied for determining the reaction pathways of a molecule undergoing unimolecular decomposition at temperature  $T$ . Furthermore, by calculating the

unimolecular constants  $k_i(T)$  for the individual paths versus  $T$ , the Arrhenius parameters  $A$  and  $E_a$  for the paths may be determined. The unimolecular reactions investigated are those for decomposition of  $\text{CH}_3\text{—NH—CH=CH—CH}_3$  (Figure 1), which represents an important constituent in the polymer of cross-linked epoxy resins.<sup>17,18</sup> This molecule is small enough that the ab initio TSs may be determined, and thus, Arrhenius parameters may be determined for these TSs using transition-state theory (TST) and compared with the simulation values. This provides a test of the simulation methodology. In the following, the theoretical model and the computational methodology are first described, followed by a presentation of the computational results.

## 2. THEORETICAL MODEL

The theoretical approach used here is based on the recent finding<sup>19</sup> that for a molecule consisting of  $s$  harmonic oscillators the classical Rice–Ramsperger–Kassel–Marcus (RRKM) rate constant:

$$k(E) = \nu \left( \frac{E - E_0}{E} \right)^{s-1} \quad (1)$$

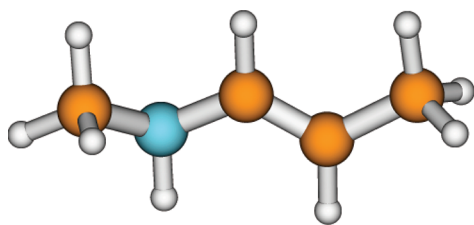
and classical TST rate constant:

$$k(T) = \nu \exp(-E_0/k_B T) \quad (2)$$

become equivalent for large  $E$  with  $E_0/E \ll 1$  and large  $s \approx s - 1$ . For  $s$  classical harmonic oscillators, the energy and the temperature are related by  $E = sk_B T$ . Thus, a simulation of the unimolecular decomposition of a microcanonical ensemble at

Received: July 2, 2011

Published: September 19, 2011



**Figure 1.** Structure of  $\text{CH}_3\text{-NH-CH=CH-CH}_3$ , the molecule investigated for the unimolecular decomposition studies.

energy  $E$  may be used to determine the thermal unimolecular rate constant  $k(T)$ . This relationship is consistent with the understanding that, for a molecule with large  $s$  and large  $E$ , the populations of the vibrational states of the individual oscillators are given by a Boltzmann distribution<sup>20</sup> and that the fluctuations of the energy of a grand canonical ensemble are negligible, making it similar to a microcanonical ensemble.<sup>21</sup> If anharmonic effects are important, they may be included by multiplying the expressions in eqs 1 and 2 by an anharmonic correction factor,<sup>19</sup> which accounts for anharmonicity in both the reactant molecule and TS.

The trajectories comprising this microcanonical ensemble are integrated until a unimolecular reaction occurs or up to a maximum time  $t_{\text{max}}$ . The reactions observed are those important for the molecule at temperature  $T$ . The total rate constant  $k$  for unimolecular decomposition of the molecule may be found by fitting the number of molecules remaining versus time:

$$N(t)/N(0) = \exp(-kt) \quad (3)$$

where  $N(0)$  is the number of trajectories for the initial microcanonical ensemble at  $t = 0$  or more approximately from the number of trajectories remaining at  $t_{\text{max}}$ , i.e.:

$$N(t_{\text{max}})/N(0) = \exp(-kt_{\text{max}}) \quad (4)$$

The total number of products formed at  $t_{\text{max}}$  is  $P(t_{\text{max}}) = N(0) - N(t_{\text{max}})$ . The total unimolecular rate constant is a sum of the rate constants for the individual unimolecular pathways, i.e.,  $k = \sum k_i$ , and the rate constant for an individual pathway is simply:

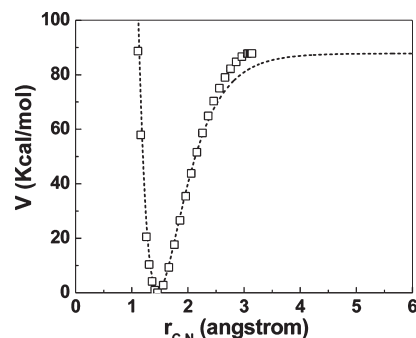
$$k_i = [P_i(t_{\text{max}})/P(t_{\text{max}})]k \quad (5)$$

where  $P_i(t_{\text{max}})$  is the number of sets of products formed for path  $i$  at  $t_{\text{max}}$ . The Arrhenius parameters  $A_i$  and  $E_{a,i}$  for path  $i$  are found from the Arrhenius equation  $k_i(T) = A_i \exp(-E_{a,i}/k_B T)$  by determining  $k_i$  as a function of  $T$ .

The above theoretical model assumes the unimolecular dynamics are intrinsically RRKM,<sup>22</sup> so that a microcanonical ensemble of states is maintained as the molecule decomposes. As shown by the simulation results, these dynamics are observed for the molecule studied here  $\text{CH}_3\text{-NH-CH=CH-CH}_3$ .

### 3. COMPUTATIONAL METHODOLOGY

**3.1. Direct Dynamics Simulations.** The simulations are performed using direct dynamics<sup>8,9</sup> in which the technology of classical trajectory calculations is coupled with electronic structure theory. In this manner, the gradient, energy, and possibly Hessian<sup>23,24</sup> needed to calculate the trajectory come directly from an electronic structure theory, without the need for an analytic potential energy function. To initialize the trajectories, the phase space of the reactant molecule is excited randomly at



**Figure 2.** MP2/6-31+G\* potential energy curve ( $\square$ ) for  $\text{CH}_3\text{-NH-CH=CH-CH}_3 \rightarrow \cdot\text{CH}_3 + \cdot\text{NH-CH=CH-CH}_3$  dissociation. The dashed line is the Morse potential energy curve.

fixed energy  $E$  to form a microcanonical ensemble.<sup>25</sup> The temperature associated with this energy may be identified from the average kinetic energy of the molecule's  $N$  atoms, i.e.:

$$\sum_{i=1}^N m_i \langle v_i^2 \rangle / 2 = 3Nk_B T / 2 \quad (6)$$

or by equating  $E$  to the average thermal energy of the molecule's  $s = 3N - 6$  classical oscillators, i.e.,  $E = sk_B T$ . Procedures for forming this microcanonical ensemble are well established.<sup>22,25-29</sup>

The simulations of  $\text{CH}_3\text{-NH-CH=CH-CH}_3$  unimolecular decomposition were performed by exciting microcanonical ensembles of molecules so that their temperatures ( $E = sk_B T$ ) were 3500, 4000, 4500, 5000, and 5500 K. The software package consisting of the chemical dynamics computer program VENUS<sup>30,31</sup> interfaced with the NWChem<sup>32</sup> electronic structure computer program was used for the simulations. A total of 100 trajectories were calculated for each temperature. The simulations were performed by direct dynamics using the MP2/6-31+G\* electronic structure theory. The trajectories with  $T$  of 3500–5000 K were integrated for  $t_{\text{max}}$  of 5.1 ps or until a unimolecular reaction occurred. The 5500 K trajectories were integrated with  $t_{\text{max}} = 3.3$  ps. The reported uncertainties for the rate constants and the Arrhenius parameters are standard deviations. The temperatures of the simulations were determined from  $E = sk_B T$ .

**3.2. PES and TST Calculations.** The model PES for the direct dynamics simulations is given by MP2/6-31+G\* theory. It is recognized that MP2 is quite approximate for homolytic bond rupture reactions,<sup>33</sup> due to the shortcomings of Hartree–Fock theory.<sup>34</sup> As discussed below, the most important decomposition pathway found here for  $\text{CH}_3\text{-NH-CH=CH-CH}_3$  is dissociation to  $\cdot\text{CH}_3 + \cdot\text{NH-CH=CH-CH}_3$ . The MP2 potential energy curve for this dissociation is given in Figure 2. The potential energy varies from 87.76 kcal/mol at  $r_{\text{C-N}} = 3.06$  Å to 87.78 kcal/mol at  $r_{\text{C-N}} = 3.14$  Å, with a maximum of 87.83 kcal/mol at 3.11 Å. This dissociation energy is similar to the experimental value of 85.1 kcal/mol for  $\text{CH}_3\text{-NH}_2$ .<sup>35</sup> The  $\text{CH}_3\text{-NH}$  dissociation energy, for the molecule studied here, is expected to be slightly lower than the value for  $\text{CH}_3\text{-NH}_2$ , since replacing one of the H-atoms of  $\text{-NH}_2$  by a C-atom lowers the bond energy.<sup>35</sup>

The TS for  $\text{CH}_3\text{-NH-CH=CH-CH}_3 \rightarrow \cdot\text{CH}_3 + \cdot\text{NH-CH=CH-CH}_3$  dissociation is variational, and its structure is determined by the shape of the potential energy curve.<sup>36-38</sup> As shown by earlier work,<sup>39</sup> MP2 theory does not give the correct

potential energy curve for such a dissociation. Using the MP2 dissociation energy of 87.83 kcal/mol and MP2 quadratic force constant for C–N stretching of 5.20 mdyn/Å, the Morse potential energy curve,  $V = D\{1 - \exp[-\beta(r - r_o)]\}^2$ , is given in Figure 2, where it is compared with the MP2 potential energy curve. (This C–N stretching force constant is only slightly larger than the experimental value of 5.12 mdyn/Å for CH<sub>3</sub>NH<sub>2</sub>.)<sup>35</sup> The MP2 curve rises more steeply and may be fit by a  $r$ -dependent Morse  $\beta$  parameter.<sup>40,41</sup> Such a curve has been referred to as a “stiff” Morse potential.<sup>40,41</sup> The variational TS<sup>42</sup> for a stiff Morse potential is at a much shorter bond length than for the standard Morse potential, resulting in a much “tighter” TS.<sup>42</sup> This leads to a significantly smaller  $A$  factor for the stiff Morse potential. For CH<sub>4</sub> → H + CH<sub>3</sub> dissociation, the  $A$  factor is an order of magnitude smaller for the stiff Morse potential.<sup>40–42</sup> The Morse potential more accurately models the actual potential energy curve than does the “stiff” Morse potential.<sup>39</sup> Thus, the MP2 direct dynamics is expected to give an  $A$  factor for CH<sub>3</sub>–NH–CH=CH–CH<sub>3</sub> → ·CH<sub>3</sub> + ·NH–CH=CH–CH<sub>3</sub> dissociation, which is significantly smaller than the assumed experimental value of  $10^{16}$ – $10^{17}$  s<sup>-1</sup>.<sup>43</sup>

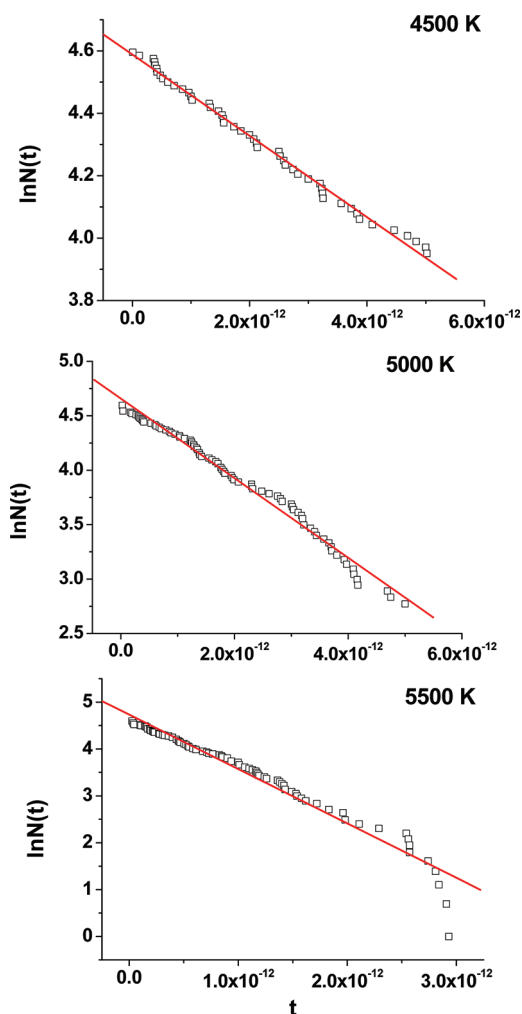
Classical TST calculations were performed to calculate Arrhenius parameters for the MP2/6-31+G\* PES, to compare with the Arrhenius parameters determined from the MP2/6-31+G\* direct dynamics simulations. As discussed below, three pathways were considered for this comparison, i.e., the decomposition to ·CH<sub>3</sub> + ·NH–CH=CH–CH<sub>3</sub> and two isomerization reactions. For the harmonic oscillator model, the classical TST rate constant in eq 2 is given by

$$k(T) = \frac{\prod_{i=1}^s \nu_i}{\prod_{i=1}^{s-1} \nu_i^{\ddagger}} \exp(-E_o/k_B T) \quad (7)$$

where the  $\nu_i$  are the vibrational frequencies for the unimolecular reaction, the  $\nu_i^{\ddagger}$  the vibrational frequencies for the TS, and  $E_o$  the classical potential energy barrier for the MP2 PES. The ratio of the products of vibrational frequencies is the Arrhenius  $A$  factor and  $E_o$  is the activation energy.

The above harmonic oscillator model is expected to be applicable for the two isomerization reactions, which have “tight” transition-state structures located at the isomerization barriers. As discussed above, the C–N dissociation reaction to form ·CH<sub>3</sub> + ·NH–CH=CH–CH<sub>3</sub> has a variational TS located at the minimum in  $k(T)$  along the dissociation path.<sup>37,38</sup> For bond dissociations, such as CH<sub>4</sub> → H + CH<sub>3</sub>, a harmonic oscillator variational TST model is accurate,<sup>44–46</sup> and this model is used here for C–N bond dissociation. Vibrational frequencies and potential energies are found along the intrinsic reaction coordinate (IRC),<sup>47,48</sup> and this information is used to find the minimum in  $k(T)$ , i.e., eq 7. These calculations were performed with the GAMESS computer program.<sup>49</sup>

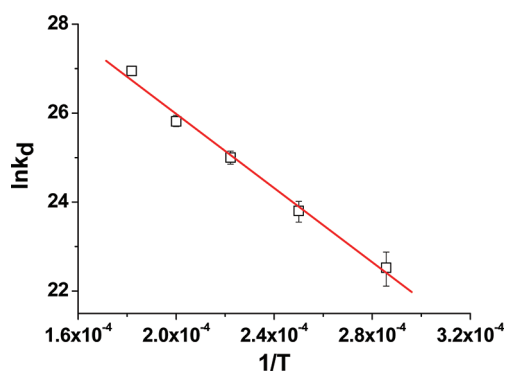
The IRC frequency for the CH<sub>3</sub> torsion about the C–N bond was found to be unstable, and its value was found by interpolation between the frequency of 208 cm<sup>-1</sup> at the CH<sub>3</sub>–NH–CH=CH–CH<sub>3</sub> potential energy minimum to 29.8 cm<sup>-1</sup> at the energy maximum at  $R_{C-N} = 3.11$  Å (see above). In previous work, the frequency for such a mode was found to decay approximately exponentially as the bond ruptures.<sup>39–51</sup> This exponential interpolation, as well as linear interpolation, was considered here and found to give similar results.



**Figure 3.** Plots of  $\ln N(t)$  for the number of excited CH<sub>3</sub>–NH–CH=CH–CH<sub>3</sub> molecules remaining versus time for simulations at 4500, 5000, and 5500 K. The fits are to eq 4. The unit for  $t$  is sec.

A more accurate representation of the C–N dissociation reaction may require a flexible variational TST model<sup>52–54</sup> which treats the transitional vibrational modes, that are transformed into product rotations, as hindered rotational degrees of freedom. An important property of the model is that it includes anharmonicity for the transitional modes.<sup>55,56</sup> A possible shortcoming is the assumed separability between the transitional and remaining modes, which becomes more approximate as the length of the rupturing bond is shortened. This could be a problem for the MP2/6-31+G\* PES, which has a “stiff” Morse potential for the C–N bond (see above). However, in future analyses of this C–N bond dissociation, it would be of interest to consider the flexible variational TST model.

The TST activation energies and  $A$  factors for the MP2/6-31+G\* PES were compared with these parameters found from the chemical dynamics simulations. The TSs for the isomerization reactions were placed at their potential barriers to give their  $E_o$  values and the  $A$  factors were calculated from eq 7. For C–N dissociation, both the bond dissociation energy and the potential energy at the IRC variational TS were considered as possible MP2 values for  $E_o$ . The dissociation  $A$  factor is that for the variational TS.



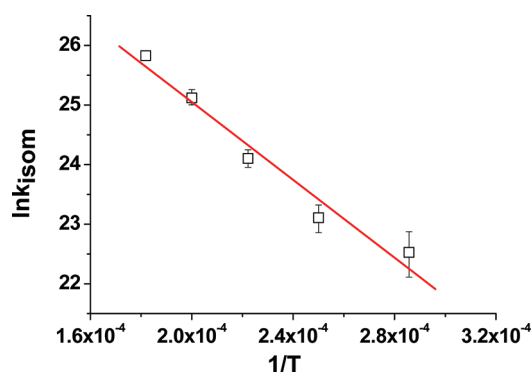
**Figure 4.** Plot of  $\ln k_d$  versus  $1/T$  for  $\text{CH}_3\text{-NH-CH=CH-CH}_3 \rightarrow \cdot\text{CH}_3 + \cdot\text{NH-CH=CH-CH}_3$  dissociation. The linear fit yields the Arrhenius parameters  $E_a = 82.8 \pm 4.3$  kcal/mol and  $A = 8.01 + 5.24/-3.17 \times 10^{14} \text{ s}^{-1}$ . The  $k_d$  is in units of  $\text{s}^{-1}$ , and  $T$  is in K.

#### 4. SIMULATION RESULTS

A total of 100 trajectories were calculated for each of the temperatures 3500, 4000, 4500, 5000, and 5500 K. The number of trajectories which reacted within  $t_{\text{max}}$  for the respective temperatures is 6, 19, 48, 84, and 100. A total of 33 different primary and secondary decomposition pathways were observed, and a mechanistic analyses of these pathways and their relationship to experiment will be presented elsewhere.<sup>57</sup> Of interest here is the kinetics analysis discussed in Section 2. The dominant pathway is C–N bond rupture to form the radicals  $\cdot\text{CH}_3$  and  $\cdot\text{NH-CH=CH-CH}_3$ . The isomerizations via H-atom transfer, to form  $\text{CH}_3\text{-N=CH-CH}_2\text{-CH}_3$  and  $\text{CH}_3\text{-NH-CH}_2\text{-CH=CH}_2$  are also important pathways. This bond rupture and these isomerizations are considered here.

Using the number of dissociations which occurred within  $t_{\text{max}}$  eq 4, the rate constant  $k$  for decomposition of  $\text{CH}_3\text{-NH-CH=CH-CH}_3$  is  $(1.21 + 0.49/-0.50) \times 10^{10}$ ,  $(4.13 + 0.93/-0.97) \times 10^{10}$ ,  $(1.28 + 0.18/-0.20) \times 10^{11}$ , and  $(3.59 + 0.41/-0.51) \times 10^{11} \text{ s}^{-1}$  for  $T$  of 3500, 4000, 4500, and 5000 K, respectively. Equation 4 may not be used to calculate a rate constant for 5500 K, since all the trajectories dissociated within  $t_{\text{max}}$ . For the calculations at 4500, 5000, and 5500 K, there is a sufficient number of reactions to find the rate constant from a plot of  $\ln N(t)$  versus  $t$ ; i.e., eq 3. The plots are given in Figure 3, and the fitted rate constants for the respective temperatures are  $1.30 \pm 0.02 \times 10^{11}$ ,  $3.65 \pm 0.05 \times 10^{11}$ , and  $1.16 \pm 0.03 \times 10^{12} \text{ s}^{-1}$ . Values for these rate constant found from nonlinear fits of eq 3 are nearly the same and  $1.37 \times 10^{10}$ ,  $3.41 \times 10^{11}$ , and  $9.75 \times 10^{11} \text{ s}^{-1}$ , respectively. The rate constants from plots of  $\ln N(t)$  and  $N(t)/N(0)$  are in excellent agreement with those found from the single point  $N(t_{\text{max}})/N(0)$ . In the following analyses, the rate constants for 3500–5000 K are from  $N(t_{\text{max}})/N(0)$ , while the 5500 K value is from the  $\ln N(t)$  plot.

**4.1.  $\text{CH}_3\text{-NH-CH=CH-CH}_3 \rightarrow \cdot\text{CH}_3 + \cdot\text{NH-CH=CH-CH}_3$  Dissociation.** The rate constant versus temperature for C–N bond rupture to form  $\cdot\text{CH}_3 + \cdot\text{NH-CH=CH-CH}_3$  was determined using eq 5, the total rate constant versus temperature and the  $P_i(t_{\text{max}})/P(t_{\text{max}})$  ratio for this decomposition pathway. The resulting  $k_d$  rate constants are given in Figure 4 as a plot of  $\ln k_d$  versus  $1/T$ . The linear fit to this plot gives  $E_a = 82.8 \pm 4.3$  kcal/mol and  $A = 8.0 + 5.2/-3.2 \times 10^{14} \text{ s}^{-1}$ . The  $E_a$  value of 82.8 kcal/mol is similar to but somewhat lower than the MP2 dissociation energy of 87.83 kcal/mol. This difference is expected from the variational



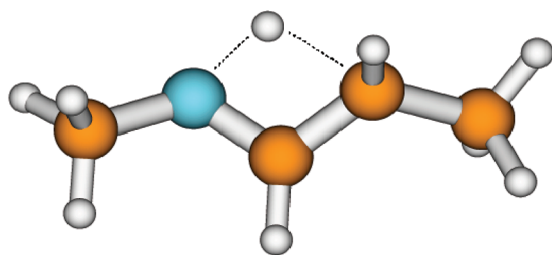
**Figure 5.** Plot of  $\ln k_{\text{isom}}$  versus  $1/T$  for forming the isomerization products  $\text{CH}_3\text{-N=CH-CH}_2\text{-CH}_3$  and  $\text{CH}_3\text{-NH-CH}_2\text{-CH=CH}_2$ . The linear fit yields the Arrhenius parameters  $E_a = 64.9 \pm 7.0$  kcal/mol and  $A = 5.2 + 4.5/-2.8 \times 10^{15} \text{ s}^{-1}$ . The  $k_{\text{isom}}$  is in units of  $\text{s}^{-1}$ , and  $T$  is in K.

nature of the TS for the C–N bond dissociation pathway.<sup>37,38</sup> As the temperature is increased, the free energy barrier for C–N bond rupture moves to a shorter C–N bond length, resulting in activation energies that are decreased as the temperature is increased<sup>37,38</sup> and are less than the dissociation energy. The direct dynamics  $A$  factor of  $\sim 10^{15} \text{ s}^{-1}$  is much smaller than the expected experimental value of  $10^{16}\text{--}10^{17} \text{ s}^{-1}$  for  $\text{CH}_3$  dissociation.<sup>42,43</sup> Such a result is expected. The MP2/6-31+G\* potential for C–N bond rupture, Figure 2, is not sufficiently attractive and, thus, gives rise to a variational TS structure that is “too tight”. The result is a small  $A$  factor.

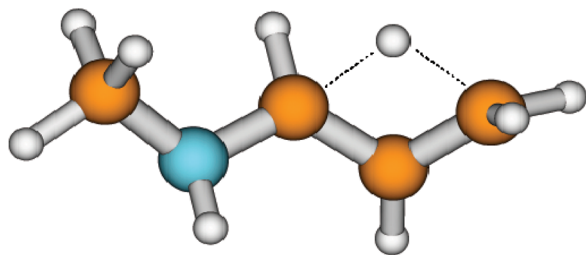
As discussed in Section 3.2, a harmonic oscillator model, based on the IRC, was used to find a variational TS for C–N bond rupture at 4000 K. The resulting TS is located at a C–N distance of 2.41 Å, giving rise to a potential energy of 73 kcal/mol and a  $A$  factor of  $8 \times 10^{14} \text{ s}^{-1}$ . The  $A$  factor is the same as the chemical dynamics value, but the  $E_0$  is smaller than that from the dynamics. The harmonic IRC variational TST rate constant is  $8.5 \times 10^{10} \text{ s}^{-1}$  and  $\sim 4$  times larger than the dynamics value of  $2.2 \times 10^{10} \text{ s}^{-1}$ . Overall, the TST Arrhenius  $E_0$  and  $A$  parameters for the MP2/6-31+G\* PES are consistent with the values found from the MP2/6-31+G\* direct dynamics simulation.

It is worth noting that there are ambiguities in the above variational TST calculation due to substantial changes in vibrational modes as the C–N bond ruptures. The  $\text{CH}_3$  torsion about this bond may be a vibration for the reactant, but for a sufficient bond extension, it will become a free rotor. Here it is treated as a vibration for both the reactant and the variational TS. The torsion’s partition function at the minimum is 13.4 and 40.6 as a vibration and free rotor, respectively. For the C–N bond dissociation maximum at 3.11 Å (Section 3.2), these respective values are 93.3 and 42.3. Thus, treating the torsion as a free rotor at the variational TS would decrease the TST rate constant. In addition, there are four rocking/bending modes whose frequencies go to zero as the C–N bond breaks. It may be better to treat these modes as hindered rotors, instead of vibrations, as is done by the flexible variational TST.<sup>52–56</sup> However, there are approximations in separating these modes from the remaining modes of the dissociating molecule, and also, the uncertainty in treating the  $\text{CH}_3$  torsion of the reactant remains.

**4.2.  $\text{CH}_3\text{-NH-CH=CH-CH}_3$  Isomerization Reactions.** Rate constants versus temperature were determined, as described above, for the isomerization reactions to form the products



$$E_o = 63.2 \text{ kcal/mol}$$



$$E_o = 82.9 \text{ kcal/mol}$$

**Figure 6.** MP2/6-31+G\* TS structures and potential energies for forming the isomerization products  $\text{CH}_3\text{-N=CH-CH}_2\text{-CH}_3$  and  $\text{CH}_3\text{-NH-CH}_2\text{-CH=CH}_2$ .

$\text{CH}_3\text{-N=CH-CH}_2\text{-CH}_3$  and  $\text{CH}_3\text{-NH-CH}_2\text{-CH=CH}_2$ . To obtain better statistics, the rate constants for these two pathways were combined to give a total isomerization rate constant  $k_{\text{isom}}$  versus  $T$ . The resulting Arrhenius plot of  $\ln k_{\text{isom}}$  versus  $1/T$  is given in Figure 5. The fitted  $E_a$  and  $A$  are  $64.9 \pm 7.01 \text{ kcal/mol}$  and  $5.2 + 4.5/-2.8 \times 10^{13} \text{ s}^{-1}$ .

To compare with these fitted Arrhenius parameters, the ab initio TSs were found for these two isomerizations. Their structures and energies are given in Figure 6. The fitted  $E_a$  is intermediate of the two ab initio barrier heights but more akin to  $E_o = 63.2 \text{ kcal/mol}$  for the  $\text{CH}_3\text{-N=CH-CH}_2\text{-CH}_3$  product. The classical TST  $A$  factor, calculated from the reactant molecule and TSs' vibrational frequencies, is  $8.8 \times 10^{13} \text{ s}^{-1}$  for the  $\text{CH}_3\text{-N=CH-CH}_2\text{-CH}_3$  product and  $6.9 \times 10^{13} \text{ s}^{-1}$  for  $\text{CH}_3\text{-NH-CH}_2\text{-CH=CH}_2$  product (the latter pathway has a reaction path degeneracy of 3). The fitted  $A$  factor is similar to these values. To make a direct comparison with the fitted  $E_a$  and  $A$  from the chemical dynamics simulation, the classical TST rate constants for the two isomerization paths were summed at each temperature to give a composite rate constant and plotted as  $\ln k_{\text{isom}}$  versus  $1/T$ . The resulting plot gives  $E_a = 62.8 \text{ kcal/mol}$  and  $A = 9.2 \times 10^{13} \text{ s}^{-1}$ , values in overall good agreement with those from the chemical dynamics simulation. That the composite  $A$  and  $E_a$  are closer to those for the  $\text{CH}_3\text{-N=CH-CH}_2\text{-CH}_3$  product is consistent with the lower  $E_a$  and larger  $A$  for this product.

The smaller  $A$  factor found from the simulations, as compared to the harmonic TST value, may be the result of anharmonic effects for the chemical dynamics on the MP2/6-31+G\* PES. This is consistent with the smaller simulation total isomerization rate constants as compared to those for the TST calculations. For

the  $T$  of 3500, 4000, 4500, 5000, and 5500 K, the respective simulation rate constants are  $5.2 \times 10^9$ ,  $1.9 \times 10^{10}$ ,  $5.0 \times 10^{10}$ ,  $1.2 \times 10^{11}$ , and  $1.5 \times 10^{11} \text{ s}^{-1}$ . The TST values are approximately a factor of 2 larger and  $1.1 \times 10^{10}$ ,  $3.4 \times 10^{10}$ ,  $8.2 \times 10^{10}$ ,  $1.7 \times 10^{11}$ , and  $3.0 \times 10^{11} \text{ s}^{-1}$ .

## 5. CONCLUSIONS

For the work presented here a classical trajectory direct chemical dynamics simulation approach is described for determining unimolecular reaction paths and Arrhenius parameters. This method is expected to be particularly useful for large molecules with many decomposition paths and whose TSs may be difficult to determine by standard electronic structure theory methods. The simulations are performed by coupling the methodology of chemical dynamics simulations with electronic structure theory<sup>8,9</sup> and involve studying the unimolecular decomposition of microcanonical ensembles of molecules.

The molecule studied here is  $\text{CH}_3\text{-NH-CH=CH-CH}_3$ , an important constituent in the polymer of cross-linked epoxy resins.<sup>17,18</sup> This epoxy resin is a component in flame resistant nanocomposites,<sup>17,58</sup> and it is important to understand its decomposition kinetics at high temperatures, as is done here. Arrhenius parameters for decomposition of  $\text{CH}_3\text{-NH-CH=CH-CH}_3$  to  $\cdot\text{CH}_3 + \cdot\text{NH-CH=CH-CH}_3$  and isomerization to  $\text{CH}_3\text{-N=CH-CH}_2\text{-CH}_3$  and  $\text{CH}_3\text{-NH-CH}_2\text{-CH=CH}_2$  were determined from direct dynamics simulations at the MP2/6-31+G\* level of theory. The Arrhenius activation energies determined from the simulation are in good agreement with the isomerization potential energy barriers and C–N bond dissociation energy for the MP2/6-31+G\* PES. TST is used to calculate Arrhenius  $A$  factors for the MP2/6-31+G\* PES, and they are in good agreement with the values found from the simulations. Overall, the TST Arrhenius parameters for the MP2/6-31+G\* PES are consistent with the values obtained from the MP2/6-31+G\* direct dynamics simulation.

It is pointed out that there are uncertainties and ambiguities in the TST calculations as a result of anharmonic effects and the treatment of modes in the variational TST calculations for C–N bond rupture. A strength of a direct dynamic simulation is that these effects are accurately represented in determining the Arrhenius parameters.

Additional applications and tests of the computational chemistry approach described and applied here, for determining unimolecular decomposition pathways and Arrhenius parameters, are expected. It may be particularly useful for studying nanomaterials and biological molecules.

## ■ AUTHOR INFORMATION

### Corresponding Author

\*E-mail: bill.hase@ttu.edu.

## ■ ACKNOWLEDGMENT

This material is based on work supported by the Office of Naval Research under award no. N00014-09-1-0626 and by the Robert A. Welch Foundation under grant no. D-0005. The Hrothgar computer cluster at Texas Tech University, under the direction of the High Performance Computing Center and Philip W. Smith, was used for the simulations reported here.

## ■ REFERENCES

- (1) Hase, W. L.; Schinke, R. Role of Computational Chemistry in the Development of Unimolecular Rate Theory. In *Theory and Applications of Computational Chemistry: The First 40 Years*; Dykstra, C. E., Frenking, G., Kim, K. S., Scuseria, G., Eds.; Elsevier: New York, 2005, pp 397–423.
- (2) Lourderaj, U.; Hase, W. L. *J. Phys. Chem. A* **2009**, *113*, 2236.
- (3) Baer, T.; Hase, W. L. *Unimolecular Reaction Dynamics. Theory and Experiments*; Oxford: New York, 1996.
- (4) Simons, J. J. *Phys. Chem.* **1991**, *95*, 1017.
- (5) Bunker, D. L. *Methods Comput. Phys.* **1971**, *10*, 287.
- (6) Hase, W. L.; Mrowka, G.; Brudzynski, R. J.; Sloane, C. S. *J. Chem. Phys.* **1978**, *69*, 3548.
- (7) Vande Linde, S. R.; Hase, W. L. *J. Phys. Chem.* **1990**, *94*, 2778.
- (8) Bolton, K.; Hase, W. L.; Peslherbe, G. H. Direct Dynamics Simulations of Reactive Systems. In *Multidimensional Molecular Dynamics Methods*; Thompson, D. L., Ed. World Scientific Publishing, Inc.: London, 1998; pp 143–189.
- (9) Sun, L.; Hase, W. L. *Rev. Comput. Chem.* **2003**, *19*, 79.
- (10) Minton, T.; Garton, D. J. Dynamics of Atomic Oxygen Induced Polymer Degradation in Low Earth Orbit. In *Chemical Dynamics in Extreme Environments*, Advanced Series in Physical Chemistry; Dressler, R. A., Ed.; World Scientific Publishing, Inc.: London, 2001; Vol. 11, pp 420–489.
- (11) Troya, D.; Schatz, G. C. *Int. Rev. Phys. Chem.* **2004**, *23*, 341.
- (12) Yan, T.-Y.; Doubleday, C.; Hase, W. L. *J. Phys. Chem. A* **2004**, *108*, 9863.
- (13) Meroueh, S. O.; Wang, Y.; Hase, W. L. *J. Phys. Chem. A* **2002**, *106*, 9983.
- (14) Laskin, J.; Futrell, J. H. *J. Am. Soc. Mass Spectrom.* **2003**, *14*, 1340.
- (15) Wang, Y.; Hase, W. L.; Song, K. *J. Am. Soc. Mass Spectrom.* **2003**, *14*, 1402.
- (16) Park, K.; Deb, B.; Song, K.; Hase, W. L. *J. Am. Soc. Mass Spectrom.* **2009**, *20*, 939.
- (17) Tseng, C.-H.; Hsueh, H.-B.; Chen, C.-Y. *Compos. Sci. Technol.* **2007**, *67*, 2350.
- (18) Lin, P.-H.; Khare, R. *Macromolecules* **2009**, *42*, 4319.
- (19) Lourderaj, U.; McAfee, J. L.; Hase, W. L. *J. Chem. Phys.* **2008**, *129*, 094701.
- (20) (a) Baer, T.; Hase, W. L. *Unimolecular Reaction Dynamics. Theory and Experiments*; Oxford: New York, 1996; p 328. (b) Park, K.; Engelkemier, J.; Persico, M.; Manikandan, P.; Hase, W. L. *J. Phys. Chem. A* **2011**, *115*, 6603.
- (21) McQuarrie, D. A. *Statistical Thermodynamics*; Harper: New York, 1973.
- (22) Bunker, D. L.; Hase, W. L. *J. Chem. Phys.* **1973**, *54*, 4621.
- (23) Lourderaj, U.; Song, K.; Windus, T. L.; Zhuang, Y.; Hase, W. L. *J. Chem. Phys.* **2007**, *126*, 044105.
- (24) Wu, H.; Rahman, M.; Wang, J.; Lourderaj, U.; Hase, W. L.; Zhuang, Y. *J. Chem. Phys.* **2010**, *133*, 074101.
- (25) Peslherbe, G. H.; Wang, H.; Hase, W. L. *Adv. Chem. Phys.* **1999**, *105*, 171.
- (26) Bunker, D. L. *J. Chem. Phys.* **1962**, *37*, 393.
- (27) Hase, W. L.; Buckowski, D. G. *Chem. Phys. Lett.* **1980**, *74*, 284.
- (28) Schranz, H. W.; Nordholm, S.; Nyman, G. *J. Chem. Phys.* **1991**, *94*, 1487.
- (29) Haile, J. M. *Molecular Dynamics Simulation*; Wiley: New York, 1992.
- (30) Hase, W. L.; Duchovic, R. J.; Hu, X.; Domornicki, A.; Lim, K. F.; Lu, D. H.; Peslherbe, G. H.; Swamy, S. R.; Vande Linde, S. R.; Varandas, A.; Wolfe, R. J. *QCPE Bull.* **1996**, *16*, 671.
- (31) Hu, X.; Hase, W. L.; Pirraglia, T. *J. Comput. Chem.* **1992**, *12*, 1014–1024.
- (32) Bylaska, E. J.; de Jong, W. A.; Govind, N.; Kowalski, K.; Straatsma, T. P.; Valiev, M.; Wang, D.; Apra, E.; Windus, T. L.; Hammond, J.; Nichols, P.; Hirata, S.; Hackler, M. T.; Zhao, Y.; Fan, P.-D.; Harrison, R. J.; Dupuis, M.; Smith, D. M. A.; Nieplocha, J.; Tipparaju, V.; Krishnan, M.; Wu, Q.; Van Voorhis, T.; Auer, A. A.; Nooijen, M.; Brown, E.; Cisneros, G.; Fann, G. L.; Fruchtl, H.; Garza, J.; Hirao, K.; Kendall, R.; Nichols, J. A.; Tsemekhman, K.; Wolinski, K.; Anchell, J.; Bernholdt, D.; Borowski, P.; Clark, T.; Clerc, D.; Dachsel, H.; Deegan, M.; Dyal, K.; Elwood, D.; Glendening, E.; Gutowski, M.; Hess, A.; Jaffe, J.; Johnson, B.; Ju, J.; Kobayashi, R.; Kutteh, R.; Lin, Z.; Littlefield, R.; Long, X.; Meng, B.; Nakajima, T.; Niu, S.; Pollack, L.; Rosing, M.; Sandrone, G.; Stave, M.; Taylor, H.; Thomas, G.; van Lenthe, J.; Wong, A.; Zhang, Z. *NWChem, A Computational Chemistry Package for Parallel Computers*, version 5.1; Pacific Northwest National Laboratory: Richland, WA, 2007.
- (33) Duchovic, R. J.; Hase, W. L.; Schlegel, H. B.; Frisch, M. J.; Raghavachari, K. *Chem. Phys. Lett.* **1982**, *89*, 120.
- (34) Levine, I. N. *Quantum Chemistry*, 3<sup>rd</sup> ed.; Allyn and Bacon: Boston, MA, 1973.
- (35) *CRC Handbook of Chemistry and Physics*, 89<sup>th</sup> ed.; Lide, D. R., Ed.; Taylor & Francis Group: Boca Raton, FL, 2008–2009.
- (36) Hase, W. L. *J. Chem. Phys.* **1972**, *57*, 730.
- (37) Hase, W. L. *J. Chem. Phys.* **1976**, *64*, 2442.
- (38) Hase, W. L. *Acc. Chem. Res.* **1983**, *16*, 258.
- (39) Hase, W. L.; Mondro, S. L.; Duchovic, R. J.; Hirst, D. M. *J. Am. Chem. Soc.* **1987**, *109*, 2916.
- (40) Duchovic, R. J.; Hase, W. L. *Chem. Phys. Lett.* **1984**, *110*, 474.
- (41) Duchovic, R. J.; Hase, W. L. *J. Chem. Phys.* **1985**, *82*, 3599.
- (42) Hase, W. L.; Duchovic, R. J. *J. Chem. Phys.* **1985**, *83*, 3448.
- (43) Johnson, R. L.; Hase, W. L.; Simons, J. W. *J. Chem. Phys.* **1970**, *52*, 3911.
- (44) Hu, X.; Hase, W. L. *J. Chem. Phys.* **1991**, *95*, 8073.
- (45) Wang, H.; Zhu, L.; Hase, W. L. *J. Phys. Chem.* **1994**, *98*, 1608.
- (46) Song, K.; de Sainte Claire, P.; Hase, W. L.; Hass, K. C. *Phys. Rev. B* **1995**, *52*, 2949.
- (47) Fukui, K. *J. Phys. Chem.* **1970**, *74*, 4161.
- (48) Miller, W. H.; Handy, N. C.; Adams, J. E. *J. Chem. Phys.* **1980**, *72*, 99.
- (49) Schmidt, M. W.; Baldrige, K. K.; Boatz, J. A.; Elbert, S. T.; Gordon, M. S.; Jensen, J. H.; Koseki, S.; Matsunaga, N.; Nguyen, K. A.; Su, S.; Windus, T. L.; Dupuis, M.; Montgomery, J. A. *J. Comput. Chem.* **1993**, *14*, 1347.
- (50) Hase, W. L. *Chem. Phys. Lett.* **1987**, *139*, 389.
- (51) Hu, X.; Hase, W. L. *J. Phys. Chem.* **1989**, *93*, 4029.
- (52) Wardlaw, D. M.; Marcus, R. A. *Adv. Chem. Phys.* **1987**, *70*, 231.
- (53) Hase, W. L.; Wardlaw, D. M. *Bimolecular Collisions*; Ashfold, M. N. R.; Baggott, J. E., Eds.; Royal Society of Chemistry: London, 1989, p 171.
- (54) Klippenstein, S. J. *J. Chem. Phys.* **1992**, *96*, 367.
- (55) Klippenstein, S. J.; Georgievskii, Y.; Harding, L. B. *Phys. Chem. Chem. Phys.* **2006**, *8*, 1133.
- (56) Harding, L. B.; Georgievskii, Y.; Klippenstein, S. J. *J. Phys. Chem. A* **2005**, *109*, 4646.
- (57) Yang, L.; Hase, W. L., to be submitted for publication.
- (58) Levchik, S. V.; Weil, E. D. *Polym. Int.* **2004**, *53*, 1901.

# Bipolar Reaction Path Hamiltonian Approach for Reactive Scattering Problems

Jeremy B. Maddox<sup>\*,†</sup> and Bill Poirier<sup>‡</sup>

<sup>†</sup>Department of Chemistry, Western Kentucky University, Bowling Green, Kentucky, 42101-1079, United States

<sup>‡</sup>Department of Chemistry and Biochemistry, Texas Tech University, Lubbock, Texas, 79409-1061, United States

**ABSTRACT:** In this work we present a method for calculating the stationary state wave functions and reaction probabilities of a multidimensional reactive scattering system. Our approach builds upon the counter-propagating wave methodology (CPWM) developed by Poirier and co-workers for calculating one-dimensional stationary state wave functions. The method involves the formulation of a bipolar decomposition for multidimensional stationary scattering wave functions within the context of a reaction path Hamiltonian, so we refer to this work as the bipolar reaction path Hamiltonian (BRPH) approach. Benchmark calculations are presented for several 2D model scattering systems with linear reaction coordinates. We show that the BRPH approach is competitive with conventional calculations based on discrete variable representation (DVR) methods.

## 1. INTRODUCTION

In this work we formulate a computational approach for calculating the stationary state wave functions and state-to-state reaction probabilities of a multidimensional (multi-D) reactive scattering system.<sup>1</sup> For such problems, the total energy is a continuous quantity and the system exhibits some unbound motion along at least one spatial coordinate. By definition, the stationary state wave functions of the system must simultaneously satisfy both the time-dependent Schrödinger equation (TDSE)

$$i\hbar \frac{\partial \Phi_E}{\partial t} = \hat{H} \Phi_E \quad (1)$$

and the time-independent Schrödinger equation (TISE)

$$\hat{H} \Phi_E = E \Phi_E \quad (2)$$

where  $\hat{H}$  is the Hamiltonian operator corresponding to the total energy of the system and  $\Phi_E$  is a wave function representing the stationary scattering state with energy  $E$ . The spatial and temporal components of  $\Phi_E$  are formally factorizable. For example, in one-dimensional (1D) space we have

$$\Phi_E(x, t) = \phi_E(x) e^{-iEt/\hbar} \quad (3)$$

where  $\phi_E(x)$  is an amplitude that depends only on the system's spatial coordinate  $x$  and also satisfies the TISE.<sup>2</sup> For unbound motion, it is well-known that stationary state wave functions are not square-integrable, and therefore, may not represent a physically realizable state.<sup>3</sup> Nevertheless, stationary scattering states are conceptually important as a formal tool for constructing normalizable wave packets, and they are a useful idealization for the limiting case of a wave packet that has a very narrow profile in momentum space and a spatial width that is much larger than the dimensions of the scattering problem.

In the context of chemical reactions, the “scattering coordinate” (or “reaction coordinate”) represents a set of pathways through the molecular configuration space that yield a transformation from reactants to products. The “reaction path” then refers to one such

pathway, generally the minimum energy pathway, which connects asymptotic minima in the reactant and product potential valleys via a saddle point corresponding to the transition state. Due to coupling with the other “perpendicular” coordinates or degrees of freedom, motion along the reaction coordinate leads to the rearrangement of chemical bonds and to the many ways in which the internal energy of the system can be redistributed among rotational, vibrational, and electronic degrees of freedom of the reactants and products, thus complicating the details of the scattering process. This is especially true for problems with many atoms and low-lying excited electronic states. At the same time, however, when all is said and done, and the quantum reactive scattering event is completed, there remains the fundamental concept that the molecular collision can be represented by a superposition of incident, transmitted, and reflected waves along the reaction coordinate. Moreover, the amplitude of the transmitted wave is directly related to the reaction probability for a particular scattering channel, which is defined by the asymptotic perpendicular quantum states of the reactants and products. In turn, the amplitudes associated with all energetically accessible channels are related to the reaction cross-section and, ultimately, the overall reaction rate.<sup>4,5</sup>

To compute the various reactive scattering quantities above in accurate quantum dynamical detail, the use of discrete variable representations (DVRs)<sup>6–9</sup> with absorbing boundary conditions (ABCs),<sup>10,11</sup> i.e., complex absorbing potentials, has over the years proven to be a very effective, general approach.<sup>12,13</sup> Variations of the DVR-ABC method have been applied to reactive scattering,<sup>14–17</sup> electron scattering,<sup>18</sup> isomerization reactions,<sup>19</sup> photo-reactions,<sup>20–24</sup> nonadiabatic systems,<sup>25</sup> and molecule-surface scattering.<sup>26</sup> Though robust and accurate, the DVR-ABC methods suffer from the well-known limitation of exponential scaling of computational effort with system size, thus limiting such calculations in practice to small molecules. Another major

Received: August 12, 2011

Published: September 21, 2011



difficulty is associated with ABCs, in particular, which necessarily expand the required reaction coordinate ranges substantially, especially in the vicinity of the channel threshold energies. By now, ABCs have been well-developed for the purpose of truncating the computational domain of a given problem as much as possible;<sup>27–29</sup> however, it can still be a nontrivial task to minimize artificial reflections, especially at energies just above threshold.

In this paper, we develop an exact computational method for determining the stationary states and reaction probabilities for a certain class of multi-D reactive scattering problems with linear and quasilinear reaction coordinates (future work will address the curvilinear case). As we shall see, this method naturally incorporates the boundary conditions of the physical scattering problem, so that there is no need for ABCs and the attendant expansion of reaction coordinate space. The method is also designed to scale well with increasing system size, provided some approximations (albeit fairly minor and reasonable) are introduced. Our approach builds upon the bipolar counter-propagating wave methodology (CPWM) of Poirier and co-workers, which for scattering problems involves the decomposition of stationary states into so-called “bipolar” traveling wave components. Over the last several years, various CPWM schemes have been developed and applied to bound stationary states,<sup>30</sup> scattering states in 1D,<sup>31–35</sup> non-adiabatic dynamics,<sup>34</sup> multi-D scattering,<sup>35</sup> and nonstationary state dynamics.<sup>36–38</sup> Although these represent a great improvement over traditional unipolar quantum trajectory methods<sup>39</sup> and have enabled the first-ever accurate synthetic quantum trajectory calculations to be performed for a system with substantial reflection interference,<sup>37,38</sup> they still can exhibit certain practical difficulties when applied to real molecular systems, such as occasional numerical instabilities or inaccuracies, thus motivating the development of new bipolar approaches.

The present work involves the formulation of a bipolar CPWM for multi-D systems in terms of adiabatic vibrational eigenstates associated with bound motion in the perpendicular degrees of freedom and a corresponding Hamiltonian that varies parametrically along a suitably defined reaction path. Some time ago, Miller and co-workers developed an approximate representation for the molecular Hamiltonian along the reaction path of a reactive system that could be constructed using a reasonable number of accurate electronic structure calculations.<sup>40</sup> By now, the reaction path Hamiltonian (RPH) approach is essentially a cornerstone of both classical and quantum mechanical theories of kinetic rates constants, and extensions of the basic RPH notion continue to be actively developed.<sup>41–44</sup> The RPH is sufficiently general that it can be applied within the context of many different computational schemes, such as self-consistent field calculations<sup>45,46</sup> and, in a more recent example, diffusion Monte Carlo.<sup>47</sup> In our work, we apply the bipolar CPWM approach to scattering problems within a framework that will be suitable for the RPH approximation; hence, we refer to this method as the bipolar reaction path Hamiltonian (BRPH) approach.

The organization of the rest of this paper is as follows. In section 2 we highlight some key points of bipolar CPWMs for 1D scattering problems. This background information is important for understanding the theoretical and numerical developments associated with the BRPH approach that are described in section 3. Some additional theoretical details are provided in Appendices A and B. We present and discuss several benchmark numerical calculations in section 4, where the state-to-state reaction probabilities are determined for model problems involving scattering motion across Eckart-type barriers with coupling to harmonic

vibrational motion. The BRPH results for these problems are quantitatively compared with analytical theory and corresponding DVR-ABC calculations. Appendix C describes our implementation of the DVR-ABC method. Finally, in section 5 we conclude with a brief summary and outlook for future studies involving more realistic problems with curvilinear reaction coordinates and larger dimensionalities.

## 2. THE 1D COUNTER-PROPAGATING WAVE METHODOLOGY (CPWM)

A thorough discussion of the CPWM for scattering problems can be found in the literature,<sup>31–33,35</sup> and we will not repeat all of those details here. However, we must summarize a few key concepts that are necessary to understand the new developments presented in section 3 for multi-D scattering problems. In 1D scattering problems, the Hamiltonian is given by

$$\hat{H} = -\frac{\hbar^2}{2m} \frac{d^2}{dx^2} + V(x) \quad (4)$$

where  $x$  represents the scattering coordinate for a particle with mass  $m$  that traverses the barrier potential  $V$ . Furthermore, we take  $V$  to be an asymptotically convergent function, i.e.,  $V(x \rightarrow \pm\infty) = \text{constant}(s)$ .<sup>48</sup>

In the CPWM approach, the stationary scattering states of the system are represented by an appropriate superposition of counter-propagating traveling waves, such as

$$\Phi(x, t) = \Phi_+(x, t) + \Phi_-(x, t) \quad (5)$$

We refer to this as a bipolar decomposition, and  $\Phi_{\pm}$  are the bipolar components. Note that we have now omitted the explicit reference to  $E$  in our notation for  $\Phi$  and  $\Phi_{\pm}$ ; henceforth, this energy dependence is implied. If we assume a left-incident scattering convention, then the  $\Phi_+$  component represents a traveling wave that moves with positive momentum (to the right) and asymptotically corresponds to the incident and transmitted plane wave portions of the total wave function. Conversely, the  $\Phi_-$  component moves with negative momentum (to the left) and is associated with a reflected plane wave in the left asymptote and approaches zero in the right asymptote. The interference between the two bipolar components determines the form of the total stationary state wave function.

For the 1D scattering problems described above, the bipolar components may be expressed as

$$\Phi_{\pm}(x, t) = \alpha_{\pm}(x) \exp\left(\pm \frac{i}{\hbar} \int p(x) dx - \frac{i}{\hbar} Et\right) \quad (6)$$

where  $\alpha_{\pm}$  are a pair of amplitudes that vary over the same region of space as the scattering potential. Asymptotically, these amplitudes should become constant and will be related to the overall transmission and reflection probabilities at a given energy. In eq 6, notice that the  $\alpha_{\pm}$  amplitudes have been formally separated from the oscillatory parts of the wave function, which are represented by the complex exponential factors. The integral term in the exponent represents the classical action of a particle with momentum  $p = (2m(E - V_{\text{eff}}))^{1/2}$  moving in an effective potential field  $V_{\text{eff}}(x)$ . Different bipolar decompositions may be specified by choosing different forms for  $V_{\text{eff}}$  provided that  $E > V_{\text{eff}}(x)$  for all  $x$  and that  $V_{\text{eff}}(x \rightarrow \pm\infty) = V(x \rightarrow \pm\infty)$ . Aside from these constraints, which guarantee that  $\Phi_{\pm}$  have the correct asymptotic

behavior, the effective potential is more or less arbitrary.<sup>49</sup> Note that the requirements on  $V_{\text{eff}}(x)$  do not imply that the true potential,  $V(x)$ , must be less than the energy; i.e., tunneling, even deep tunneling, is in principle allowed and treated exactly.

Even for a given  $V_{\text{eff}}$  the bipolar decomposition described above does not provide a unique specification of the  $\Phi_{\pm}$ , although the allowed form of these is greatly constrained. To obtain a unique decomposition, leading to slowly varying  $\alpha_{\pm}(x)$ , we must impose an additional relation

$$\Phi' = -\frac{p'}{2p}\Phi + \frac{i}{\hbar}p(\Phi_+ - \Phi_-) \quad (7)$$

which was originally introduced by Fröman and Fröman (FF) in the context of a generalized semiclassical theory for tunneling phenomena.<sup>50</sup> In eq 7 and hereafter we use a prime to denote the derivative of a 1D function with respect to the scattering coordinate, e.g.,  $f' = df(x)/dx$ .

Starting from an essentially arbitrary initial guess for  $\Phi_{\pm}$ , the exact FF decomposition solution  $\Phi_{\pm}$  are obtained in the long-time limit by solving a pair of time-dependent equations of motion (see eq 10 below), involving the total (hydrodynamic) time derivatives of the bipolar components

$$d_t \Phi_{\pm} = \partial_t \Phi_{\pm} \pm \frac{p}{m} \Phi'_{\pm} \quad (8)$$

The notation  $d_t$  signifies the total time derivative and implies that  $\Phi_{\pm}$  are evolving on a pair of counter-propagating Lagrangian-type reference frames  $x_{\pm}(t)$  that satisfy the following auxiliary equations of motion, defining the left- and right-traveling trajectories

$$d_t x_{\pm} = \pm \frac{p}{m} \quad (9)$$

Using the FF condition and the TISE to expand the convective term  $\Phi'_{\pm}$  leads to the following coupled equations of motion

$$d_t \Phi_{\pm} = F_{\pm} \Phi_{\pm} + G(\Phi_+ + \Phi_-) \quad (10)$$

where the factors  $F_{\pm}$  and  $G$  depend on  $V_{\text{eff}}$  according to

$$F_{\pm} = \frac{i}{\hbar} \left( E - 2V_{\text{eff}} \mp \frac{i\hbar V'_{\text{eff}}}{\sqrt{8m(E - V_{\text{eff}})}} \right) \quad (11a)$$

$$G = -\frac{i}{\hbar} \left( V - V_{\text{eff}} - \frac{\hbar^2}{8m} \left[ \frac{V''_{\text{eff}}}{E - V_{\text{eff}}} + \frac{5}{4} \left( \frac{V'_{\text{eff}}}{E - V_{\text{eff}}} \right)^2 \right] \right) \quad (11b)$$

The solutions of eq 10 are subject to the boundary conditions

$$\Phi_+(x \rightarrow -\infty, t) = \exp \left[ \frac{i}{\hbar} \int p \, dx - \frac{i}{\hbar} Et \right] \quad (12a)$$

$$\Phi_-(x \rightarrow +\infty, t) = 0 \quad (12b)$$

that serve to reinforce the left-incident scattering convention. Given an initial guess for  $\Phi_{\pm}(x, t=0)$ , the dynamics that follows from eq 10 can be viewed as a pseudo-time-dependent relaxation brought about through the cooperative influence of  $G$ , which couples the evolution of  $\Phi_{\pm}$ , and the application of the boundary conditions.

To help illustrate this, we suppose that the “initial” ( $t = 0$ ) bipolar components are represented by a  $\Phi_+$  plane-wave with

constant positive momentum  $p_0 = (2m(E - V(-\infty)))^{1/2}$ , and a  $\Phi_-$  reflected wave with zero amplitude:

$$\Phi_+(x, t = 0) = \exp \left( \frac{i}{\hbar} p_0 x \right) \quad (13a)$$

$$\Phi_-(x, t = 0) = 0 \quad (13b)$$

Clearly, the superposition of these components is not a stationary state solution of the TISE for  $V(x) \neq 0$ . Asymptotically,  $G$  is negligible, and the coupling between the bipolar components vanishes; therefore,  $\Phi_{\pm}$  evolve like free-particle wave functions in those regions. This evolution is enforced by the boundary conditions and the motion of the counter-propagating reference frames, which carry the  $\Phi_{\pm}$  amplitude away from the interaction region in opposite directions. Within the interaction region, however, where  $G \neq 0$ , there is coupling that leads to a transfer of amplitude between the  $\Phi_{\pm}$  components, and over time a nonzero reflected wave builds up. Eventually a steady-state is reached between the flux associated with the boundary conditions and the flux induced by the coupling. The long-time limit solutions of the CPWM equations of motion provide a pair of bipolar components whose superposition is the desired stationary scattering state wave function.

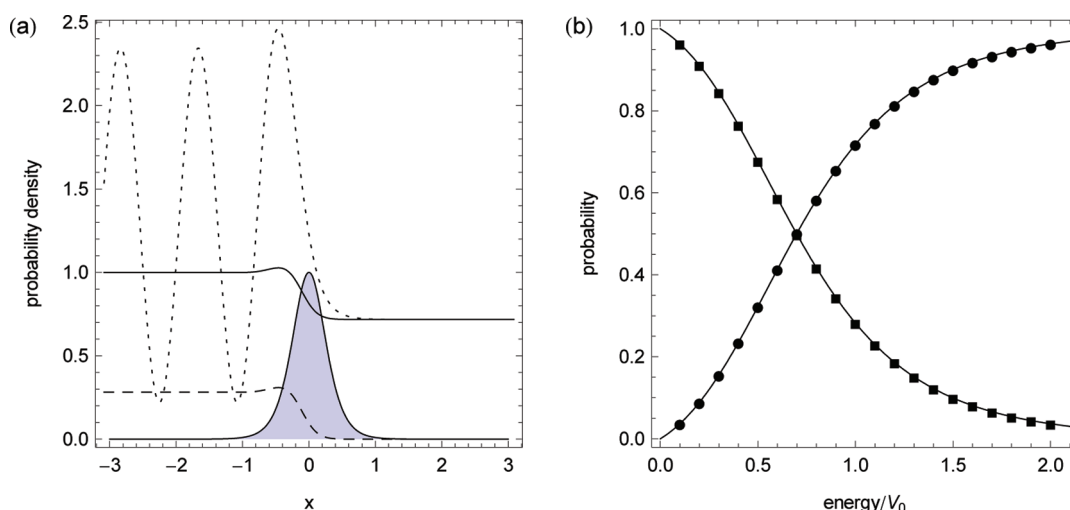
Figures 1 and 2 illustrate typical CPWM results for the case of a particle with mass  $m = 2000$  au scattering across a 1D Eckart barrier defined by

$$V(x) = V_0 \operatorname{sech}^2(\alpha x) \quad (14)$$

where  $V_0 = 0.0018$  hartree and  $\alpha = 3.0 \text{ b}^{-1}$ . This parametrization of the Eckart potential has been referred to as the Eckart A barrier in previous work.<sup>32,33</sup> The potential is plotted in Figure 1a as a function of  $x$  in units of  $V_0$ , and the area beneath the curve has been shaded.

We have numerically integrated the so-called “constant velocity” form of the CPWM equations,<sup>31–33</sup> where  $V_{\text{eff}}(x) = 0$  for all  $x$ , starting from the initial conditions given in eq 13. The bipolar solutions are represented on a grid of equally spaced points with grid spacing  $\Delta x$ , and the time step is defined by  $\Delta t = \Delta x(m/2E)^{1/2}$ . The solutions are propagated in time using a second-order Runge–Kutta method until the probability densities of the bipolar components  $\rho_{\pm} = |\Phi_{\pm}|^2$  are essentially fully converged. For this work, we have calculated the bipolar components using  $\Delta x = 0.02, 0.01, 0.005$ , and  $0.0025 \text{ b}$  for 20 stationary states with energies in the range of 10–200% of the barrier height.

The  $\rho_{\pm}$  densities for the stationary state with energy  $E = V_0$  (and  $\Delta x = 0.0025 \text{ b}$ ) are plotted as the solid and dashed lines in Figure 1a along with the total probability density  $\rho = |\Phi|^2 = |\Phi_+ + \Phi_-|^2$ , which is represented by the dotted line. Clearly, the  $\rho_{\pm}$  densities are constant in the asymptotic regions and vary more or less smoothly across the range of the potential barrier. The total density  $\rho$  oscillates in the region to the left of the barrier because of interference between the  $\Phi_{\pm}$  components. On the right-side of the barrier the reflected wave vanishes and there is no interference. Here, the value  $\rho_+(x \rightarrow \infty)$  is related to the transmission probability  $P_T$ , and the value  $\rho_-(x \rightarrow -\infty)$  is related to the reflection probability  $P_R$ . In practice, these quantities are estimated at the edges of the numerical grid ( $x = \pm 3.085 \text{ b}$  in this case). In principle, the 1D CPWM above is exact, so these quantities as well as  $\Phi_{\pm}$  may be computed to arbitrary precision, via a suitable choice of the numerical parameters (grid edges, grid spacing, time step, etc.).



**Figure 1.** (a) The (solid)  $\rho_+$ , (dashed)  $\rho_-$ , and (dotted)  $\rho$  probability densities corresponding to the  $E = V_0$  stationary scattering state of the 1D Eckart barrier are plotted as a function of position  $x$  in units of bohr. The barrier height is  $V_0 = 0.0018$  hartree, and the shaded area shows the potential function in units of  $V_0$ . (b) Benchmark CPWM results for the transmission ( $P_T$ , circles) and reflectance ( $P_R$ , squares) probabilities as a function of energy for 20 stationary states of the 1D Eckart barrier. The grid spacing for these results is  $\Delta x = 0.0025$  b. The solid curves corresponds to the exact result for this problem.

In Figure 1b we plot the calculated  $P_T$  and  $P_R$  (circles and squares, respectively) as a function of energy in units of the barrier height, and the solid lines correspond with analytical results.<sup>51</sup> As expected,  $P_T$  increases ( $P_R$  decreases) with increasing energy, and qualitatively, there is good agreement between the present CPWM results and exact theory. In Figure 2a we show the fractional error in  $P_T$  as a function of energy for several values of the CPWM grid spacing. This error is computed according to the formula

$$(P_T)_{\text{error}} = \frac{1}{P_T} |P_T^{\text{calc}} - P_T| \quad (15)$$

where  $P_T^{\text{calc}}$  is the calculated transmission probability and  $P_T$  is the corresponding exact result. The error is generally larger at lower energies, where the time steps are also larger, and the sporadically low errors, e.g., at  $E = 0.5V_0$ , are likely due to a fortuitous cancellation that can occur when the error changes sign. The fractional errors in  $P_T$  are less than 0.1% across the given energy range for all grid spacings, and at a given energy, the error is reduced by about half an order of magnitude as  $\Delta x$  is decreased by a factor of one-half.

One advantage of working with stationary scattering states is that the asymptotic form of the wave functions and the energy of the system are known in advance. This fact can be used in different ways to estimate the numerical error in calculations for which analytical results are unavailable. Panels b–d of Figure 2 illustrate three additional fractional error measures as a function of energy for different grid spacings. If the exact values of  $P_T$  are not known, then the error can be estimated by calculating  $P_T$  and  $P_R$  from  $\rho_+(x_{\text{max}})$  and  $\rho_-(x_{\text{min}})$ , respectively, and comparing with the exact relationship  $P_T + P_R = 1$ . The normalization errors for the present calculations are shown in Figure 2b, where it is seen that this error estimate increases very regularly with increasing energy. Also, the normalization error for this problem drops by a full order of magnitude as the grid spacing decreases by a factor of one-half.

For the second error estimate, we note that the bipolar superposition should constitute a solution of the TISE, and we can use an independent numerical method to compute the expectation

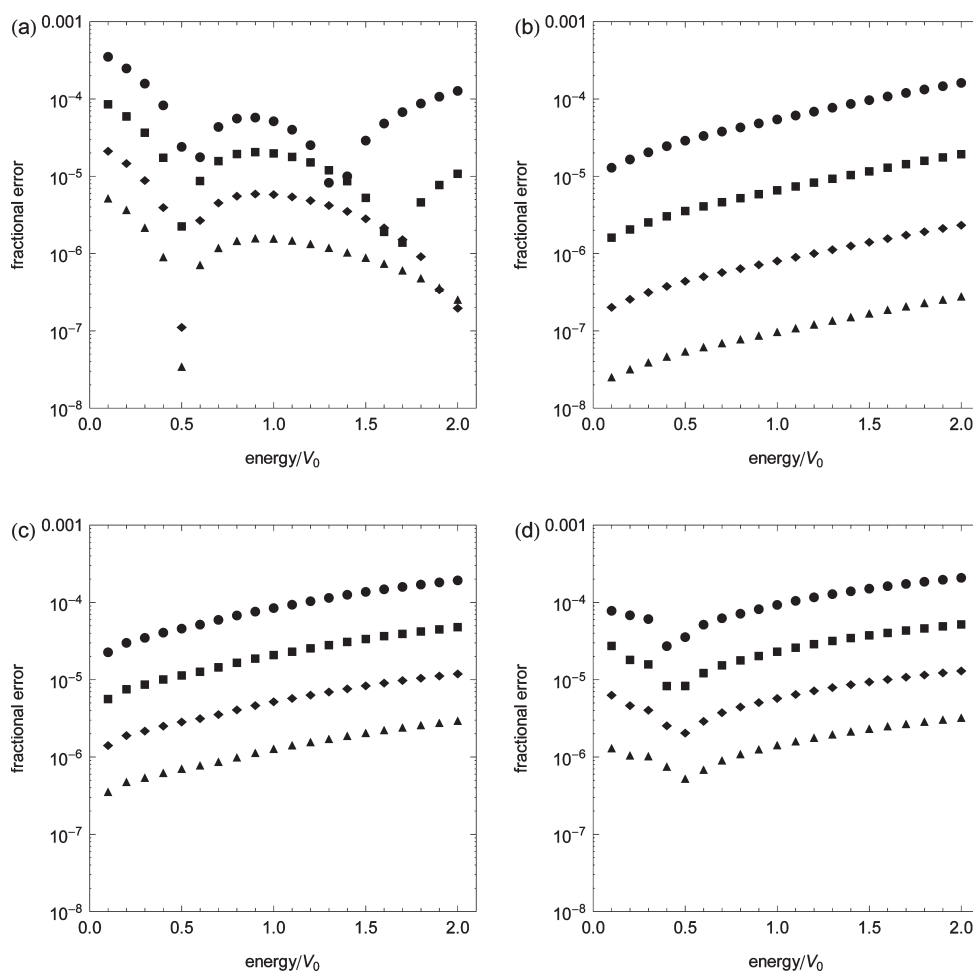
value of the energy from the CPWM stationary state. The points shown in Figure 2c represent the error associated with the energy expectation value for different grid spacings. These errors are computed according to the formula

$$\langle \hat{H} \rangle_{\text{error}} = \frac{1}{E} \left| E - \frac{\int \Phi^*(x) \hat{H} \Phi(x) dx}{\int \Phi^*(x) \Phi(x) dx} \right| \quad (16)$$

where the operation of the Hamiltonian is evaluated using fourth-order finite difference derivatives for the kinetic energy term and Boole's method for the numerical integrations.<sup>52</sup> In this way, the error estimate is averaged over the entire stationary state via numerical integration. These calculations do introduce additional errors; however, because the grid spacings are relatively small, such errors are expected to be less significant than the error in the CPWM calculations. For the given energy range, we see that the fractional error is less than 0.1% and decreases with the grid spacing in the same fashion as the fractional error in  $P_T$ . For the third error estimate, we note that the FF condition must also be satisfied at long times. Equation 7 can be rearranged and combined with  $p = (2mE)^{1/2}$  to solve for the system's energy, and this provides another independent method to estimate the error. The points shown in Figure 2d represent the error associated with how well the FF condition is satisfied according to the equation

$$\text{FF}_{\text{error}} = \frac{1}{E} \left| E + \frac{\hbar^2}{2m} \frac{\int \Phi(x) \Phi'(x) dx}{\int \Phi^*(x) (\Phi_+(x) - \Phi_-(x)) dx} \right|^2 \quad (17)$$

where we have used finite differences and numerical integration to perform the error calculation. These fractional errors follow the same trend as the error in the energy expectation value; however, these results show that the calculated bipolar CPWM solutions are internally consistent with the FF condition.



**Figure 2.** Various fractional error measures for the 1D Eckart A problem as a function of energy and for several different CPWM grid spacings  $\Delta x$ : (circles) 0.02, (squares) 0.01, (diamonds) 0.005, and (up-triangles) 0.0025 in units of bohr. (a) Fractional error measure in the calculated transmission probability  $P_T$ . (b) Fractional error associated with the normalization condition  $P_T + P_R = 1$ . (c) Fractional error estimate associated with the expectation value of the Hamiltonian. (d) Fractional error estimate associated with the FF condition.

Finally, we note that all of these errors can be reduced by employing a wider grid, smaller grid-point spacings, or smaller time steps in order to achieve higher accuracy, and it is not necessary to adjust these in a nonlinear way to improve the convergence. This is not the case for ABCs, where there are subtleties associated with finessing the onset of the absorbing potential, its width, height, and general shape. Also, of course, the coordinate range needed here for a given level of accuracy is always much smaller than if an ABC were used.

### 3. THE BIPOLAR REACTION PATH HAMILTONIAN (BRPH) APPROACH

In this section we formulate the BRPH approach for the simplest class of multi-D reactive scattering problems; namely, a 2D system with a linear reaction coordinate. We also present a discussion of our numerical implementation for the calculation of stationary state wave functions involving multiple scattering channels and state-to-state reaction probabilities, i.e., the probabilities associated with elements of the scattering matrix ( $S$ -matrix). Our general strategy is to transform the 2D scattering problem into a corresponding 1D multichannel scattering

problem, involving individual channel scattering amplitudes defined along a suitably chosen linear reaction path. We invoke a bipolar representation for each channel scattering amplitude, and this leads to a set of coupled equations of motion that are formally similar to the 1D CPWM equations described in section 2. Hence, we can then exploit the same numerical algorithms for 1D scattering problems to determine the stationary state wave functions of the 2D problem.

**3.1. Theory.** We consider a two-dimensional (2D) scattering problem described in some appropriately chosen mass-weighted Cartesian (MWC) coordinate system. In this case, the Hamiltonian operator has the form

$$\hat{H} = -\frac{\hbar^2}{2m} \left( \frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} \right) + V(x, y) \quad (18)$$

where  $m$  is the reduced mass and  $V$  is a 2D potential energy surface. For realistic problems, such as the collinear  $A + B - C \rightarrow A - B + C$  exchange reaction, the potential will necessarily exhibit a bend in the MWC space.<sup>5</sup> However, in the present work, we limit our consideration to the simplified case, where the reaction

path is both linear and parallel to the  $x$ -axis. The  $x$  (reaction) coordinate is then directly associated with some unbound scattering motion across a potential barrier, and the  $y$  (perpendicular) coordinate is associated with bound vibrational motion. Also, we must require that  $V$  approach a function that is independent of  $x$  in each asymptotic limit, and we further assume (for this paper) that this scattering potential is asymptotically symmetric, i.e.,  $V(x \rightarrow \pm\infty, y) = V_{\text{asympt}}(y)$ . As discussed in section 2 and below, the latter assumption allows us to utilize the constant velocity CPWM approach (where  $V_{\text{eff}} = 0$ ), which leads to a more natural approach for the benchmark problems described in section 4. Extension of the following formulation to the more general case of an asymptotically asymmetric problem is relatively straightforward but involves details that have already been addressed in previous work.<sup>33</sup>

Next, we introduce a set of vibrational states that vary slowly, i.e., adiabatically, along the reaction coordinate. Note that this does not imply that we are assuming that the vibrational quantum number of the stationary state be conserved as a function of  $x$ . The adiabatic vibrational states are the solutions of a corresponding adiabatic TISE

$$\left[ -\frac{\hbar^2}{2m} \frac{\partial^2}{\partial y^2} + V(x, y) \right] \phi_i(x, y) = \varepsilon_i(x) \phi_i(x, y) \quad (19)$$

where both the eigenvalues  $\varepsilon_i$  and the wave functions  $\phi_i$  depend parametrically on  $x$ .<sup>53</sup> Our assumptions on the 2D potential imply that  $\varepsilon_i(x) \rightarrow E_i$  as  $x \rightarrow \pm\infty$ , where  $E_i$  are the eigenenergies of the asymptotic system with potential energy  $V_{\text{asympt}}$ . It is further assumed that these adiabatic eigenstates can be determined exactly or by some appropriate approximation, and it is this last point where we can make a connection to the RPH method. For example, if only the potential energy, force, and Hessian were known along the reaction path, then the adiabatic eigenstates could be approximately represented with harmonic oscillator wave functions.

The key ansatz within the BRPH approach is that the total stationary scattering state for a given energy  $E$  can be written as an infinite sum over the adiabatic eigenfunctions

$$\Phi(x, y, t) = \sum_i a_i(x, t) \phi_i(x, y) \quad (20)$$

and we refer to the set of 1D functions  $a_i$  as channel scattering amplitudes. Each channel scattering amplitude is then expressed as a superposition of bipolar components:

$$a_i(x, t) = a_{i+}(x, t) + a_{i-}(x, t) \quad (21)$$

where  $a_{i\pm}$  take the form of traveling waves that move in opposite directions along the reaction coordinate  $x$ . If the potential energy meets the criteria discussed above, then it is sensible to set  $V_{\text{eff}} = 0$  and invoke the constant velocity form for the bipolar components

$$a_{i\pm}(x, t) = \alpha_{i\pm}(x) \exp\left[\frac{i}{\hbar} (\pm p_i x - Et)\right] \quad (22)$$

where  $p_i = (2m(E - E_i))^{1/2}$  is the momentum of a particle with mass  $m$  and kinetic energy  $E - E_i$ . Similar to the 1D case, the oscillatory components of eq 22 are formally separated from the  $\alpha_{i\pm}$  amplitudes, and these are expected to vary over the

interaction region and asymptotically converge to values that may then be related to the elements of the  $S$ -matrix.

There are several points here that merit further discussion. First, we have presumed that the adiabatic energy eigenvalues form a discrete spectrum; however, the adiabatic Hamiltonian may also possess states corresponding to a continuous range of eigenenergies if the perpendicular degrees of freedom are semi-bound. Such considerations will be saved for our future work, and here we will only consider the discrete case. Second, in practice, we cannot numerically represent an infinite number of adiabatic states, so we will have to truncate the sum in eq 20 at some point in order to carry out feasible calculations. We expect the number of terms included will affect both the computational effort and the convergence associated with numerical results and that this may lead to various approximation strategies for different types of scattering problems. Finally, for a given total energy  $E$ , both open channel  $E_i < E$  and closed channel  $E_i > E$  scattering amplitudes contribute to the sum in eq 20. For open channels, the momenta  $p_i$  are real-valued and the bipolar components  $a_{i\pm}$  contain complex exponentials that oscillate with respect to position. For closed channels the  $p_i$  are imaginary and  $a_{i\pm}$  contain real-valued exponentials. Both cases must be included if the total stationary state is to be represented exactly, and it is clear that the spatial amplitudes  $\alpha_{i\pm}$  must vanish asymptotically for closed channels. Next, we develop the equations of motion for the bipolar components of the channel scattering amplitudes.

First, the FF condition is applied as a separate condition for each channel scattering amplitude

$$a'_i = \frac{i}{\hbar} p_i (a_{i+} - a_{i-}) \quad (23)$$

which along with the TISE, provide the necessary relationships to develop a set of coupled time-dependent equations of motion for the bipolar components. To accomplish this, we appeal to an adiabatic representation of the 2D Hamiltonian. The details of the derivation are provided in Appendix A and the resulting equations of motion are given by

$$d_t a_{i\pm} = F_i a_{i\pm} + G_i a_i + H_i \quad (24)$$

Before discussing the individual terms here, we consider the general structure of these equations, and in particular, the total time derivative:

$$d_t a_{i\pm} = \partial_t a_{i\pm} \pm \frac{p_i}{m} a'_{i\pm} \quad (25)$$

Like the 1D CPWM equations, this expression implies that the bipolar components  $a_{i\pm}$  are evolving within a set of corresponding Lagrangian reference frames  $x_{i\pm}(t)$  that satisfy the ancillary equations of motion

$$d_t x_{i\pm} = \pm \frac{p_i}{m} \quad (26)$$

whose solutions move in opposite directions for the  $\pm$  components and with different constant velocities  $v_i = p_i/m = (2(E - E_i)/m)^{1/2}$  for different channel scattering amplitudes. For open channels, the Lagrangian dynamics is more or less similar to the 1D case. For the closed channels, however, where  $E_i > E$ , there is an apparent problem that the velocities are imaginary. The solutions of eq 26 for the closed channel components can be expressed as

$$x_{i\pm}(t) = x_{i\pm}(0) \pm i|v_i|t \quad (27)$$

which are complex for real values of  $t$ . This may be avoided, however, if we impose the requirement that the closed channel components evolve in imaginary time, i.e., if  $t = -i|t|$ , then the closed channel trajectories will be real provided that the initial conditions  $x_{i\pm}(0)$  are also real. This and other important issues will be discussed further in section 3.2, where our numerical strategy is described in detail. A similar strategy was employed, and justified, in a previous study.<sup>31</sup> For now, we will continue with an analysis of the terms in eq 24.

The various factors in the BRPH equations of motion are given by

$$F_i = \frac{i}{\hbar}(E - 2E_i) \quad (28a)$$

$$G_i = -\frac{i}{\hbar}(\varepsilon_i - E_i) \quad (28b)$$

$$H_i = \frac{i\hbar}{2m} \sum_j I_{ij}^{(1)} a_j' + I_{ij}^{(0)} a_j \quad (28c)$$

The  $F_i$  terms are associated with free particle motion and provide a flux of amplitude both into and out of the scattering region. The  $G_i$  terms also play a role in determining how the bipolar components evolve within the scattering region. Most importantly, they provide a coupling between left- and right-traveling components of a given scattering channel,  $i$ , leading to intrachannel reflection. This coupling becomes significant in the region of space where the  $\varepsilon_i(x)$  deviates significantly from  $E_i$ . Finally, the  $H_i$  terms provide nonadiabatic coupling across different scattering channels, leading to the redistribution of vibrational energy along the scattering coordinate. Analogous terms are found within certain implementations of the 1D CPWM for applications to scattering systems involving multiple diabatic electronic states.<sup>34</sup> The nonadiabatic coupling terms presented here involve the functions

$$I_{ij}^{(1)}(x) = 2 \int \phi_i^*(x, y) \phi_j^{(1,0)}(x, y) dy \quad (29a)$$

$$I_{ij}^{(0)}(x) = \int \phi_i^*(x, y) \phi_j^{(2,0)}(x, y) dy \quad (29b)$$

that depend on various overlap integrals between the adiabatic wave functions and their spatial derivatives with respect to the reaction coordinate. Within the integrands, the notation  $\phi_j^{(m,n)}(x, y)$  represents a mixed partial derivative with respect to  $x$  and  $y$

$$\phi_j^{(m,n)}(x, y) = \left( \frac{\partial^m}{\partial x^m} \left[ \left( \frac{\partial^n \phi_j(x, y)}{\partial y^n} \right)_x \right] \right)_y \quad (30)$$

and the integrals are labeled with superscripts according to a correspondent spatial derivative of the channel scattering amplitudes. For example, in eq 28c the  $I_{ij}^{(1)}$  factor is paired with  $a_j'$  and the  $I_{ij}^{(0)}$  factor is paired with  $a_j$ . These integrals are significant in regions of space where the adiabatic eigenfunctions (but not necessarily eigenenergies) are changing, i.e., regions where motion along  $x$  and  $y$  are strongly coupled through the scattering potential,  $V(x, y)$ .

There are several important and interesting limiting cases to consider. The first of these involves the asymptotic behavior of the BRPH equations and the boundary conditions for the bipolar

components. For a given energy  $E$ , there will be  $i_{\max}$  open scattering channels that satisfy  $E > E_i$  and therefore a total of  $i_{\max}$  degenerate left-incident solutions, i.e., one per open channel. The numerical solution for each such degenerate state requires a separate BRPH calculation with distinct initial and boundary conditions. For the solution left-incident on channel  $n$ , the boundary conditions are given by

$$a_{i+}(x \rightarrow -\infty, t) = \delta_{in} \exp \left[ \frac{i}{\hbar}(p_n x - Et) \right] \quad (31a)$$

$$a_{i-}(x \rightarrow +\infty, t) = 0 \quad (31b)$$

and a set of working initial conditions may be defined as

$$a_{i+}(x, t = 0) = \delta_{in} \exp(ip_n x / \hbar) \quad (32a)$$

$$a_{i-}(x, t = 0) = 0 \quad (32b)$$

As  $x \rightarrow \pm\infty$ , the adiabatic eigenenergies and eigenfunctions, by definition, become constant with respect to  $x$ . Hence, the  $G_i$  and  $H_i$  terms will vanish, and the resulting equations of motion in  $x$ , i.e., eq 24, are consistent with free-particle evolution. In the long-time limit, the  $a_{i\pm}$  solutions reach a steady state, and the  $S$ -matrix elements can be calculated from the asymptotic values. The transmission probability from state  $n$  to state  $j$  (state-to-state reaction probability) is given by

$$P_{n \rightarrow j} = \sqrt{\frac{E - E_j}{E - E_n}} |a_{j+}(x \rightarrow +\infty)|^2 \quad (33)$$

Partially state resolved and cumulative reaction probabilities can be calculated by taking different summations of the state-to-state reaction probabilities.

A second interesting limit involves the reduction of the system from a 2D problem to a 1D problem, such as for the case where the system becomes highly confined along  $y$ . Here, the adiabatic eigenenergies will be vastly separated from one another so that the only open scattering channel will be the ground state ( $i = 0$ ) of the asymptotic system. The asymptotic zero-point energy  $E_0$  is an arbitrary quantity and may be set to zero without consequence, and the asymptotic excited state energies become infinite by comparison  $E_{i \neq 0} \rightarrow \infty$ . The closed channel scattering amplitudes have imaginary momenta  $p_{i \neq 0} = i|p_i|$ , where  $|p_{i \neq 0}| = (2m|E - E_i|)^{1/2} \rightarrow \infty$  in the highly confined limit. According to eq 22,  $a_{i \neq 0, +} \rightarrow 0$  for all  $x$ , and the  $a_{i \neq 0, -}$  will either be divergent if  $\alpha_{i \neq 0, -} \neq 0$  or vanish if  $\alpha_{i \neq 0, -} = 0$ . The former case is unphysical; therefore, all closed channel scattering amplitudes must vanish as the system is reduced to 1D. In this case, the nonadiabatic term  $H_0$  also reduces to

$$H_0 = \frac{i\hbar}{2m} I_{00}^{(1)} a_0' + I_{00}^{(0)} a_0 \quad (34)$$

By expanding the adiabatic eigenfunctions in eq 29a using any complete and orthonormal set of 1D basis functions along  $y$ , it can be shown in general that the function  $I_{ii}^{(1)} = 0$ , for all  $i$  and for all  $x$ ; however, the function  $I_{ii}^{(0)}$  is generally nonzero. The quantity  $\varepsilon_0 - I_{00}^{(0)}$  takes the role of a 1D scattering potential, and the BRPH equations reduce to

$$d_t a_{0\pm} = -\frac{i}{\hbar} E a_{0\pm} - \frac{i}{\hbar} (\varepsilon_i - I_{00}^{(0)}) a_0 \quad (35)$$

which are formally equivalent to the 1D constant velocity CPWM equations.

A third interesting limit is that of a separable potential,  $V(x,y) = V_x(x) + V_y(y)$ , for which the adiabatic eigenfunctions,  $\phi_i(x,y)$ , are independent of  $x$ , but the eigenenergies,  $\varepsilon_i(x) = E_i + V_x(x)$ , need not be. In this case, the  $H_i$  terms obviously vanish completely, so there are no nonadiabatic transitions, as is appropriate.

**3.2. Numerical Implementation.** Next we present our numerical implementation of the BRPH approach. Before solving the BRPH equations, we must first determine the eigenenergies and eigenfunctions of the adiabatic Hamiltonian as a function of the reaction path. For model problems, such as the harmonic and Morse oscillators, this can be done analytically, and relevant formulas for the harmonic case are given in Appendix B. For more general problems, however, one must resort to approximations or numerical calculations. In our development work we have employed DVRs to calculate eigenstate wave functions of the adiabatic Hamiltonian. Here a set of grid points is defined along the reaction path  $x$  and we use Colbert and Miller's "universal DVR" with equally spaced grid points to define the matrix elements of the kinetic energy operator for motion along the  $y$ -coordinate

$$T_{\alpha\beta}^{(y)} = \frac{(-1)^{\alpha-\beta}}{2m\Delta y^2} \begin{cases} \frac{\pi^2}{3} & \alpha = \beta \\ \frac{2}{(\alpha-\beta)^2} & \alpha \neq \beta \end{cases} \quad (36)$$

where  $\Delta y$  is the grid spacing and the indices  $\alpha$  and  $\beta$  run over interior DVR grid points.<sup>54</sup> The potential energy is approximated as a diagonal matrix over the DVR grid points  $V_{\alpha\beta}^{(y)} = V(x,y_\alpha)\delta_{\alpha\beta}$  and varies parametrically with the reaction coordinate. For each grid point along  $x$  the perpendicular DVR Hamiltonian, with elements  $H_{\alpha\beta}^{(y)} = T_{\alpha\beta}^{(y)} + V_{\alpha\beta}^{(y)}$ , is diagonalized to give a finite set of eigenvalues  $\varepsilon_i$  and DVR eigenvectors  $\phi_i^{\text{DVR}}$ . To obtain the desired eigenfunctions, the eigenvectors are multiplied by the appropriate weights, which in this case are related to the grid spacing

$$\phi_i(x, y_\alpha) = \frac{[\phi_i^{\text{DVR}}]_\alpha}{\sqrt{\Delta y}} \quad (37)$$

Numerically speaking, the sign of the DVR eigenvectors is irrelevant, such that  $\pm\phi_i^{\text{DVR}}$  are equivalent, and typical numeric algorithms will give results with arbitrary sign at different points along the reaction coordinate. We must be careful to correct this sign mismatch before postprocessing these quantities. The integrands of eq 29 involve the first and second derivative of the adiabatic wave functions with respect to  $x$ . These can be efficiently calculated using the DVR representation of the appropriate derivative operators and the inherent quadrature properties of the DVR approach or by using finite difference derivatives and common numerical integration formulas.

The BRPH solutions are represented over a uniform spatial grid with spacing that covers the range over which the scattering potential is numerically significant. This grid should be large enough that the calculated reaction probabilities are converged, and the grid spacing  $\Delta x$  should be small enough to achieve the desired accuracy. This is also important for calculating the spatial derivatives in the nonadiabatic coupling term  $H_i$ . Notably, these derivatives involve the total channel scattering amplitudes  $a'_i = a'_{i+} + a'_{i-}$ , which can be evaluated using either finite differences or the FF conditions. We have implemented both strategies and found that using the FF conditions leads to a more stable convergence of the BRPH solutions compared to finite differences,

especially for energy values close to the onset of a channel threshold. For the results discussed in section 4 we have used the FF conditions to evaluate the derivative terms in  $H_i$ .

As we have noted, the BRPH equations of motion are expressed in a Lagrangian reference frame, and this means that the  $\Phi_{i\pm}$  components are moving in opposite directions and with different velocities  $v_i = (2(E - E_i)/m)^{1/2}$  for different channels. The grid points over which these functions are represented move with the flow of the BRPH components so that at different times the grid points for different  $a_{i\pm}$  will no longer be coincident with one another. Moreover, for closed channel amplitudes, the velocity is imaginary and we must invoke imaginary time propagation to keep the associated reference frames on the real axis. For an arbitrary universal time step  $\Delta t$  (and  $-i|\Delta t|$ ), we would need to employ interpolation methods to properly evaluate the Lagrangian time derivatives for each component. Such methods have previously been implemented for 1D CPWM numerical algorithms,<sup>33</sup> and have also been carried out for 1D problems involving multiple diabatic potential surfaces.<sup>34</sup> Certainly this scheme is a viable approach here and we note that this would offer much greater control over the time step and accuracy of the BRPH calculations. However, the interpolation codes are somewhat awkward to work with and create some extra computational overhead. At the present stage of development we have avoided the interpolation issue by using a unique time step for each channel

$$\Delta t_i = \frac{\Delta x}{v_i} = \Delta x \sqrt{\frac{m}{2(E - E_i)}} \quad (38)$$

so that the + (or -) BRPH grid points associated with the different channels will be coincident with one another after each time step. This choice automatically yields imaginary time steps for the closed channel amplitudes. Of course, this then implies that the solutions will not be coincident with respect to time, and we will come back to this important point momentarily. For individual time steps, we employ a second-order Runge–Kutta integration scheme,<sup>55</sup> where the BRPH solutions are propagated using the formula

$$a_{i\pm}(x \pm \Delta x, t + \Delta t_i) = a_{i\pm}(x, t) + K_{i\pm}^{(2)} \quad (39)$$

where

$$K_{i\pm}^{(1)} = \Delta t_i d_t(a_{i\pm}(x, t)) \quad (40a)$$

$$K_{i\pm}^{(2)} = \Delta t_i d_t(a_{i\pm}(x, t) + K_{i\pm}^{(1)}/2) \quad (40b)$$

The factor enclosed by parentheses in eq 40b represents a first-order half-step, i.e.,  $\Delta t_i/2$ , while eq 39 is the second-order full-step. The error for each Runge–Kutta step is third-order in  $|\Delta t_i|$ , which is different for different channel amplitudes, and the cumulative error is second-order.

Note that in the application of eqs 39 and 40 we must take care of the fact that  $\pm$  grid points at the half-step will be offset from one another. This can be efficiently handled by applying a shift function, as necessary, to evaluate all of the  $\Phi_{i\pm}$  components at the same point in space. For example, consider a discrete numerical representation of the channel scattering amplitudes:

$$a_{i\pm}(t) = \{a_{i\pm}(x_1), \dots, a_{i\pm}(x_{k-1}), a_{i\pm}(x_k), a_{i\pm}(x_{k+1}), \dots, a_{i\pm}(x_N)\} \quad (41)$$

where the set of points  $\{x_1, \dots, x_N\}$  are the numerical grid points along the reaction coordinate. Note that  $x_{i+1} = x_i + \Delta x$ , so that after a full first-order time step  $\Delta t_i$ , the components will be given by

$$a_{i+}(t + \Delta t_i) = \{a_{i+}(x_2), \dots, a_{i+}(x_k), a_{i+}(x_{k+1}), a_{i+}(x_{k+2}), \dots, a_{i+}(x_{N+1})\} \quad (42a)$$

$$a_{i-}(t + \Delta t_i) = \{a_{i-}(x_0), \dots, a_{i-}(x_{k-2}), a_{i-}(x_{k-1}), a_{i-}(x_k), \dots, a_{i-}(x_{N-1})\} \quad (42b)$$

where it is clear that  $a_{i\pm}$  components are not properly aligned with one another in space. Equation 42 is also valid for first-order half-steps provided that we make the replacement  $x_i \rightarrow x'_i = x_i + \Delta x/2$ . To evaluate the time derivative of the components for the next time step, the  $\pm$  components need to be calculated at the same point in space, so we apply a pair of shift functions that yield

$$\text{Shift}_+[a_{i+}(t + \Delta t_i)] = \{a_{i+}(x_{N+1}), \dots, a_{i+}(x_{k-1}), a_{i+}(x_k), a_{i+}(x_{k+1}), \dots, a_{i+}(x_N)\} \quad (43a)$$

$$\text{Shift}_-[a_{i-}(t + \Delta t_i)] = \{a_{i-}(x_1), \dots, a_{i-}(x_{k-1}), a_{i-}(x_k), a_{i-}(x_{k+1}), \dots, a_{i-}(x_0)\} \quad (43b)$$

Here the interior positions are now coincident in space and the points at the edges of the BRPH grid are meaningless and potentially dangerous. This is not a problem, however, because between time steps the edge points are replaced with the appropriate boundary conditions for the desired solutions

$$a_{i+}(t + \Delta t_i) = \{e^{i/\hbar(p_i x_1 - E \Delta t_i)}, \dots, a_{i+}(x_{k-1}), a_{i+}(x_k), a_{i+}(x_{k+1}), \dots, a_{i+}(x_N)\} \quad (44a)$$

$$a_{i-}(t + \Delta t_i) = \{a_{i-}(x_1), \dots, a_{i-}(x_{k-1}), a_{i-}(x_k), a_{i-}(x_{k+1}), \dots, 0\} \quad (44b)$$

These shifts are applied after both the half and full time step so that the BRPH solutions are always coincident in space, thus eliminating the need for interpolation schemes.

Certainly, the mismatch in the treatment of time between different scattering amplitudes is somewhat counterintuitive and warrants some concern. To justify this, we must recognize that the time-dependent dynamics encapsulated by both the CPWM and BRPH equations do not represent the actual physical dynamics in the deterministic sense. By this we mean that the BRPH solutions at intermediate times have no physical importance apart from the fact that the (exact FF solution) bipolar components possess a dynamic phase factor,  $e^{-iEt/\hbar}$ , that is consistent with the true stationary state. However, this time dependence is clearly known *a priori*, and the value of  $t$  is completely arbitrary for stationary states. In our numerical applications we have often found that it is useful to reset the phase of the CPWM or BRPH solutions between time steps  $\Delta t_i$  by applying a conjugate phase factor of  $e^{+iE\Delta t_i/\hbar}$ . This works for both real and imaginary time steps and has no consequential effect in the long-time limit, changing only the intermediate-time dynamics, which are unimportant. Thus, neither the precise form of the initial conditions for the BRPH propagation nor the slight phase shifts between time-steps as discussed above prevent the method from reaching the desired steady-state solution in the long-time limit, at least not in principle.

That said, CPWM algorithms are somewhat similar to numerical optimization problems, such as the Newton–Raphson method, or other iterative methods for self-consistently solving nonlinear equations, like the Hartree–Fock equations. Nonconvergent, or worse yet divergent, solutions will often occur when the initial guess is too far removed, or perhaps even completely isolated (in solution space) from the desired solution. We have observed similar behavior in some BRPH calculations, especially for energies very close to the onset of a scattering channel. From experience, we expect to encounter difficulties with convergence when the quantities  $E - E_i$  are small and the troublesome scattering amplitude has a very broad de Broglie wavelength compared to the other components and the size of the interaction region. Related issues are also encountered in 1D applications, as the kinetic energy approaches zero. The obvious solution is to add points to the numerical grid. Ideally it would be desirable to have enough grid points outside the interaction region to guarantee that  $\varepsilon_i(x)$  and the adiabatic couplings are numerically converged. At the same time, the grid spacing must be chosen to accurately represent the scattering potential. At some point, usually when the energy is an extremely small fraction of the barrier height, one cannot afford to add enough points to maintain an adequate representation of both the potential and the scattering state simultaneously. For multi-D calculations, the problem is more serious because the same issues will occur over an energy range where the S-matrix elements are nontrivial and do exhibit interesting features. Note that this situation is not unique to BRPH and is, in fact, problematic for virtually all accurate quantum scattering methods (especially DVR-ABC). In any case, it is clear that, like the 1D case, the BRPH will ultimately fail for some small enough value of  $E - E_i$ ; however, we would like to be able to push the method as far as possible.

In this context, there are several avenues that could lead to enhanced numerical stability. One idea is to take the final BRPH components from a presumably stable calculation at higher energy and use these as the initial conditions for lower kinetic energies, where stability is problematic. To compute the S-matrix as a function of energy, one would scan the energy backward from higher to lower values. We could also mix different solutions from completed calculations above and below the channel onset. Another approach may be to slowly “switch on” the nonadiabatic coupling terms with an appropriately defined scaling factor. In our work we have employed a very simple formula

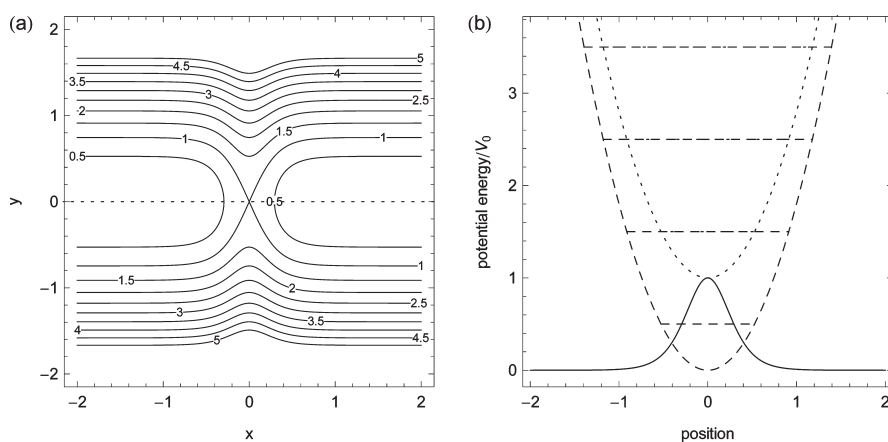
$$a_{i\pm}(\text{new}) = a_{i\pm}(\text{old}, t) + \eta[a_{i\pm}(\text{old}, t + \Delta t_i) - a_{i\pm}(\text{old}, t)] \quad (45)$$

which mixes the BRPH solutions between time steps, and the mixing parameter  $\eta$  is a number between (0,1).<sup>56</sup> We have found that this does enhance numerical stability for some cases, but ultimately fails as the energy gets even closer to the onset of a scattering channel.

## 4. RESULTS

In this section we present benchmark BRPH results for several simple 2D scattering problems. The first example involves the trivial case of a separable system defined by an Eckart barrier along  $x$  and a harmonic oscillator (HO) along  $y$ . We refer to this system as the uncoupled Eckart+HO problem. Since there is no coupling between  $x$  and  $y$ , the adiabatic eigenfunctions are constant with respect to the reaction coordinate, such that the





**Figure 3.** (a) Contour lines illustrating the 2D potential energy surface for the uncoupled Eckart+HO problem. The isovalues are reported in units of the barrier height  $V_0 = 0.0018$  hartree. The dotted line corresponds to the minimum energy path. (b) Various 1D slices through the potential energy surface are shown. The solid curve illustrates the Eckart barrier along the linear reaction path  $V(x,0)$ . The dashed curves correspond to the asymptotic harmonic oscillator potential  $V(\pm\infty,y)$  and its four lowest energy eigenvalues. The dotted curves are associated with the harmonic potential and eigenvalues along the dividing surface  $V(0,y)$ .

nonadiabatic coupling terms vanish and the interchannel reaction probabilities  $P_{n \rightarrow i \neq n} = 0$  for all energies. Our results here will serve as a reference for a nonseparable problem that we consider later. The scattering potential for the uncoupled Eckart+HO problem is given by

$$V(x,y) = V_0 \operatorname{sech}^2(\alpha x) + \frac{1}{2} m \omega_0^2 y^2 \quad (46)$$

where the parameters of the problem are given as follows:  $m = 2000$  au,  $V_0 = \hbar \omega_0 = 0.0018$  hartree, and  $\alpha = 3 \text{ b}^{-1}$ . Diagrams illustrating the 2D structure and various 1D slices of the scattering potential are shown in Figure 3.

In Figure 4 we plot 1D probability densities for the channel scattering square amplitudes,  $\rho_i = |a_i|^2$ , and corresponding BRPH components,  $\rho_{i\pm} = |a_{i\pm}|^2$ , for the two degenerate left-incident stationary scattering state solutions of the uncoupled Eckart+HO problem with energy  $E = 2V_0$ . In each case, only a single channel is involved, i.e., the incident channel,  $i = n$ , because there is no nonadiabatic coupling. These functions are quite similar to the 1D example in section 2, except that there are multiple solutions to consider when there is more than one open scattering channel. Panels a and b show the densities for the  $n = 0$  and  $n = 1$  incident states, respectively. Clearly, the  $n = 0$  state has a larger kinetic energy and the scattering amplitude exhibits a shorter wavelength and greater transmission than the  $n = 1$  state.

The BRPH components and scattering amplitudes can be combined with the adiabatic eigenfunctions to generate 2D wave functions and probability densities:

$$\rho_+(x,y) = \left| \sum_i a_{i+}(x) \phi_i(x,y) \right|^2 \quad (47a)$$

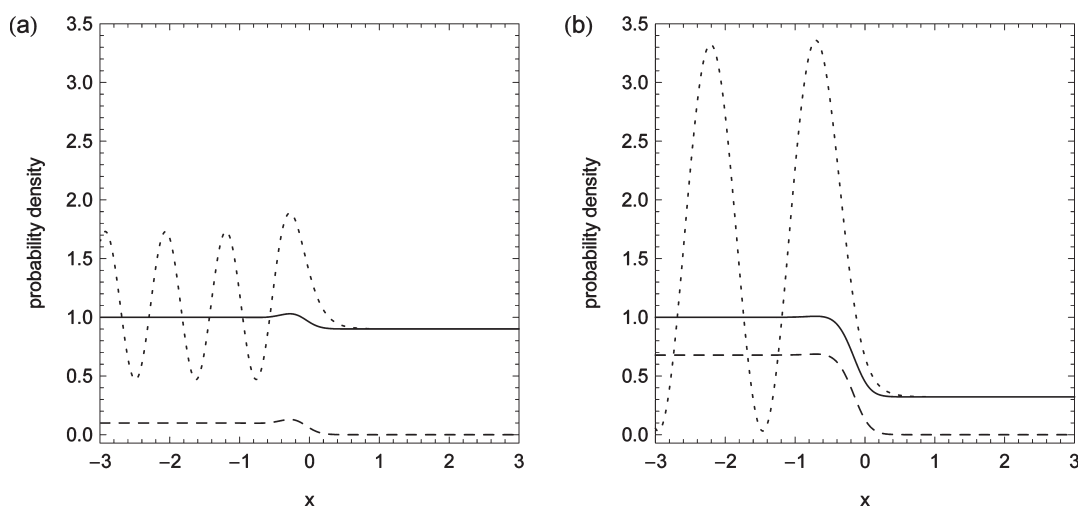
$$\rho_-(x,y) = \left| \sum_i a_{i-}(x) \phi_i(x,y) \right|^2 \quad (47b)$$

$$\rho(x,y) = \left| \sum_i a_i(x) \phi_i(x,y) \right|^2 \quad (47c)$$

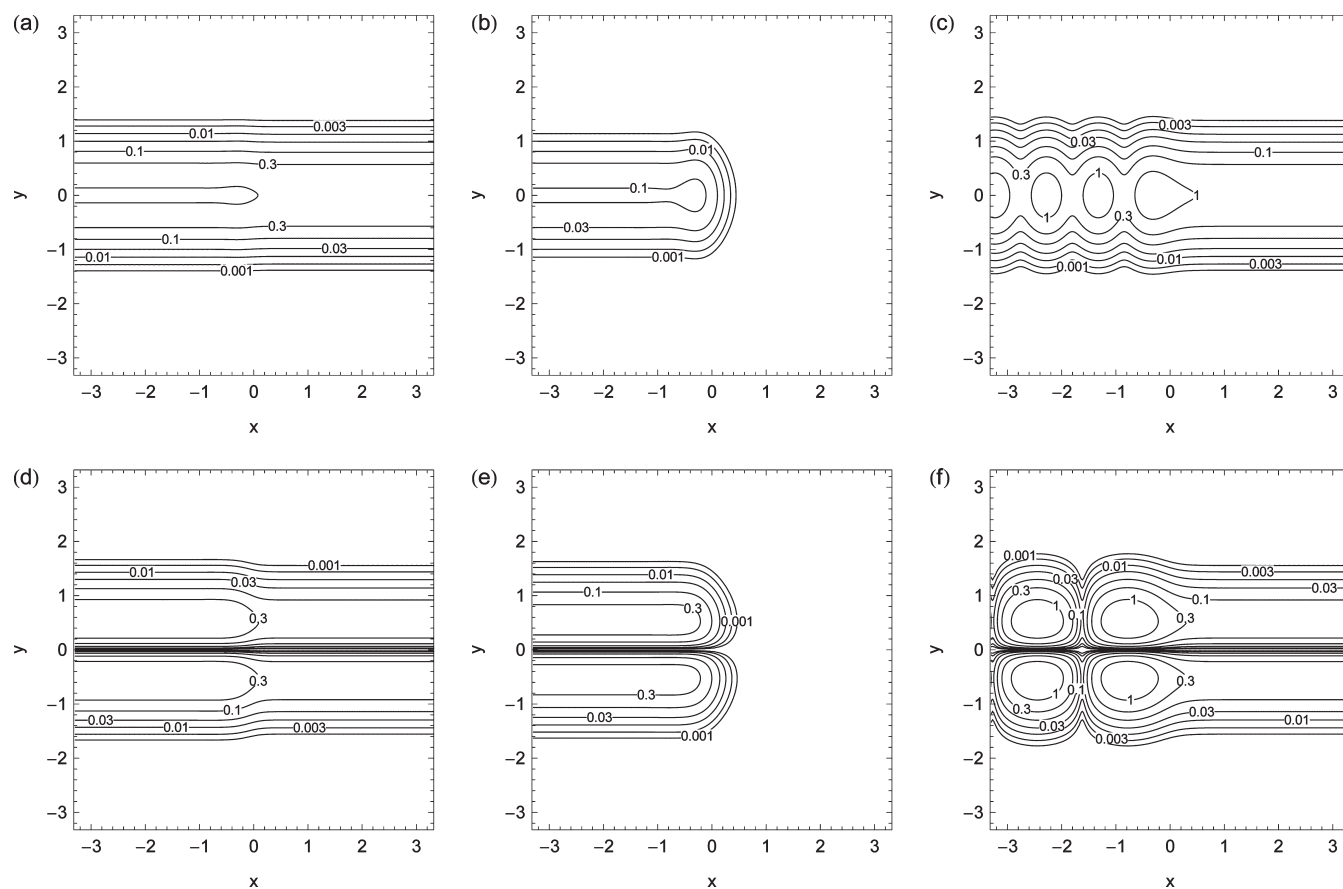
where  $\rho_+$  represents the density of the total incident and transmitted wave,  $\rho_-$  is the density of the total reflected wave, and  $\rho$  is the total probability density of the stationary state. In Figure 5, we plot these densities for the two degenerate left-incident states at  $E = 2V_0$ . Because there is no interference between channels in this case, the wave functions involved are simply the product of a single 1D scattering wave function (single-channel scattering amplitude) and an HO eigenfunction. Consequently, these densities are not particularly interesting for the uncoupled case; however, they do serve as a useful point of reference.

For the uncoupled Eckart+HO problem, the intrachannel reaction probabilities are identical to the transmission probability for the 1D Eckart barrier. In Figure 6a we plot the intrachannel reaction probabilities  $P_{n \rightarrow n}$  as a function of energy for the first three scattering channels with  $E < 3.5V_0$ . The calculated values for the  $n = 0, 1,$  and  $2$  channels are represented by circles, squares, and diamonds, respectively. The vertical dashed lines in the figure indicate the onset of a scattering channel, i.e.,  $E = E_j$  and the exact results are shown as solid lines. We have calculated these results for several different grid spacings; however, only the most accurate results using  $\Delta x = 0.0025$  b are presented. Figure 6b shows the fractional error of the BRPH transmission probabilities. These errors are all less than 0.1% across the energy range and are quite similar compared to those for the 1D case; however, for the 2D problem, we have degenerate states at higher energies, and a separate error is given for each one. Generally, the error is larger for states where the kinetic energy is small and the time step is large. We note the anomalously low error in the data point for the  $n = 2$  channel at  $E = 2.6 V_0$ . This is the most challenging state for this set of calculations, and we suspect that there may be a small unconverted error leading to a coincidental cancellation in favor of a seemingly more accurate transmission probability.

The error estimates introduced in section 2 for CPWM calculations can also be generalized for BRPH calculations. In panels c and d of Figure 6 we plot the normalization error and  $\langle \hat{H} \rangle_{\text{error}}$  respectively. For multi-D problems, we define the normalization as the sum of the partially state resolved transmission and reflection probabilities, which should be unity for all open channels. The magnitude of both the normalization and energy expectation errors is less than 0.1% across the given energy range and exhibit similar trends with respect to the grid spacing (not shown) as



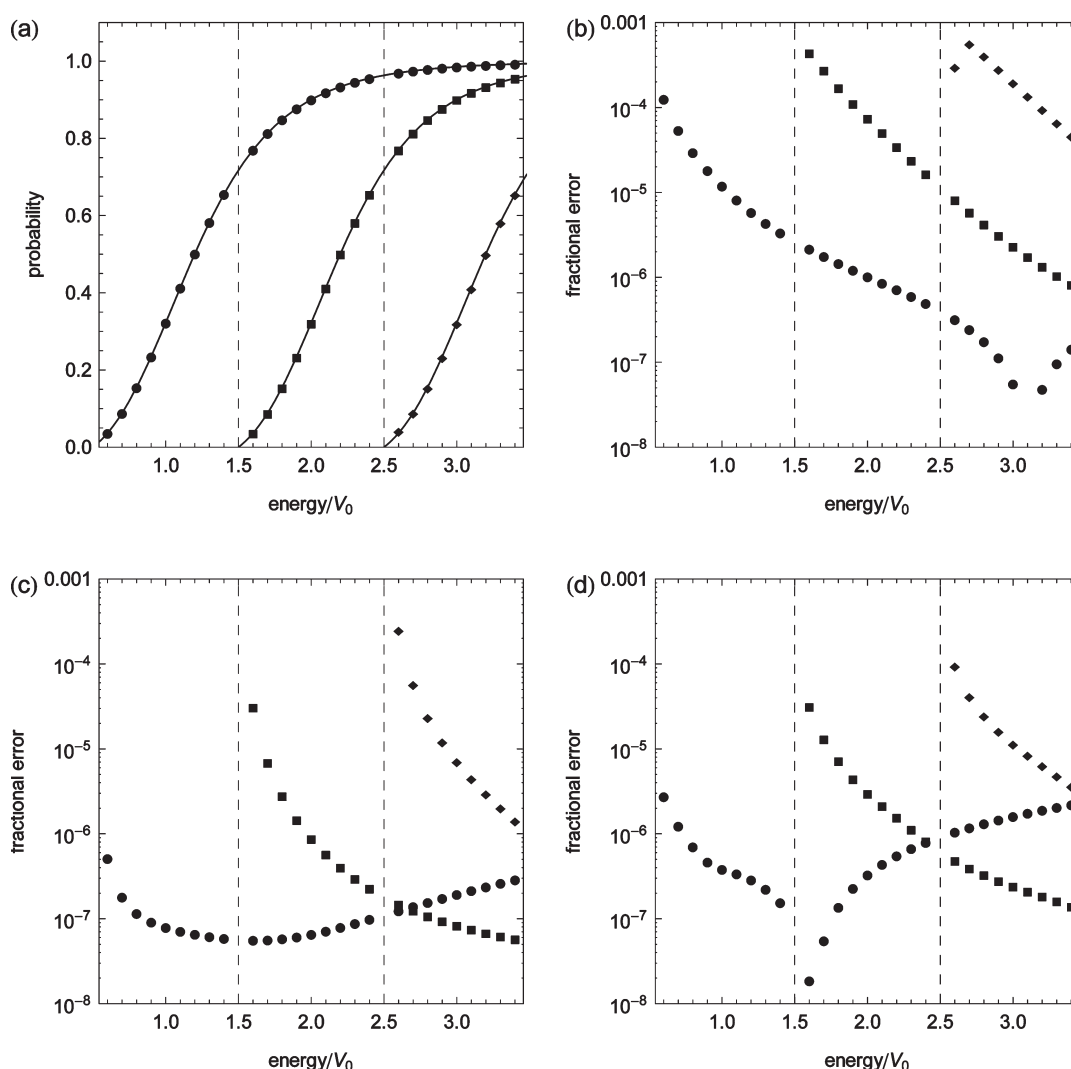
**Figure 4.** The BRPH densities for two degenerate stationary states of the uncoupled Eckart+HO problem are plotted as a function of the reaction coordinate: (solid)  $\rho_{n+}$ , (dashed)  $\rho_{n-}$ , and (dotted)  $\rho_n$ . These states correspond to the two open scattering channels with  $E = 2V_0$ . Panels a and b correspond to the  $n = 0$  and  $n = 1$  incident channels, respectively.



**Figure 5.** Contour plots illustrating 2D BRPH densities for the two degenerate stationary states of the uncoupled Eckart+HO problem with  $E = 2V_0$ . For the  $n = 0$  incident channel, we have (a)  $\rho_{+}$ , (b)  $\rho_{-}$ , and (c)  $\rho$ . Similarly, panels d–f are the corresponding densities for the  $n = 1$  incident channel.

compared to the CPWM results for the 1D Eckart A problem. Interestingly, however, the energy dependence of these errors is quite different compared to the errors shown in Figure 2. For the 1D problem, the errors increase regularly with energy, whereas in panels c and d of Figure 6 they generally decrease, although, the errors for the  $n = 0$  incident channel do increase

with energy, but only after the total energy has exceeded the  $n = 1$  channel onset. We speculate that the difference between the multi-D and 1D errors is related to the fact that the BRPH equations of motion contain constant terms, i.e.,  $E_i$ , that are not present in the 1D CPWM equations. Certainly this is an interesting, although very subtle, feature of our results that we will continue



**Figure 6.** Transmission probabilities and fractional error estimates as a function of energy in units of the barrier height  $V_0 = 0.0018$  hartree for the uncoupled Eckart+HO problem. Circles, squares, and diamonds represent the  $n = 0, 1,$  and  $2$  incident channels, respectively. Vertical dashed lines indicate the onset of a scattering channel. (a) Intrachannel transmission probabilities; solid lines are the corresponding exact results. (b) Fractional error in the calculated transmission probabilities. (c) Fractional error of the normalization condition for the 2D stationary states. (d) Fractional error in the energy expectation value of the stationary states.

to examine in future studies; however, the important point here is that BRPH approach reproduces analytical results extremely well for this uncoupled problem and that the error can be controllably reduced to arbitrary precision by decreasing the grid spacing.

In the next example, we consider a coupled Eckart+HO problem where the harmonic oscillator potential is displaced along the reaction coordinate. The potential energy surface is defined by

$$V(x, y) = V_0 \operatorname{sech}^2(ax) + \frac{1}{2} m\omega_0^2 (y - Y(x))^2 \quad (48)$$

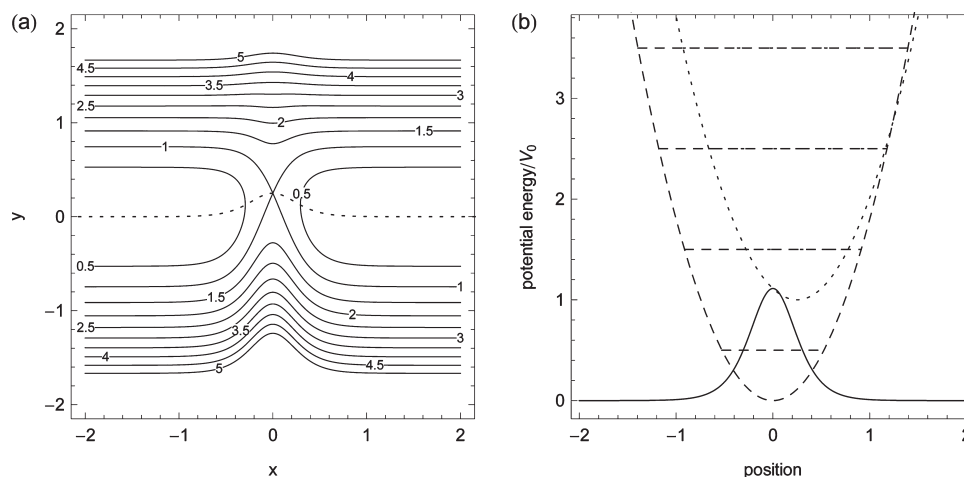
where the function

$$Y(x) = Y_0 \operatorname{sech}^2(ax) \quad (49)$$

provides a displacement that couples motion along the  $x$  and  $y$  coordinates. The displacement is zero as  $x \rightarrow \pm\infty$  and has a maximum value of  $Y_0 = 0.25$  b at the dividing surface  $x = 0$ . In principle, we could use this curve to define the reaction coordinate; however, this is not necessarily required because the reactant and product valleys of the potential are coincident with the  $x$ -axis. In this

sense, the problem is quasilinear and we will take the  $x$ -axis to be the reaction coordinate in our calculations. The other parameters of the coupled problem are identical to those in the uncoupled Eckart+HO example. Figure 7 shows the 2D potential energy surface and several 1D slices for the coupled Eckart+HO problem.

The fact that we have included a displacement of the oscillator, as opposed to varying only the harmonic frequency along the reaction coordinate, is important. In the case of the coupled Eckart+HO with no displacement<sup>35,57</sup> (i.e., symmetric about  $y$ ), the nonadiabatic coupling between even and odd adiabatic eigenstates vanishes due to symmetry, and one would need to probe higher energies in order to observe nonzero state-to-state transitions. However, the intrachannel reaction probabilities at high energies will far outweigh the interchannel probabilities for the symmetric Eckart+HO problem, i.e.,  $P_{0 \rightarrow 0} \approx 1 \gg P_{0 \rightarrow 2}$ , which makes it difficult to assess whether the nonadiabatic coupling terms are treated correctly. Breaking the symmetry of the problem by simply displacing the oscillator leads to much more obvious and interesting nonadiabatic effects at lower energies, even without



**Figure 7.** (a) Contour lines illustrating the coupled Eckart+HO potential. The isovalues are reported in units of the barrier height  $V_0 = 0.0018$  hartree. The dotted line corresponds to the minimum energy path. (b) Various 1D slices of the 2D potential are shown. The solid curve illustrates  $V(x,0)$ . The dashed curves correspond to  $V(\pm\infty,y)$  and the four lowest energy asymptotic eigenvalues  $E_i$ . The dotted curves represent  $V(0,y)$  and the four lowest adiabatic eigenenergies  $\varepsilon_i(0)$ . The maximum displacement of the oscillator is 0.25 b at the dividing surface.

attenuating the harmonic frequency. In Appendix B we derive expressions for the nonadiabatic BRPH integrals (eq 29) for the case where the harmonic oscillator is both displaced and scaled along the reaction coordinate. For the displaced only oscillator, we have the following result

$$I_{ij}^{(1)} = -\sqrt{2}\beta Y' \sqrt{(i+1)}\delta_{i,j-1} + \sqrt{2}\beta Y' \sqrt{i}\delta_{i,j+1} \quad (50a)$$

$$I_{ij}^{(0)} = \frac{1}{2}\beta^2 Y'^2 (\sqrt{(i+2)(i+1)}\delta_{i,j-2} + \sqrt{i(i-1)}\delta_{i,j+2}) - \frac{\sqrt{2}}{2}\beta Y'' (\sqrt{(i+1)}\delta_{i,j-1} - \sqrt{i}\delta_{i,j+1}) - \frac{1}{2}\beta^2 Y'^2 (2i+1)\delta_{i,j} \quad (50b)$$

where  $\beta = (m\omega_0/\hbar)^{1/2}$ . These terms introduce coupling between specific scattering channel amplitudes, and the magnitude of this coupling is scaled by functions of the reaction coordinate that depend on the derivatives of  $Y$ .

We have calculated the stationary states and reaction probabilities for the coupled Eckart+HO problem over a range of energies that includes up to three open scattering channels. Recall that the BRPH ansatz (eq 20) involves an infinite sum over both the open and closed scattering channels and that for coupled problems this sum must be artificially truncated. We have performed calculations including up to three additional closed channel amplitudes, and we found that with inclusion of only two closed channel terms, the accuracy in the results shown below is limited by the grid spacing, which is  $\Delta x = 0.0025$  b for our most accurate BRPH calculations. This is certainly not a general result, and we expect that more terms will be required to achieve convergence for different types of scattering problems, especially those with significant nonadiabatic coupling and also anharmonic character in the perpendicular degrees of freedom.

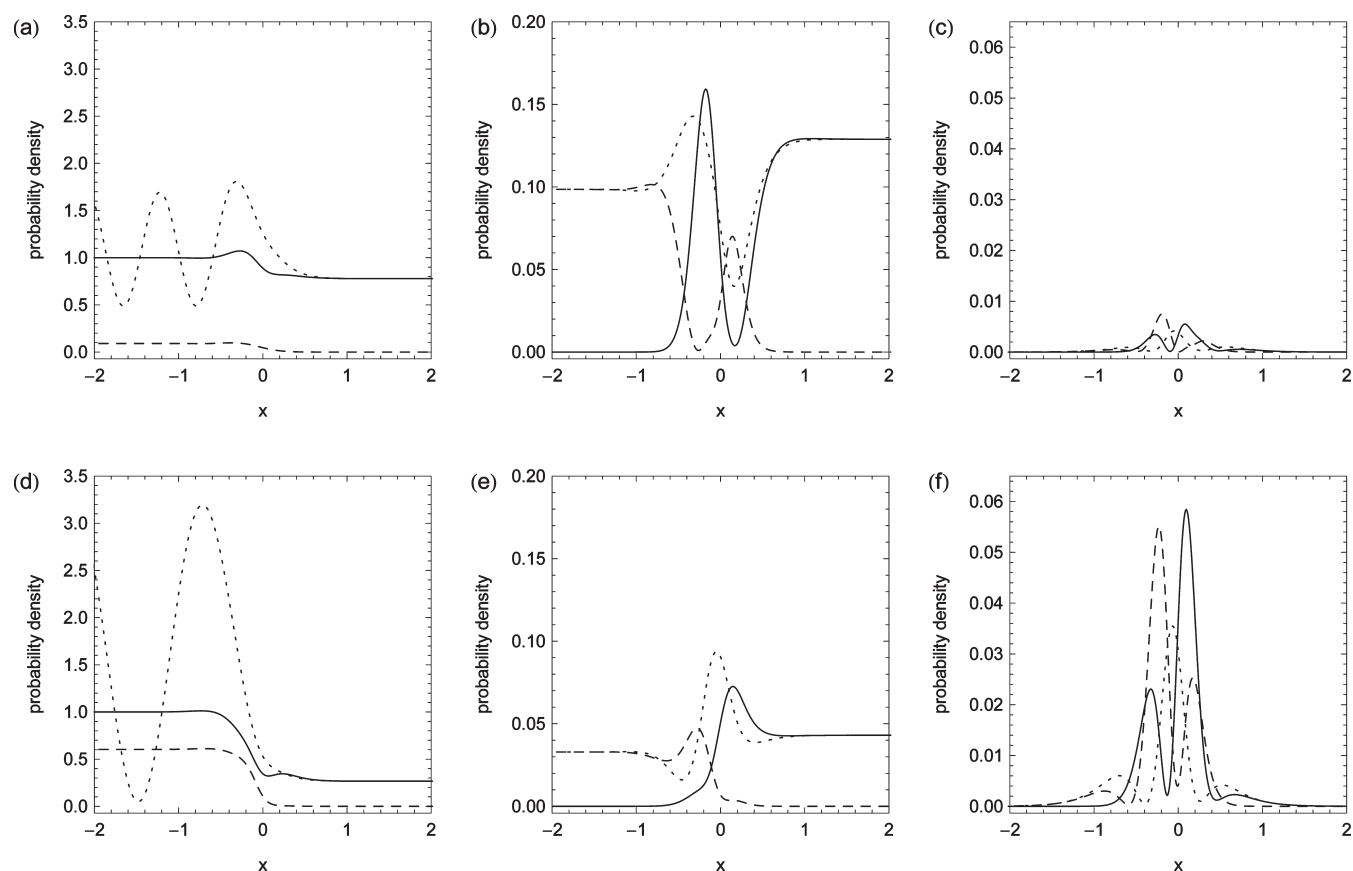
Figure 8 shows the BRPH densities as a function of  $x$  for the two degenerate left-incident stationary state solutions with energy  $E = 2V_0$ . In each panel, the solid line corresponds to a  $\rho_{i+}$  density, the dashed line to  $\rho_{i-}$ , and the dotted line to  $\rho_i$ . Panels a and d show the incident channel densities corresponding to the  $n = 0$  and  $n = 1$  incident channel stationary states, respectively. Qualitatively, the curves are very similar to the uncoupled case;

however, the intrachannel transmission is decreased. The densities shown in panels b and e are associated, respectively, with the non-incident channels corresponding to panels a and d; these would be formally zero if the problem were uncoupled. Thus, Figure 8b illustrates the  $\rho_{1\pm}$  and  $\rho_1$  densities for the  $n = 0$  incident state, and similarly, Figure 8d shows the  $\rho_{0\pm}$  and  $\rho_0$  densities for the  $n = 1$  incident state. In both cases  $\rho_{i \neq n \pm} \rightarrow 0$  as  $x \rightarrow \mp\infty$ , which is consistent with the required boundary conditions. Panels c and f show the closed channel densities  $\rho_{2\pm}$  and  $\rho_2$  for the  $n = 0$  and  $n = 1$  incident states, respectively. For both cases, the closed channel densities all vanish as  $x \rightarrow \pm\infty$ , so that there is no net transmission or reflectance probability associated with the closed channel. While the magnitude of the closed channel densities are smaller compared to the open channels, they are clearly not negligible. The 2D densities  $\rho_{\pm}$  and  $\rho$  for the degenerate states with  $E = 2V_0$  are shown in Figure 9. The patterns are similar to those for the uncoupled case; however, there are distortions attributed to the nonadiabatic coupling between scattering amplitudes. This is most clearly seen in Figure 9c, where the irregularities in the total  $\rho$  indicate the presence of nontrivial interference effects.

As we have discussed previously, the asymptotic values of the open channel  $\rho_{j+}$  as  $x \rightarrow \infty$  for a given left-incident channel  $n$  are related to the state-to-state transmission probabilities  $P_{n \rightarrow j}$ . Also, the cumulative reaction probability is given by the sum

$$N(E) = \sum_{n,j} P_{n \rightarrow j}(E) \quad (51)$$

In Figure 10 we compare the calculated BRPH reaction probabilities with those obtained using the DVR-ABC method.<sup>12–14,54</sup> See Appendix C for a discussion of the latter calculations. Figure 10a shows the state-to-state probabilities for the coupled Eckart+HO problem on a logarithmic scale as a function of energy. Note that the energy scale here is measured in units of the barrier height  $V_0$  and that the vertical dashed lines indicate the onset of a scattering channel. The BRPH probabilities for individual transitions  $n \rightarrow j$  are given by circles, squares, etc. and the corresponding DVR-ABC results are indicated with the symbol  $\times$  for all transitions. At low energies, only the ground state ( $n = 0$ ) scattering channel is open. The BRPH  $0 \rightarrow 0$  transmission probabilities are represented by filled circles, and these increase more or less steadily across the given energy range; however, there



**Figure 8.** BRPH densities as a function of  $x$  for two degenerate ( $E = 2V_0$ ) left-incident stationary states of the coupled Eckart+HO problem: (solid)  $\rho_{i+}$ , (dashed)  $\rho_{i-}$ , and (dotted)  $\rho_i$ . Panels a and d show the intrachannel densities for the  $n = 0$  and  $n = 1$  incident scattering channels, respectively, while panels b and e show the nonincident channel densities, i.e.,  $i = 1$  and  $i = 0$ , respectively. Panels c and f are the closed channel  $i = 2$  scattering amplitudes contributing to the  $n = 0$  and  $n = 1$  incident channels, respectively.

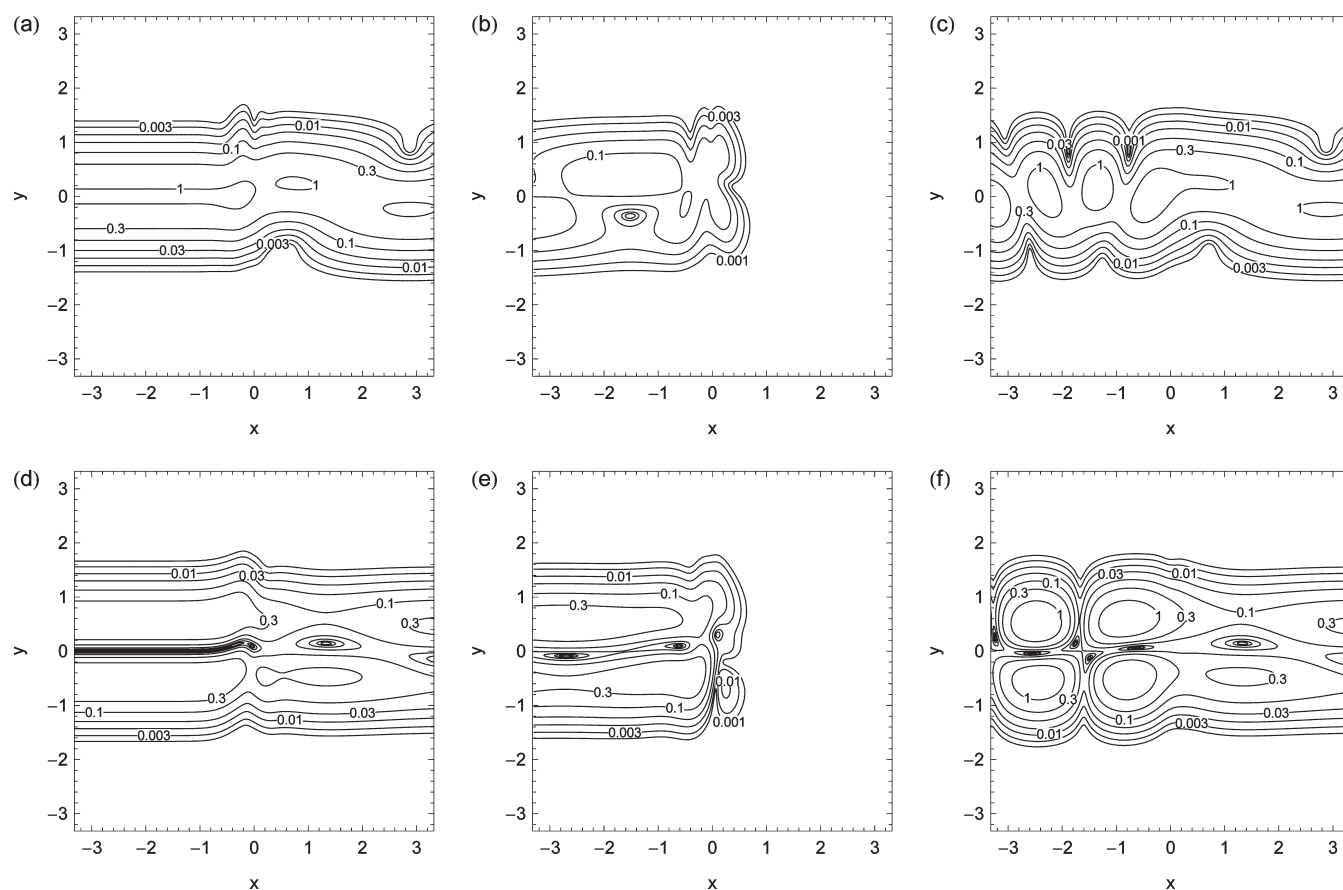
are subtle variations at higher energies, where interstate transitions occur. The  $1 \rightarrow 1$  and  $2 \rightarrow 2$  transmission probabilities are depicted as up-triangles and empty circles, respectively. These are qualitatively similar to the  $0 \rightarrow 0$  transmission, although the variation at higher energies is visibly larger for the  $1 \rightarrow 1$  transmission probability compared to the  $0 \rightarrow 0$  case. This makes sense in light of the fact that the scaling factors in eq 50 are proportional to the incident state quantum number and the fact that the  $n = 0$  incident state can only couple to higher energy scattering states, whereas the  $n = 1$  incident state is coupled to higher and lower energy states. The filled squares, diamonds, and down-triangles in Figure 10a represent the  $0 \rightarrow 1$ ,  $0 \rightarrow 2$ , and  $1 \rightarrow 2$  transmission probabilities, respectively. These begin to increase after the onset of their respective scattering channels. The  $0 \rightarrow 1$  and  $0 \rightarrow 2$  transmission probabilities are observed to decrease with increasing energy, and we speculate the same behavior would also be observed for  $1 \rightarrow 2$  at higher energies.

Figure 10a shows that there is very good qualitative agreement between the BRPH and DVR-ABC calculations. In Figure 10b we plot the fractional (or relative) differences between the BRPH and DVR-ABC transmission probabilities, which are calculated according to the formula

$$\text{fractional difference} = \frac{2|P_{\text{BRPH}} - P_{\text{DVR}}|}{|P_{\text{BRPH}}| + |P_{\text{DVR}}|} \quad (52)$$

Panel b illustrates that the two sets of calculations are also, generally, quantitatively consistent with one other across the given energy range. Note that the largest differences, roughly 1% and 5%, are found for the two energy values just above the channel thresholds at  $E = 1.5V_0$  and  $E = 2.5V_0$ , respectively. For reasons that we have discussed previously, it is difficult to obtain convergence for states with very low kinetic energies, and we attribute the larger differences here to this issue, which affects the accuracy in both the BRPH and DVR-ABC methods. Figure 10c compares the cumulative reaction probability as a function of energy for our BRPH and DVR-ABC calculations. Qualitatively speaking, the agreement is excellent, and Figure 10d examines the fractional differences between the two methods. Generally, the differences are all much less than 0.1% over the given energy range with the exception of just above the channel thresholds, where they are roughly equal to 0.1%; i.e., there are at least two significant figures in common, which is still fairly good. The convergence of the cumulative reaction probability with respect to the grid parameters in our DVR calculations has been monitored closely, and the differences here between the BRPH and DVR-ABC results are of the same order of magnitude as the individual convergence in these numbers with respect to the BRPH and DVR grid parameters; hence, we can conclude that the two methods are quantitatively consistent with one another.

As an independent assessment of the BRPH performance, we also examine the fractional errors in the normalization and energy



**Figure 9.** Contour plots illustrating 2D BRPH densities for the two degenerate stationary states ( $E = 2V_0$ ) of the coupled Eckart+HO problem. For the  $n = 0$  incident states, we have (a)  $\rho_+$ , (b)  $\rho_-$ , and (c)  $\rho$ . Similarly, panels d–f are the corresponding densities for the  $n = 1$  incident state.

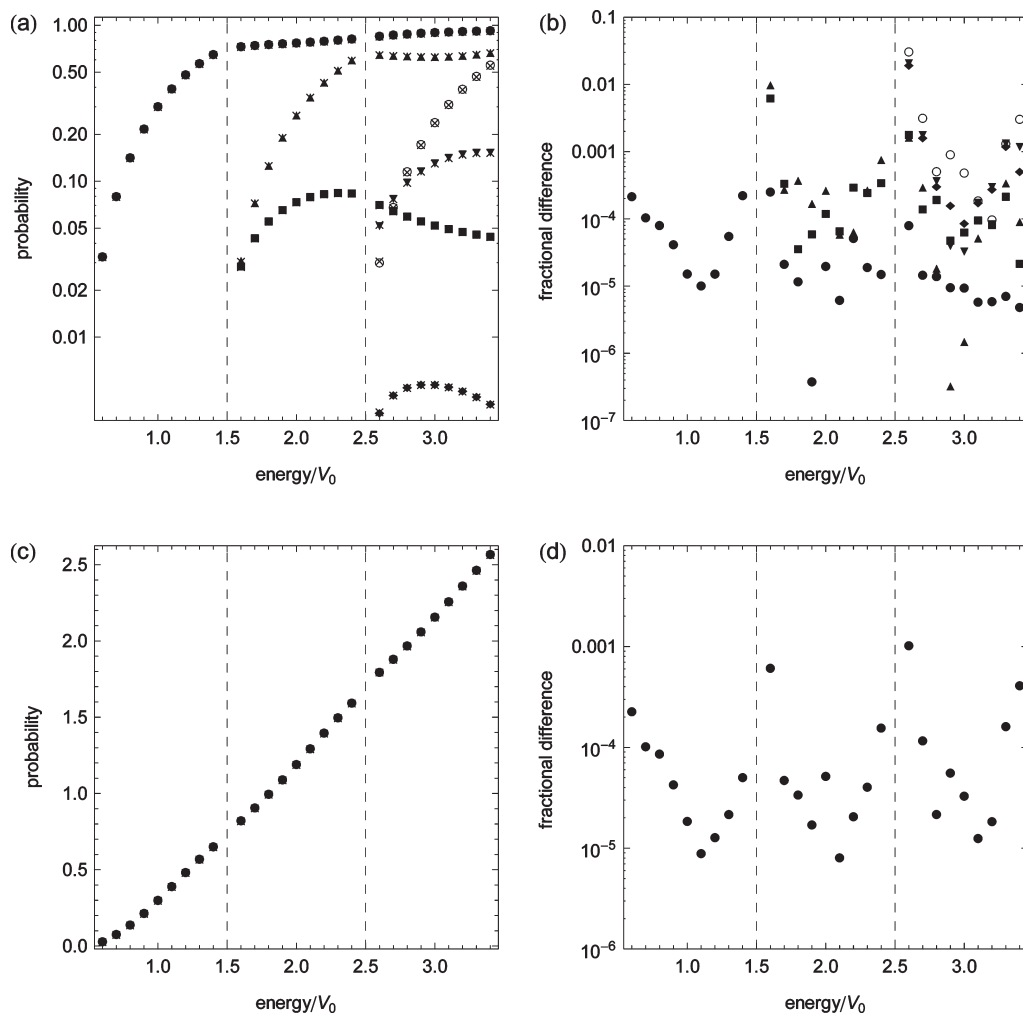
expectation value, which are plotted in panels a and b of Figure 11, respectively. Here each incident channel has its own error measure and there are up to three channels for the highest energy values shown. Generally, the errors for the coupled problem are comparable to those for the uncoupled Eckart+HO potential; however, the energy dependence is qualitatively distinct. Of course, we should expect some differences between the uncoupled and coupled problems; the latter are nontrivial and it is reasonable to expect larger errors. In any case, both fractional error measures are less than 0.1% across the given energy range, so that we may conclude the BRPH solutions for the coupled problem are also quantitatively consistent with both the normalization requirements of the stationary states and the TISE.

Taken together, the differences and errors shown in Figures 10 and 11, respectively, suggest that the BRPH results are likely more accurate than the DVR-ABC calculations at the energies just above the scattering channel thresholds. One concrete indication of this is that the energy expectation errors shown in Figure 11b are on the order of  $10^{-4}$  or better across the energy range; however, the consistency between the BRPH and DVR-ABC calculations is only  $10^{-2}$  in the state-to-state probabilities near the channel threshold and  $10^{-4}$  or better away from threshold (see Figure 10b). As discussed in Appendix C, the discrepancy can be attributed to the fact that we could not fully converge our DVR-ABC calculations with respect to the size and density of the DVR grid at the energies just above the thresholds. This suggests that the BRPH approach may offer a computational advantage over the DVR-ABC method for calculating reaction probabilities at near-threshold energies.

The time complexity for both the DVR-ABC and BRPH approaches scales as  $O(N^3)$ , where  $N$  is the total number of grid points used in the calculations. For DVR-ABC, the value of  $N$  required to achieve a certain level of precision increases as the total energy gets closer to a channel threshold, while for the BRPH, the precision is more or less constant with respect to energy for a fixed  $N$ , assuming a fixed number of scattering channels. This issue also affects the memory requirements of the two methods, and it seems that BRPH has an advantage near a channel threshold. We intend to continue exploring these issues in future work.

## 5. BRIEF SUMMARY AND OUTLOOK

In this work we have described the development of a computational methodology, the BRPH approach, for the calculation of multi-D stationary scattering state wave functions and reaction probabilities in reactive scattering problems. We have presented benchmark results for the simplest class of 2D scattering problems with linear reaction coordinates. The BRPH approach utilizes an adiabatic representation of the system's Hamiltonian to recast the multi-D problem into a set of coupled 1D scattering problems, and we can then exploit the same numerical algorithms used in 1D CPWM applications to calculate the stationary states and state-to-state reaction probabilities of 2D problems. In our numerical applications, we have demonstrated that BRPH calculations are both qualitatively and quantitatively consistent with conventional methods based upon the DVR-ABC approach. Importantly, the BRPH method does not require the use of ABCs



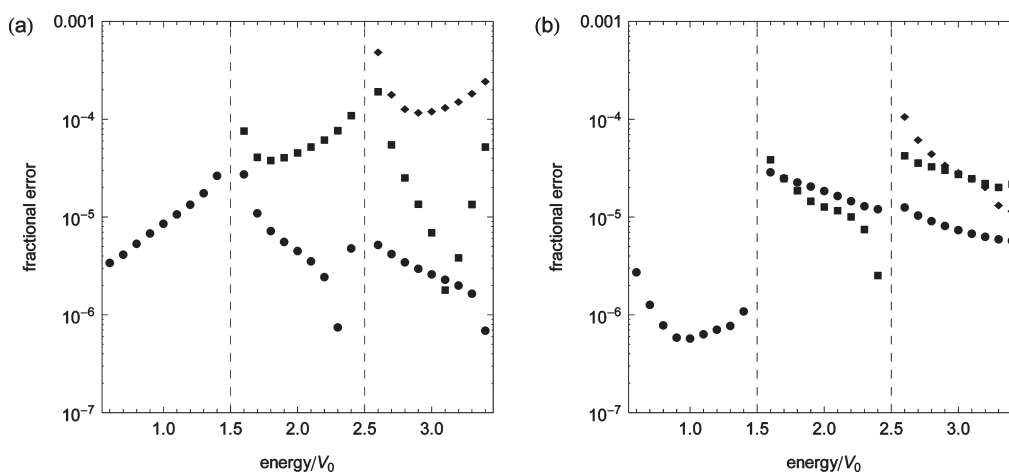
**Figure 10.** Reaction probabilities and fractional differences as a function of energy in units of the barrier height  $V_0 = 0.0018$  for the coupled Eckart+HO problem. Vertical lines represent the onset of a scattering channel. (a) The BRPH state-to-state transmission probabilities are represented as follows: (filled circles)  $0 \rightarrow 0$ , (squares)  $0 \rightarrow 1$ , (diamonds)  $0 \rightarrow 2$ , (up-triangles)  $1 \rightarrow 1$ , (down-triangles)  $1 \rightarrow 2$ , and (empty circles)  $2 \rightarrow 2$ . The corresponding DVR-ABC transmission probabilities are represented with the symbol  $\times$  for all probabilities. (b) Fractional differences between the BRPH and DVR-ABC state-to-state reaction probabilities. The differences for the various probabilities are represented according to the same scheme used in panel a for the BRPH probabilities. (c) Comparison of (circles) BRPH and ( $\times$ ) DVR-ABC cumulative reaction probabilities. (d) Fractional difference between the BRPH and DVR-ABC cumulative reaction probabilities.

so that the range of the computational grid needed in BRPH calculations is much smaller than that for DVR-ABC.

In future work, we plan to extend the theoretical and numerical implementation of the BRPH approach to include asymptotically asymmetric potentials, larger dimensionalities, and curvilinear reaction paths. Asymptotically asymmetric problems do not really present a major challenge, since we can exploit techniques developed in previous work for dealing with such problems.<sup>33</sup> Curvilinear reaction paths and higher dimensionalities, on the other hand, present a much more interesting and challenging class of problems. To address the issue of curvature, in 2D for example, we speculate that the BRPH approach could be applied within a set of orthogonal natural collision coordinates  $s = s(x, y)$  and  $q = q(x, y)$ , where the channel scattering amplitudes  $a_i(s)$  would be defined along some curved reaction path, which is a function of MWC coordinates  $x$  and  $y$ . Likewise, the adiabatic eigenfunctions would vary as  $\phi_i(s, q)$ , where  $q$  is associated with bound motion perpendicular to the reaction coordinate, e.g., a bead oscillating along a curved wire. Generally, the kinetic energy

operator expressed in terms of  $s$  and  $q$  will not have a simple form, so we expect that the adiabatic eigensystem will have to be represented numerically (possibly with 1D DVRs for bound states) in order to obtain exact results. It is also expected that the BRPH integrals will depend upon the metric tensor for the curvilinear coordinate system, and the BRPH equations of motion and the nonadiabatic coupling terms will be correspondingly more complicated.

Another area of interest is how the BRPH approach, compared with conventional methods, will scale with respect to the number of physical dimensions in the scattering problem. It is well-known that the computational effort in DVR calculations, measured by the size of the DVR grid, scales exponentially with the number of degrees of freedom. In BRPH calculations the grid size scales linearly with the number of scattering channels included in the calculation; however, the number of channels also scales exponentially with the number of dimensions. The question will then be, can the BRPH offer a lower pre-exponential factor compared to DVR? The present work indicates that BRPH is advantageous



**Figure 11.** Fractional error estimates for the coupled Eckart+HO problem. Circles, squares, and diamonds represent the  $n = 0, 1,$  and  $2$  incident channel states, respectively. Vertical dashed lines indicate the onset of a scattering channel. (a) Fractional error associated with normalization of the partially state-resolved transmission and reflectance probabilities. (b) Fractional error associated with the energy expectation value.

for energies near threshold, and it will be important to establish whether this transfers to curvilinear problems. In any case, however, the primary utility of the BRPH scheme is that it can readily be used in tandem with various approximate methods for describing the bound degrees of freedom. As dimensionality increases, we must ultimately invoke some level of approximation for the solution of the adiabatic eigenvalue problem. These may include harmonic or anharmonic model representations and perturbation theory. For more accurate work, we could also incorporate quantum Monte Carlo representations, supersymmetric quantum mechanics,<sup>58–61</sup> and massively parallel numerical schemes.<sup>62–64</sup> The most appropriate method will likely be dictated by the details of the problem at hand, so it will be useful to carefully explore the benefits and limitations associated with a variety of different approaches.

## APPENDIX A. DERIVATION OF BRPH EQUATIONS OF MOTION

In this section, we derive the BRPH equations of motion. We begin by constructing the total time derivative of  $a_{i\pm}$ :

$$d_t a_{i\pm} = \partial_t a_{i\pm} \pm \frac{p_i}{m} a'_{i\pm} \quad (53)$$

The partial time derivative is simply given by

$$\partial_t a_{i\pm} = -\frac{i}{\hbar} E a_{i\pm} \quad (54)$$

The convective term contains a spatial derivative of the bipolar components and is evaluated using the expression

$$a'_{i\pm} = \frac{1}{2}(a'_{i+} + a'_{i-}) \pm \frac{1}{2}(a'_{i+} - a'_{i-}) \quad (55)$$

The first term is determined by the FF condition (eq 23). The second term is found by taking the second derivative of the FF condition to give

$$a''_i = \frac{i}{\hbar} p_i (a'_{i+} - a'_{i-}) \quad (56)$$

and then solving for the difference. Substituting these results back into the total time derivative and rearranging leads to the

following equations of motion

$$d_t a_{i\pm} = \frac{i}{\hbar} (E - 2E_i) a_{i\pm} - \frac{i}{\hbar} (E - E_i) a_i + \frac{i}{\hbar} \left[ -\frac{\hbar^2}{2m} a''_i \right] \quad (57)$$

where we note that the last term now contains the second derivative of the scattering amplitude. We appeal to an adiabatic representation of the TISE to evaluate this term. The total stationary state is then expressed as a vector of scattering amplitudes

$$\phi_E = \{a_1, a_2, \dots\}^T \quad (58)$$

and the TISE can be recast as matrix vector product  $\hat{H} \cdot \Phi = E\Phi$ , whose elements are given by

$$\sum_j \hat{H}_{ij} a_j = E a_i \quad (59)$$

The matrix elements of the Hamiltonian are functions of  $x$  and are formally defined by

$$\hat{H}_{ij} = \int \phi_i^*(x, y) \hat{H} \phi_j(x, y) dy \quad (60)$$

To calculate these more explicitly, we first operate with  $\hat{H}$  on the quantity  $\phi_j(x, y) a(x)$ , where the function  $a(x)$  serves as a book-keeping factor to help track the order of derivative operators. Applying the Hamiltonian operator and using the adiabatic TISE (eq 19) we obtain

$$\hat{H}(\phi_j a) = -\frac{\hbar^2}{2m} (\phi_j a'' + 2\phi_j^{(1,0)} a' + \phi_j^{(2,0)} a) + \varepsilon_j \phi_j a \quad (61)$$

where we have used the product rule to expand the second derivative in  $x$ . Next, we multiply both sides of the equation by  $\phi_i^*(x, y)$  and integrate over the  $y$ -coordinate to yield

$$\hat{H}_{ij} a = -\frac{\hbar^2}{2m} [\delta_{ij} a'' + I_{ij}^{(1)} a' + I_{ij}^{(0)} a] + \delta_{ij} \varepsilon_i a \quad (62)$$

where we have used the fact that the adiabatic eigenstates are orthonormal. The BRPH integrals  $I_{ij}^{(1)}$  and  $I_{ij}^{(0)}$  are defined in



eq 29. It is important to note that these quantities are functions of the reaction coordinate that vanish as  $x \rightarrow \pm\infty$ . Similarly,  $\varepsilon_i$  are functions of  $x$  that converge to a set of fixed eigenvalues that are associated with the vibrational states of the asymptotic reactants or products. Removing the book-keeping factor from the last equation gives form for the matrix elements of the Hamiltonian:

$$\hat{H}_{ij} = -\frac{\hbar^2}{2m} \left[ \delta_{ij} \frac{d^2}{dx^2} + I_{ij}^{(1)} \frac{d}{dx} + I_{ij}^{(0)} \right] + \varepsilon_i \delta_{ij} \quad (63)$$

Substituting eq 63 into eq 59 and rearranging the expression for  $a''_i$  gives

$$-\frac{\hbar^2}{2m} a''_i = (E - \varepsilon_i) a_i + \frac{\hbar^2}{2m} \sum_j (I_{ij}^{(1)} a'_j + I_{ij}^{(0)} a_j) \quad (64)$$

Finally, substituting this last expression into eq 57 and rearranging leads to the BRPH equations of motion given in eqs 24 and 28.

## APPENDIX B. ANALYTIC FORMULAS FOR BRPH INTEGRALS

In this section we derive analytic expressions for eq 29 in the special case where the potential function  $V(x,y)$  along the  $y$ -coordinate is described by a displaced and scaled harmonic oscillator. Both the equilibrium position (with respect to the reaction path  $y=0$ ) and the harmonic frequency (force constant) may vary as a function of  $x$ . The potential energy function is given by

$$V(x,y) = V_{\text{scatter}}(x) + \frac{1}{2} m \omega(x)^2 (y - Y(x))^2 \quad (65)$$

where  $V_{\text{scatter}}$  is some 1D barrier potential that depends only on  $x$ ,  $m$  is the reduced mass of the oscillator,  $\omega(x)$  is the harmonic frequency, and  $Y(x)$  is the equilibrium position. It is assumed that  $\omega(x) \rightarrow \omega_0$  and  $Y(x) \rightarrow 0$  as  $x \rightarrow \pm\infty$ , such that the asymptotic potential function is independent of the reaction coordinate. The adiabatic eigenenergies will be given by

$$\varepsilon_i(x) = V_{\text{scatter}}(x) + \hbar \omega(x) \left( i + \frac{1}{2} \right) \quad (66)$$

and asymptotically we have  $\varepsilon_i(x \rightarrow \pm\infty) = E_i = \hbar \omega_0 (i + 1/2)$ .

The adiabatic eigenfunctions for this problem depend parametrically on  $x$  and are given by a set of Gauss–Hermite polynomials

$$\phi_j(x,y) = A_j \sqrt{\beta(x)} H_j(u(x,y)) G(u(x,y)) \quad (67)$$

where  $A_j = (2^j j!)^{-1/2} \pi^{-1/4}$  is a normalization factor,  $\beta(x) = (m\omega(x)/\hbar)^{1/2}$  is a function with units of inverse length,  $H_j(u)$  is the  $j$ th Hermite polynomial, and  $G(u) = \exp(-u^2/2)$  is a Gaussian function. The Gauss–Hermite polynomials are expressed in terms of the dimensionless function  $u(x,y) = \beta(x)(y - Y(x))$ .

Expressions for eq 29 are obtained by taking derivatives of eq 67 with respect to  $x$  and using well-known results for the moments of the harmonic oscillator eigenfunctions:<sup>65</sup>

$$\int \phi_i^* u^0 \phi_n dy = \delta_{i,n} \quad (68a)$$

$$\int \phi_i^* u^1 \phi_n dy = \frac{\sqrt{2}}{2} [\sqrt{i+1} \delta_{i,n-1} + \sqrt{i} \delta_{i,n+1}] \quad (68b)$$

$$\int \phi_i^* u^2 \phi_n dy = \frac{1}{2} [\sqrt{(i+2)(i+1)} \delta_{i,n-2} + (2i+1) \delta_{i,n} + \sqrt{i(i-1)} \delta_{i,n+2}] \quad (68c)$$

$$\int \phi_i^* u^3 \phi_n dy = \frac{\sqrt{2}}{4} [\sqrt{(i+3)(i+2)(i+1)} \delta_{i,n-3} + (3i+3) \sqrt{(i+1)} \delta_{i,n-1} + (3i) \sqrt{i} \delta_{i,n+1} + \sqrt{i(i-1)(i-2)} \delta_{i,n+3}] \quad (68d)$$

$$\int \phi_i^* u^4 \phi_n dy = \frac{1}{4} [\sqrt{(i+4)(i+3)(i+2)(i+1)} \delta_{i,n-4} + (4i+6) \sqrt{(i+2)(i+1)} \delta_{i,n-2} + (6i^2+6i+3) \delta_{i,n} + (4i-2) \sqrt{i(i-1)} \delta_{i,n+2} + \sqrt{i(i-1)(i-2)(i-3)} \delta_{i,n+4}] \quad (68e)$$

Applying the chain rule to eq 67 and using the properties of Gauss–Hermite polynomials one can show that

$$\phi_j^{(1,0)} = \frac{\beta'}{\beta} \left( -u^2 \phi_j + \frac{1}{2} \phi_j + \sqrt{2j} u \phi_{j-1} \right) + \beta Y' (u \phi_j - \sqrt{2j} \phi_{j-1}) \quad (69a)$$

$$\begin{aligned} \phi_j^{(2,0)} = & \frac{\beta''}{\beta} \left( -u^2 \phi_j + \frac{1}{2} \phi_j + \sqrt{2j} u \phi_{j-1} \right) \\ & + \frac{\beta'^2}{\beta^2} \left( u^4 \phi_j - 2u^2 \phi_j - \frac{1}{4} \phi_j - \sqrt{8j} u^3 \phi_{j-1} + \sqrt{2j} u \phi_{j-1} \right. \\ & + \left. \sqrt{4j(j-1)} u^2 \phi_{j-2} \right) + \beta^2 Y'^2 (u^2 \phi_j - \phi_j - \sqrt{8j} u \phi_{j-1} \\ & + \sqrt{4j(j-1)} \phi_{j-2}) + \beta Y'' (u \phi_j - \sqrt{2j} \phi_{j-1}) \\ & + \beta' Y' (-2u^3 \phi_j + 5u \phi_j + \sqrt{32j} u^2 \phi_{j-1} - \sqrt{18j} \phi_{j-1} \\ & - \sqrt{16j(j-1)} u \phi_{j-2}) \end{aligned} \quad (69b)$$

Substituting eq 69 into eq 29 and using eq 68 we find that

$$I_{ij}^{(1)} = \frac{\beta'}{\beta} (\sqrt{(i+2)(i+1)} \delta_{i,j-2} - \sqrt{i(i-1)} \delta_{i,j+2}) - \sqrt{2} \beta Y' (\sqrt{(i+1)} \delta_{i,j-1} - \sqrt{i} \delta_{i,j+1}) \quad (70a)$$

$$\begin{aligned} I_{ij}^{(0)} = & \frac{1}{4} \frac{\beta'^2}{\beta^2} (\sqrt{(i+4)(i+3)(i+2)(i+1)} \delta_{i,j-4} \\ & + \sqrt{i(i-1)(i-2)(i-3)} \delta_{i,j+4}) + \frac{1}{2} \left( \frac{\beta''}{\beta} - \frac{\beta'^2}{\beta^2} \right) \\ & \times (\sqrt{(i+2)(i+1)} \delta_{i,j-2} - \sqrt{i(i-1)} \delta_{i,j+2}) \\ & - \frac{1}{2} \frac{\beta'^2}{\beta^2} (i^2 + i + 1) \delta_{i,j} + \frac{1}{2} \beta^2 Y'^2 (\sqrt{(i+2)(i+1)} \delta_{i,j-2} \\ & + \sqrt{i(i-1)} \delta_{i,j+2}) - \frac{\sqrt{2}}{2} \beta Y'' (\sqrt{(i+1)} \delta_{i,j-1} - \sqrt{i} \delta_{i,j+1}) \\ & - \frac{1}{2} \beta^2 Y'^2 (2i+1) \delta_{i,j} - \frac{\sqrt{2}}{2} \beta' Y' (\sqrt{(i+3)(i+2)(i+1)} \delta_{i,j-3} \\ & + \sqrt{i(i-1)(i-2)} \delta_{i,j+3}) + \frac{\sqrt{2}}{2} \beta' Y' (i \sqrt{(i+1)} \delta_{i,j-1} \\ & + (i+1) \sqrt{i} \delta_{i,j+1}) \end{aligned} \quad (70b)$$

There are two interesting subcases. First, if only the harmonic frequency is scaled along the reaction coordinate, then the integrals simplify to

$$I_{ij}^{(1)} = \frac{\beta'}{\beta} (\sqrt{(i+2)(i+1)}\delta_{i,j-2} - \sqrt{i(i-1)}\delta_{i,j+2}) \quad (71a)$$

$$I_{ij}^{(0)} = \frac{1}{4} \frac{\beta^2}{\beta^2} (\sqrt{(i+4)(i+3)(i+2)(i+1)}\delta_{i,j-4} + \sqrt{i(i-1)(i-2)(i-3)}\delta_{i,j+4}) + \frac{1}{2} \left( \frac{\beta''}{\beta} - \frac{\beta^2}{\beta^2} \right) \times (\sqrt{(i+2)(i+1)}\delta_{i,j-2} - \sqrt{i(i-1)}\delta_{i,j+2}) - \frac{1}{2} \frac{\beta^2}{\beta^2} (i^2 + i + 1)\delta_{i,j} \quad (71b)$$

where we note that there is no coupling between states with different parity. If the oscillator is only displaced along the reaction coordinate, then the integrals reduce to eq 50 of section 4.

### APPENDIX C. DVR-ABC CALCULATIONS

In this section, we describe our implementation of the DVR-ABC approach. For the present work, we utilize the universal DVR developed by Miller and Colbert.<sup>54</sup> In 2D problems, this involves a rectangular set of DVR grid points in  $x$  and  $y$ :

$$q_\alpha = (x_\alpha, y_\alpha) = (\alpha_x \Delta x, \alpha_y \Delta y) \quad (72)$$

where the index  $\alpha = 1, 2, \dots$ , runs over the total number of DVR points and the two subindices  $\alpha_x = 0, \pm 1, \dots$ , and  $\alpha_y = 0, \pm 1, \dots$ , are used to label the  $x$  and  $y$  components of the DVR points within the respective  $x$  and  $y$  subspaces. It is assumed that the problem is more or less centered around the point  $(x, y) = (0, 0)$  and, for simplicity, that the DVR grid is symmetric along  $x$  and along  $y$ .

Miller and Seideman<sup>12,13</sup> analyzed flux-correlation functions to develop simple and very useful DVR expressions for computing the cumulative reaction probability and elements of the  $S$ -matrix; see eqs 76 and 78 below. For a given energy  $E$ , these formulas involve the system's Greens function, which in the DVR-ABC approach is constructed via a matrix inversion

$$G = (EI - H + i(\Gamma_+ + \Gamma_-)/2)^{-1} \quad (73)$$

where  $I$  is an identity matrix over the DVR grid points and  $H$  is the 2D DVR Hamiltonian:

$$H_{\alpha\beta} = T_{\alpha_x, \beta_x}^{(x)} + T_{\alpha_y, \beta_y}^{(y)} + \delta_{\alpha\beta} V(x_\alpha, y_\alpha) \quad (74)$$

The DVR kinetic energy matrix elements for motion along  $y$  are explicitly defined in eq 36, and the expression for  $x$  is similar. The quantities  $\Gamma_\pm$  in eq 73 define a pair of complex absorbing potentials (CAPs) associated with the asymptotic regions of the scattering system. Like the physical potential, the CAPs are approximated by diagonal matrices over the DVR points. In our work, we use a fourth-order polynomial form

$$(\Gamma_\pm)_{\alpha\beta} = \delta_{\alpha\beta} \begin{cases} Z(x_\alpha \mp X_0)^4 / W_x^4 & |x_\alpha| > X_0 \\ 0 & |x_\alpha| \leq X_0 \end{cases} \quad (75)$$

where  $Z$ ,  $W_x$ , and  $X_0$  are the CAP height, width, and onset, respectively. Note that the total spatial extent of the DVR grid along  $x$  is given by  $2(X_0 + W_x)$ .

State-to-state reaction probabilities are calculated according to

$$P_{n \rightarrow j} = |S_{jn}|^2 = \frac{1}{16\hbar^2} (\phi_j^{\text{DVR}})^* \cdot \Gamma_+ \cdot G \cdot \Gamma_- \cdot \phi_n^{\text{DVR}} \quad (76)$$

where  $\phi_j^{\text{DVR}}$  (and  $\phi_n^{\text{DVR}}$ ) is the DVR representation of the asymptotic scattering state associated with channel  $j$  (and  $n$ ). The elements of these vectors are given by

$$(\phi_j^{\text{DVR}})_\alpha = w_\alpha^{1/2} (m/p_j)^{1/2} e^{(i/\hbar)p_j x_\alpha} \phi_j(y_\alpha) \quad (77)$$

where  $w_\alpha = \Delta x \Delta y$  is the 2D DVR weight,  $p_j = (2m(E - E_j))^{1/2}$  is the momentum of a particle with kinetic energy  $E - E_j$ , and  $\phi_j(y)$  is a 1D eigenstate of the asymptotic system with eigenenergy  $E_j$ . The cumulative reaction probability can be calculated by taking the sum in eq 51 or by using the so-called "direct" expression

$$N(E) = \text{Tr}[\Gamma_- \cdot G \cdot \Gamma_+ \cdot G^*] \quad (78)$$

We have used eqs 76 and 78 to generate the DVR-ABC results for comparison with our BRPH calculations.

DVR-ABC calculations involve a number of parameters that affect the accuracy and convergence of the calculated reaction probabilities. The CAP parameters  $Z$ ,  $W_x$ , and  $X_0$  introduced above are important, and ideally, these should be tuned to completely absorb the outgoing flux over as short a length as possible while minimizing artificial reflections. Similarly, the total spatial extent of the DVR grid along  $y$ , which we introduce as the length  $2W_y$ , must be large enough so that  $E_j$  and  $\phi_j(y)$  are numerically well-represented. The DVR grid spacings  $\Delta x$  and  $\Delta y$  together with the spatial widths  $X_0$ ,  $W_x$ , and  $W_y$  determine the number of DVR grid point and the computational effort of the DVR-ABC calculations. The accuracy of the DVR-ABC calculation may be limited by any one or more of these parameters, and it is generally useful to know the extent to which the calculated reaction probabilities are converged.

To establish the precision of our results, we introduce a length scale defined by the following de Broglie wavelength

$$\lambda = \frac{h}{p} = \frac{2\pi\hbar}{\sqrt{2mE}} \quad (79)$$

Note that we are not implying that  $\lambda$  is associated with any specific physical wave function; rather, we are using it as length scale that can be loosely associated with both the bound and scattering components of the true stationary state. Furthermore, we use  $\lambda$  in different ways to set the 2D DVR grid and CAP parameters. For example, we define the grid spacings as

$$\Delta x = \lambda/n_x \quad (80a)$$

$$\Delta y = \lambda/n_y \quad (80b)$$

where  $n_x$  and  $n_y$  are integers that give the number of points per de Broglie wavelength. The total energy  $E$  represents an upper bound for both the kinetic energy of the scattering component and the eigenenergies of the bound states contributing to the true stationary state. Consequently,  $\lambda$  is a lower bound on the physical wavelengths of these components, and the spacings defined by eq 80 are quite useful because the appropriate energy dependence needed to represent the various components is built in. The CAP

parameters are determined by the formulas

$$X_0 = N_0\lambda \quad (81a)$$

$$W_x = N_x\lambda \quad (81b)$$

$$Z = 10E \quad (81c)$$

where  $N_0$  and  $N_x$  are integers that determine the width of the spatial grid along  $x$  in units of  $\lambda$ .

The spatial extent of the DVR grid along  $y$  should cover, at least, the space between the classical turning points defined by the minimum and maximum roots of  $V(x) = E$ . In fact, the grid should be extended beyond these points to account for the tails of the bound state wave functions associated with tunneling. We have found that  $\lambda$  is not very well-suited for defining  $W_y$ ;  $\lambda$  is far too large at low energies and becomes much too small as  $E$  increases. Instead, we use the following formula

$$W_y = y_{cl} + N_y A^{-1/3} \quad (82a)$$

where  $N_y$  is an integer similar to  $N_x$  and  $y_{cl}$  is the maximum classical turning point along  $y$ . Note again, that we are assuming the problem is symmetric about  $y$ ; if it were not, then one would have an analogous expression for the minimum classical turning point. The parameter  $A$  comes from the approximate representation of a 1D wave function in the vicinity of the turning point:

$$\phi(y) \approx \text{Ai}\left(\frac{B + A(y - y_{cl})}{A^{2/3}}\right) \quad (83)$$

where  $\text{Ai}$  is an Airy function, and the parameters  $A$  and  $B$  are given by

$$A = 2mV_1/\hbar^2 \quad (84a)$$

$$B = 2m(V_0 - E)/\hbar^2 \quad (84b)$$

Here,  $V_0$  and  $V_1$  are the values of the potential and its first derivative, respectively, evaluated at  $y_{cl}$ . The wave function in eq 83 is the well-known solution for a particle evolving under the influence of a linear potential,  $V(y) = V_0 + V_1(y - y_{cl})$ . The length scale  $A^{-1/3}$  is inversely proportional to  $V_1$ ; therefore, when the force is weak,  $A^{-1/3}$  will be large, and the DVR grid will extend further into the forbidden region.

We have established the integers  $n_x$ ,  $n_y$ ,  $N_x$ ,  $N_y$ , and  $N_0$  as the tunable parameters for controlling the accuracy of our DVR-ABC calculations, and we have thoroughly investigated the convergence of our DVR-ABC calculations for the Eckart+HO problem. The following protocol has been used to generate converged reaction probabilities. One begins by calculating  $N(E)$  using the parameters  $n_x = 3$ ,  $n_y = 3$ ,  $N_x = 1$ ,  $N_y = 1$ , and  $N_0 = 0$ . Typically, this gives a very inaccurate result for  $N(E)$ ; nevertheless, this value is saved as a reference. Next, one performs a set of five distinct calculations for  $N(E)$ , where each integer parameter is increased by 1. The relative difference between the original  $N(E)$  and the five updated values is used to estimate the convergence with respect to these parameters. For each parameter a decision is made as to whether to accept the new parameter or continue with the original value. Once the updated parameters are selected, the reference calculation is repeated and the new  $N(E)$  is saved. The parameters are sequentially updated again and the convergence is retested. This process is repeated iteratively until the desired precision has been reached

or some maximum number of grid points, determined by the computer system's memory, has been exceeded.

In our work, we have set the convergence tolerance to  $10^{-5}$  and the convergence iterations will stop once the reference calculations exceed 4000 DVR points. For the results shown in section 4, the cumulative reaction probabilities converged to a relative difference of  $10^{-5}$  with respect to the integer parameters. The energy points just above the onset of the scattering channels are the exception to this, and these calculations were stopped before this level of convergence could be achieved.

## AUTHOR INFORMATION

### Corresponding Author

\*E-mail: jeremy.maddox@wku.edu.

## ACKNOWLEDGMENT

J.B.M. gratefully acknowledges financial support from the Applied Research and Technology Program (ARTP), Western Kentucky University. B.P. gratefully acknowledges support from the Robert A. Welch Foundation (D-1523) and the National Science Foundation (CHE-0741321).

## REFERENCES

- Maddox, J. B.; Poirier, B. In *Quantum Trajectories*; Hughes, K., Parlant, G., Eds.; Collaborative Computational Project on Molecular Quantum Dynamics (CCP6): Daresbury, UK, 2011; pp 9–11.
- Cohen-Tannoudji, C.; Diu, B.; Laloë, F. *Quantum Mechanics*; John Wiley: New York, 1977; pp 32–34.
- Tannor, D. J. *Introduction to Quantum Mechanics: A Time-Dependent Perspective*; University Science Books: Sausalito, CA, 2007; pp xv–xvi.
- Steinfeld, J. I.; Francisco, J. S.; Hase, W. L. *Chemical Kinetics and Dynamics*; Prentice-Hall: Englewood Cliffs, NJ, 1989; pp 268–272.
- Levine, R. D.; Bernstein, R. B. *Molecular Reaction Dynamics and Chemical Reactivity*; Oxford University Press: New York, 1987; pp 276–289.
- Lill, J. V.; Parker, G. A.; Light, J. C. *Chem. Phys. Lett.* **1982**, *89*, 483–489.
- Lill, J. V.; Parker, G. A.; Light, J. C. *J. Chem. Phys.* **1986**, *85*, 900–910.
- Light, J. C.; Hamilton, I. P.; Lill, J. V. *J. Chem. Phys.* **1985**, *82*, 1400–1409.
- Whitnell, R. M.; Light, J. C. *J. Chem. Phys.* **1987**, *86*, 2007–2019.
- Engquist, B.; Majda, A. *Proc. Natl. Acad. Sci. U. S. A.* **1977**, *74*, 1765–1766.
- Neuhasuer, D.; Baer, M. J. *Chem. Phys.* **1989**, *90*, 4351–4355.
- Seideman, T.; Miller, W. H. *J. Chem. Phys.* **1992**, *96*, 4412–4422.
- Seideman, T.; Miller, W. H. *J. Chem. Phys.* **1992**, *97*, 2499–2514.
- Thompson, W. H.; Miller, W. H. *Chem. Phys. Lett.* **1993**, *206*, 123–129.
- Auerbach, S. M.; Miller, W. H. *J. Chem. Phys.* **1994**, *100*, 1103–1112.
- Wang, H.; Thompson, W. H.; Miller, W. H. *J. Chem. Phys.* **1997**, *107*, 7194–7201.
- Poirier, B.; Miller, W. H. *Chem. Phys. Lett.* **1997**, *265*, 77–83.
- Edlund, A.; Peskin, U. *Int. J. Quantum Chem.* **1998**, *69*, 167–173.
- Gezelter, J. D.; Miller, W. H. *J. Chem. Phys.* **1995**, *103*, 7868–7876.
- Seideman, T. *J. Chem. Phys.* **1993**, *98*, 1989–1998.

- (21) Thompson, W. H.; Miller, W. H. *J. Chem. Phys.* **1994**, *101*, 8620–8627.
- (22) Peskin, U.; Miller, W. H. *J. Chem. Phys.* **1995**, *102*, 4084–4092.
- (23) Thompson, W. H.; Karlsson, H. O.; Miller, W. H. *J. Chem. Phys.* **1996**, *105*, 5387–5396.
- (24) Gezelter, J. D.; Miller, W. H. *J. Chem. Phys.* **1996**, *104*, 3546–3554.
- (25) Qi, J.; Bowman, J. M. *J. Chem. Phys.* **1996**, *104*, 7545–7553.
- (26) Saalfrank, P.; Miller, W. H. *Surf. Sci.* **1994**, *303*, 206–230.
- (27) Poirier, B.; Carrington, J.; Tucker J. *J. Chem. Phys.* **2003**, *118*, 17–28.
- (28) Poirier, B.; Carrington, J.; Tucker J. *J. Chem. Phys.* **2003**, *119*, 77–89.
- (29) Muga, J. G.; Palao, J. P.; Navarro, B.; Egusquiza, I. L. *Phys. Rep.* **2004**, *395*, 357–426.
- (30) Poirier, B. *J. Chem. Phys.* **2004**, *121*, 4501–4515.
- (31) Trahan, C.; Poirier, B. *J. Chem. Phys.* **2006**, *124*, 034115/1–034115/18.
- (32) Trahan, C.; Poirier, B. *J. Chem. Phys.* **2006**, *124*, 034116/1–034116/14.
- (33) Poirier, B. *J. Theor. Comput. Chem.* **2007**, *6*, 99–125.
- (34) Poirier, B.; Parlant, G. *J. Phys. Chem. A* **2007**, *111*, 10400–10408.
- (35) Poirier, B. *J. Chem. Phys.* **2008**, *129*, 084103/1–084103/18.
- (36) Poirier, B. *J. Chem. Phys.* **2008**, *128*, 164115/1–164115/15.
- (37) Park, K.; Poirier, B.; Parlant, G. *J. Chem. Phys.* **2008**, *129*, 194112/1–194112/16.
- (38) Park, K.; Poirier, B. *J. Theor. Comput. Chem.* **2010**, *9*, 711–734.
- (39) Wyatt, R. E. *Quantum Dynamics with Trajectories: Introduction to Quantum Hydrodynamics*; Springer: New York, 2005; pp 11–21.
- (40) Miller, W. H.; Handy, N. C.; Adams, J. E. *J. Chem. Phys.* **1980**, *72*, 99–112.
- (41) Fast, P. L.; Truhlar, D. G. *J. Chem. Phys.* **1998**, *109*, 3721–3729.
- (42) Gonzalez, J.; Gimenez, X.; Bofill, J. M. *Theor. Chem. Acc.* **2004**, *112*, 75–83.
- (43) Fernandez-Ramos, A.; Ellingson, B. A.; Garrett, B. C.; Truhlar, D. G. *Rev. Comp. Chem.* **2007**, *23*, 125–232.
- (44) Gonzalez, J.; Gimenez, X.; Bofill, J. M. *J. Chem. Phys.* **2009**, *131*, 054108/1–054108/16.
- (45) Fang, J.-Y.; Hammes-Schiffer, S. *J. Chem. Phys.* **1998**, *108*, 7085–7099.
- (46) Fang, J.-Y.; Hammes-Schiffer, S. *J. Chem. Phys.* **1998**, *109*, 7051–7063.
- (47) Hinkle, C. E.; McCoy, A. B. *J. Phys. Chem. Lett.* **2010**, *1*, 562–567.
- (48) Presently, all CPWMs require the scattering potential to converge asymptotically. For example, asymptotically periodic potentials are beyond the scope of existing CPWMs; however, we see no reason why the CPWM approach could not be extended to such problems, for example, by using Bloch-type bipolar components.
- (49) The form of the effective potential does make a difference in numerical applications and should ideally be a smooth, slowly varying function connecting the asymptotic values of  $V(x)$ .
- (50) Fröman, N.; Fröman, P. O. *JWKB Approximation*; North-Holland: New York, 1965.
- (51) Ahmed, Z. *Phys. Rev. A* **1993**, *47*, 4761–4767.
- (52) Abramowitz, M.; Stegun, I. A. *Handbook of Mathematical Functions with Formulas, Graphs, and Mathematical Tables*; Dover Publications: New York, 1965; pp 877–898.
- (53) Please note the distinction between the index  $i$  and the imaginary unit  $i = (-1)^{1/2}$ .
- (54) Colbert, D. T.; Miller, W. H. *J. Chem. Phys.* **1992**, *96*, 1982–1991.
- (55) Press, W. H.; Teukolsky, S. A.; Vetterling, W. T.; Flannery, B. P. *Numerical Recipes in Fortran: The Art of Scientific Computing*, 3rd ed.; Cambridge University Press: New York, 2007; pp 155–162; 907–910
- (56) Datta, S. *Quantum Transport: Atom to Transistor*; Cambridge University Press: New York, 2005; pp 17–18.
- (57) Pauler, D. K.; Kendrick, B. K. *J. Chem. Phys.* **2004**, *120*, 603–611.
- (58) Bittner, E. R.; Maddox, J. B.; Kouri, D. J. *J. Phys. Chem. A* **2009**, *113*, 15276–15280.
- (59) Kouri, D. J.; Markovich, T.; Maxwell, N.; Bittner, E. R. *J. Phys. Chem. A* **2009**, *113*, 15257–15264.
- (60) Kouri, D. J.; Maji, K.; Markovich, T.; Bittner, E. R. *J. Phys. Chem. A* **2010**, *114*, 8202–8216.
- (61) Kouri, D. J.; Maji, K.; Markovich, T.; Bittner, E. R. *J. Phys. Chem. A* **2011**, *115*, 950.
- (62) Chen, W.; Poirier, B. *J. Theor. Comput. Chem.* **2010**, *9*, 825–846.
- (63) Chen, W.; Poirier, B. *J. Theor. Comput. Chem.* **2010**, *9*, 435–469.
- (64) Yang, B.; Chen, W.; Poirier, B. *J. Chem. Phys.* **2011**, *135*, 094306/1–094306/17.
- (65) Wilson, E. B.; Decius, J. C.; Cross, P. C. *Molecular Vibrations: The Theory of Infrared and Raman Vibrational Spectra*; McGraw-Hill Book Company, Inc.: New York, 1955; pp 289–291.

# On the Intramolecular Hydrogen Bond in Solution: Car–Parrinello and Path Integral Molecular Dynamics Perspective

Przemyslaw Dopieralski,<sup>\*,†</sup> Charles L. Perrin,<sup>\*,‡</sup> and Zdzislaw Latajka<sup>†</sup><sup>†</sup>Faculty of Chemistry, University of Wrocław, Joliot–Curie 14, 50-383 Wrocław, Poland<sup>‡</sup>Department of Chemistry and Biochemistry, University of California at San Diego, La Jolla, California 92093-0358, United States Supporting Information

**ABSTRACT:** The issue of the symmetry of short, low-barrier hydrogen bonds in solution is addressed here with advanced ab initio simulations of a hydrogen maleate anion in different environments, starting with the isolated anion, going through two crystal structures (sodium and potassium salts), then to an aqueous solution, and finally in the presence of counterions. By Car–Parrinello and path integral molecular dynamics simulations, it is demonstrated that the position of the proton in the intramolecular hydrogen bond of an aqueous hydrogen maleate anion is entirely related to the solvation pattern around the oxygen atoms of the intramolecular hydrogen bond. In particular, this anion has an asymmetric hydrogen bond, with the proton always located on the oxygen atom that is less solvated, owing to the instantaneous solvation environment. Simulations of water solutions of hydrogen maleate ion with two different counterions,  $K^+$  and  $Na^+$ , surprisingly show that the intramolecular hydrogen-bond potential in the case of the  $Na^+$  salt is always asymmetric, regardless of the hydrogen bonds to water, whereas for the  $K^+$  salt, the potential for H motion depends on the location of the  $K^+$ . It is proposed that repulsion by the larger and more hydrated  $K^+$  is weaker than that by  $Na^+$  and competitive with solvation by water.

## 1. INTRODUCTION

Today no one will deny that hydrogen bonds (H-bonds) are a key feature of molecular structure and reactivity.<sup>1–11</sup> Despite the fact that an enormous amount of experimental and theoretical work has been devoted to unveiling their peculiar character in the past decade, H-bonds are still a rich source of new and old challenging, unsolved problems.

One of those problems is the symmetry of short, low-barrier H-bonds in solution.<sup>12</sup> For many years it was hoped that the crystal structure would describe a molecule in solution. Yet crystal structures of the same H-bond can differ. For example, the potassium salt of hydrogen maleate (Hmaleate) shows a centered hydrogen,<sup>13</sup> whereas the hydrogen is asymmetrically positioned in the corresponding sodium salt.<sup>14</sup> Nevertheless, NMR studies showed that those H-bonds are invariably asymmetric in aqueous solution.<sup>15,16</sup> It was initially proposed that the asymmetry is a consequence of the polarity of water, which stabilizes the localized negative charge of  $O-H\cdots O^-$  or  $^-O\cdots H-O$  more than the delocalized one of  $(O\cdots H\cdots O)^-$ .<sup>17</sup> This rationale was supported by computer simulations.<sup>18</sup> However, further NMR studies in nonpolar organic solvents continued to show asymmetric H-bonds,<sup>19</sup> even in the NHN H-bonds of protonated 1,8-bis(dimethylamino)naphthalenes,<sup>20</sup> in non-ionic species,<sup>21</sup> in a zwitterion,<sup>22</sup> in the intermolecular H-bond of pyridine–dichloroacetic acid complexes,<sup>23</sup> and in the “strongest” of H-bonds.<sup>24</sup> Therefore it was concluded that although the environments around the two carboxyl groups can be identical in a crystal, a solution is disorganized, with one of the carboxyls instantaneously solvated better than the other,<sup>25</sup> leading to the presence of equilibrating solvatomers (isomers, stereoisomers, or tautomers that differ in solvation).<sup>22</sup> QM/MM calculations (AM1-SRP/AMBER) on hydrogen phthalate anion in solution

support this interpretation.<sup>26</sup> Moreover, the absence of symmetric H-bonds in solution suggests that they have no special stabilization.<sup>27</sup>

An early theoretical study on the simplest case, the isolated Hmaleate anion, found the potential for the proton transfer to be a double minimum ( $C_s$  symmetry), with a barrier height of 1.4 kcal/mol and with a structure of  $C_{2v}$  symmetry being a transition state.<sup>28</sup> Later studies confirmed the asymmetric structure as more stable, although the barrier varies from 0.1 kcal/mol up to a few kcal/mol, depending on the method used.<sup>29–33</sup> Plane-wave DFT calculations found a broad, flat potential energy surface for the unit cell but returned to a shallow double-well surface when the crystal packing forces were removed.<sup>34,35</sup> Thus, questions dealing with the symmetry of Hmaleate anion are still discussed.

Here we use both ab initio Car–Parrinello molecular dynamics (CPMD) simulations and fully quantum mechanical path integral molecular dynamics (PIMD) on aqueous Hmaleate ion at 298 K to provide a theoretical understanding of how the solvent and the counterions  $K^+$  and  $Na^+$  modify the symmetry of the intramolecular H-bond. Based on a new procedure, which distinguishes separate free-energy profiles for those periods during the simulation when oxygen atoms of the intramolecular H-bond are solvated to different or similar extents, we can discover the determinants of the H-bond symmetry and confirm some intriguing experimental findings by Perrin et al.<sup>19</sup>

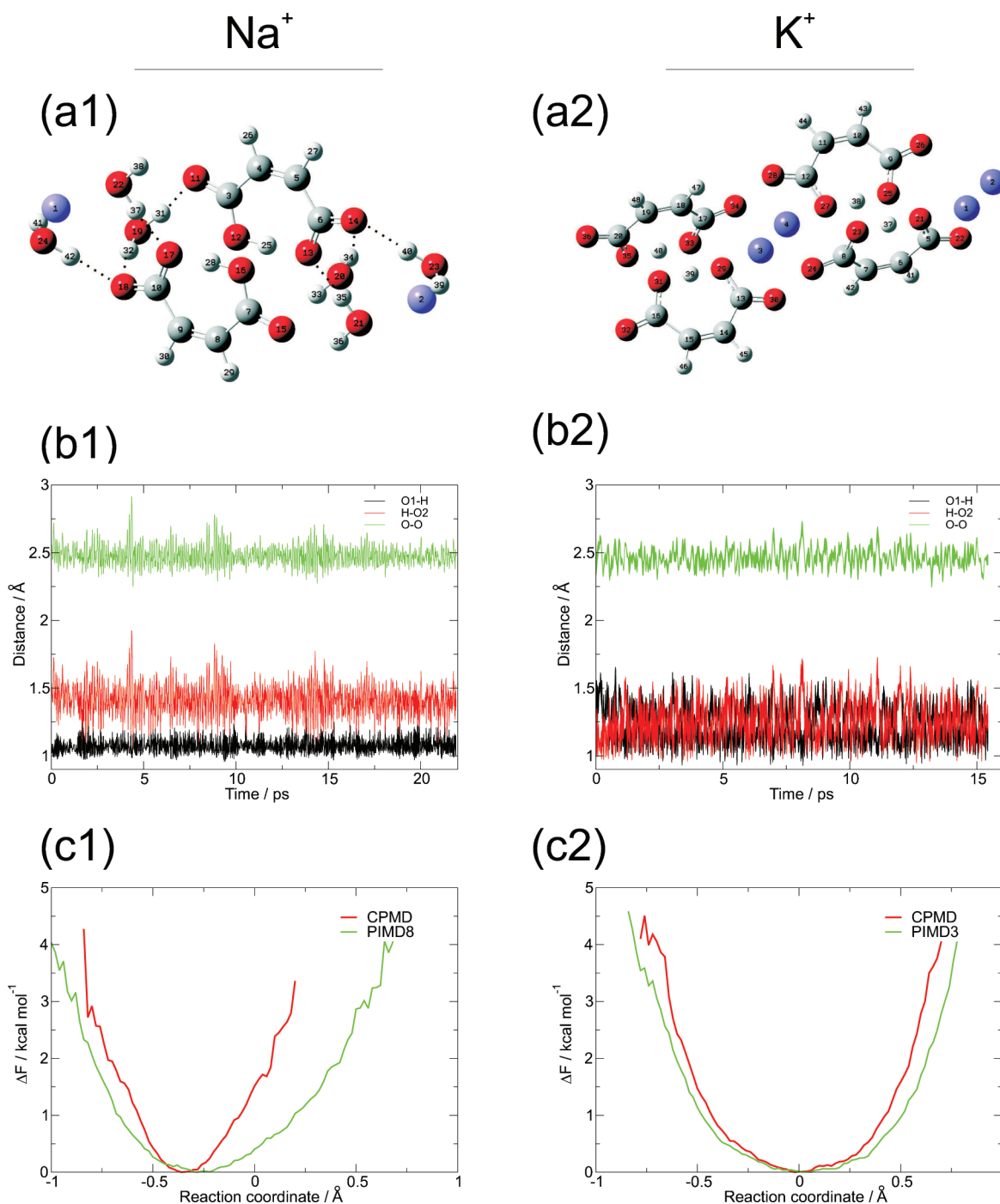
## 2. RESULTS AND DISCUSSION

### 2.1. Isolated Hydrogen Maleate Ion and Its Potassium and Sodium Salts.

Before addressing the main theme, one must

Received: August 19, 2011

Published: October 10, 2011



**Figure 1.** Two crystal structures of Hmaleate salts from CPMD and PIMD simulations: sodium (left panel: a1, b1, and c1) and potassium (right panel: a2, b2, and c2). (a1,b1) Representative snapshot of the crystal unit cell. (b1,b2) Time evolution of the distances involved in the intramolecular H-bond. (c1,c2) Free-energy profiles for proton transfer within the H-bond.

consider the isolated Hmaleate ion and the crystal structures of its salts. To confirm that the isolated ion represents a situation with a truly symmetric-centered single, well potential or, equivalently, a double well with a low barrier, CPMD and PIMD room temperature simulations have been performed. Additionally we have performed CPMD and PIMD simulations for sodium and potassium Hmaleate crystals.

The results for the isolated ion show the O···O distance to be 2.457 Å, and both O···H distances are the same, 1.235 Å. Calculated

structures of sodium and potassium Hmaleate crystals, shown in Figure 1a1 and a2, reproduce the experimental structures. In particular, the sodium salt shows a low symmetry and an asymmetric H-bond, whereas the potassium salt is a highly symmetric crystal, with a centered proton. However, it must be noted that although the potassium crystal is truly a high-symmetry one, in the sodium crystal the presence of water molecules reduces the symmetry and is responsible for different environments of the two oxygen atoms involved in the intramolecular H-bond of Hmaleate.

**Table 1.** Time-Averaged Intramolecular H-bond Distances from Calculations on Sodium and Potassium Hmaleate Salts

bond	NaHMal			KHMal	
	CPMD	PIMD8	expt <sup>14</sup>	CPMD	expt <sup>36</sup>
O–H	1.08	1.13	1.079	1.22	–
H···O	1.41	1.37	1.369	1.23	–
O···O	2.48	2.47	2.445	2.44	2.434

Time-averaged H-bond distances for these two simulations have been collected in Table 1, and instantaneous distances are displayed in Figure 1b1 and b2.

The graph in Figure 1b1 shows the results from simulations on the crystal structure of sodium Hmaleate. As expected, the O–O distance shows little variation, and this behavior is maintained in all the subsequent simulations. In this crystal one O–H distance remains short, near 1.0 Å, while the other is longer, near 1.4 Å, but with greater variability. This leads to a highly asymmetric single-well potential for proton motion within the H-bond, as shown in Figure 1c1. In this crystal the proton is not able to jump from one oxygen atom to the other within the intramolecular H-bond because the two oxygens are in different molecular environments. By inspection of the structure, it can be seen that one of the oxygens is closer to its nearest sodium than the other oxygen is to its closest sodium. The difference is approximately 1 Å. Additionally, the oxygen that has the closer sodium ion participates in a H-bond to the H of a water molecule, whereas the other oxygen does not coordinate any water molecules. Thus, because of crystal packing, which produces an inequivalence of the two oxygens, an asymmetric single-well potential for proton motion within the H-bond is induced, as presented in Figure 1c1, with the proton located on the oxygen that is farther from the sodium ion and that lacks a water of solvation. Inclusion of quantum effects by the PIMD treatment does not qualitatively change the results, and the potential remains an asymmetric single well but with larger fluctuations of the proton motion, as included in Figure 1c1.

For the potassium Hmaleate crystal, the situation is different. The graph in Figure 1b2 shows the results of the simulations. In this crystal both O–H distances vary rapidly but stay near  $1.2 \pm 0.1$  Å. This leads to a broad and symmetric single-well potential for proton motion within the H-bond, as shown in Figure 1c2. Both the X-ray crystal structure and our simulations show that a chain,  $\cdots K^+ \cdots O1-H-O2 \cdots K^+ \cdots$ , is observed, with the potassium ions placed symmetrically, at equal distances, with respect to both oxygens. Because each oxygen atom has exactly the same neighborhood, the proton is located most of the time in the center of the H-bond, as shown in Figure 1b2, and the potential is a symmetric single-well one, as depicted in Figure 1c2. The result from PIMD simulations is also shown in Figure 1c2. The inclusion of quantum effects does not change the picture, as expected.

**2.2. Hydrogen Maleate Ion in Water.** In Figure 2 the results from CPMD and PIMD simulations on aqueous Hmaleate anion are collected. The graph in Figure 2b1 shows that according to CPMD, the proton jumps from one oxygen to the other, with an average residence time of around 1 ps but sometimes longer. The proton is found more often on O1. These proton positions correspond to an asymmetric double-well potential (Figure 2c1) with a barrier height around 1 kcal/mol. This is likely to be a consequence of the instantaneous solvation environment, which favors one solvatomer over another.

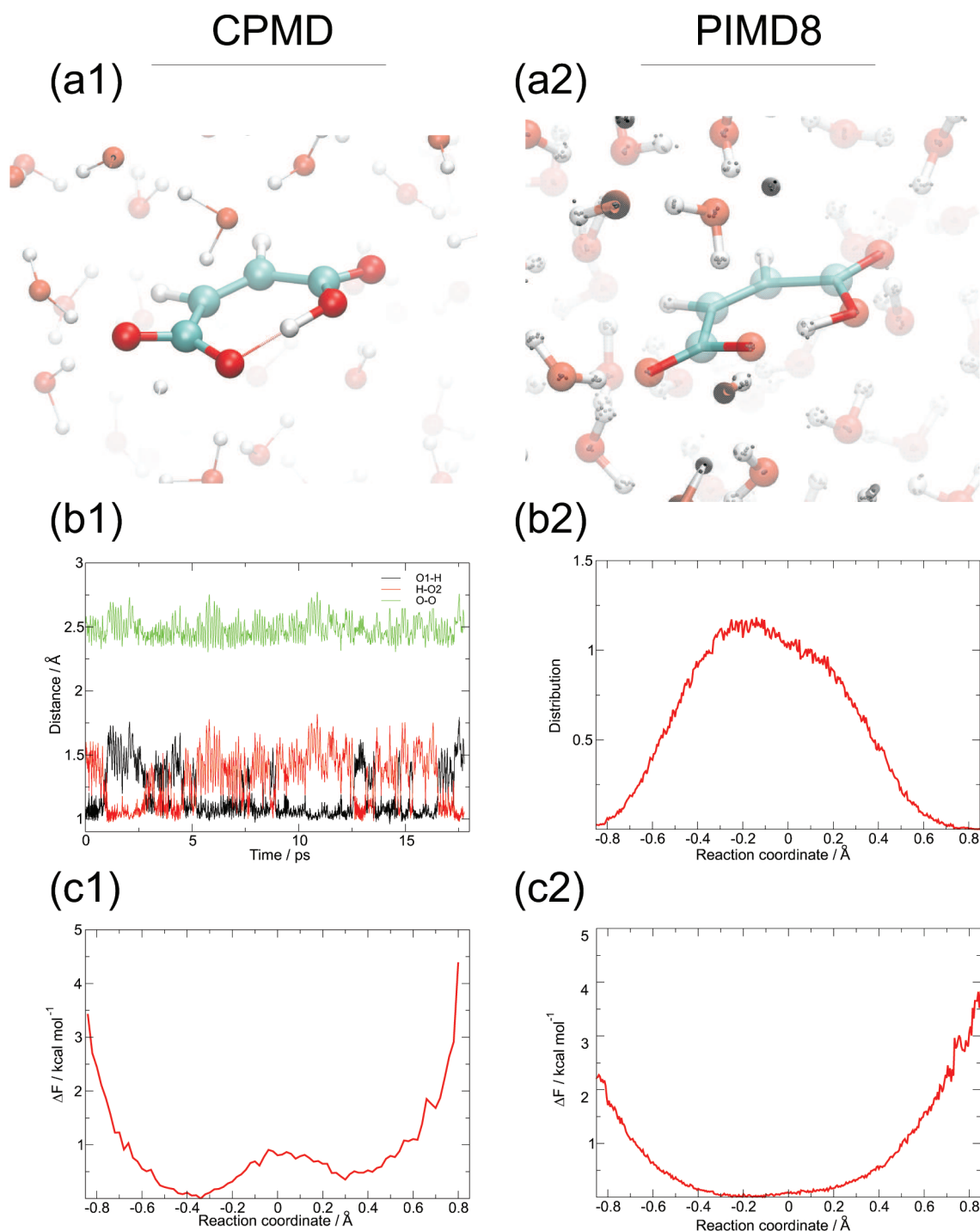
We have probed solvation in more detail by summing the number of H-bonds between the hydrogen atoms of the surrounding waters and each of the two oxygen atoms in the intramolecular H-bond of the Hmaleate anion. The distance dependence for deciding to include a H-bond was investigated and was finally taken as 2.5 Å. Based on this routine, we find that oxygen atom O1 has been less solvated than O2 (with fewer H-bonds from water molecules) for approximately 48.2% of the simulation time, O2 has been less solvated than O1 for 15.1% of the time, and both oxygen atoms have been solvated similarly during approximately 26.7% of the time.

Moreover, the position of the proton in the intramolecular H-bond is correlated with the relative solvation of the two Hmaleate oxygen atoms. The remarkable result is that the proton is always located on the oxygen that is less well solvated.

This result, and other considerations presented here, bring to mind a picture of donor and acceptor solvation theory similar to the Marcus theory of electron transfer.<sup>37</sup> Proton or electron movement (but not the proton or electron transfer itself) is much faster than reorganization of the solvent, assumed in Marcus theory. However, we cannot assume that the donor and acceptor moieties are only weakly coupled—this assumption of Marcus would not be true in Hmaleate anion. Marcus theory depicts electron transfer as made possible by prearrangement of fluctuating solvent molecules. When, by chance, such a correct arrangement happens, the transfer can take place. Our observation is very similar. The proton is always located on the oxygen that is less well solvated; using the Marcus idea, we can reverse this statement and say that the less well solvated oxygen atom is preferred to possess the proton, and the better solvated oxygen is preferred as the bearer of more negative charge, making it a future acceptor. When by chance this solvation pattern is reversed by solvent fluctuations, the proton can jump. These ideas were exploited in studies by Borgis and Hynes<sup>38</sup> and by Mavri et al.<sup>39</sup> However, the time frame of the CPMD simulations and the nature of delocalized plane-wave basis set do not allow us to use directly the methodologies described there.

These results suggest that the asymmetric double-well potential and the apparent barrier of 1 kcal/mol in Figure 2c1 are a consequence of averaging over instantaneous solvation environments that sometimes favor one solvatomer and then the other. Indeed, if the averaging is restricted to shorter time intervals, a different picture emerges. The separate free-energy profiles for the 48.2% of the simulation time when O1 was less solvated than O2, for the 15.1% of the time when O2 was less solvated than O1, and for the 26.7% of time when both oxygen atoms were solvated similarly are shown in Figure 3a1. For the first two cases, where the two oxygens were solvated to different extents, the free-energy profile becomes a single-well potential, with its minimum at the oxygen that is less solvated. There is no longer a minimum at the other, better-solvated oxygen. For the third case, where both oxygen atoms were solvated similarly, the free-energy profile remains a double-well potential, as in the earlier Figure 2c1, with nearly the same difference in well depths, but the barrier has become indistinct. Thus these profiles in Figure 3a1 show how the instantaneous potential responds to the solvation environment. An important conclusion is that the apparent barrier in the CPMD simulations is an artifact of the long-term averaging over different solvation environments.

One must be aware of the fact that the free energies presented in Figures 3 and 5 below are not in a sense the most rigorously defined ones. Definition of a free energy profile requires a



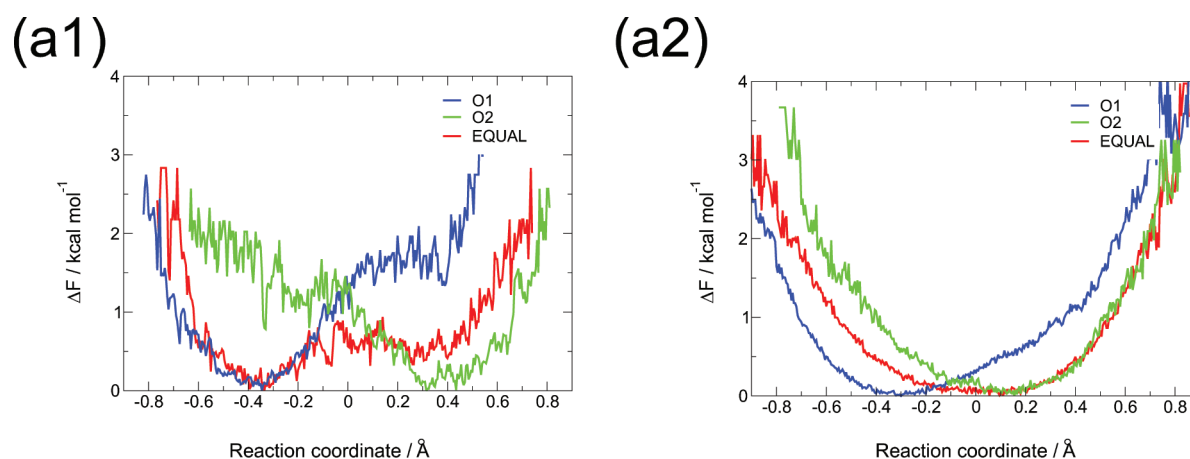
**Figure 2.** Hmaleate ion in water from CPMD (left column: a1, b1, and c1) and from PIMD (right column: a2, b2, and c2). (a1,a2) Representative snapshot from simulations, for PIMD the replicas of each atom as small gray spheres are marked. (b1) Time evolution of the distances involved in the intramolecular H-bond. (b2) Distribution function for the reaction coordinate (c1,c2) Free energy profiles for proton transfer within the H-bond.

complete phase-space average over all orthogonal coordinates, so our averaging cannot yield a proper free energy. Nevertheless our projections/averaging of “free energies” are still reasonable approximations and are able to capture the correlation between proton localization in the H-bond and the local solvation pattern. However, it is not possible to analyze energetics based on these figures. A more rigorous way to define the free energy is to introduce a second reaction coordinate (as was done by Tuckerman

for  $\text{OH}^-$  solvation),<sup>40</sup> which in our case is the difference in hydration number of oxygen atoms involved in the H-bond. For figures and more details see Supporting Information.

According to PIMD simulation, the situation is slightly different. Figure 2b2 shows the distribution of the reaction coordinate, and Figure 2b3 shows the corresponding free-energy profile. Because of the quantum character of the proton there is no energy barrier associated with proton transfer, and a slightly





**Figure 3.** Free-energy profiles for H motion in aqueous Hmaleate ion from CPMD (a1) and from PIMD (a2) separately for the simulation time when O1 was less solvated than O2, for the time when O2 was less solvated than O1, and for the time when both oxygen atoms were solvated similarly.

asymmetric single-well potential is observed, with the proton located predominantly on one oxygen atom of Hmaleate.

As previously, we have quantified the solvation and find that for the first 6 ps of simulation time one oxygen atom (designated O1) continued to be less solvated than the other, thereby favoring one solvatomer. Afterward O2 is sometimes less solvated, so that the other solvatomer becomes favored. But there is no possibility of a single centered potential. More specifically, population analysis shows that for approximately 44.6% of the simulation time oxygen atom O1 was less solvated than oxygen O2 and for 16.5% of the time oxygen atom O2 was less solvated. For approximately 38.9% of the time both oxygen atoms in the H-bond were solvated more or less the same.

However, the slightly asymmetric single-well potential of Figure 2c2, with the proton apparently located on oxygen that is less well solvated, is again a consequence of averaging over instantaneous solvation environments. If the averaging is restricted to shorter time intervals, a different picture emerges. The separate free-energy profiles for the 44.6% of the simulation time when O1 was less solvated than O2, for the 16.5% of the time when O2 was less solvated than O1, and for the 38.9% of time when both oxygen atoms were solvated similarly are shown in Figure 3a2. For the first two cases, where one oxygen was less solvated than the other, the slightly asymmetric single-well potential is no longer restricted to one in which the proton is located on O1. Instead, the proton is located on the oxygen that is less solvated. Moreover, for the 38.9% of time when both oxygen atoms were solvated to nearly the same extent, the single well is no longer asymmetric but is symmetric, within the accuracy of the sampling. Thus these profiles show how the instantaneous potential responds to the solvation environment. A further conclusion is that the strong asymmetry in the distribution function and the free-energy profile of Figure 2b2, c2 is a result of the long-term averaging over solvation environments, during most of which O1 happens to be less well solvated.

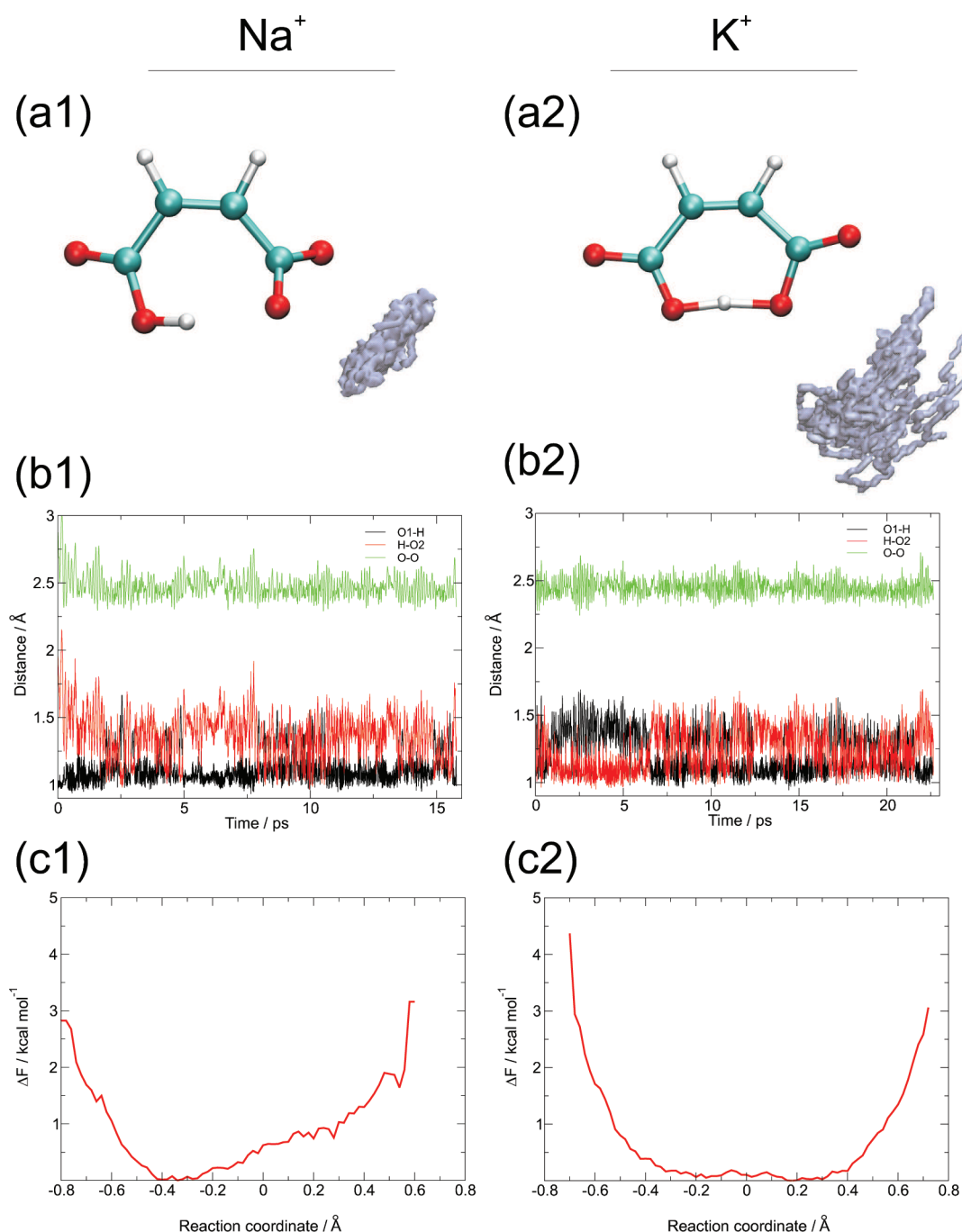
Thus we have shown that the intramolecular H-bond in Hmaleate ion in water has an asymmetric H-bond owing to the instantaneous solvation environment. In principle, with sampling over a longer time, the oxygens must become equivalent, restoring an apparent symmetry between them. Under such circumstances, the total time during which O1 is better solvated than O2 must be equal to the total time during which O2 is better solvated than O1, along with times during which they are solvated equally.

However, the simulation time that we used was too short to permit the oxygens to become equivalent. The instantaneous asymmetry of the H-bond then reflects the instantaneous solvation environment.

**2.3. Counterion Effects.** We next consider the influence of counterions on H-bond symmetry. We have examined two counterions, sodium and potassium, for which crystal structures have been discussed in the Isolated Hydrogen Maleate Ion and Its Potassium and Sodium Salts Section. The presence of a counterion is a major factor that can stabilize an asymmetric structure. According to Lluch's calculations on the QM/MM level (AM1-SRP/AMBER) for the H-bond in the potassium salt of Hphthalate anion in chloroform, the energy profile for the intramolecular proton transfer along the H-bond is a double well with two equivalent asymmetric minima.<sup>26</sup> Those calculations show further that a transition state with a centered position of the proton is observed when potassium ion is equidistant to both oxygen atoms of the H-bond and that an energy minimum is observed when the potassium ion is equidistant to the two oxygen atoms of one carboxyl group.

Results of our CPMD and PIMD simulations are presented in Figure 4. At the beginning of the simulations a bare  $\text{Na}^+$  or  $\text{K}^+$  was placed equidistant to the two oxygens of one of the carboxylates. During the equilibration period, 4.0 water molecules hydrate the  $\text{Na}^+$  and 6.4 waters hydrate the  $\text{K}^+$ . However, there are no waters directly between either  $\text{M}^+$  and the anion, because the equilibration and the simulation times are too short to overcome the barrier to separating the ions sufficiently to permit water to insert between them. In principle, the simulations could have been extended to much longer times or could have been started with fully hydrated  $\text{Na}^+(\text{H}_2\text{O})_5$  or  $\text{K}^+(\text{H}_2\text{O})_6$  near the Hmaleate anion, but this would be less interesting because the influence of the  $\text{M}^+$  would be smaller and would exert less control on the H-bond.

Figure 4a1 and a2 shows the spatial distribution functions of the  $\text{Na}^+$  and  $\text{K}^+$  in the vicinity of Hmaleate anion. The figures reveal differences in the positioning of sodium and potassium ions around the carboxyl group. Most of the time the sodium ion was located equidistant to the two oxygen atoms of one carboxyl, whereas the potassium ion tended to stay closer to the oxygen atom that is involved in the intramolecular H-bond. Figure 4b1 and b2 shows further how the position of the hydrogen in the H-bond varies during the simulations. Figure 4c1 and c2 show the resulting free-energy profiles for motion of the hydrogen.

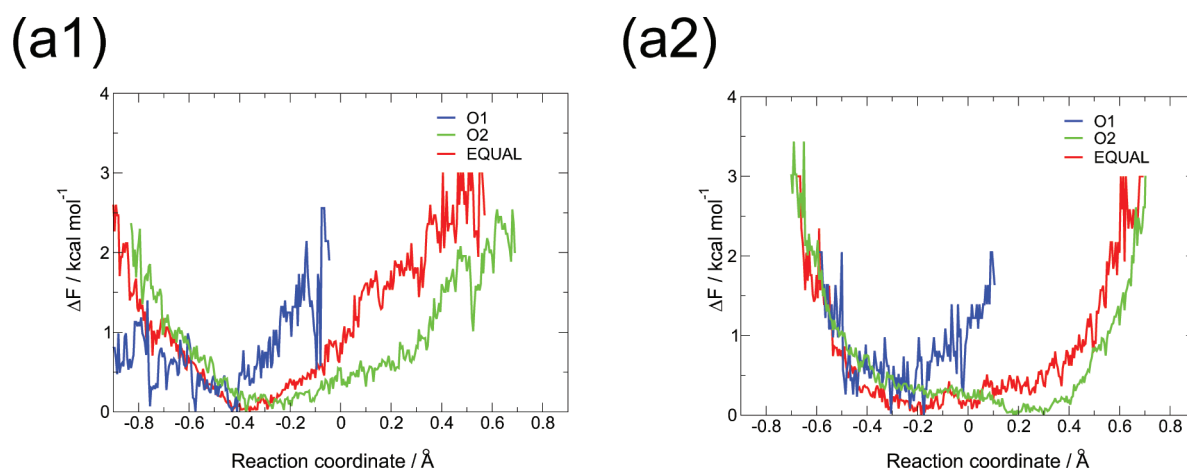


**Figure 4.** Hmaleate ion in water from CPMD simulations with two different counterions: sodium (left panel: a1, b1, and c1) and potassium (right panel: a2, b2, and c2). (a1,a2) Spatial distribution function of the counterion around Hmaleate ion. (b1,b2) Time evolution of the O–H distances involved in the intramolecular H-bond. (c1,c2) Free-energy profile for proton transfer within the H-bond.

Sodium and potassium ions affect the symmetry of the H-bond in Hmaleate ion in different ways. Potassium ion results in a centered single-well potential energy curve, whereas sodium ion results in an asymmetric single-well potential. The strong asymmetry of the H-bond in the presence of Na<sup>+</sup> is not surprising, because the Na<sup>+</sup> repels the proton, and thus a position nearer to O1 is preferred. The surprising result is the single-well potential centered at the midpoint of the H-bond in the presence of K<sup>+</sup>, which might have shown a greater repulsion because it is closer to the OHO.

To resolve this puzzle, the potential for H motion can again be separated into instantaneous solvation environments that favor one solvatomer or the other. We have quantified that solvation by summing the number of H-bonds between O1 or O2 of the OHO and the neighboring water hydrogens within 2.5 Å. We can thus distinguish which O is less solvated by water.

Population analysis for sodium Hmaleate in water shows that during approximately 58.3% of the simulation time oxygen atom O2 was less solvated by water than oxygen O1 and 6.8% of the time oxygen atom O1 was less solvated. For approximately 34.9%



**Figure 5.** CPMD free-energy profiles for Hmaleate anion with  $\text{Na}^+$  (a1) and  $\text{K}^+$  (a2) separately for the simulation time when O1 was less solvated than O2, for the time when O2 was less solvated than O1, and for the time when both oxygen atoms were solvated similarly.

of the time, both oxygen atoms in the intramolecular H-bond were solvated more or less the same. In comparison, population analysis for potassium Hmaleate in water shows that approximately 67.9% of the simulation time oxygen atom O2 was less solvated by water than oxygen O1 and 2.3% of time oxygen atom O1 was less solvated. For approximately 29.8% of time, both oxygen atoms in the intermolecular H-bond were solvated more or less the same.

Figure 5 shows free-energy profiles for H motion in Hmaleate anion in the presence of  $\text{Na}^+$  (a1) and  $\text{K}^+$  (a2) from CPMD, separated into instantaneous solvation environments. For  $\text{Na}^+$  the potential is always asymmetric, with a minimum at O1, which is farther from the sodium, regardless of the H-bonds from water. We suggest that this is because the repulsion by  $\text{Na}^+$  is so strong that a H position nearer to O1 was always preferred. For  $\text{K}^+$  the potential depends on the solvation by water. For the rare times when O1 was less solvated, a position for the H of the H-bond nearer to O1 was strongly preferred. In contrast, for the 67.9% of the time when O2 was less solvated, a position for the H nearer to O2 was slightly preferred, but with an energy cost of only 0.5 kcal/mol, the H could be nearer the other O, resulting in an almost symmetric potential. For the remaining time when both oxygens were solvated similarly, the potential is effectively symmetric. We suggest that these nearly symmetric potentials arise because the repulsion by the larger, more distant, and more hydrated  $\text{K}^+$  is weaker than for  $\text{Na}^+$  and competitive with solvation by water, so that the potential for H motion depends on the location of the  $\text{K}^+$ . Because the distribution of the  $\text{K}^+$  is more diffuse than that of the  $\text{Na}^+$ , which is more localized between the carboxylate oxygens, the  $\text{K}^+$  exerts a greater repulsion for the H when it is close to the OHO but a lesser repulsion when it is distant. The variable repulsion by the  $\text{K}^+$  balances the repulsion by the waters of hydration, which depends on which O is less solvated.

Because these results for potassium ion are contrary to previous results by Lluch,<sup>26</sup> one should understand why? A partial explanation is that our simulations never reached a situation where either potassium or sodium ion is equidistant to both oxygen atoms of the intramolecular H-bond, which is a transition state according to Lluch's calculations. Thus we should not have been able to observe a centered proton. That is the case with sodium, but not with potassium, according to Figure 5a2, where we suggest that the

centered proton arises because repulsion by the potassium ion near one carboxyl balances the repulsion by the waters of solvation on the other carboxyl.

### 3. CONCLUSIONS AND OUTLOOK

We have studied theoretically by means of *ab initio* Car–Parrinello molecular dynamics the symmetry of the intramolecular H-bond of Hmaleate anion as an isolated structure, in the crystals of its sodium and potassium salts, as a hydrated ion in water, and in water with counterions ( $\text{K}^+$  and  $\text{Na}^+$ ). The results confirm and clarify experimental findings by Perrin et al.<sup>19</sup>

When maleate ion is an isolated structure, molecular dynamics at 298 K predicts a truly symmetric potential with a centered proton. For two crystal structures (sodium and potassium salts) two different situations are observed. For the highly symmetric case of potassium Hmaleate, the potential is similar to that of the isolated ion, with a single-well centered proton. For sodium Hmaleate trihydrate crystal, which is of low symmetry, the different environments of the two oxygen atoms of the intramolecular H-bond of the Hmaleate result in an asymmetric single-well potential with a proton located on one oxygen atom, in agreement with the neutron diffraction study by Olovsson.<sup>14</sup>

Thanks to a new procedure that can produce separate free-energy profiles from periods during the simulation when one oxygen atom of the intramolecular H-bond was less well solvated than the other one and periods where both oxygen atoms were solvated to a similar extent, we have shown that the position of the proton in aqueous Hmaleate ion is entirely dependent on the solvation pattern around the oxygen atoms in the intramolecular H-bond. It is shown that the proton is always located on the oxygen atom that is less well solvated and that there is no longer a minimum at the other, better solvated oxygen.

Additionally, separation into instantaneous solvation environments has been applied to an aqueous solution of Hmaleate ion with two different counterions, namely  $\text{K}^+$  and  $\text{Na}^+$ . Analysis of their influence on intramolecular H-bond symmetry revealed that, whereas the potential of the intramolecular H-bond in the  $\text{Na}^+$  salt is always asymmetric, owing to strong repulsion by the  $\text{Na}^+$  regardless of the H-bonds to water, for the  $\text{K}^+$  salt the repulsion by this larger and more hydrated ion is weaker than for  $\text{Na}^+$  and competitive with solvation by water, so that the potential for H motion depends on the location of the  $\text{K}^+$ .

The systems described here combine several factors that make computational studies of H-bonds difficult. First, the proton potential is highly anharmonic and fluctuating, and it is necessary to use a statistical description using molecular dynamics methods, and second, the barriers are of such height that quantum effects can be influential, as our path integral results indicate. It is very challenging to reproduce spectra of such systems. It is possible to reuse the trajectories used in this work, for example, to extract snapshots and construct their proton potential energy surfaces with the aim of further solving the vibrational Schrödinger equation and obtaining a statistically averaged vibrational spectrum.<sup>41–43</sup> If NMR parameters are also calculated for each point of the potential energy surface, the expectation value of the NMR shift of the proton can be obtained—such a technique was proven successful for a highly anharmonic short H-bond.<sup>44</sup> However, such calculations are beyond the scope of the current study.

#### 4. METHODS

All calculations are based on ab initio molecular dynamics<sup>45</sup> using the efficient Car–Parrinello<sup>46</sup> propagation scheme as implemented in the CPMD program package.<sup>47</sup> These pseudo-potential calculations have been carried out using the PBE<sup>48</sup> exchange–correlation functional within the spin-restricted Kohn–Sham formalism together with a plane-wave basis set at a kinetic energy cutoff of 100 Ry,  $\Gamma$ -point sampling of the Brillouin zone, and Troullier–Martins<sup>49</sup> norm-conserving pseudopotentials. The supercell for all these calculations was a cubic box 15 Å in length with periodic boundary conditions. All dynamic simulations were performed in the canonical ensemble at 298 K using Nosé–Hoover chain thermostats<sup>50</sup> in order to control the kinetic energy of the nuclei (as well as the fictitious kinetic energy of the orbitals). For the path integral case, a separate thermostat was used for each degree of freedom.<sup>51</sup> To account for the fact that the PBE functional tends to overstructure water compared to experiment, some researchers have conducted simulations at 400 K.<sup>52</sup> Nevertheless in all our simulations, we have used a proper temperature of 298 K together with long enough trajectories to avoid overly rapid proton transfer inside the intramolecular H-bond.

We have adopted the same approach as Miura et al.,<sup>53</sup> in which the positions of the atoms initially evolve according to the classical equations of motion. Then we proceed with PIMD simulation,<sup>54–56</sup> which explores the quantum behavior of both the nuclear and electronic degrees of freedom. It maps the problem of a quantum particle into one of a classical ring polymer model with beads that interact through temperature- and mass-dependent spring forces. Such mapping is known in the literature as a quantum classical isomorphism.<sup>57–59</sup> It should be underlined that “real” properties of the quantum systems are recovered only when the number of beads is extrapolated to infinity. The path integral simulations in the present study used eight beads and the normal mode variable transformation.<sup>56</sup>

A molecular dynamics time step of  $\delta t = 3$  au ( $\approx 0.073$  fs) was used for the integration of the Car–Parrinello equations of motion using a fictitious mass parameter for the orbitals of 400 au together with the proper atomic masses. The initial configurations were generated with classical molecular dynamics simulations of 1 ns. After this initial equilibration period (ca. 30 000 steps), the Car–Parrinello molecular dynamics simulations were performed, and the data were collected over trajectories spanning 300 000 steps (ca. 22 ps) for the sodium crystal

(NaHMAL), 200 000 integration steps for the potassium crystal (KHMAL) (ca. 16 ps), 280 000 steps (ca. 20 ps) for Hmaleate ion with 103 water molecules, 250 000 steps (ca. 18 ps) for sodium Hmaleate with 102 water molecules, and 300 000 steps (ca. 22 ps) for potassium Hmaleate with 101 water molecules. Path integral runs were performed for similar time periods as the Car–Parrinello simulations.

The reaction coordinate  $\delta$  is defined in eq 1 as the difference between O1–H and O2–H distances, where O1 labels the oxygen that bears the H at the beginning of the simulation and O2 labels the oxygen that is initially H-bonded to the H:

$$\delta = R_{\text{O1-H}} - R_{\text{O2-H}} \quad (1)$$

The reaction coordinate is a measure of the degree of proton transfer in the H-bond, with zero corresponding to the midpoint of the H-bond. Also, for the studies of counterions, the oxygen atoms involved in the H-bond are labeled as O2 for the one closer to the  $M^+$  ion and O1 for the farther ( $\text{O1} \cdots \text{H} \cdots \text{O2} \cdots M^+$ ). The free-energy profiles were obtained from eq 2, where  $k$  is the Boltzmann constant,  $N_A$  is the Avogadro number,  $T$  is the simulation temperature, and  $P$  is the proton distribution as a function of the reaction coordinate.

$$\Delta F = -k \cdot N_A \cdot T \ln(P(\delta)) \quad (2)$$

The visualize molecular dynamics (VMD)<sup>60</sup> program has been used for data visualization.

#### ■ ASSOCIATED CONTENT

**S Supporting Information.** Figures 1–4 compile free energy profiles generated for all four studied cases: standard Car–Parrinello (Figure 1) and path integral (Figure 2) simulations of hydrogen maleate ion in water and hydrogen maleate ion in water with potassium (Figure 3) and sodium (Figure 4) counterions. This material is available free of charge via the Internet at <http://pubs.acs.org>.

#### ■ AUTHOR INFORMATION

##### Corresponding Author

\*E-mail: [mclar@elrond.chem.uni.wroc.pl](mailto:mclar@elrond.chem.uni.wroc.pl); [cperrin@ucsd.edu](mailto:cperrin@ucsd.edu).

#### ■ ACKNOWLEDGMENT

We are grateful to Jaroslaw Panek for useful discussions and to University of Wroclaw (Internal Grant for Young Scientist no. 105/10/E-344/M/2011 to P.D.) and to the U.S. National Science Foundation (grant CHE07-42801 to C.L.P.) for financial support. The calculations were carried out using resources from Wroclaw Supercomputer Center (WCSS) and the GALERATION Cluster and the Academic Computer Center in Gdańsk (CI TASK).

#### ■ REFERENCES

- (1) Pimentel, G. C.; McClellan, A. L. *The Hydrogen Bond*; Freeman: San Francisco, CA, 1960.
- (2) Vinogradov, S.; Linnel, R. *The Hydrogen Bond*; Van Nostrand-Reinhold: New York, 1971.
- (3) *The Hydrogen Bond: Recent Developments in Theory and Experiments*; Schuster, P.; Zundel, G., Sandorfy, C., Eds.; North-Holland: Amsterdam, The Netherlands, 1976.
- (4) Warshel, A. *Biochemistry* **1981**, *20*, 3167–3177.

- (5) Jeffrey, G. A.; Saenger, W. *Hydrogen Bonding in Biological Structures*; Springer: Berlin, Germany, 1991.
- (6) Warshel, A. *Computer Modeling of Chemical Reactions in Enzymes and Solutions*; John Wiley & Sons: New York, 1997.
- (7) Jeffrey, G. *An Introduction to Hydrogen Bonding*; Oxford University Press: Oxford, U.K., 1997.
- (8) Scheiner, S. *Hydrogen Bonding: A Theoretical Perspective*; Oxford University Press: New York, 1997.
- (9) Hadzi, D. *The Hydrogen Bond*; Wiley: Chichester, U.K., 1997.
- (10) Desiraju, G.; Steiner, T. *The Weak Hydrogen Bond in Structural Chemistry and Biology*; Oxford University Press: Oxford, U.K., 1999.
- (11) Warshel, A. *Acc. Chem. Res.* **2002**, *35*, 385–395.
- (12) Perrin, C. L. *Pure Appl. Chem.* **2009**, *81*, 571–583.
- (13) Darlow, S.; Cochran, W. *Acta Crystallogr.* **1961**, *14*, 1250–1257.
- (14) Olovsson, G.; Olovsson, I.; Lehmann, M. S. *Acta Crystallogr.* **1984**, *C40*, 1521–1526.
- (15) Perrin, C. L.; Thoburn, J. D. *J. Am. Chem. Soc.* **1989**, *111*, 8010–8012.
- (16) Perrin, C. L.; Thoburn, J. D. *J. Am. Chem. Soc.* **1992**, *114*, 8559–8565.
- (17) Perrin, C. L. *Science* **1994**, *266*, 1665–1668.
- (18) Mavri, J.; Hodoscek, M.; Hadzi, D. *J. Mol. Struct. (Theochem)* **1990**, *209*, 421–431.
- (19) Perrin, C. L.; Nielson, J. B. *J. Am. Chem. Soc.* **1997**, *119*, 12734–12741.
- (20) Perrin, C. L.; Ohta, B. K. *J. Am. Chem. Soc.* **2001**, *123*, 6520–6526.
- (21) Perrin, C. L.; Ohta, B. K. *Bioorg. Chem.* **2002**, *30*, 3–15.
- (22) Perrin, C. L.; Lau, J. S. *J. Am. Chem. Soc.* **2006**, *128*, 11820–11824.
- (23) Perrin, C. L.; Karri, P. *Chem. Commun. (Cambridge, U. K.)* **2010**, *46*, 481–483.
- (24) Perrin, C. L.; Lau, J. S.; Kim, Y.-J.; Karri, P.; Moore, C.; Rheingold, A. L. *J. Am. Chem. Soc.* **2009**, *131*, 13548–13554.
- (25) Perrin, C. L.; Ohta, B. K. *J. Mol. Struct.* **2003**, *644*, 1–12.
- (26) Garcia-Viloca, M.; Gonzalez-Lafont, A.; Lluch, J. M. *J. Am. Chem. Soc.* **1999**, *121*, 9198–9207.
- (27) Perrin, C. L. *Acc. Chem. Res.* **2010**, *43*, 1550–1557.
- (28) George, P.; Bock, C. W.; Trachtman, M. *J. Phys. Chem.* **1983**, *87*, 1839–1841.
- (29) Hodoscek, M.; Hadzi, D. *J. Mol. Struct. (Theochem)* **1990**, *209*, 411–419.
- (30) Garcia-Viloca, M.; Gonzalez-Lafont, n.; Lluch, J. M. *J. Am. Chem. Soc.* **1997**, *119*, 1081–1086.
- (31) Bach, R. D.; Dmitrenko, O.; Glukhovtsev, M. N. *J. Am. Chem. Soc.* **2001**, *123*, 7134–7145.
- (32) Woo, H.-K.; Wang, X.-B.; Wang, L.-S.; Lau, K.-C. *J. Phys. Chem. A* **2005**, *109*, 10633–10637.
- (33) Ratajczak, H.; Barnes, A.; Baran, J.; Yaremko, A.; Latajka, Z.; Dopieralski, P. *J. Mol. Struct. (Theochem)* **2008**, *887*, 9–19.
- (34) Wilson, C. C.; Thomas, L. H.; Morrison, C. A. *Chem. Phys. Lett.* **2003**, *381*, 102–108.
- (35) Wilson, C. C.; Thomas, L. H.; Morrison, C. A. *Chem. Phys. Lett.* **2004**, *399*, 292–293.
- (36) Darlow, S. *Acta Crystallogr.* **1961**, *14*, 1257–1259.
- (37) Marcus, R. A. *J. Chem. Phys.* **1956**, *24*, 979–989.
- (38) Borgis, D.; Hynes, J. T. *J. Chem. Phys.* **1991**, *94*, 3619–3628.
- (39) Mavri, J.; Berendsen, H. J. C.; van Gunsteren, W. F. *J. Phys. Chem.* **1993**, *97*, 13469–13476.
- (40) Tuckerman, M. E.; Marx, D.; Parrinello, M. *Nature* **2002**, *417*, 925–929.
- (41) Pejov, L.; Spangberg, D.; Hermansson, K. *J. Phys. Chem. A* **2005**, *109*, 5144–5152.
- (42) Jezierska, A.; Panek, J.; Borstnik, U.; Mavri, J.; Janezic, D. *J. Phys. Chem. B* **2007**, *111*, 5243–5248.
- (43) Stare, J.; Panek, J.; Eckert, J.; Grdadolnik, J.; Mavri, J.; Hadzi, D. *J. Phys. Chem. A* **2008**, *112*, 1576–1586.
- (44) Stare, J.; Jezierska, A.; Ambrozic, G.; Kosir, I. J.; Kidric, J.; Koll, A.; Mavri, J.; Hadzi, D. *J. Am. Chem. Soc.* **2004**, *126*, 4437–4443.
- (45) Marx, D.; Hutter, J. *Ab Initio Molecular Dynamics: Basic Theory and Advanced Methods*; Cambridge University Press: Cambridge, U.K., 2009.
- (46) Car, R.; Parrinello, M. *Phys. Rev. Lett.* **1985**, *55*, 2471–2474.
- (47) Hutter, J.; et al. *CPMD Program Package*; IBM Corporation and Max-Planck Institut: Stuttgart, Germany, 1990; <http://www.cpm.org>, (date accessed January 2, 2010).
- (48) Perdew, J. P.; Burke, K.; Ernzerhof, M. *Phys. Rev. Lett.* **1996**, *77*, 3865–3868.
- (49) Troullier, N.; Martins, J. L. *Phys. Rev. B* **1991**, *43*, 1993–2006.
- (50) Martyna, G. J.; Klein, M. L.; Tuckerman, M. *J. Chem. Phys.* **1992**, *97*, 2635–2643.
- (51) Tuckerman, M.; Laasonen, K.; Sprik, M.; Parrinello, M. *J. Chem. Phys.* **1995**, *99*, 5749–5752.
- (52) Sit, P. H.-L.; Marziari, N. *J. Chem. Phys.* **2005**, *122*, 204510–204518.
- (53) Miura, S.; Tuckerman, M.; Klein, M. *J. Chem. Phys.* **1998**, *109*, 5290–5299.
- (54) Marx, D.; Parrinello, M. *Z. Phys. B* **1994**, *95*, 143–144.
- (55) Marx, D.; Parrinello, M. *J. Chem. Phys.* **1996**, *104*, 4077–4082.
- (56) Tuckerman, M.; Marx, D.; Klein, M.; Parrinello, M. *J. Chem. Phys.* **1996**, *104*, 5579–5588.
- (57) Feynman, R.; Hibbs, A. *Quantum Mechanics and Path Integrals*; McGraw-Hill: New York, 1965.
- (58) Schweizer, K.; Stratt, R.; Chandler, D.; Wolynes, P. *J. Chem. Phys.* **1981**, *75*, 1347–1369.
- (59) Chandler, D.; Wolynes, P. *J. Chem. Phys.* **1981**, *74*, 4078–4095.
- (60) Humphrey, W.; Dalke, A.; Schulten, K. *J. Mol. Graphics* **1996**, *14*, 33–38.

# Thermodynamic Properties of Water Molecules at a Protein–Protein Interaction Surface

David J. Huggins,<sup>\*,†,‡,§</sup> May Marsh,<sup>†,⊥</sup> and Mike C. Payne<sup>†,§</sup>

<sup>†</sup>Cambridge Molecular Therapeutics Programme, Hutchison/MRC Research Centre, University of Cambridge, Hills Road, Cambridge, CB2 0XZ, United Kingdom

<sup>‡</sup>Department of Chemistry, University of Cambridge, Lensfield Road, Cambridge, UK CB2 1EW, United Kingdom

<sup>§</sup>TCM Group, Cavendish Laboratory, University of Cambridge, 19 J J Thomson Avenue, Cambridge CB3 0HE, United Kingdom

<sup>⊥</sup>Department of Biochemistry, University of Cambridge, Tennis Court Road, Cambridge, UK CB2 1QW, United Kingdom

**ABSTRACT:** Protein–protein interactions (PPIs) have been identified as a vital regulator of cellular pathways and networks. However, the determinants that control binding affinity and specificity at protein surfaces are incompletely characterized and thus unable to be exploited for the purpose of developing PPI inhibitors to control cellular pathways in disease states. One of the key factors in intermolecular interactions that remains poorly understood is the role of water molecules and in particular the importance of solvent entropy. This factor is expected to be particularly important at protein surfaces, and the release of water molecules from hydrophobic regions is one of the most important drivers of PPIs. In this work, we have studied the protein surface of a mutant of the protein RadA to quantify the thermodynamics of surface water molecules. RadA and its human homologue RAD51 function as recombinases in the process of homologous recombination. RadA binds to itself to form oligomeric structures and thus contains a well-characterized protein–protein binding surface. Similarly, RAD51 binds either to itself to form oligomers or to the protein BRCA2 to form filaments. X-ray crystallography has determined that the same interface functions in both interactions. Work in our group has generated a partially humanized mutant of RadA, termed MAYM, which has been crystallized in the apo form. We studied this apo form of MAYM using a combination of molecular dynamics (MD) simulations and inhomogeneous fluid solvation theory (IFST). The method locates a number of the hydration sites observed in the crystal structure and locates hydrophobic sites where hydrophobic species are known to bind experimentally. The simulations also highlight the importance of the restraints placed on the protein in determining the results. Finally, the results identify a correlation between the predicted entropy of water molecules at a given site and the solvent-accessible surface area and suggest that correlations between water molecules only need to be considered for water molecules separated by less than 3.2 Å. The combination of MD and IFST has been used previously to study PPIs and represents one of the few existing methods to quantify solvent thermodynamics. This is a vital aspect of molecular recognition and one which we believe must be developed.

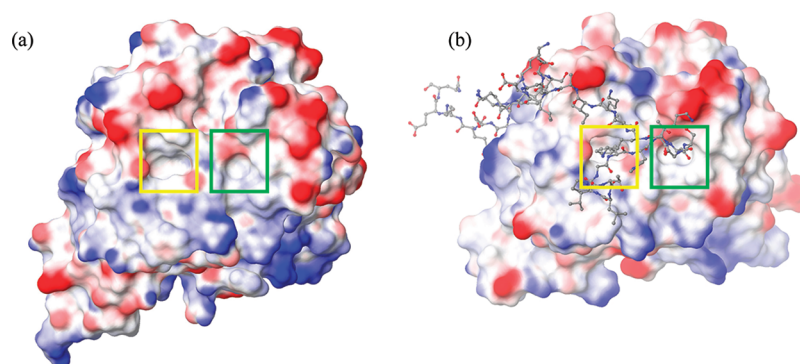
## INTRODUCTION

Protein–protein interactions (PPIs) are essential in controlling cellular networks and play an important role in many disease states.<sup>1</sup> Significant efforts are now being focused on understanding the nature of the intermolecular interactions in PPIs, and computational methods are a key aspect of increasing our understanding.<sup>2,3</sup> In addition, PPIs are now increasingly being targeted for drug development, and computational methods are commonly combined with structural data in virtual screening and lead optimization for PPI targets.<sup>4</sup> One aspect of molecular interactions that is particularly important for understanding PPIs is hydrophobic association driven by desolvation of nonpolar protein surfaces. Water molecules form significant hydrogen bonding interactions in bulk water and are somewhat ordered. Conversely, water molecules at a hydrophobic surface have reduced hydrogen bonding interactions and have differing levels of order, dependent upon the environment. The balance of these components is one of the key factors that controls the thermodynamics of binding. This has been proposed as the principal driving force for binding in a number of systems and also impacts protein folding and stability.<sup>5</sup> In this study, we apply solvation thermodynamics to a prototypical PPI surface.

**Recombinase Biology.** Recombinases such as RadA and RAD51 are key factors in the process of homologous recombination (HR) to repair broken double strand breaks (DSBs) in DNA.<sup>6</sup> The human RAD51 recombinase is known to form an oligomeric structure in the cell, where it is sequestered until needed for HR. Shortly after DNA replication, RAD51 is loaded onto DNA around DSBs by association with the so-called BRC repeats of the regulatory BRCA2 protein.<sup>7</sup> RadA, the archaeal homologue of RAD51, is sequestered in oligomeric structure in the cell but appears to bind DNA as a helical filament without the presence of a regulatory protein.<sup>8</sup> The interface for oligomerization has been identified in RadA and RAD51 by crystallography.<sup>9,10</sup> The key determinant of binding is the presence of a hydrophobic pocket on the surface that binds a phenylalanine residue.<sup>11</sup> Another smaller pocket is found in close proximity and binds an alanine residue. These pockets are termed the phenylalanine pocket and the alanine pocket. RadA and RAD51 oligomerize by bringing together their hydrophobic surfaces with an FMRA and an FTTA sequence, respectively. The BRC repeats of BRCA2 also exploit these pockets to bind RAD51 with a

Received: July 13, 2011

Published: September 20, 2011



**Figure 1.** (a) The molecular surface of RAD51 in complex with the BRC4 peptide from PDBID 1N0W. (b) The molecular surface of MAYM. RAD51 and MAYM are colored by electrostatic potential and BRC4 is displayed as atom colored balls and sticks. The phenylalanine and alanine pockets are boxed in yellow and green, respectively.

conserved FXXA motif.<sup>12</sup> In addition, these pockets are surrounded by a surface dotted with hydrophobic patches, as shown in Figure 1a for MAYM and Figure 1b for RAD51. This surface is thus typical of a PPI and provides a good test case to explore the thermodynamics of solvation and how it contributes to protein–protein association.

**Inhomogeneous Fluid Solvation Theory.** Inhomogeneous fluid solvation theory (IFST) was developed by Lazaridis as a method to study hydrophobic hydration<sup>13</sup> by calculating interactions and correlations between water molecules through an analysis of molecular dynamics (MD) or Monte Carlo (MC) simulations. IFST was initially used to study pure water,<sup>14</sup> but the theory was then extended to consider small hydrophobic solutes<sup>15</sup> and then to consider protein binding sites.<sup>16</sup> In IFST, bulk water is considered as a reference state, and other molecules perturb this state, resulting in a change in enthalpy and entropy.<sup>17</sup> This is quantified by calculating the interaction energies and the correlation functions between the water molecules and the solute.<sup>15</sup> Regions of high water density are identified and then analyzed to compare the enthalpy and entropy with water molecules in bulk solvent. The methodology is described in detail below. IFST has been used to analyze a number of ligand binding sites to elucidate the role of water molecules.<sup>16,18,19</sup> IFST has also shown success in predicting binding affinities and has recently been implemented in Schrodinger's WaterMap software.<sup>20,21</sup> WaterMap has also been applied to explain binding affinities and specificities for PDZ domain<sup>22</sup> and for the polo-box domain of the mitotic kinase PLK1.<sup>23</sup> It has also been employed recently by Zielkiewicz to study water molecules around simple polypeptides.<sup>24</sup>

Here, we apply IFST to the protein surface of the RadA MAYM mutant and explore the thermodynamic properties of water molecules at a PPI interface. This analysis quantifies the intermolecular interactions that underlie PPIs and allows the identification of potential binding hotspot regions.

## MATERIALS AND METHODS

We performed MD simulations of bulk water and of the apo MAYM protein using NAMD<sup>25</sup> using a number of simulation protocols.

**Crystallography.** The crystal structure of RadA was taken from a protein construct of *Pyrococcus furiosus* RadA (accession number AF052597) containing residues 108–349 (Marsh et al., unpublished). Residues 288–300 in the L2 loop were replaced by a single Asn residue, and residues 108–286, 304–329, and

336–349 have assigned density. The MAYM form of RadA has four humanizing mutations: I169M, Y201A, V202Y, K221M. The crystal structure contains one DMSO solvent molecule and one phosphate group. This protein construct lacks an N-terminal domain and thus does not oligomerize. However, the N-terminal domain is located over 15 Å from the phenylalanine and alanine pockets<sup>9,26</sup> and is thus unlikely to affect the properties of this surface.

**Structure Preparation.** The protein structure was initially prepared as follows. Atom coordinates for the protein and the water molecules were taken from the X-ray crystal structure. The DMSO solvent molecule and the phosphate group were deleted from the structure. The hydrogen-atom positions for the protein and the water molecules were then built using the PSFGEN mode of VMD<sup>27</sup> with the CHARMM27 energy function.<sup>28,29</sup> Histidine residues were then manually checked for protonation state. His210, His243, and His269 were assigned as epsilon protonated. All remaining histidines were assigned as delta protonated. The residues lysine, arginine, aspartate, glutamate, cysteine, and tyrosine were also analyzed to check their protonation state. There was no evidence of any unusual protonation states, and thus all lysine and arginine residues were assigned as positively charged, all aspartate and glutamate residues were assigned as negatively charged, and all cysteine and tyrosine residues were assigned as neutral. The terminal residues 304 and 336 were patched with an *N*-acetyl group, and the terminal residues 286 and 329 were patched with an *N*-methyl amide group. The atomic charges were assigned from the CHARMM27 forcefield.<sup>28,29</sup> All water molecules were modeled with the TIP4P/2005 water model.<sup>30</sup> The next stage was to solvate the protein with water molecules. All the water molecules observed in the crystal structure were retained. Solvation was performed with the SOLVATE program<sup>31</sup> version 1.0 from the Max Planck Institute to generate a solvation sphere of radius 50 Å around the center of the protein. No ions were included in the solution, as the protein has a net charge of zero. The system was then cut to form a rhombic dodecahedron (RHDO) with an edge length of 60 Å using the CHARMM program (version 34b1).<sup>32</sup>

**Equilibration.** During all simulations with the RHDO, the protein atoms were fixed, the RHDO was treated using periodic boundary conditions, and the electrostatics were modeled using the particle mesh Ewald method.<sup>33</sup> The water molecules in the RHDO were first subjected to energy minimization for 10 000 steps using NAMD. This was followed by MD equilibration for

100 ps in an NPT ensemble and then MD equilibration for 100 ps in an NVT ensemble. This stage of preparation was undertaken to equilibrate the density of the water molecules at the surface. The density of the water molecules plays an important role in IFST and is thus important to converge accurately. We ensured that the system was brought to equilibrium before continuing our simulations by verifying that the system reached a point where the energy fluctuations were stable. In the next stage, the RHDO was cut to form a sphere of water molecules around the binding pocket of interest using the CHARMM program. The solvent sphere of radius 20 Å was centered at the coordinates of the CA atom of Ala201. This is defined as the centroid of the solvent sphere. The resulting system containing the protein and a sphere of water molecules was then treated with three protocols. For each protocol, the system was subjected to MD equilibration for 100 ps using NAMD with spherical boundary conditions.<sup>34</sup> Again, we ensured that the system was brought to equilibrium before beginning the MD simulation by verifying that the system reached a point where the energy fluctuations were stable for each protocol. The three protocols are as follows:

- (1) Fixed: All protein atoms were kept fixed.
- (2) Restrained: All atoms of any residue partially or completely outside the 20 Å sphere were fixed in place. All heavy atoms of any residue completely inside the 20 Å sphere were restrained using a 1.0 kcal/mol/Å<sup>2</sup> harmonic force.
- (3) Free: All atoms of any residue partially or completely outside the 20 Å sphere were fixed in place. All atoms of any residue completely inside the 20 Å sphere were not constrained.

**Molecular Dynamics.** Production simulations were performed for 10.0 ns at 300 K. All MD simulations were performed using the NAMD program version 2.7b3<sup>32</sup> with the CHARMM27 force field<sup>28,29</sup> using an MD time step of 2.0 fs. Electrostatic interactions were modeled with a uniform dielectric and a dielectric constant of 1.0 throughout the setup and production runs. Van der Waals interactions were truncated at 12.0 Å with switching from 8.0 Å. Bulk solvent was simulated as a periodic box of edge length 25 Å for a period of 8 ns using the same methods, parameters, and equilibration procedures detailed above.

**Clustering.** The 10.0 ns MD runs were first analyzed to cluster the water molecules into distinct spherical regions of high number density. These regions have been termed hydration sites in previous work using IFST,<sup>20</sup> and we retain this terminology here. We employed a radius of 1.2 Å for these hydration sites, in line with prior work.<sup>18</sup> The hydration sites were selected by sampling 1000 snapshots from the MD trajectory. All 1000 snapshots were superposed to generate a profile of the water density. Within the complete water density profile, we identified the oxygen atom of the water molecule with the largest number of water molecules within a 1.2 Å radius. The 1.2 Å sphere around the position of this oxygen atom was defined as a hydration site. This water molecule and all of its neighboring water molecules within 1.2 Å from any snapshot were excluded from further consideration. The process was then repeated to identify more hydration sites, allowing no new hydration sites within 1.2 Å of a previously defined hydration site. This iteration was terminated once when the density of an identified hydration sites fell below 1.5 times the number density of bulk water, which corresponds to an occupancy of 0.36 in the sphere of radius 1.2 Å. Only hydration sites within 12.0 Å of the solvation sphere center were considered. The resultant set of hydration sites was then subjected to energy and entropy calculations using IFST.

**Energy Calculations.** The interaction energy of each hydration site was calculated by sampling 5000 snapshots, taken every 2 ps from the 10.0 ns simulation. For each snapshot, we computed the average interaction energy with both the protein and all the other water molecules with VMD version 1.8.7 using the namdenergy plugin. This was then compared with the interaction energy of a water molecule determined from the bulk water simulation (−23.62 kcal/mol) to calculate the energy difference  $\Delta E$  shown in eq 1.

$$\Delta E = \bar{E}_{\text{water/protein}}^{\text{surface}} + \bar{E}_{\text{water/water}}^{\text{surface}} - \bar{E}_{\text{water/water}}^{\text{bulk}} \quad (1)$$

In this equation,  $\Delta E$  is the energy difference,  $\bar{E}_{\text{water/protein}}^{\text{surface}}$  is the mean interaction energy between a water molecule in the hydration site and the protein,  $\bar{E}_{\text{water/water}}^{\text{surface}}$  is the mean interaction energy between a water molecule in the hydration site and all of the other water molecules, and  $\bar{E}_{\text{water/water}}^{\text{bulk}}$  is the mean total interaction energy of a water molecule in bulk.

**Entropy Calculations.** The entropy of each hydration site was calculated by sampling 100 000 snapshots, taken every 100 fs from the 10.0 ns simulation. The entropy difference between a water molecule at a hydration site and in bulk was calculated from the contributions of the protein–water term ( $S_{pw}$ ), the water–water reorganization term ( $S_{ww}$ ), and a term arising from the change in density ( $S_{\text{density}}$ ).<sup>35</sup> These terms can be calculated by integrating over the protein–water  $g_{pw}(r, \omega)$  and water–water  $g_{ww}(r, \omega, r', \omega')$  correlation functions, where the variable  $r$  represents the position of the water molecule with respect to the center of the hydration site, and the Euler angles  $\omega$  represent the orientation of the water molecule in the fixed protein reference frame. As in previous work, only correlations between two species were considered.<sup>18,20</sup> The protein–water correlations function were calculated using a bin size of 0.06 Å for the radial component and 18° for the angular components. The protein–water and contribution to the entropy of changing the number density<sup>35</sup> can be calculated for each hydration site using eqs 2 and 3, where  $k$  is Boltzmann's constant,  $\rho$  is the number density of bulk solvent,  $\rho_{\text{site}}$  is the number density of the hydration site being considered, and  $\Omega$  is the integral over the Euler angles  $\omega$ .

$$S_{pw} = -k\rho/\Omega \int g_{sw}(r, \omega) \ln g_{sw}(r, \omega) dr d\omega \quad (2)$$

$$S_{\text{density}} = k \ln \left[ \frac{\rho}{\rho_{\text{site}}} \right] \quad (3)$$

As in previous work, the protein–water term was separated into translational,  $S_{pw}^{\text{trans}}$ , and orientational,  $S_{pw}^{\text{orient}}$ , entropic contributions, and the orientational distributions were assumed to be independent of the position of the water molecules within the sites.<sup>18</sup> The entropies were calculated using eqs 4 and 5, where  $g_{pw}^{\text{trans}}(r)$  and  $g_{pw}^{\text{orient}}(\omega)$  are the translational and orientational correlation functions.

$$S_{pw}^{\text{trans}} = -k\rho \int g_{pw}^{\text{trans}}(r) \ln g_{pw}^{\text{trans}}(r) dr \quad (4)$$

$$S_{pw}^{\text{orient}} = -k\rho/\Omega \int g_{pw}^{\text{trans}}(r) dr \int g_{pw}^{\text{orient}}(\omega) \ln g_{pw}^{\text{orient}}(\omega) d\omega \quad (5)$$

The water–water reorganization term was calculated for each pair of hydration sites within a distance of 3.5 Å. This distance corresponds to water molecules in the first solvation shell of a



water molecule in bulk. The water–water correlation functions were calculated using a bin size of 0.1 Å for the radial component and 18° for the angular components. For a given hydration site, the total reorganization entropy was calculated as the sum of the pairs of proximal sites. This term was then compared with the entropy of a water molecule from the bulk water simulation due to other water molecules within 3.5 Å (11.24 cal/mol/K). The entropies were calculated using eq 6.

$$\Delta S_{ww} = \sum S_{w,w'} - S_{w,w'}^{\text{bulk}} \quad (6)$$

$\Delta S_{ww}$  is the water–water entropy change,  $S_{w,w'}$  is the pair entropy between a water molecule in the hydration site and a water molecule in another hydration site and  $S_{w,w'}^{\text{bulk}}$  is the pair entropy of a water molecule in bulk. The contribution to the enthalpy from water–water correlations was also split into translational and orientational contributions. However, because of the vast amount of data required to accurately calculate the multidimensional water–water correlation functions, we employed two approximations first proposed by Li and Lazaridis.<sup>18</sup> The first is that the water–water correlation functions can be treated as dependent only on the relative orientation of the two water molecules and the distance between the centers of the two hydration sites. This correlation function can, in turn, be separated into translational and orientational contributions.

$$g_{ww}(r,r',\omega,\omega') = g_{ww}(R,\omega^{\text{rel}}) \quad (7)$$

$$g_{ww}(R,\omega^{\text{rel}}) = g_{ww}(R)g_{ww}(\omega^{\text{rel}}|R) \quad (8)$$

In these equations,  $g_{ww}$  are the water–water correlation functions,  $r'$  represents the position of the second water molecule with respect to the center of its hydration site,  $\omega'$  represents the orientation of the second water molecule in the fixed protein reference frame, the variable  $R$  is the distance between the centers of the two hydration sites, and  $\omega^{\text{rel}}|R$  is the relative orientation of two water molecules at a distance  $R$ . The second approximation is that the water–water correlation functions for the bound waters are the same as the water–water correlation functions in bulk water. This leads to eqs 9, 10, and 11, where the variables  $\theta_1$ ,  $\theta_2$ ,  $\chi_1$ ,  $\chi_2$ , and  $\varphi$  are the five angles that specify the relative orientation of two water molecules.<sup>14</sup>

$$g_{ww}(R) = g_{ww}^{\text{bulk}}(R) \quad (9)$$

$$g_{ww}(\omega^{\text{rel}}|R) = g_{ww}^{\text{bulk}}(\omega^{\text{rel}}|R) \quad (10)$$

$$g_{ww}^{\text{bulk}}(\omega^{\text{rel}}|R)g_{ww}^{\text{bulk}}(\theta_1,\theta_2,\chi_1,\chi_2,\varphi|R) \quad (11)$$

Application of these approximations leads to eqs 12 and 13.

$$S_{ww}^{\text{trans}} = -\frac{1}{2}k\rho^2 \int g_{pw(a)}^{\text{trans}}(r)g_{pw(b)}^{\text{trans}}(r')\{g_{ww}^{\text{bulk}}(R)\ln g_{ww}^{\text{bulk}}(R) - g_{ww}^{\text{bulk}}(R) + 1\}drdr' \quad (12)$$

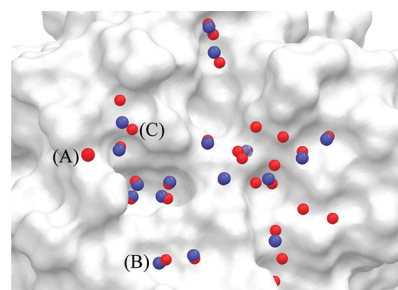
$$S_{ww}^{\text{orient}} = -\frac{1}{2}k\rho^2 \int g_{pw(a)}^{\text{trans}}(r)g_{pw(b)}^{\text{trans}}(r')\{g_{ww}^{\text{bulk}}(R)dR \times \int g_{ww}^{\text{orient}}(\omega)g_{ww}^{\text{orient}}(\omega')\{g_{ww}^{\text{bulk}}(\omega^{\text{rel}}|R)\ln g_{ww}^{\text{bulk}}(\omega^{\text{rel}}|R)\}d\omega d\omega' \quad (13)$$

The water–water correlation functions were calculated from the 8 ns simulation of bulk water, using all available water pairs.

**Table 1.** Effect of the MD Protocol on the Predictions<sup>a</sup>

MD Scheme	free	restrained	fixed
total sites predicted	52	65	78
crystal waters matched (within 1.2 Å)	18	20	21
percentage of predictions correct (%)	34.62	30.77	26.92
percentage crystal waters matched (%)	47.37	52.63	55.26
rmsd of matches (Å)	0.76	0.64	0.62

<sup>a</sup>The effect of the MD protocol on the hydration site clustering and the accuracy with respect to the crystal structure water molecules. The percentage of predictions correct is the percentage of predictions made that are correct. The percentage crystal waters matched is the percentage of the crystal water molecules that were correctly identified.



**Figure 2.** The molecular surface of the MAYM mutant showing the positions of water molecules in the crystal structure and the predicted hydration sites from the restrained protein simulation. The oxygen atoms of the crystal structure water molecules are colored red and the correctly predicted hydration sites are colored blue.

All calculations were performed using the Darwin Supercomputer of the University of Cambridge High Performance Computing Service (<http://www.hpc.cam.ac.uk/>) and were funded by the EPSRC under grant EP/F032773/1. All MD simulations were performed using NAMD compiled for use with CUDA-accelerated GPUs.

## RESULTS

The initial stage of the analysis was to cluster the water molecules from the MD trajectories to identify the hydration sites. To assess the predictions, we compared the positions of the hydration sites to the experimental positions of the oxygen atoms of water molecules from the crystal structure. The experimental sites should represent regions of high water density. We counted the number of predictions where the hydration sites were within 1.2 Å of the crystal structure oxygen atom position. Density was assigned to 38 water molecules in the crystal structure of apo MAYM within 12 Å of the site centroid. Each MD methodology produced a different number of hydration sites. This data can be found in Table 1. The fixed protein simulation predicts the largest number of hydration sites (78) and identifies the largest number of water molecules from the crystal structure.<sup>21</sup> The sites are predicted with an rmsd of 0.62 Å from the crystal structure positions. However, the restrained simulation also performs well, identifying 65 hydration site and 20 water molecules from the crystal structure with an rmsd of 0.64 Å. The correctly predicted hydration sites (blue) and crystal structure water molecules (red) for the restrained simulation are shown in Figure 2. Some water molecules and some hydration sites lie under the surface and thus do not appear in the figure. The water molecules labeled A, B, and

Table 2. Effect of the MD Protocol on Specific Hydration Sites<sup>a</sup>

site	occupancy			$\Delta E$ (kcal mol <sup>-1</sup> )			$-T\Delta S$ (kcal mol <sup>-1</sup> )			$\Delta F$ (kcal mol <sup>-1</sup> )		
	free	rest	fix	free	rest	fix	free	rest	fix	free	rest	fix
A	0.79	0.93	0.98	-0.05	0.94	1.49	0.92	-0.16	1.09	0.87	0.78	2.58
B	NA	0.73	0.85	NA	-1.37	-0.99	NA	0.50	0.77	NA	-0.87	-0.22
C	NA	0.95	0.99	NA	3.65	3.57	NA	0.98	0.41	NA	4.62	3.98
D	NA	0.98	0.97	NA	3.51	5.2	NA	0.68	2.10	NA	4.19	7.30
E	NA	0.94	0.94	NA	2.41	1.01	NA	1.09	2.27	NA	3.50	3.28
F	0.66	0.90	0.96	-1.05	0.06	-0.01	-1.00	-0.16	0.19	-2.05	-0.10	0.18
G	0.53	0.91	0.98	5.40	6.27	7.69	-1.65	0.07	0.33	3.75	6.33	8.02
H	0.68	0.77	0.85	-0.52	-0.14	-0.73	-0.80	-0.50	-0.81	-1.32	-0.65	-1.54
I	0.68	0.82	0.88	1.10	1.65	1.22	0.12	0.99	0.91	1.22	2.64	2.13
J	0.70	0.85	0.95	-0.12	0.27	-0.17	-0.99	1.52	1.79	-1.11	1.79	1.62

<sup>a</sup> The effect of the MD protocol on ten hydration sites on the surface of MAYM for the free, restrained (rest), and fixed (fix) schemes.  $E$  is the interaction energy, and  $F$  is the free energy.

C are in close proximity to neighboring crystal units in the X-ray structure (3.60, 5.27, and 4.74 Å to the closest heavy atoms, respectively), and their positions may thus be affected. The free simulation compares less favorably with the crystal structure, identifying 18 water molecules from the crystal structure with an rmsd of 0.76 Å and 52 hydration sites in total. It is important to note that the two metrics of the number of crystal structure water molecules identified and the rmsd of the water molecules are reliant on assigning X-ray density to specific points, which is an artifact of crystallography.

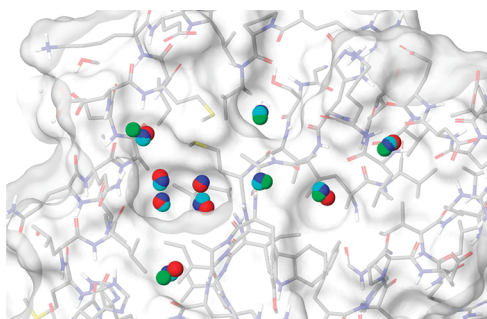
In addition to comparing the positions of the hydration sites with the crystal structure, we calculated the effect of the three schemes on the calculated occupancy and thermodynamic properties of the hydration sites. The results of this analysis can be seen in Table 2, which details the calculated properties of five hydration sites. In general, despite small differences in the number of predicted sites and their position and occupancy, the fixed and restrained schemes agree reasonably well on the majority of the hydration sites. However, the free scheme yields quite different results, with markedly lower occupancies for all the hydration sites. There is also a key disparity that it is interesting to note. When restraints on the protein are removed, the hydrophobic phenylalanine pocket is filled by two methionine residues for a significant portion of the simulation. These two methionine residues form one side of the phenylalanine pocket. This reduces the apparent occupancy of the four water molecules within the pocket to an average of 0.19 in the free simulation. This low occupancy means that they are not identified as hydration sites under the clustering protocol. These four sites have appreciable occupancies of 0.94 and 0.90 on average from the fixed and restrained simulations. This prediction is not completely unexpected, as the opening and closing of hydrophobic pockets on protein surfaces has been observed.<sup>36</sup> Furthermore, these two methionines have relatively high average  $B$ -factors of 15.99 Å<sup>2</sup> and 11.93 Å<sup>2</sup>, suggesting high mobility. Because of the limitations of MD and of crystallography, it is difficult to assess whether the phenylalanine pocket spends an appreciable time in a closed conformation. However, as this clearly affects the MD simulations and the subsequent IFST analysis, it is a very important consideration. If the protein structure is treated as fully flexible, the energy function must be accurate or the predictions of IFST will be misleading. Previous implementations of this

Table 3. Calculated Thermodynamic Properties for the Hydration Site Lying within the Alanine Pocket<sup>a</sup>

MD scheme	free (kcal mol <sup>-1</sup> )	restrained (kcal mol <sup>-1</sup> )	fixed (kcal mol <sup>-1</sup> )
occupancy	0.77	0.82	0.88
$E$ (pw)	-13.41	-13.98	-13.79
$E$ (ww)	-9.11	-7.99	-8.61
$E$ (total)	-22.52	-21.97	-22.40
$\Delta E$	+1.10	+1.65	+1.22
TS (density)	0.68	0.72	0.76
TS (pw, trans)	0.12	0.16	0.26
TS (pw, orient)	1.83	2.03	2.28
TS (pw)	1.94	2.19	2.54
TS (ww, trans)	0.01	0.01	0.01
TS (ww, orient)	0.84	1.42	0.95
TS (ww)	0.85	1.43	0.96
TS (total)	3.47	4.34	4.26
$-T\Delta S$	+0.12	+0.99	+0.91
$\Delta F$	+1.22	+2.64	+2.13

<sup>a</sup> Details on the thermodynamic properties for the hydration site lying within the alanine pocket, calculated using the restrained MD scheme. The protein–water terms are denoted pw, and the water–water terms are denoted ww. The translational contributions are denoted trans, and the orientational contributions are denoted orient.  $E$  is the interaction energy, and  $F$  is the free energy.

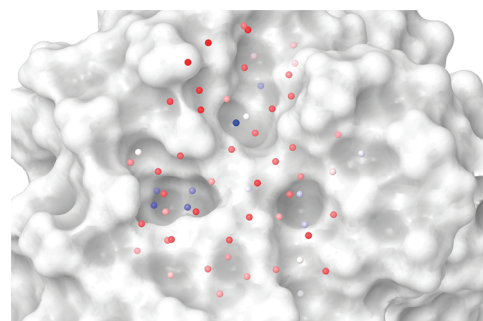
methodology have treated the protein as fixed<sup>18</sup> or as restrained.<sup>21</sup> We predict that this can have a significant effect on the location and occupancies of hydration sites. It also has a significant effect on the calculated thermodynamic properties, as can be seen in Table 2 and Table 3. Table 2 details the interaction energy, entropy, and free energy for the three different MD protocols for ten hydration sites. For many of the hydration sites, the three schemes agree both qualitatively and quantitatively. However, some hydration sites are predicted to have different thermodynamic properties in the three schemes. This is true for sites A and F in Table 2, where the predictions for the free energies vary by 1.80 and 2.23 kcal/mol, respectively. Such a difference impacts the conclusions of the modeling and would affect any quantitative treatment of the results. Table 2 shows that the hydrophobic sites C, D, and G have a free



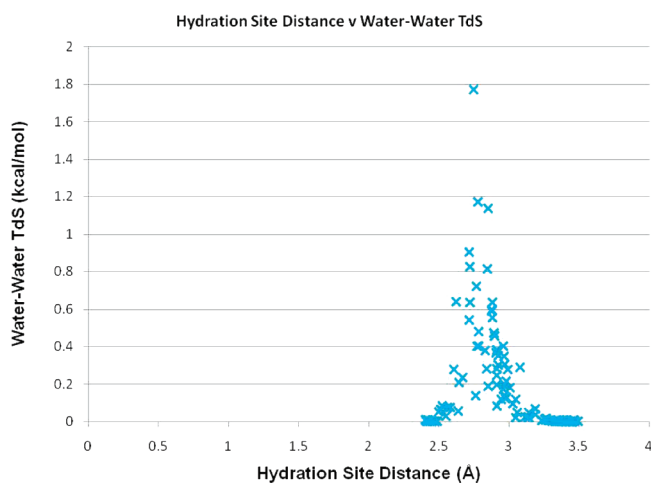
**Figure 3.** The molecular surface of the MAYM mutant showing the positions of ten water molecules in the crystal structure and the predicted hydration sites from the three simulation schemes. The oxygen atoms of the crystal structure water molecules are colored red, the hydration sites from the free simulation are colored green, the hydration sites from the restrained simulation are colored dark blue, and the hydration sites from the fixed simulation are colored cyan.

energy with respect to bulk of +4.62, +4.19, and +6.33 kcal/mol. This agrees very well with previous applications of IFST to hydrophobic sites, where the maximum free energy with respect to bulk was approximately 5 kcal/mol.<sup>20,21</sup> Table 3 provides more specific details on the thermodynamic properties for the hydration site lying within the alanine pocket. The protein–water entropy decreases from the free scheme to the restrained scheme and then to the fixed scheme. This trend occurs throughout the results. Fixing or restraining the protein also restrains the surrounding water molecules, and this has a direct effect on the entropies.

The ten hydration sites shown in Table 2 are illustrated in Figure 3. For the hydration site labeled A, the three schemes agree closely with one another in position and also agree with the crystal structure position. However, the fixed scheme has a markedly different thermodynamic profile from the other schemes. This is due to the increased order in the fixed scheme at this hydration site, with the resulting decreased entropy leading to a less favorable free energy with respect to bulk. Hydration sites B, C, D, and E lie in the phenylalanine pocket and form a conserved square network with few hydrogen bonds per water. This is most marked for hydration sites C and D at the base of the pocket, which have very reduced interaction energies with respect to bulk. However, these hydration sites do not have a high overall entropy with respect to bulk water because of the reduction in water–water correlations in the pocket. Hydration site G lies on the surface of the protein in the same location as the DMSO solvent molecule in the crystal structure. The highly unfavorable free energy for this hydration site may explain why a DMSO molecule is found there in the apo state. Hydration site H also lies on the protein surface above a backbone amide group but is mostly exposed to solvent. It has a more favorable interaction energy than in bulk due to hydrogen bonding, and the reduced water–water correlations at the surface also lead to a favorable entropy with respect to bulk. Displacement of a water molecule from this hydration site by a ligand is predicted to contribute unfavorably to the binding free energy. Formation of a strong hydrogen bond between the ligand and the backbone amide group at this site could lead to a net favorable contribution to the binding free energy whereas a hydrophobic group would lead to a net unfavorable contribution. Hydration site I lies in the alanine pocket, and water molecules within this site have a strong degree of orientational ordering due to the formation of hydrogen bonds with two backbone carbonyls. Hydration site J is on a flat hydrophobic



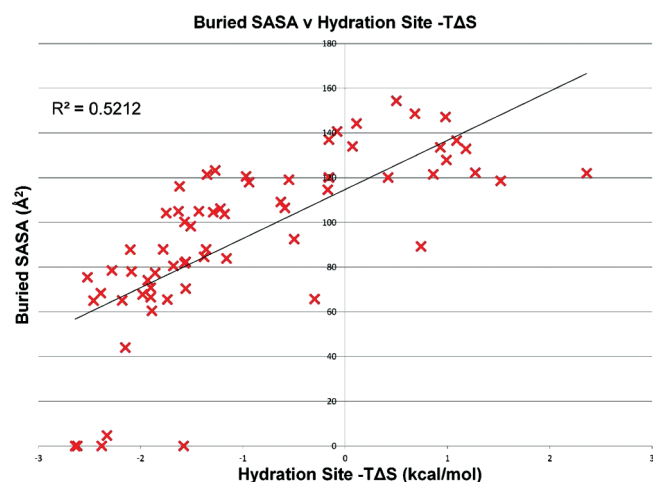
**Figure 4.** The molecular surface of the MAYM mutant showing the predicted hydration sites from the restricted simulation. The hydration sites are colored by free energy with respect to bulk water from more positive in blue to more negative in red.



**Figure 5.** A plot of the distance between two hydration sites against the calculated water–water entropic contribution to the free energy of that site ( $T\Delta S$ ), predicted by the restrained simulation.

surface and makes weak interactions with the protein. However, its overall interaction energy is only 0.27 kcal/mol higher than in bulk water due to favorable interactions with other water molecules. However, these interactions lead to a strong degree of order and unfavorable protein–water entropy (+1.79 kcal/mol) and water–water entropy (+1.77 kcal/mol) terms. The property of increased ordering around hydrophobic solutes to yield favorable interactions has been likened to the formation clathrate cages and has been used previously to explain the hydrophobic effect.<sup>37</sup> The surface of RadA along with the predicted hydration sites from the restricted simulation can be seen in Figure 4. The sites are colored by hydrophobicity from hydrophobic in blue to hydrophilic in red. Such a view has been used previously to study protein binding sites and to explain binding affinity and selectivity.<sup>22,23</sup> Here it can be used to identify binding hotspot regions and provide a quantitative comparison. The phenylalanine and alanine pockets are clearly visible with blue hydrophobic sites on the left- and right-hand sides, respectively.

As well as studying the effect of the three simulation schemes, we have also considered the effect of other computational parameters in the IFST methodology. In this study we only considered water–water entropies for pairs of hydration sites up to 3.5 Å apart, because of the high computational cost of considering a large number of pairs. We thus looked at the correlation in the water–water pair distance and the water–



**Figure 6.** A plot of the change in SASA when a carbon atom is placed at each of the 56 hydration sites predicted by the restrained simulation against the calculated total  $-T\Delta S$  of that site with respect to bulk water.

water pair entropy. A graph of the water–water pair distance against the water–water pair entropy for the restrained scheme can be seen in Figure 5. Because of the dependence on the radial distribution function in bulk, the significant pair entropies are found when the distance between the hydration sites is similar to the maximum in the radial distribution function (2.7 Å). No significant pair entropies are found for hydration sites separated by more than 3.2 Å using this methodology. The majority of the pair entropies result from the orientational term, with the largest translational term being only 0.006 kcal/mol. With sufficient data, it would be very instructive to repeat this calculation without the approximations to the correlation functions.

As a final test, we also calculated the change in solvent-accessible surface area ( $\Delta S_{ASA}$ ) of a carbon atom placed at the centroid of each hydration site. The  $\Delta S_{ASA}$  upon binding is commonly employed as an estimate of the contribution of the hydrophobic effect to binding, so we were interested in how it correlates with the thermodynamic properties of the hydration sites. Figure 6 shows the plot of  $\Delta S_{ASA}$  against the entropic contribution to the free energy ( $-T\Delta S$ ) for all 65 hydration sites in the restrained simulation. The coefficient of determination between  $\Delta S_{ASA}$  and  $-T\Delta S$  is 0.52, suggesting a reasonable correlation, with buried sites tending to have more negative entropies and thus more unfavorable contributions to the free energies. The coefficients of determination for  $\Delta S_{ASA}$  with the interaction energy (0.06) and the total free energy (0.31) were not as high. The  $\Delta S_{ASA}$  for a shape comprised of all 65 hydration spheres was 2167.46, and the sum of the entropic contributions to the free energies for the 65 sites was 62.14 kcal/mol. This corresponds to a value of 28.67 cal/mol/Å<sup>2</sup>, which is consistent with previous estimates used in MMPBSA (38) and MMGBSA (39) of between 5.0 and 50.0 cal/mol/Å<sup>2</sup>.

In summary, the results of this study highlight the importance of the molecular dynamics scheme on the results of IFST and illustrate how the predictions from IFST can be used to understand the thermodynamics of hydration at a protein surface.

## DISCUSSION

This paper describes the application of IFST to a prototypical PPI surface. In particular, we studied the effect of freezing or

restraining the protein structure during the simulation. This approximation has been applied previously, and we were interested in the effect. The free, fixed, and restrained schemes perform comparably in terms of correctly predicting the location of water molecules in the crystal structure. The fixed and restrained schemes identify the primary hotspot in the phenylalanine binding site as three hydration sites that are entropically unfavorable and strongly enthalpically unfavorable. However, these sites are not identified in the free simulation, as the protein shifts to close the pocket with two methionine residues. This may be due to inaccuracies in the forcefield, but it may, however, represent a lowly populated state of the apo protein that is incorrectly scored and thus overly populated. It may also be due to incorrect pressure in the MD simulation. Creation of the spherical boundary region and simulation in an NVT ensemble are likely to affect the pressure and the density of the water, which could lead to cavitation. All three schemes predict a secondary hotspot in the alanine binding site and also locate a third hotspot, which is filled by a solvent DMSO molecule in the crystal structure.

In general, the locations of the hydration sites are very similar with the three schemes. However, the results predict that fixing the protein significantly restricts movement of water molecules at the surface, and this impacts the predicted density and thermodynamic properties of the hydration sites. In particular, the protein–water entropies decrease when the protein is frozen, and this leads to less favorable free energies with respect to bulk water. Incorporating at least some protein flexibility into the simulation seems to be very important, and this is consistent with recent implementations of IFST.<sup>20,21</sup> However, the effect of the degree and nature of the restraints have not been fully explored, and this remains as an important task for future work. In particular, quantifying hydration thermodynamics in highly flexible protein regions is a significant challenge but a very important one. The findings of our study also suggest that water–water pair entropies need only be calculated for pairs that are less than 3.5 Å apart for this implementation of IFST, as contributions from more distant pairs were found to be negligible. However, due to the dependence on the radial distribution function in bulk, this may not be true in a more complete treatment of water–water pair correlations and should be investigated in further work. It is also interesting that the degree of burial of a hydration site correlates to some degree with the entropy but not with the interaction energy. This suggests that the surface area term of MMGBSA and MMPBSA approaches to calculating binding free energy captures some aspects of solvent entropy changes.

IFST is one of the most important methods to quantify solvent thermodynamics, and it has numerous important potential applications. As shown here, it is ideally suited to scanning a protein surface to locate binding hotspots, and it can also be used to predict PPI surfaces on proteins of unknown function. When combined with a scoring function to compute protein–ligand interactions, it can also be applied to molecular docking and the computation of protein–ligand binding affinities.<sup>21,35</sup> This also allows it to be applied to molecular design algorithms for increasing binding affinities. However, in common with other methodologies that utilize MD, this method is highly sensitive to implementation details. This work details one aspect of the implementation that is very important and suggests a number of others. The utility of the method depends on using accurate forcefields, water models, restraints, and simulation parameters. However, the potential of IFST to greatly improve prediction of protein–ligand binding affinities makes the development of this method a very important goal of computational modeling.

## AUTHOR INFORMATION

## Corresponding Author

\*E-mail: djh210@cam.ac.uk.

## ACKNOWLEDGMENT

The authors thank Marko Hyvonen, Tom Blundell, Bracken King, Nate Silver, Duncan Scott, Chris Abell, Ashok Venkitaraman, and John Skidmore for helpful discussions. We also thank Stuart Rankin for technical help and the Cambridge HPCS for use of the CUDA-accelerated GPUs. We are grateful for financial support from the MRC, Wellcome Trust, and EPSRC. We also acknowledge financial support from the Wellcome Trust Translation Award GR080083 (2006-2010), as the structural work from this paper builds upon work from that project.

## REFERENCES

- (1) Shoemaker, B. A.; Panchenko, A. R. Deciphering protein-protein interactions. Part I. Experimental techniques and databases. *PLoS Comput. Biol.* **2007**, *3* (3), 337–344.
- (2) Massova, I.; Kollman, P. A. Computational alanine scanning to probe protein–protein interactions: A novel approach to evaluate binding free energies. *J. Am. Chem. Soc.* **1999**, *121* (36), 8133–8143.
- (3) Shoemaker, B. A.; Panchenko, A. R. Deciphering protein-protein interactions. Part II. Computational methods to predict protein and domain interaction partners. *PLoS Comput. Biol.* **2007**, *3* (4), 595–601.
- (4) Wells, J. A.; McClendon, C. L. Reaching for high-hanging fruit in drug discovery at protein-protein interfaces. *Nature* **2007**, *450* (7172), 1001–1009.
- (5) Tsai, C. J.; Lin, S. L.; Wolfson, H. J.; Nussinov, R. Studies of protein-protein interfaces: A statistical analysis of the hydrophobic effect. *Protein Sci.* **1997**, *6* (1), 53–64.
- (6) West, S. C. Molecular views of recombination proteins and their control. *Nat. Rev. Mol. Cell Biol.* **2003**, *4* (6), 435–445.
- (7) Davies, A. A.; Masson, J. Y.; Mcllwraith, M. J.; Stasiak, A. Z.; Stasiak, A.; Venkitaraman, A. R.; West, S. C. Role of BRCA2 in control of the RAD51 recombination and DNA repair protein. *Mol. Cell* **2001**, *7* (2), 273–282.
- (8) Yang, S. X.; Yu, X.; Seitz, E. M.; Kowalczykowski, S. C.; Egelman, E. H. Archaeal RadA protein binds DNA as both helical filaments and octameric rings. *J. Mol. Biol.* **2001**, *314* (5), 1077–1085.
- (9) Shin, D. S.; Pellegrini, L.; Daniels, D. S.; Yelent, B.; Craig, L.; Bates, D.; Yu, D. S.; Shivji, M. K.; Hitomi, C.; Arvai, A. S.; Volkman, N.; Tsuruta, H.; Blundell, T. L.; Venkitaraman, A. R.; Tainer, J. A. Full-length archaeal Rad51 structure and mutants: mechanisms for RAD51 assembly and control by BRCA2. *EMBO J.* **2003**, *22* (17), 4566–4576.
- (10) Conway, A. B.; Lynch, T. W.; Zhang, Y.; Fortin, G. S.; Fung, C. W.; Symington, L. S.; Rice, P. A. Crystal structure of a Rad51 filament. *Nat. Struct. Mol. Biol.* **2004**, *11* (8), 791–796.
- (11) Rajendra, E.; Venkitaraman, A. R. Two modules in the BRC repeats of BRCA2 mediate structural and functional interactions with the RAD51 recombinase. *Nucleic Acids Res.* **2010**, *38* (1), 82–96.
- (12) Pellegrini, L.; Yu, D. S.; Lo, T.; Anand, S.; Lee, M.; Blundell, T. L.; Venkitaraman, A. R. Insights into DNA recombination from the structure of a RAD51-BRCA2 complex. *Nature* **2002**, *420* (6913), 287–93.
- (13) Lazaridis, T.; Paulaitis, M. E. Entropy of Hydrophobic Hydration - a New Statistical Mechanical Formulation. *Fluid Phase Equilib.* **1993**, *83*, 43–49.
- (14) Lazaridis, T.; Karplus, M. Orientational correlations and entropy in liquid water. *J. Chem. Phys.* **1996**, *105* (10), 4294–4316.
- (15) Lazaridis, T. Solvent reorganization energy and entropy in hydrophobic hydration. *J. Phys. Chem. B* **2000**, *104* (20), 4964–4979.
- (16) Li, Z.; Lazaridis, T. Thermodynamic contributions of the ordered water molecule in HIV-1 protease. *J. Am. Chem. Soc.* **2003**, *125* (22), 6636–6637.
- (17) Lazaridis, T. Inhomogeneous fluid approach to solvation thermodynamics. I. Theory. *J. Phys. Chem. B* **1998**, *102* (18), 3531–3541.
- (18) Li, Z.; Lazaridis, T. Thermodynamics of buried water clusters at a protein–ligand binding interface. *J. Phys. Chem. B* **2006**, *110* (3), 1464–1475.
- (19) Li, Z.; Lazaridis, T. The effect of water displacement on binding thermodynamics: concanavalin A. *J. Phys. Chem. B* **2005**, *109* (1), 662–70.
- (20) Young, T.; Abel, R.; Kim, B.; Berne, B. J.; Friesner, R. A. Motifs for molecular recognition exploiting hydrophobic enclosure in protein–ligand binding. *Proc. Natl. Acad. Sci. U.S.A.* **2007**, *104* (3), 808–13.
- (21) Abel, R.; Young, T.; Farid, R.; Berne, B. J.; Friesner, R. A. Role of the active-site solvent in the thermodynamics of factor Xa ligand binding. *J. Am. Chem. Soc.* **2008**, *130* (9), 2817–31.
- (22) Beuming, T.; Farid, R.; Sherman, W. High-energy water sites determine peptide binding affinity and specificity of PDZ domains. *Protein Sci.* **2009**, *18* (8), 1609–19.
- (23) Huggins, D. J.; McKenzie, G.; Robinson, D.; Narváez, A.; Hardwick, B.; Roberts-Thomson, M.; Venkitaraman, A.; Grant, G.; Payne, M., Computational Analysis of Phosphopeptide Binding to the Polo-Box Domain of the Mitotic Kinase PLK1 Using Molecular Dynamics Simulation. *PLoS Comput. Biol.* **2010**, *6* (8).
- (24) Czapiewski, D.; Zielkiewicz, J. Structural properties of hydration shell around various conformations of simple polypeptides. *J. Phys. Chem. B* **2010**, *114* (13), 4536–50.
- (25) Phillips, J. C.; Braun, R.; Wang, W.; Gumbart, J.; Tajkhorshid, E.; Villa, E.; Chipot, C.; Skeel, R. D.; Kale, L.; Schulten, K. Scalable molecular dynamics with NAMD. *J. Comput. Chem.* **2005**, *26* (16), 1781–1802.
- (26) Chen, L. T.; Ko, T. P.; Chang, Y. C.; Lin, K. A.; Chang, C. S.; Wang, A. H. J.; Wang, T. F. Crystal structure of the left-handed archaeal RadA helical filament: identification of a functional motif for controlling quaternary structures and enzymatic functions of RecA family proteins. *Nucleic Acids Res.* **2007**, *35* (6), 1787.
- (27) Humphrey, W.; Dalke, A.; Schulten, K. VMD: Visual molecular dynamics. *J. Mol. Graph.* **1996**, *14* (1), 33–8.
- (28) MacKerell, A. D.; Bashford, D.; Bellott, M.; Dunbrack, R. L.; Evanseck, J. D.; Field, M. J.; Fischer, S.; Gao, J.; Guo, H.; Ha, S.; Joseph-McCarthy, D.; Kuchnir, L.; Kuczera, K.; Lau, F. T. K.; Mattos, C.; Michnick, S.; Ngo, T.; Nguyen, D. T.; Prodhom, B.; Reiher, W. E.; Roux, B.; Schlenkrich, M.; Smith, J. C.; Stote, R.; Straub, J.; Watanabe, M.; Wiorkiewicz-Kuczera, J.; Yin, D.; Karplus, M. All-atom empirical potential for molecular modeling and dynamics studies of proteins. *J. Phys. Chem. B* **1998**, *102* (18), 3586–3616.
- (29) Mackerell, A. D.; Feig, M.; Brooks, C. L. Extending the treatment of backbone energetics in protein force fields: Limitations of gas-phase quantum mechanics in reproducing protein conformational distributions in molecular dynamics simulations. *J. Comput. Chem.* **2004**, *25* (11), 1400–1415.
- (30) Abascal, J. L. F.; Vega, C., A general purpose model for the condensed phases of water: TIP4P/2005. *J. Chem. Phys.* **2005**, *123* (23), -.
- (31) Grubmüller, H. *Solvate: A Program to Create Atomic Solvent Models*. 1996.
- (32) Brooks, B. R.; Bruccoleri, R. E.; Olafson, B. D.; States, D. J.; Swaminathan, S.; Karplus, M. Charmm - a Program for Macromolecular Energy, Minimization, and Dynamics Calculations. *J. Comput. Chem.* **1983**, *4* (2), 187–217.
- (33) Darden, T.; York, D.; Pedersen, L. Particle Mesh Ewald - an N. Log(N) Method for Ewald Sums in Large Systems. *J. Chem. Phys.* **1993**, *98* (12), 10089–10092.
- (34) Brooks, C. L., III; Brunger, A.; Karplus, M. Active site dynamics in protein molecules: a stochastic boundary molecular-dynamics approach. *Biopolymers* **1985**, *24* (5), 843–65.
- (35) Abel, R.; Wang, L.; Friesner, R. A.; Berne, B. J. A Displaced-Solvent Functional Analysis of Model Hydrophobic Enclosures. *J. Chem. Theory Comput.* **2010**, *6* (9), 2924–2934.
- (36) Eyrich, S.; Helms, V. Transient pockets on protein surfaces involved in protein–protein interaction. *J. Med. Chem.* **2007**, *50* (15), 3457–3464.

(37) Head-Gordon, T. Is Water-Structure around Hydrophobic Groups Clathrate-Like. *Proc. Natl. Acad. Sci. U.S.A.* **1995**, *92* (18), 8308–8312.

(38) Huggins, D. J.; Altman, M. D.; Tidor, B. Evaluation of an inverse molecular design algorithm in a model binding site. *Proteins* **2009**, *75* (1), 168–186.

(39) Fogolari, F.; Brigo, A.; Molinari, H. Protocol for MM/PBSA molecular dynamics simulations of proteins. *Biophys. J.* **2003**, *85* (1), 159–66.

#### ■ NOTE ADDED AFTER ASAP PUBLICATION

This article was published ASAP on October 5, 2011. Changes have been made to Figure 6 and its caption, and to the penultimate paragraph of the Results section. The correct version was published on November 8, 2011.

# Magnetic Coupling in Transition-Metal Binuclear Complexes by Spin-Flip Time-Dependent Density Functional Theory

Rosendo Valero,<sup>\*,†</sup> Francesc Illas,<sup>‡</sup> and Donald G. Truhlar<sup>§</sup><sup>†</sup>Research Unit "Molecular Physical Chemistry", University of Coimbra, Rua Larga, 3004-535 Coimbra, Portugal<sup>‡</sup>Departament de Química Física and Institut de Química Teòrica i Computacional (IQTCUB), Universitat de Barcelona, C/Martí i Franquès 1, E-08028 Barcelona, Spain<sup>§</sup>Department of Chemistry and Supercomputing Institute, University of Minnesota, Minneapolis, Minnesota 55455-0431, United States

**ABSTRACT:** Spin-flip time-dependent density functional theory (SF-TDDFT) has been applied to predict magnetic coupling constants for a database of 12 spin-1/2 homobinuclear transition-metal complexes previously studied by Phillips and Peralta employing spin-projected broken-symmetry density functional theory (Phillips, J. J.; Peralta, J. E. *J. Chem. Phys.* **2011**, *134*, 034108). Several global hybrid density functionals with a range of percentages of Hartree–Fock exchange from 20% to 100% have been employed within the collinear-spin formalism, and we find that both the high-spin reference state and low-spin state produced by SF-TDDFT are generally well adapted to spin symmetry. The magnetic coupling constants are calculated from singlet–triplet energy differences and compared to values arising from the popular broken-symmetry approach. On average, for the density functionals that provide the best comparison with experiment, the SF-TDDFT approach performs as well as or better than the spin-projected broken-symmetry strategy. The constrained density functional approach also performs quite well. The SF-TDDFT magnetic coupling constants show a much larger dependence on the percentage of Hartree–Fock exchange than on the other details of the exchange functionals or the nature of the correlation functionals. In general, SF-TDDFT calculations not only avoid the ambiguities associated with the broken-symmetry approach, but also show a considerably reduced systematic deviation with respect to experiment and a larger antiferromagnetic character. We recommend MPW1K as a well-validated hybrid density functional to calculate magnetic couplings with SF-TDDFT.

## 1. INTRODUCTION

The synthesis and study of bi- and polynuclear transition-metal complexes has been motivated in part by the remarkable magnetic properties they often exhibit, ultimately leading to what has been called single-molecule magnets.<sup>1–10</sup> Magnetic molecules have a nonzero total spin and other properties that make them suitable for potential technological applications such as high-density information storage and quantum computing.<sup>11–13</sup> In transition-metal complexes, the metal atoms may act as paramagnetic centers with effective localized spin moments,  $S_i$ , where  $i$  identifies the atom on which the spin is localized, that interact with each other ferro-, ferri-, or antiferromagnetically. Experimental measurements of magnetic susceptibilities versus temperature or neutron diffraction, among other techniques, permit one to study the lower lying electronic states of magnetic systems. In many cases, the experimental data can be interpreted by describing the magnetic interactions with the isotropic Heisenberg–Dirac–Van Vleck (HDV) Hamiltonian, which for a binuclear complex takes the form<sup>14,15</sup>

$$\hat{H} = -J \mathbf{S}_1 \cdot \mathbf{S}_2 \quad (1)$$

where  $J$  represents the phenomenological magnetic exchange coupling between the two magnetic centers. A positive sign for  $J$  corresponds to ferromagnetic coupling and a negative sign to antiferromagnetic coupling. The HDV Hamiltonian is appropriate for the physical description of magnetic coupling in a wide variety of systems, including some organic biradicals, transition-metal complexes, and ionic solids.<sup>15</sup>

The prediction of magnetic couplings from first principles is a very important albeit difficult task, since this property depends strongly on a balanced treatment of electron exchange and electron correlation effects. Furthermore, both nondynamical and dynamical correlation have to be described accurately, as evidenced in a number of studies<sup>16–21</sup> that have employed high-level wave function methods such as difference-dedicated configuration interaction (DDCI)<sup>22</sup> and multiconfigurational second-order perturbation theory (e.g., CASPT2<sup>23,24</sup>). However, the size and complexity of transition-metal clusters of chemical interest preclude in most cases the use of these potentially accurate but computationally demanding electronic structure methods. In recent years, density functional theory (DFT) has emerged as a robust and practical electronic structure method in quantum chemistry and solid-state physics. DFT is formally a theory designed for the ground electronic state, which in the Kohn–Sham formulation is represented by a single Slater determinant formed by orbitals of a fictitious noninteracting system obtained by solving pseudoeigenvalue equations.

For two spin-1/2 centers, the HDV Hamiltonian in eq 1 has one triplet and one singlet eigenstate, with eigenvalues (energies) equal to  $-J/4$  and  $+3J/4$ , respectively. In this case, the magnetic coupling can be obtained by simply mapping the lowest triplet and singlet electronic states to the HDV eigenstates, and it is given by the difference between the energies of these states. For wave function methods, this mapping can be carried out by

Received: June 9, 2011

Published: September 09, 2011

expressing the electronic state functions as expansions in configuration state functions having a well-defined spin symmetry.<sup>15</sup> In DFT, one way to calculate the splitting of two levels is to take the difference in energy of separate self-consistent-field (SCF) calculations on the two spin states, which is called the  $\Delta$ SCF approach. In the unrestricted Kohn–Sham formalism, the lowest triplet state, with energy  $E(\text{HS})$ , is in most cases approximately well represented by a single Slater determinant, but if we make an analogy with wave function theory, the lowest singlet state would require a spin-adapted linear combination of at least two determinants.<sup>25</sup> Noodleman<sup>26–28</sup> advocated a workaround to this problem that involves converging to a broken-symmetry (BS) solution of the Kohn–Sham equations such that the two spins are localized at the two centers. The BS determinant has neither singlet nor triplet spin symmetry, and a strategy must be adopted to relate its energy,  $E(\text{BS})$ , to that of the relevant singlet state. At this point there are two extreme approaches that can be called spin-unprojected and spin-projected. In the spin-unprojected approach,<sup>29,30</sup> one assumes that the energy of the BS state is an approximation of the energy of the open-shell singlet state, in which case the magnetic coupling would be obtained as

$$J = E(\text{BS}) - E(\text{HS}) \quad (2)$$

whereas in the spin-projected approach,<sup>26–28</sup> one assumes that the BS state is a weighted average of spin states; in this case that would be an equal mixture of singlet and triplet, which yields

$$J = 2(E(\text{BS}) - E(\text{HS})) \quad (3)$$

These two equations can be seen as limiting cases of the weighted-average formula proposed by Yamaguchi and co-workers,<sup>31–33</sup> in particular

$$J = \frac{2[E(\text{BS}) - E(\text{HS})]}{\langle S^2 \rangle_{\text{HS}} - \langle S^2 \rangle_{\text{BS}}} \quad (4)$$

where  $\langle S^2 \rangle_{\text{HS}}$  and  $\langle S^2 \rangle_{\text{BS}}$  are the expectations of the square of the spin angular momentum for the HS and BS solutions, respectively. If we assume that the high-spin unrestricted DFT determinant is a good approximation of the triplet, then  $\langle S^2 \rangle_{\text{HS}} = 2$ , and if the BS determinant is a good approximation of the singlet state, as assumed in the spin-unprojected approach, then  $\langle S^2 \rangle_{\text{BS}} \approx 0$ , and eq 4 reduces to eq 2. In practice, this would be achieved if the two centers were strongly coupled, so that the  $\alpha$  and  $\beta$  spins are covalently shared (not localized on the two centers at all); hence, this may be called the strong interaction limit. Conversely, in the weak interaction limit in which the two spin orbitals are completely localized, the BS determinant is a 50:50 mixture of pure singlet and triplet states,  $\langle S^2 \rangle_{\text{BS}} \approx 1$ , and eq 4 reduces to eq 3.<sup>34</sup> For unrestricted Hartree–Fock wave functions, this argument is straightforward, but in the case of DFT,  $\langle S^2 \rangle$  is not rigorously defined, and one can argue that spin symmetry does not have to be respected, although this argument presents conceptual complications.<sup>35</sup> In the present study we will compare the weighted-average BS approach to spin-flip time-dependent density functional theory<sup>36–39</sup> (SF-TDDFT), where the spin symmetry is less ambiguous. To anticipate the results, we will find that spin projection is required for consistency between the BS approach and the SF-TDDFT approach.

In section 2 we introduce the computational methods and motivate their use for the calculation of magnetic couplings. In section 3 we present the details of the database of transition-metal complexes employed. Section 4 contains the main results of the study and their relation to previous work. Finally, section 5 draws conclusions.

## 2. COMPUTATIONAL METHODOLOGY

Note that only singlet–triplet and doublet–quartet pairs of states can be studied in the present implementation of SF-TDDFT, and here we restrict our attention to the case of singlet–triplet splittings of systems with two spin-1/2 centers, each of which is an identical transition-metal atom, but with its own set of nonmagnetic ligands. Symmetry-adapting localized states leads to a spatially symmetric state and a spatially anti-symmetric state,  $\phi^+$  and  $\phi^-$ , respectively, and these are simply the sum and difference of singly occupied particle or hole states centered on the transition-metal atoms. These states can be used to form three singlet states and three triplet states:<sup>25,40</sup>

$${}^1\Gamma_1 \propto \left( \lambda \phi^+ \bar{\phi}^+ - \sqrt{1 - \lambda^2} \phi^- \bar{\phi}^- \right) \quad (5)$$

$${}^1\Gamma_2 \propto \left( \lambda \phi^+ \bar{\phi}^+ + \sqrt{1 - \lambda^2} \phi^- \bar{\phi}^- \right) \quad (6)$$

$${}^1\Gamma_3 \propto (\phi^+ \bar{\phi}^- - \phi^- \bar{\phi}^+) \quad (7)$$

$${}^3\Gamma_1 \propto (\phi^+ \phi^-) \quad (8)$$

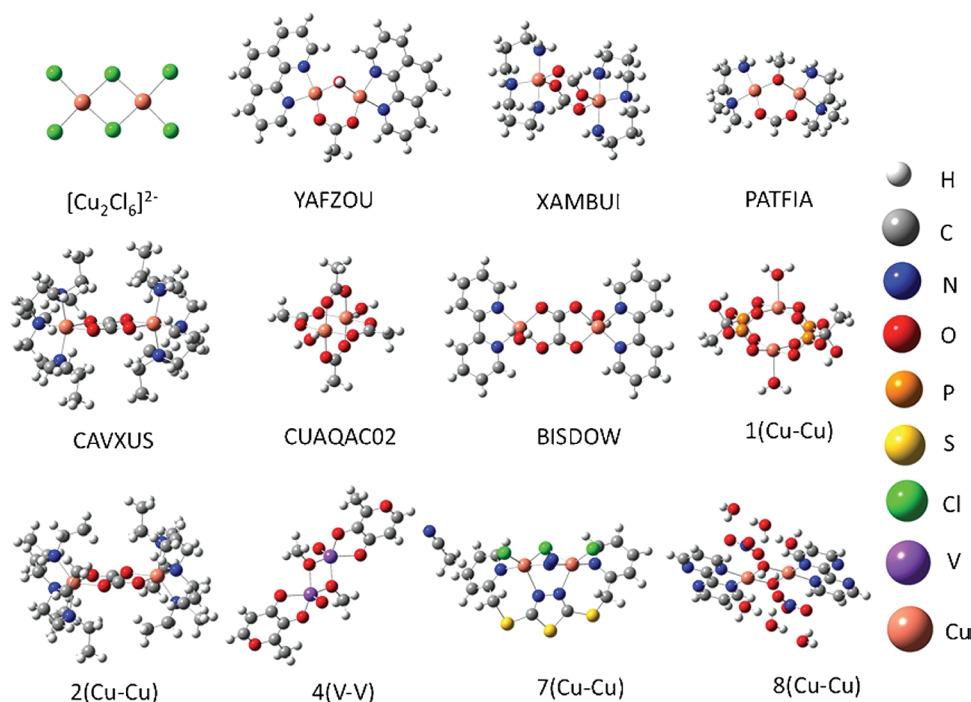
$${}^3\Gamma_2 \propto (\bar{\phi}^+ \bar{\phi}^-) \quad (9)$$

$${}^3\Gamma_3 \propto (\phi^+ \bar{\phi}^- + \phi^- \bar{\phi}^+) \quad (10)$$

where a spin orbital without or with an overbar has  $\alpha$  or  $\beta$  spin, respectively, and each product of spin orbitals should be interpreted as shorthand for a determinant. The three singlet spin states are each composed of two determinants; the first two triplet wave functions, eqs 8 and 9, which have  $M_S = 1$  and  $-1$ , respectively, are singly determinantal, and the third triplet spin state, eq 10 with  $M_S = 0$ , consists of two determinants. The lowest singlet state of two weakly coupled centers is eq 5, and the 3-fold-degenerate lowest triplet state is given by eq 8, 9, or 10.

In time-dependent DFT (TDDFT),<sup>41–47</sup> a time-dependent perturbation is added to the ground-state Hamiltonian and the poles of the response function are the frequencies of the allowed excitations (here “excitations” include both excitations and de-excitations, if any), thereby yielding the energies of the excited states. The perturbation is assumed small enough for the response function to be in the linear regime; furthermore, the dependence of the exchange-correlation potential on the frequency of the excitation is ignored—which is called the adiabatic approximation and entails using the ground-state exchange-correlation potentials. In the conventional formulation of TDDFT, which may be called the low-spin formulation or LS-TDDFT, the state before the perturbation is a closed-shell singlet. With these approximations, only single excitations from the closed-shell reference state can be captured by the formalism, and in particular the state of eq 5 cannot be obtained in LS-TDDFT.<sup>41</sup> However, in SF-TDDFT,<sup>36–39</sup> which may also be called high-spin TDDFT, the triplet wave function in eq 8 is taken as a reference, and both an excitation and a spin flip are applied to obtain both determinants in eq 5. The same kind of process also yields the triplet state with  $M_S = 0$  in eq 10. Thus, to represent the triplet state, there would be, in principle, two options, namely, to take the reference state with  $M_S = 1$  or to take the  $M_S = 0$  component of the triplet state generated in the spin-flip excitation process; this is a choice that also shows up when





**Figure 1.** Binuclear complexes studied in this investigation.

using configuration interaction (CI) wave functions. In principle, these two states should be exactly degenerate (i.e., the excitation energy would be zero). In practical calculations, there is an energy difference between these two components of the triplet state, and this energy difference has been called the self-splitting test.<sup>48</sup> It is supposed to be a measure of the consistency of the SF-TDDFT method. We follow here the usual procedure of calculating the singlet–triplet splitting from the  $M_S = 0$  component of the triplet because eq 10 generated by the response is more consistent than the reference (eq 8) for comparing with the generated eq 5.<sup>49,50</sup> In general, one expects more accurate energy splittings when one compares parallel calculations than when one compares disparate ones. In fact, we find that computing magnetic couplings using the  $M_S = 1$  results for the triplet leads to meaningless values. Because the use of the  $M_S = 0$  state is the standard approach in the literature, results obtained by this method are labeled SF-TDDFT, as usual.

Applications of the spin-flip strategy in the literature include the study of the singlet–triplet splitting and diradical character of organic systems,<sup>51–55</sup> bioinorganic chemistry,<sup>56–59</sup> conical intersections,<sup>60,61</sup> and electron transfer couplings.<sup>62–64</sup> To our knowledge, there are only two previous studies of magnetic splittings like those considered here with a formalism equivalent to SF-TDDFT, namely, the recent work of Ziegler and co-workers,<sup>50</sup> where they apply their spin-flip constricted variational DFT formalism to the study of trinuclear copper complexes, and the even more recent work of Yang et al.<sup>49</sup> on low-spin–high-spin splittings in p-block atoms.

The calculations in this study have been carried out with a collinear formulation<sup>36</sup> of SF-TDDFT and within the Tamm–Dancoff<sup>61,65,66</sup> approximation, as implemented in the Q-Chem program.<sup>67</sup> Within the collinear approach, only the Hartree–Fock exchange part of the exchange–correlation functional contributes to the SF coupling.<sup>36</sup> Therefore, only hybrid functionals can be employed in this formulation. For the triplet reference

state, the calculations were performed with a grid composed of 120 radial points and 302 Lebedev angular points.

The density functionals used in the present calculations and in those included in the tabular comparisons made in section 4) are all global hybrids (in a global hybrid, the percentage  $X$  of Hartree–Fock exchange is the same for all interelectronic distances) of the hybrid generalized gradient approximation (GGA) type and of the hybrid meta-GGA type. In particular, the functionals studied, in order of increasing  $X$ , are as follows:

- the popular B3LYP<sup>68–70</sup> hybrid GGA functional (20%)
- the Minnesota M06<sup>71,72</sup> hybrid meta-GGA (27%)
- PBE35 (a modification, with  $X = 35$ , of the PBE0<sup>73,74</sup> functional, which itself is a hybrid version ( $X = 25$ ) of the Perdew–Burke–Ernzerhof (PBE)<sup>75</sup> GGA)
- B3LYP40 (B3LYP with  $X = 40$ )
- B1LYP40 (the B1LYP<sup>76</sup> one-parameter hybrid GGA with the percentage of Hartree–Fock exchange raised from 25% to 40%)
- B1PW40 (the B1PW91<sup>76</sup> one-parameter hybrid GGA with the percentage of Hartree–Fock exchange raised from 25% to 40%)
- BMK<sup>77</sup> (the Boese–Martin model for kinetics hybrid meta-GGA with 42% Hartree–Fock exchange)
- MPW1K<sup>78</sup> (the modified Perdew–Wang one-parameter model for kinetics hybrid GGA with  $X = 42.8$ )
- B3LYP54 (B3LYP with 54% Hartree–Fock exchange)
- the Minnesota M06-2X<sup>71,72</sup> (54%) and M06-HF<sup>72,79</sup> (100%) hybrid meta-GGA functionals

In addition, a few functionals to which we compare will be explained in section 4.

### 3. TRANSITION-METAL COMPLEX DATABASE

We have chosen a database of 12 bimolecular transition-metal complexes, each containing two spin-1/2 metal centers,

Table 1. Magnetic Couplings ( $\text{cm}^{-1}$ ) for the 12 Transition-Metal Complexes Studied at the SF-TDDFT Level

system	B3LYP ( $X = 20$ )	M06 ( $X = 27$ )	B3LYP40 <sup>a</sup> ( $X = 40$ )	B1LYP40 ( $X = 40$ )	B1PW40 ( $X = 40$ )	BMK ( $X = 42$ )	MPW1K ( $X = 42.8$ )	B3LYP54 ( $X = 54$ )	M06-2X ( $X = 54$ )	M06-HF ( $X = 100$ )	exptl
$\text{Cu}_2\text{Cl}_6^{2-}$	-342	-210	-94 (-44)	-94	-88	-59	-69	-15	-13	35	0 to -94
YAFZOU	96	91	76 (73)	75	76	78	73	60	65	90	111
XAMBUI	-1	1	1 (-1)	1	2	1	1	1	0	3	2
PATFIA <sup>b</sup>	-399	-198	-98 (-98)	-98	-98	-71	-81	-35	-22	32	-11
	[-420]	[-140]	[-100]	[-101]	[-105]	[-86]	[-86]	[-37]	[-24]	[26]	
CAVXUS	-65	-37	-15 (-19)	-15	-15	-14	-12	-7	-8	-2	-19
CUAQAC02	-721	-523	-245 (-245)	-244	-249	-214	-215	-123	-129	-31	-286
BISDOW	-1126	-743	-349 (-343)	-347	-354	-306	-306	-173	-181	-40	-382
1(Cu-Cu)	-426	-226	-74 (-57)	-74	-74	-59	-60	-27	-29	-8	-62
2(Cu-Cu)	-397	-280	-115 (-100)	-114	-115	-105	-98	-53	-56	-8	-75
4(V-V)	-444	-290	-159 (-154)	-159	-166	-156	-147	-85	-81	-3	-214
7(Cu-Cu)	-421	-145	177 (164)	183	181	159	198	198	194	143	168
8(Cu-Cu)	35	79	111 (102)	110	112	108	112	103	104	72	114
MSE <sup>c</sup>	-292	-147	-3 (-2)	-3	-4	7	11	47	46	82	
MUE <sup>c</sup>	292	147	29 (29)	30	28	32	35	62	59	98	
RMSE <sup>c</sup>	374	190	39 (39)	39	37	43	46	93	89	147	
MURE <sup>c</sup>	5.3	2.6	0.9 (1.0)	0.9	0.9	0.7	0.8	0.6	0.5	1.0	

<sup>a</sup>The values in parentheses are the results obtained with the def2-QZVPPD basis set for the metallic centers. <sup>b</sup>The values in brackets are the results obtained with the full PATFIA model. <sup>c</sup>The mean errors are calculated over rows YAFZOU through 8(Cu-Cu), as explained in the second paragraph of section 4.

to study the performance of SF-TDDFT in the calculation of magnetic couplings. The motivations for choosing this database are twofold: First, since this database was used by Phillips and Peralta<sup>80</sup> for a study of the performance of range-separated hybrid functionals by the spin-projected broken-symmetry approach, it allows us to make a precise comparison to that approach for a relatively extended set of complexes. Second, this database provides a particularly straightforward way to examine magnetic exchange constants because only singlet-triplet energy differences need to be calculated to obtain  $J$  for this database.

The 12 complexes are illustrated in Figure 1. Eleven of them have Cu(II)  $d^9$  atoms as metal centers, and one has a V(IV)  $d^1$  metal center. The first seven complexes are Cu(II) complexes, to be called here  $\text{Cu}_2\text{Cl}_6^{2-}$ , YAFZOU, XAMBUI, PATFIA, CAVXUS, CUAQAC02, and BISDOW, where the names of the latter six complexes correspond to their Cambridge Structural Database reference codes. Note that these seven cases have been used before to test various approaches to the calculation of magnetic couplings.<sup>81–85</sup> As in previous work, the XAMBUI and PATFIA systems are simplified as compared to the systems for which the spin splittings were measured experimentally; in particular, the ferrocenecarboxylate groups were replaced by formate groups. To test the accuracy of this approximation, the full PATFIA complex has also been used to compute magnetic couplings using the crystallographic data from López et al.<sup>86</sup> The last five complexes studied are four Cu(II) complexes and one V(IV) complex, to be called 1(Cu-Cu), 2(Cu-Cu), 4(V-V), 7(Cu-Cu), and 8(Cu-Cu), as in the work of Phillips and Peralta.<sup>80</sup> These complexes are the complexes of spin-1/2 centers from a larger database used by Rudra et al.<sup>35</sup> and later by Peralta and Melo<sup>87</sup> in their magnetic coupling studies (their database also included complexes with higher spin states). The geometries of all the complexes are taken from their crystallographic structures. The counterions are neglected in all cases.

The Gaussian basis sets employed for the first seven complexes are the same as in ref 83, and those for the last five complexes are taken from ref 87. For the full PATFIA complex, the Los Alamos ECP double- $\zeta$ -type basis set LANL2DZ<sup>88</sup> was employed for the Fe atoms. Additional SF-TDDFT calculations were carried out for the B3LYP40 functional using the larger def2-QZVPPD<sup>89,90</sup> basis set for the transition-metal atoms and the same basis sets as before for the rest of the atoms.

## 4. RESULTS AND DISCUSSION

The main results of the present work are presented in Table 1, where the calculated magnetic couplings for the 12 transition-metal complexes with each of the 10 density functionals are compared with experiment. The magnetic couplings are calculated simply as the difference between the singlet and triplet energies, as one would do in wave function theory. This is justified by the values of the spin-squared operator  $\langle S^2 \rangle$  of the two spin-flip states, eqs 5 and 10, which are presented in Table 2. It can be seen that the values of  $\langle S^2 \rangle$  are in most cases close to the theoretical values of 0.0 and 2.0 for the singlet and the triplet, respectively. The exceptions are the XAMBUI and the 7(Cu-Cu) complexes and for the M06-HF density functional also YAFZOU and PATFIA. In these cases we have taken a pragmatic approach and considered the singlet and triplet to be the states with the lower and higher values of  $\langle S^2 \rangle$ , respectively. The  $\langle S^2 \rangle$  values that show significant deviations from the nominally correct values may be an indication of the inadequacy of some presently available functionals for particular systems, but they might also indicate that  $\langle S^2 \rangle$  cannot always be used as a reliable indicator of the success of a given calculation. A more fundamental reason for this behavior might be a larger multiconfigurational character of a given complex with the Kohn-Sham orbitals differing from the magnetic orbitals.

We have computed four statistical measures of accuracy, namely, the mean signed error (MSE), the mean unsigned error (MUE), the root mean squared error (RMSE), and—following

Table 2. Values of  $\langle S^2 \rangle$  for the Magnetic States Obtained at the SF-TDDFT Level for the 12 Complexes Studied<sup>a</sup>

system	B3LYP	M06	B3LYP40	B1LYP40	B1PW40	MPW1K	BMK	B3LYP54	M06-2X	M06-HF
Cu <sub>2</sub> Cl <sub>6</sub> <sup>2-</sup>	0.03	0.01	0.02	0.02	0.02	0.02	0.01	0.02	0.01	2.01
	1.99	2.01	2.00	2.00	2.00	2.01	2.01	2.01	2.02	0.02
YAFZOU	1.71	2.00	2.00	2.00	2.01	2.01	1.92	2.01	2.00	1.45
	0.30	0.02	0.02	0.02	0.02	0.02	0.10	0.02	0.02	0.58
XAMBUI	0.98	1.73	1.52	1.51	1.49	1.47	1.50	1.52	1.85	1.05
	1.03	0.29	0.50	0.51	0.53	0.55	0.52	0.50	0.17	0.97
PATFIA	0.02	0.01	0.02	0.02	0.02	0.02	0.02	0.02	0.04	1.69
	1.99	2.01	2.01	2.01	2.01	2.01	2.00	2.01	1.98	0.34
CAVXUS	0.03	0.01	0.03	0.03	0.03	0.03	0.02	0.02	0.01	0.45
	1.99	2.01	2.00	2.00	2.00	2.00	2.00	2.01	2.01	1.58
CUAQAC02	0.03	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01
	1.98	2.00	2.00	2.00	2.00	2.00	2.01	2.01	2.01	2.01
BISDOW	0.03	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01
	1.99	2.01	2.01	2.01	2.01	2.01	2.01	2.01	2.01	2.01
1(Cu–Cu)	0.04	0.01	0.02	0.02	0.02	0.02	0.01	0.01	0.01	0.01
	1.98	2.01	2.00	2.00	2.00	2.01	2.01	2.01	2.01	2.01
2(Cu–Cu)	0.03	0.01	0.02	0.02	0.02	0.02	0.01	0.01	0.01	0.02
	1.98	2.01	2.01	2.00	2.00	2.01	2.01	2.01	2.02	2.02
4(V–V)	0.03	0.05	0.05	0.05	0.05	0.06	0.02	0.06	0.05	0.08
	2.03	2.06	2.07	2.07	2.07	2.08	2.03	2.10	2.08	2.13
7(Cu–Cu)	0.09	0.21	1.62	1.60	1.62	1.73	1.95	1.97	2.02	1.62
	1.94	1.83	0.45	0.46	0.44	0.34	0.10	0.20	0.04	0.43
8(Cu–Cu)	1.99	2.01	2.01	2.01	2.01	2.01	2.01	2.01	2.01	2.01
	0.02	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01

<sup>a</sup>The first and second values correspond to the lower and higher energy states, respectively.

Phillips and Peralta<sup>80</sup>—also the mean unsigned relative error (MURE). The definitions of these quantities are as follows:

$$\text{MSE} = \frac{1}{N} \sum_{i=1}^N [J_{\text{calcd},i} - J_{\text{exptl},i}] \quad (11)$$

$$\text{MUE} = \frac{1}{N} \sum_{i=1}^N \text{Abs}[J_{\text{calcd},i} - J_{\text{exptl},i}] \quad (12)$$

$$\text{RMSE} = \sqrt{\frac{\sum_{i=1}^N [J_{\text{calcd},i} - J_{\text{exptl},i}]^2}{N}} \quad (13)$$

$$\text{MURE} = \frac{1}{N} \sum_{i=1}^N \text{Abs} \left[ \frac{J_{\text{calcd},i} - J_{\text{exptl},i}}{J_{\text{exptl},i}} \right] \quad (14)$$

First, note that there is some uncertainty as to the experimental value of  $J$  for the Cu<sub>2</sub>Cl<sub>6</sub><sup>2-</sup> complex. One experimental reference quotes an antiferromagnetic coupling of between 0 and –40 cm<sup>-1</sup>,<sup>91</sup> whereas another experiment suggests a significantly larger value of –94 cm<sup>-1</sup>.<sup>92</sup> There is also a spread in high-level ab initio results for this splitting. Thus, studies carried out with the DDCI approach<sup>21,93,94</sup> have found a magnetic coupling for the experimental structure in agreement with the first experimental interval and with more recent CASPT2 calculations.<sup>17</sup> In contrast, a recent application of state-specific multireference coupled cluster theory with single and double excitations yielded<sup>95</sup> values between –66 and –84 cm<sup>-1</sup>, in better agreement with the

second experimental value (–94 cm<sup>-1</sup>). The B3LYP40 result in Table 1 shows that the magnetic coupling for Cu<sub>2</sub>Cl<sub>6</sub><sup>2-</sup> depends strongly on the basis set used for the Cu atoms. Because of the larger experimental and theoretical uncertainties in this case, we omitted the magnetic couplings of Cu<sub>2</sub>Cl<sub>6</sub><sup>2-</sup> from all calculations of mean errors.

The most conspicuous trend in Table 1 is the strong dependence of the magnetic coupling on the percentage of Hartree–Fock exchange. In general, local functionals (i.e., those without Hartree–Fock exchange) tend to favor the low-spin states, and high- $X$  hybrid functionals tend to favor the high-spin states, probably because Hartree–Fock exchange correlates electrons with the same spin by enforcing the Fermi hole, but this does not apply to electrons with opposing spins. We find, for example, that B3LYP and M06 with  $X = 20$  and 27, respectively, are in most cases strongly antiferromagnetic, and they show a large systematic deviation from experiment. The optimal value of  $X$  is found to be about 40, as in B3LYP40, B1LYP40, and B1PW40, or a little higher, as in MPW1K and BMK. The MSE for all these functionals is very small, on the order of only a few wavenumbers. The best results are found with the B1PW40 density functional, and we will use this as a reference for subsequent comparisons. However, note that the MURE is still relatively large, mainly due to the large deviations (in terms of MURE) found for the PATFIA complex. The MURE becomes less appropriate in cases like this where some of the couplings are very small because these small quantities appear in the denominator of the relative error.

Table 1 also shows that the dependence of  $J$  on  $X$  is stronger for the antiferromagnetic than for the ferromagnetic complexes. Another interesting feature of the results is the small dependence

**Table 3.** Comparison of the Optimal Results of the Present Study with Those of Other Methods To Obtain Magnetic Couplings ( $\text{cm}^{-1}$ )

	SF-TDDFT, <sup>a</sup> B1PW40	BS- $\Delta$ SCF, <sup>b</sup> PBE35	BS- $\Delta$ SCF, <sup>c</sup> M06	BS- $\Delta$ SCF, <sup>c</sup> M06-2X	BS- $\Delta$ SCF, <sup>d</sup> B2-PLYP	BS- $\Delta$ SCF, <sup>d</sup> B2GP-PLYP	SF-TDDFT, <sup>a</sup> B1PW40	REKS, <sup>e</sup> B3LYP	REKS, <sup>e</sup> BH&HLYP	C-DFT, <sup>f</sup> B3LYP	exptl
Cu <sub>2</sub> Cl <sub>6</sub> <sup>2-</sup>	-88	-9	5	0.1	-121	-61					0 to -94
YAFZOU	76	132	294	75	164	123	76	264	87		111
XAMBUI	2	1.5	3	0.8	-15	-11	2	6.2	0.75		2
PATFIA	-98	-7.5	-15	-19	15	19	-98	139	32		-11
CAVXUS	-15	-10.5	-28	-6	-17	-14	-15	3.3	-3.4		-19
CUAQAC02	-249	-233	-436	-143	-262	-177	-249	-285	-91		-286
BISDOW	-354	-308	-632	-177	-336	-122	-354	-429	-135		-382
MSE <sup>g</sup>	-9	27	-38	53	22	51	-9	47	79		
MUE <sup>g</sup>	32	27	100	68	28	55	32	63	88		
RMSE <sup>g</sup>	43	38	141	103	33	80	43	90	130		
MURE <sup>g</sup>	1.4	0.3	0.7	0.6	1.9	1.7	1.4	3.1	1.1		
1(Cu–Cu)	-74	-71								-32	-62
2(Cu–Cu)	-115	-89								-88	-75
4(V–V)	-166	-129								-166	-214
7(Cu–Cu)	181	277								224	168
8(Cu–Cu)	112	179								114	114
MSE <sup>h</sup>	1.5	47									24
MUE <sup>h</sup>	23	56									29
RMSE <sup>h</sup>	29	69									36
MURE <sup>h</sup>	0.2	0.4									0.2

<sup>a</sup> Present work. <sup>b</sup> Reference 80. <sup>c</sup> Reference 83. <sup>d</sup> Reference 85. <sup>e</sup> Reference 81. The results included here correspond to those computed using eq 14 of that work. <sup>f</sup> Reference 35. <sup>g</sup> Mean errors for rows YAFZOU though BISDOW. <sup>h</sup> Mean errors for rows 1(Cu–Cu) though 8(Cu–Cu).

found on the pure exchange and correlation functionals employed, as evidenced by the set of five density functionals mentioned above that have nearly the same accuracy for the prediction of magnetic couplings, despite significant differences in the functional forms and/or parameters of the exchange and correlation functionals.

The observation that magnetic couplings and, more generally, spin-state energy differences, depend strongly on  $X$  has been discussed at length in the literature for transition-metal-containing compounds<sup>96–110</sup> and p-block atoms.<sup>49</sup> In general,  $X$  in the range of 40–60 is often necessary to obtain good agreement with experiment for spin energy differences<sup>36,49,52,53,60,64</sup> (note, however, that B3LYP was found to perform well for bioinorganic copper complexes).<sup>56–59</sup> The optimum  $X$  value of about 40 found here is roughly in agreement with the observation made many years ago that  $X = 35$  is optimal for strongly correlated solids such as NiO<sup>96–98</sup> and others.<sup>111</sup> However,  $X$  of about 15 was found to be optimum for single-center Fe(II) complexes,<sup>99</sup> and local functionals ( $X = 0$ ) such as OPBE<sup>75,112</sup> and OLYP<sup>112,113</sup> were also found to perform well in several cases.<sup>100,108,110</sup> Thus, the calculation of magnetic couplings in transition-metal complexes could be seen in the wider context of the study of multiplicity-changing transitions in transition-metal chemistry, crucial to the understanding of, e.g., reaction mechanisms<sup>114</sup> and spin-crossover complexes.<sup>115–117</sup>

An important consideration to keep in mind in seeking generalizations is that the HDV model was designed for when the magnetic coupling is due to a pure spin flip without changes in the spatial orbitals, as is most likely to occur for weakly coupled centers, whereas many of the cases just mentioned involve spin states of different orbital parentage or orbitals on the same center that are not weakly coupled. The recent study of single-center

splittings in p-block atoms<sup>49</sup> showed that splittings depend strongly on  $X$  but also that the dependence on  $X$  and the optimum value of  $X$  depend on the system studied and on the method, being different for  $\Delta$ SCF, LS-TDDFT, and SF-TDDFT. In this context, it is interesting to compare the present SF-TDDFT results with previous studies, and we have prepared Table 3 to facilitate such a comparison. In this table, we will compare the present results with weighted-average BS- $\Delta$ SCF results based on spin projection (eq 3) and with those of two further strategies to obtaining magnetic couplings: the constrained DFT (C-DFT) method<sup>118</sup> as implemented by Wu and Van Voorhis<sup>119</sup> and the spin-restricted ensemble-referenced Kohn–Sham (REKS) method of Filatov and Shaik.<sup>120,121</sup>

The first set of results in Table 3 is the optimal set of results of the present study, which—on the basis of the values of MSE, MUE, and RMSE—are the B1PW40 results. The next column has the optimal spin-projected weighted-average BS- $\Delta$ SCF results of Phillips and Peralta,<sup>80</sup> namely, their PBE35 results. Note that these authors employed three density functionals in their study: PBEX, HSE $\Omega$ , and LC- $\omega$ PBE $\Omega$ , where  $X$  denotes a variable percentage of Hartree–Fock exchange and  $\Omega$  represents a variable range parameter  $\omega$  for the range-separated HSE and LC- $\omega$ PBE density functionals. The rest of their parameters were taken from the standard PBE0,<sup>73,74</sup> HSE,<sup>122</sup> and LC- $\omega$ PBE<sup>123</sup> functionals. Their results showed that one obtained the best results with HSE0 (HSE $\Omega$  with  $\Omega = 0.0 a_0^{-1}$ ) but that one obtains similar results for any  $\Omega$  between 0.0 and 0.2  $a_0^{-1}$ . (Note that HSE0 is the same as PBE0). However, Phillips and Peralta also pointed out that since one could expect results with the original HSE value of  $\Omega$  very similar to the best results with  $\Omega = 0.0$ , the standard HSE method with the value of  $\Omega = 0.11 a_0^{-1}$  would be advantageous

for extended systems because of its more favorable computational cost as compared to that using  $\Omega = 0.0$ . The mean errors in Table 3 are for two subsets. If we compute the mean errors of the results from weighted-average BS- $\Delta$ SCF calculations with PBE35 over the 11 complexes used for averages in Table 1, we obtain MSE, MUE, and RMSE values of 36, 40, and 54  $\text{cm}^{-1}$ , all larger than the corresponding values in Table 1 for SF-TDDFT with B1PW40; the MURE is smaller though (0.3 vs 0.9) because of the large relative error of SF-TDDFT with B1PW40 for PATFIA.

For the seven complexes  $\text{Cu}_2\text{Cl}_6^{2-}$ , YAFZOU, XAMBUI, PATFIA, CAVXUS, CUAQAC02, and BISDOW, we also compare the present results with two other sets of spin-projected weighted-average BS- $\Delta$ SCF results, the ones from our previous study<sup>83</sup> using the Minnesota M06 and M06-2X functionals<sup>71,72</sup> and those reported with the double-hybrid B2-PLYP<sup>124</sup> and B2GP-PLYP<sup>125</sup> functionals by Schwabe and Grimme.<sup>85</sup> The B2-PLYP and B2GP-PLYP density functionals are doubly hybrid functionals,<sup>126</sup> in which an SCF step is followed by a post-SCF perturbative calculation of the correlation energy; both the density functional correlation energy and the perturbative contribution are empirically scaled. For the six complexes YAFZOU, XAMBUI, PATFIA, CAVXUS, CUAQAC02, and BISDOW, we report the best results of two implementations of the REKS method as described in eqs 13 and 14 of the work in which the couplings are calculated.<sup>81</sup> Finally, for the five complexes 1(Cu–Cu), 2(Cu–Cu), 4(V–V), 7(Cu–Cu), and 8(Cu–Cu), we compare the present results with the C-DFT results of Rudra et al.<sup>35</sup>

The results in Table 3 show several interesting features. First, note that only SF-TDDFT with the present  $X = 40$  functional (or any other functional in Table 1 with  $X \cong 40$  since they all give similar results), PBE35 with the spin-projected weighted-average BS- $\Delta$ SCF approach, and C-DFT predict the correct sign (ferromagnetic or antiferromagnetic) for all the complexes for which they have been tested. (The only two methods for which it is certain that they predict all 12 signs correctly are the first two.) If we look at the MSEs as a measure of systematic deviations from experiment, it is clear that the best method is SF-TDDFT with  $X \cong 40$ . Furthermore, the only three methods with MUE and RMSE comparable to those of SF-TDDFT with  $X \cong 40$  are spin-projected weighted-average BS- $\Delta$ SCF with PBE35, doubly hybrid B2-PLYP, which was applied<sup>85</sup> using the spin-projected eq 3, and C-DFT. Note that, for the last two, results have been reported for only the second subset of magnetic complexes. If instead of the MSE, MUE, and RMSE we would take the MURE as a measure of accuracy, it is clear that PBE35 with the spin-projected broken-symmetry approach would be the best method overall (i.e., for the 12 complexes). The MURE obtained with SF-TDDFT is large because of the large relative error of the PATFIA complex, which has a very small splitting. The REKS method is systematically too ferromagnetic and does not provide results competitive with those of SF-TDDFT or the best spin-projected weighted-average BS ones. One can conclude that, to the extent that the database of 12 complexes studied is representative of homobinuclear transition-metal complexes with spin-1/2 centers, SF-TDDFT with  $X \cong 40$  makes more accurate predictions, on average, than any other method studied. Since all the SF-TDDFT functionals with  $X \cong 40$  make similar predictions, we recommend using MPW1K or BMK since they are standard functionals that have been well validated for a great variety of chemical properties; for example, MPW1K has been shown to

provide relatively accurate predictions for hydrogen-bonding and charge-transfer interactions.<sup>127</sup> Of the two, MPW1K is simpler (being a hybrid GGA, whereas BMK is a hybrid meta-GGA) and is therefore easiest to implement in a wide variety of programs.

The discussion above regarding the accuracy of the different approaches has focused on the comparison to experiment. From a fundamental point of view it is also interesting to compare SF-TDDFT and BS energy splitting values obtained with two different formalisms. We found that, in all cases studied, the results obtained from the BS approach compare to those arising from SF-TDDFT if and only if spin projections (eq 3) are taken into account; this is another indication that spin symmetry has to be taken into account in DFT calculations as is usually done when wave functions are used.

## 5. CONCLUSIONS

Magnetic exchange coupling constants have been computed for a database of 12 spin-1/2 homobinuclear transition-metal complexes previously studied by Phillips and Peralta<sup>80</sup> and others. In the present work, several global hybrid density functionals, with the percentage of Hartree–Fock exchange ranging from 20% to 100%, have been employed with collinear, Tamm–Dancoff spin-flip time-dependent density functional theory. The magnetic coupling constants are calculated from singlet–triplet energy differences, as one would do in wave function theory, with both spin states generally well adapted to spin symmetry.

For a given functional, the spin-state energy splitting values predicted by the SF-TDDFT formalism are consistent with those obtained from the broken-symmetry approach if and only if spin projection is taken into account in the latter. Considering all 12 complexes, we find that 40% Hartree–Fock exchange provides the best agreement with experiment and—in terms of mean signed error, mean unsigned error, and root-mean-square error—the SF-TDDFT approach performs systematically better than the spin-projected weighted-average broken-symmetry strategy, although the optimal percentage is slightly different in each case (about 40% for SF-TDDFT and 35% for spin-projected weighted-average broken symmetry). If one considers subsets of the database for which previous results are available, one finds that the spin-projected weighted-average broken-symmetry doubly hybrid functional B2-PLYP and the constrained density functional theory based on B3LYP also perform quite well.

For SF-TDDFT, the magnetic couplings show a much larger dependence on the percentage of Hartree–Fock exchange than on other aspects of the exchange and correlation density functionals employed. Given that hybrid meta-GGAs do not seem to improve the results with respect to hybrid GGAs, one can use a less computationally demanding hybrid GGA such as MPW1K for SF-TDDFT on this kind of system. Further studies on a greater variety of systems would be welcome.

In conclusion, we find that the SF-TDDFT approach provides more accurate spin-state energy splittings than the spin-projected weighted-average broken-symmetry scheme for binuclear spin-1/2 transition-metal complexes, with the added advantage of avoiding the ambiguities associated with the weighted-average broken-symmetry approach. We recommend MPW1K as a well-known, standard hybrid density functional to calculate magnetic couplings in the context of SF-TDDFT.

## AUTHOR INFORMATION

## Corresponding Author

\*E-mail: rvalero@qui.uc.pt.

## ACKNOWLEDGMENT

We thank Jordan J. Phillips for providing us with tables of magnetic coupling constants derived from his work. R.V. thanks the Portuguese Foundation for Science and Technology for Grant C2008-FCTUC UQFM v29. Financial support has been provided by the Spanish MICINN (Grant FIS2008-02238) and Generalitat de Catalunya (Grants 2009SGR1041 and XRQTC) and through the 2009 ICREA Academia Award for excellence in research granted to F.I. This work was supported in part by the U.S. National Science Foundation under Grant CHE09-56776.

## REFERENCES

- (1) Kahn, O. *Molecular Magnetism*; VCH Publishers: New York, 1993; pp 1–380.
- (2) Blundell, S. J.; Pratt, F. L. *J. Phys.: Condens. Matter* **2004**, *16*, R771.
- (3) Davidson, E. R.; Clark, A. E. *Phys. Chem. Chem. Phys.* **2007**, *9*, 1881.
- (4) Miyasaka, H.; Saitoh, A.; Abe, S. *Coord. Chem. Rev.* **2007**, *251*, 2622.
- (5) Mezei, G.; Zaleski, C. M.; Pecoraro, V. L. *Chem. Rev.* **2007**, *107*, 4933.
- (6) Bagai, R.; Christou, G. *Chem. Soc. Rev.* **2009**, *38*, 1011.
- (7) Zeng, Y.-F.; Hu, X.; Liu, F.-C.; Bu, X.-H. *Chem. Soc. Rev.* **2009**, *38*, 469.
- (8) Atanasov, M.; Comba, P.; Hausberg, S.; Martin, B. *Coord. Chem. Rev.* **2009**, *253*, 2306.
- (9) Sessoli, R.; Powell, A. K. *Coord. Chem. Rev.* **2009**, *253*, 2328.
- (10) Wang, S.; Ding, X.-H.; Zuo, J.-L.; You, X.-Z.; Huang, W. *Coord. Chem. Rev.* **2011**, *255*, 1713.
- (11) Cornia, A.; Constantino, A. F.; Zobbi, L.; Caneschi, A.; Gatteschi, D.; Mannini, M.; Sessoli, R. *Struct. Bonding (Berlin)* **2006**, *112*, 133 and references therein.
- (12) Gómez-Segura, J.; Veciana, J.; Ruiz-Molina, D. *Chem. Commun.* **2007**, *36*, 3699.
- (13) Morán-López, J. L.; Guirado-López, R. A.; Montejano-Carrizalez, J. M.; Aguilera-Granja, F.; Rodríguez-Alba, R.; Mejía-López, J.; Romero, A. H.; García, M. E. *Curr. Sci.* **2008**, *95*, 1177.
- (14) Cramer, C. J.; Truhlar, D. G. *Phys. Chem. Chem. Phys.* **2009**, *11*, 10757.
- (15) Moreira, I. d. P. R.; Illas, F. *Phys. Chem. Chem. Phys.* **2006**, *8*, 1645 and references therein.
- (16) de Loth, P.; Cassoux, P.; Daudey, J. P.; Malrieu, J. P. *J. Am. Chem. Soc.* **1981**, *103*, 4007.
- (17) de Graaf, C.; Sousa, C.; Moreira, I. d. P. R.; Illas, F. *J. Phys. Chem. A* **2001**, *105*, 11371.
- (18) Moreira, I. d. P. R.; Illas, F.; Calzado, C. J.; Sanz, J. F.; Malrieu, J. P.; Ben Amor, N.; Maynau, D. *Phys. Rev. B* **1999**, *59*, 6593.
- (19) Muñoz, D.; Illas, F.; Moreira, I. d. P. R. *Phys. Rev. Lett.* **2000**, *84*, 1579.
- (20) Calzado, C. J.; Cabrero, J.; Malrieu, J. P.; Caballol, R. *J. Chem. Phys.* **2002**, *116*, 2728.
- (21) Calzado, C. J.; Cabrero, J.; Malrieu, J. P.; Caballol, R. *J. Chem. Phys.* **2002**, *116*, 3985.
- (22) Miralles, J.; Castell, O.; Caballol, R.; Malrieu, J. P. *Chem. Phys.* **1993**, *172*, 33.
- (23) Andersson, K.; Malmqvist, P.-Å.; Roos, B. O.; Sadlej, A. J.; Wolinski, K. *J. Phys. Chem.* **1990**, *94*, 5483.
- (24) Andersson, K.; Malmqvist, P.-Å.; Roos, B. O. *J. Chem. Phys.* **1992**, *96*, 1218.
- (25) Neese, F. *Coord. Chem. Rev.* **2009**, *253*, 526.
- (26) Noodleman, L. *J. Chem. Phys.* **1981**, *74*, 5737.
- (27) Noodleman, L.; Davidson, E. R. *J. Chem. Phys.* **1986**, *109*, 131.
- (28) Noodleman, L.; Peng, C. Y.; Case, D. A.; Mouesda, J. M. *Coord. Chem. Rev.* **1995**, *144*, 199.
- (29) Ruiz, E.; Cano, J.; Alvarez, S.; Alemany, P. *J. Comput. Chem.* **1999**, *20*, 1391.
- (30) Ruiz, E. *J. Comput. Chem.* **2011**, *32*, 1998.
- (31) Yamaguchi, K.; Tsunekawa, T.; Toyoda, Y.; Fueno, T. *Chem. Phys. Lett.* **1988**, *143*, 371.
- (32) Yamaguchi, K.; Fueno, T.; Ueyama, N.; Nakamura, A.; Ozaki, M. *Chem. Phys. Lett.* **1989**, *164*, 210.
- (33) Nagao, H.; Mitani, M.; Nishino, M.; Yoshioka, Y.; Yamaguchi, K. *Int. J. Quantum Chem.* **1997**, *65*, 947.
- (34) Caballol, R.; Castell, O.; Illas, F.; Malrieu, J. P.; Moreira, I. d. P. R. *J. Phys. Chem. A* **1997**, *101*, 7860.
- (35) Rudra, I.; Wu, Q.; Van Voorhis, T. *J. Chem. Phys.* **2006**, *124*, 024103.
- (36) Shao, Y.; Head-Gordon, M.; Krylov, A. I. *J. Chem. Phys.* **2003**, *118*, 4807.
- (37) Wang, F.; Ziegler, T. *J. Chem. Phys.* **2004**, *121*, 12191.
- (38) Krylov, A. I. *J. Phys. Chem. A* **2005**, *109*, 10638.
- (39) Levine, B. G.; Ko, C.; Quenneville, J.; Martinez, T. *J. Mol. Phys.* **2006**, *104*, 1039.
- (40) Slipchenko, L. V.; Krylov, A. I. *J. Chem. Phys.* **2002**, *117*, 4694.
- (41) Casida, M. E. In *Recent Advances in Density Functional Methods, Part I*; Chong, D. P., Ed.; World Scientific: Singapore, 1995; p 155.
- (42) Runge, E.; Gross, E. K. U. *Phys. Rev. Lett.* **1984**, *52*, 997.
- (43) Petersilka, M.; Grossman, U. J.; Gross, E. K. U. *Phys. Rev. Lett.* **1996**, *76*, 1212.
- (44) Bauernschmitt, R.; Ahlrichs, R. *Chem. Phys. Lett.* **1996**, *256*, 454.
- (45) Stratmann, R. E.; Scuseria, G. E.; Frisch, M. J. *J. Chem. Phys.* **1998**, *109*, 8218.
- (46) Marques, M. A. L.; Gross, E. K. U. *Annu. Rev. Phys. Chem.* **2004**, *55*, 427.
- (47) Dreuw, A.; Head-Gordon, M. *Chem. Rev.* **2005**, *105*, 4009.
- (48) Krylov, A. I. *Chem. Phys. Lett.* **2001**, *338*, 375.
- (49) Yang, K.; Peverati, R.; Truhlar, D. G.; Valero, R. *J. Chem. Phys.* **2011**, *135*, 044118.
- (50) Zhekova, H.; Seth, M.; Ziegler, T. *J. Chem. Theory Comput.* **2011**, *7*, 1858.
- (51) Jung, Y.; Head-Gordon, M. *ChemPhysChem* **2003**, *4*, 522.
- (52) Jung, Y.; Head-Gordon, M. *J. Phys. Chem. A* **2003**, *107*, 7475.
- (53) Jung, Y.; Heine, T.; Schleyer, P. v. R.; Head-Gordon, M. *J. Am. Chem. Soc.* **2004**, *126*, 3132.
- (54) Jung, Y.; Brynda, M.; Power, P. P.; Head-Gordon, M. *J. Am. Chem. Soc.* **2006**, *128*, 7185.
- (55) Rinkevicius, Z.; Vahtras, O.; Ågren, H. *J. Chem. Phys.* **2010**, *133*, 114104.
- (56) de la Lande, a.; Moliner, V.; Parisel, O. *J. Chem. Phys.* **2007**, *126*, 035102.
- (57) de la Lande, A.; Gerard, H.; Parisel, O. *Int. J. Quantum Chem.* **2008**, *108*, 1898.
- (58) de la Lande, A.; Parisel, O.; Gerard, H.; Moliner, V.; Reinaud, O. *Chem.—Eur. J.* **2008**, *14*, 6465.
- (59) de la Lande, A.; Salahub, D.; Moliner, V.; Gerard, H.; Piquemal, J.-P.; Parisel, O. *Inorg. Chem.* **2009**, *48*, 7003.
- (60) Minezawa, N.; Gordon, M. S. *J. Phys. Chem. A* **2009**, *113*, 12749.
- (61) Huix-Rotllant, M.; Natarajan, B.; Ipatov, A.; Wawire, C. M.; Deutsch, T.; Casida, M. E. *Phys. Chem. Chem. Phys.* **2010**, *12*, 12811.
- (62) You, Z.-Q.; Shao, Y.; Hsu, C. P. *Chem. Phys. Lett.* **2004**, *390*, 116.
- (63) Yang, C. H.; Hsu, C. P. *J. Chem. Phys.* **2006**, *124*, 244507.
- (64) Zhang, W.; Zhu, W.; Liang, W.; Zhao, Y.; Nelsen, S. F. *J. Phys. Chem. B* **2008**, *112*, 11079.
- (65) Tamm, I. *J. Phys. (Moscow)* **1945**, *9*, 449.
- (66) Hirata, S.; Head-Gordon, M. *Chem. Phys. Lett.* **1999**, *314*, 291.

- (67) Shao, Y.; Fusti-Molnar, L.; Jung, Y.; Kussmann, J.; Ochsenfeld, C.; Brown, S. T.; Gilbert, A. T. B.; Slipchenko, L. V.; Levchenko, S. V.; O'Neill, D. P.; Distasio, R. A., Jr.; Lochan, R. C.; Wang, T.; Beran, G. J. O.; Beasley, N. A.; Herbert, J. M.; Lin, C. Y.; Van Voorhis, T.; Chien, S. H.; Sodt, A.; Steele, R. P.; Rassolov, V. A.; Maslen, P. E.; Korambath, P. P.; Adamson, R. D.; Austin, B.; Baker, J.; Byrd, E. F. C.; Dachsel, H.; Doerksen, R. J.; Dreuw, A.; Dunietz, B. D.; Dutoi, A. D.; Furlani, T. R.; Gwaltney, S. R.; Heyden, A.; Hirata, S.; Hsu, C.-P.; Kedziora, G.; Khalliulin, R. Z.; Klunzinger, P.; Lee, A. M.; Lee, M. S.; Liang, W.; Lotan, I.; Nair, N.; Peters, B.; Proynov, E. I.; Pieniazek, P. A.; Rhee, Y. M.; Ritchie, J.; Rosta, E.; Sherrill, C. D.; Simmonett, A. C.; Subotnik, J. E.; Woodcock, H. L., III; Zhang, W.; Bell, A. T.; Chakraborty, A. K.; Chipman, D. M.; Keil, F. J.; Warshel, A.; Hehre, W. J.; Schaefer, H. F., III; Kong, J.; Krylov, A. I.; Gill, P. M. W.; Head-Gordon, M. *Phys. Chem. Chem. Phys.* **2006**, *8*, 3172.
- (68) Becke, A. D. *J. Chem. Phys.* **1993**, *98*, 5648.
- (69) Stephens, P. J.; Devlin, F. J.; Chabalowski, C. F.; Frisch, M. J. *J. Phys. Chem.* **1994**, *98*, 11623.
- (70) Stephens, P. J.; Devlin, F. J.; Ashvar, C. S.; Bak, K. K.; Taylor, P. R.; Frisch, M. J. *ACS Symp. Ser.* **1996**, *629*, 105.
- (71) Zhao, Y.; Truhlar, D. G. *Theor. Chem. Acc.* **2008**, *120*, 215.
- (72) Zhao, Y.; Truhlar, D. G. *Acc. Chem. Res.* **2008**, *41*, 157.
- (73) Perdew, J.; Ernzerhof, M.; Burke, K. *J. Chem. Phys.* **1996**, *105*, 9982.
- (74) Adamo, C.; Barone, V. *J. Chem. Phys.* **1999**, *110*, 6158.
- (75) Perdew, J. P.; Burke, K.; Ernzerhof, M. *Phys. Rev. Lett.* **1996**, *77*, 3865.
- (76) Adamo, C.; Barone, V. *Chem. Phys. Lett.* **1997**, *274*, 242.
- (77) Boese, A. D.; Martin, J. M. L. *J. Chem. Phys.* **2004**, *121*, 3405.
- (78) Lynch, B. J.; Fast, P. L.; Harris, M.; Truhlar, D. G. *J. Phys. Chem. A* **2000**, *104*, 4811.
- (79) Zhao, Y.; Truhlar, D. G. *J. Phys. Chem. A* **2006**, *110*, 13126.
- (80) Phillips, J. J.; Peralta, J. E. *J. Chem. Phys.* **2011**, *134*, 034108.
- (81) Moreira, I. d. P. R.; Costa, R.; Filatov, M.; Illas, F. *J. Chem. Theory Comput.* **2007**, *3*, 764.
- (82) Rivero, P.; Moreira, I. d. P. R.; Illas, F.; Scuseria, G. E. *J. Chem. Phys.* **2008**, *129*, 184110.
- (83) Valero, R.; Costa, R.; Moreira, I. d. P. R.; Truhlar, D. G.; Illas, F. *J. Chem. Phys.* **2008**, *128*, 114103.
- (84) Rivero, P.; Loschen, C.; Moreira, I. d. P. R.; Illas, F. *J. Comput. Chem.* **2009**, *30*, 2316.
- (85) Schwabe, T.; Grimme, S. *J. Phys. Chem. Lett.* **2010**, *1*, 1201.
- (86) López, C.; Costa, R.; Illas, F.; De Graaf, C.; Turnbull, M. M.; Landee, C. P.; Espinosa, E.; Mata, I.; Molins, E. *Dalton Trans.* **2005**, *13*, 2322.
- (87) Peralta, J. E.; Melo, J. I. *J. Chem. Theory Comput.* **2010**, *6*, 1894.
- (88) Hay, P. J.; Wadt, W. R. *J. Chem. Phys.* **1985**, *82*, 299.
- (89) Rappoport, D.; Furche, F. *J. Chem. Phys.* **2010**, *133*, 134105.
- (90) These basis sets were obtained from the Basis Set Exchange Database, version 1.2.2, as developed and distributed by the Environmental and Molecular Sciences Laboratory, which is part of the Pacific Northwestern Laboratory, P.O. Box 999, Richland, WA 99352, and is funded by the U.S. Department of Energy. See: Feller, D. *J. Comput. Chem.* **1996**, *17*, 1571. Schuchardt, K. L.; Didier, B. T.; Elsethagen, T.; Sun, L.; Gurumoorthis, V.; Chase, J.; Li, J.; Windus, T. L. *J. Chem. Inf. Model.* **2007**, *47*, 1045.
- (91) Willet, R. D. In *Magneto Structural Correlations in Exchange Coupled Systems*; Willet, R. D., Gatteschi, D., Kahn, O., Eds.; NATO Advanced Studies Series C: Mathematical and Physical Sciences; Reidel: Dordrecht, The Netherlands, 1985; p 140.
- (92) Honda, M.; Katayama, C.; Tanaka, J.; Tanaka, M. *Acta Crystallogr., Sect. C* **1985**, *41*, 197.
- (93) Miralles, J.; Daudey, J. P.; Caballol, R. *Chem. Phys. Lett.* **1992**, *198*, 555.
- (94) Castell, O.; Miralles, J.; Caballol, R. *Chem. Phys.* **1994**, *179*, 377.
- (95) Saito, T.; Nishihara, S.; Yamanaka, S.; Kitagawa, Y.; Kawakami, T.; Okumura, M.; Yamaguchi, K. *Chem. Phys. Lett.* **2011**, *505*, 11.
- (96) Martin, R. L.; Illas, F. *Phys. Rev. Lett.* **1997**, *79*, 1539.
- (97) Illas, F.; Martin, R. L. *J. Chem. Phys.* **1998**, *108*, 2519.
- (98) Moreira, I. d. P. R.; Illas, F.; Martin, R. L. *Phys. Rev. B* **2002**, *65*, 155102.
- (99) Reiher, M.; Solomon, O.; Hess, B. A. *Theor. Chem. Acc.* **2001**, *107*, 48.
- (100) Swart, M.; Groenhof, A. R.; Ehlers, A. W.; Lammertsma, K. *J. Phys. Chem. A* **2004**, *108*, 5479.
- (101) Harvey, J. N. *Struct. Bonding (Berlin)* **2004**, *112*, 151.
- (102) Daku, L. M. L.; Vargas, A.; Hauser, A.; Fouqueau, A.; Casida, M. E. *ChemPhysChem* **2005**, *6*, 1393.
- (103) Pierloot, K.; Vancoillie, S. *J. Chem. Phys.* **2006**, *125*, 124303.
- (104) Rong, C.; Lian, S.; Yin, D.; Shen, B.; Zhong, A.; Bartolotti, L.; Liu, S. *J. Chem. Phys.* **2006**, *125*, 174102.
- (105) Brewer, G.; Olida, M. J.; Schmiedekamp, A. M.; Viragh, C.; Zavalij, P. Y. *Dalton Trans.* **2006**, *47*, 5617.
- (106) Vargas, A.; Zerara, M.; Krausz, E.; Hauser, A.; Daku, L. M. L. *J. Chem. Theory Comput.* **2006**, *2*, 1342.
- (107) Strickland, N.; Harvey, J. N. *J. Phys. Chem. B* **2007**, *111*, 841.
- (108) Pierloot, K.; Vancoillie, S. *J. Chem. Phys.* **2008**, *128*, 034104.
- (109) Oláh, J.; Harvey, J. N. *J. Phys. Chem. A* **2009**, *113*, 7338.
- (110) Vancoillie, S.; Zhao, H.; Radoń, M.; Pierloot, K. *J. Chem. Theory Comput.* **2010**, *6*, 576.
- (111) Feng, X.; Harrison, N. M. *Phys. Rev. B* **2004**, *70*, 092402.
- (112) Handy, N. C.; Cohen, A. J. *Mol. Phys.* **2001**, *99*, 403.
- (113) Lee, C.; Yang, W.; Parr, R. G. *Phys. Rev. B* **1988**, *37*, 785.
- (114) Keogh, D. W.; Poli, R. *J. Am. Chem. Soc.* **1997**, *119*, 2516.
- (115) Gülich, P.; Goodwin, H. A. *Topics in Current Chemistry*; Springer: New York, 2004; Vols. 233–235.
- (116) Hostettler, M.; Törnroos, K. W.; Chernyshov, D.; Vangdal, B.; Bürgi, H.-B. *Angew. Chem., Int. Ed.* **2004**, *43*, 4589.
- (117) Enachescu, C.; Hauser, A.; Girerd, J.-J.; Boillot, M.-L. *ChemPhysChem* **2006**, *7*, 1127.
- (118) Dederichs, P. H.; Blugel, S.; Zeller, R.; Akai, H. *Phys. Rev. Lett.* **1984**, *53*, 2512.
- (119) Wu, Q.; Van Voorhis, T. *Phys. Rev. A* **2005**, *72*, 024502.
- (120) Filatov, M.; Shaik, S. *Chem. Phys. Lett.* **1998**, *288*, 689.
- (121) Filatov, M.; Shaik, S. *Chem. Phys. Lett.* **1999**, *304*, 429.
- (122) Heyd, J.; Scuseria, G. E.; Ernzerhof, M. *J. Chem. Phys.* **2003**, *118*, 8207.
- (123) Vydrov, O. A.; Scuseria, G. E. *J. Chem. Phys.* **2006**, *125*, 234109.
- (124) Grimme, S. *J. Chem. Phys.* **2006**, *124*, 034108.
- (125) Karton, A.; Tarnopolsky, A.; Lamere, J.-F.; Schatz, G. C.; Martin, J. M. L. *J. Phys. Chem. A* **2008**, *112*, 12868.
- (126) Zhao, Y.; Lynch, B. J.; Truhlar, D. G. *J. Phys. Chem. A* **2004**, *108*, 4786.
- (127) Zhao, Y.; Truhlar, D. G. *J. Chem. Theory Comput.* **2005**, *1*, 415.

# Toward Reliable DFT Investigations of Mn-Porphyrins through CASPT2/DFT Comparison

Mikael Kepenekian,<sup>†,‡</sup> Adrian Calborean,<sup>†</sup> Valentina Vetere,<sup>†,||</sup> Boris Le Guennic,<sup>‡</sup> Vincent Robert,<sup>\*,‡,§</sup> and Pascale Maldivi<sup>\*,†</sup>

<sup>†</sup>SCIB, UMR E 3 CEA/UJF-Grenoble 1, Laboratoire de Reconnaissance Ionique et Chimie de Coordination, INAC, Grenoble, F-38054, France

<sup>‡</sup>Université de Lyon, CNRS, Institut de Chimie de Lyon, Ecole Normale Supérieure de Lyon, 15 Parvis René Descartes, 69342 Lyon Cedex 07, France

<sup>§</sup>Laboratoire de Chimie Quantique, Institut de Chimie de Strasbourg, UMR7177 CNRS/Université de Strasbourg, 4, rue Blaise Pascal, CS 90032, 67081 Strasbourg-Cedex, France

**ABSTRACT:** The low-energy spectroscopies of Mn(II) and Mn(III) porphyrin (P) complexes were investigated using complete active space and subsequent perturbative treatment (CASPT2) as well as DFT-based calculations. Starting from DFT optimizations of Mn<sup>II</sup>P and Mn<sup>III</sup>PCl using crystallographic data, the CASPT2 results show that whatever the relative position of the Mn(II) ion with respect to the porphyrin cavity, the high-spin state  $S = 5/2$  of the [MnP] unit lies much lower in energy than the intermediate  $S = 3/2$  state. Not only are these results in agreement with experimental observations but they also differ from previous theoretical conclusions. In the Mn(III) complexes,  $\sigma$  and  $\pi$  charge redistributions compete to result in a  $S = 2$  ground state. The performances of different functionals have been tested in the reproduction of the CASPT2 spin gaps. Our results confirm that the Mn(II) system is very challenging, as GGA functionals fail in the spin states ordering and in the reproduction of the gaps, unless a high percentage of exact HF exchange (55%), as in KMLYP, is incorporated. This inspection demonstrates the need for specific active space functional to investigate the low-energy spectroscopy of [MnP] units.

## 1. INTRODUCTION

The prominent role of porphyrin-based complexes in biological processes as in heme active sites has stimulated intense work from both the experimental and theoretical communities since the 1970s.<sup>1–5</sup> Much attention has been devoted to revealing the strong relationships between the electronic structures of metalloporphyrins and their molecular parameters. Such remarkable interplay has also led to intense efforts in order to take advantage of these features in widespread areas of interest such as health, catalysis,<sup>6–8</sup> and molecular materials.<sup>9</sup> In particular, metalloporphyrins have been recently investigated as possible information storage devices taking advantage of charge transfer effects.<sup>10–13</sup> The association of a redox metal center such as manganese with smart functionalization of the porphyrin has turned out to be a promising route in the design of molecular switches.<sup>14</sup> From this perspective, we got interested in gaining a better understanding of the particular relationships between the electronic structures of Mn(II) and Mn(III) porphyrins and their structural features through theoretical approaches.

Quantum chemistry descriptions of metalloporphyrins are well documented in the literature, especially after the advent of density functional theory (DFT) methods.<sup>15–23</sup> This class of coordination compounds is also considered a good candidate for quantum chemistry benchmarking,<sup>19–23</sup> thanks to the large amount of chemical and spectroscopic data available. The most difficult case is for Mn(II) ( $d^5$ ), which may give rise to two low-lying states: the high-spin one (HS,  $S = 5/2$ ) and an intermediate one (IS,  $S = 3/2$ ). Much experimental evidence (Mn–N distances, EPR, magnetic data) has unambiguously shown that

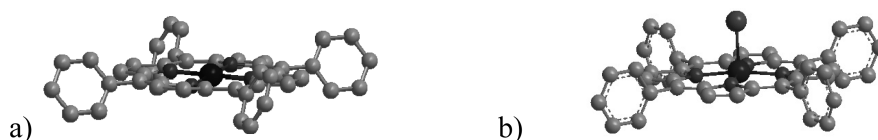
the sextet is the ground state.<sup>24–26</sup> In contrast, from the early extended Hückel calculations<sup>3</sup> to more recent DFT studies,<sup>20,27</sup> the electronic structure of Mn(II) porphyrins is predicted to exhibit a quartet ground state. In particular, it seems that GGA functionals are unable to reproduce the high-spin character of the ground state. In contrast, hybrid functionals which contain a Hartree–Fock exchange component do predict most of the time the correct spin energetics.<sup>20</sup> This result is in line with the well-known property of hybrid functionals to favor the HS state as ground state, contrarily to GGA predicting generally a low spin state.<sup>28–30</sup> It was long ago underlined<sup>1,26</sup> that the HS state of Mn(II) implies the population of a  $d_{x^2-y^2}$  orbital compared to its IS state (or to Mn(III),  $d^4$ ) thus leading to significantly longer Mn–N distances in HS Mn(II) porphyrins ( $d_{\text{Mn–N}}$  ca. 2.09 Å) compared to IS six-coordinate ones ( $d_{\text{Mn–N}}$  ca. 2.03 Å) or to HS ( $S = 2$ ) Mn(III) porphyrins ( $d_{\text{Mn–N}}$  ca. 2.02 Å).<sup>26</sup> Actually, X-ray structural studies have clearly shown that the Mn(II) ion is positioned out of the mean porphyrinic plane at a height of 0.19 Å in a tetra-coordinated Mn(II) tetraphenyl porphyrin (TPP).<sup>24,26</sup> Moreover, a recent DFT study on Mn(II) porphyrins has enlightened this difficulty: the authors have obtained a (wrong) quartet ground state. It is not until a rather large height  $h = 0.37$  Å of the Mn ion with respect to the porphyrin mean plane is reached that the HS state becomes the ground state.<sup>27</sup>

The dependence of spin states ordering with the functional is also recurrent in Fe(II) and Fe(III) porphyrins.<sup>19,20,22,31–33</sup> Even

Received: June 15, 2011

Published: September 12, 2011





**Figure 1.** Crystallographic structures of (a)  $[\text{Mn}(\text{II})\text{TPP}]^{26}$  and (b)  $[\text{Mn}(\text{III})\text{TPPCl}]^{57}$ . Hydrogen atoms are not depicted, for clarity.

with a hybrid functional such as B3LYP, the spin energetics are not always reliable.<sup>33</sup> More generally, within the last 10 years, a large amount of work in the DFT community has been devoted to the reproduction of spin states ordering within transition metal complexes.<sup>23,28–30,33–40</sup> Some of the main findings that emerged from these studies were (i) the efficiency of the OPTX exchange functional of Handy and Cohen<sup>41</sup> associated with standard local correlation functionals (OPBE, OLYP, etc.) to reproduce low spin energetics of metal complexes,<sup>34,35,39,42</sup> (ii) the design of a modified B3LYP functional, namely B3LYP\*, including a weaker exact exchange (ex. ex.) contribution (15%) than the standard one (20%),<sup>28,36</sup> and (iii) the correlation between the HS/LS ordering and the nature of the metal–ligand bonding.<sup>28,29</sup> The latter is of course a direct consequence of the role of the ligand field in the ordering of the lowest electronic states of a metal complex. In the meantime, there have been considerable advances in the design of new families of functionals—the third rung in the so-called Jacob’s ladder<sup>43</sup>—especially with the advent of meta-GGAs and hybrid meta-GGAs.<sup>44–46</sup> These functionals have been built with the intent of being as universal as possible, including all chemical elements and a wide range of properties.<sup>44,45</sup>

For all of the above-mentioned reasons, a better analysis of the performance of various functionals for describing Mn(II) and Mn(III) porphyrins became necessary. Such benchmarking may be realized against experimental data or highly accurate *ab initio* calculations as a reference.<sup>23,29,33,47</sup> In this context, explicitly correlated calculations are particularly appealing since (i) they manipulate the exact Hamiltonian and (ii) the multireference character of the wave function gives access to important information with respect to the weights of the different configurations. In a configuration interaction method, the zeroth-order wave function is formed by a linear expansion of Slater determinants. Such a description is accessible by means of Complete Active Space Self-Consistent Field (CASSCF)<sup>48</sup> calculations which incorporate qualitatively the leading electronic configurations distributing  $n$  electrons in  $m$  molecular orbitals (MOs), defining an active space referenced as  $\text{CAS}[n,m]$ . At this level of calculation, the so-called static correlation effects are taken into account variationally, provided that the active space is flexible enough. The dynamical correlation effects can be included using second-order perturbation treatment (CASPT2) to produce reference calculations. However, large basis sets and extended active spaces might be necessary to reach convergence in the spectroscopy. Thus, particular attention was paid to (i) the active space characteristics and (ii) the nature of the functional in the low-energy spectroscopy determination.

Several functionals were compared in this study, from GGAs to meta-GGAs, hybrids and hybrid meta-GGAs, including recent and/or already proven efficient functionals, as above-mentioned: (i) GGAs BP86,<sup>49,50</sup> PBE,<sup>51</sup> OPBE,<sup>41</sup> and BLYP;<sup>49,52</sup> (ii) meta-GGAs MO6-L<sup>53</sup> and TPSS;<sup>45</sup> and (iii) hybrid functionals including ex. ex. (value given in parentheses): B3LYP (20%)<sup>54</sup> and B3LYP\* (15%),<sup>36</sup> PBE0 (25%),<sup>55</sup> BHandHLYP (50%),<sup>54</sup> and KMLYP (55.7%).<sup>56</sup> The two latter choices were driven by the large exchange

energy due to the half-filled  $d^5$  shell in Mn(II). BHandLYP is based on the half-and-half approach of Becke<sup>54</sup> based on 50% Hartree–Fock exchange and 50% BLYP exchange and a correlation functional. Another recent functional, KMLYP, that shows a similar percentage of exact HF exchange (55,7%) has appeared on the basis of the exchange part of Kang and Musgrave.<sup>56</sup> This functional was originally developed to reproduce energy reaction barriers and also contains a built-in reduction of self-interaction errors, which may be too interesting for such issues. Finally, recent meta hybrid functionals have also been investigated, namely, M06 (27%) and M06-2X (54%), from the Minnesota suite<sup>44</sup> completing the local M06-L meta-GGA, and TPSSH (10%) from Staroverov et al.<sup>45</sup> M06 was designed to be efficient for transition metals, while M06-2X was designed for main group elements but incorporates a high content of ex. ex. Finally, TPSS (meta-GGA) and TPSSH<sup>45</sup> were chosen because they are based on a local part, which is the PKZB functional from Perdew et al.,<sup>46</sup> that reproduces the exchange energy to second order in expansion of the density gradient and does not contain self-interaction spurious effects.

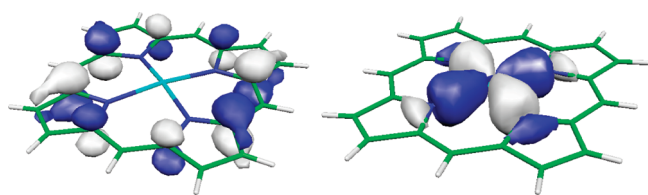
Our comparative study based on both the *ab initio* CASPT2 approach and the DFT scheme has been applied on simple Mn(II) and Mn(III) porphyrins, with H atoms in all meso and  $\beta$  positions (named porphin and abbreviated below as P), i.e.,  $\text{Mn}^{\text{II}}\text{P}$  and  $\text{Mn}^{\text{III}}\text{PCL}$ . Using the spin-dependent optimized geometries, the adiabatic energy differences for the  $\text{Mn}^{\text{II}}\text{P}$  and  $\text{Mn}^{\text{III}}\text{PCL}$  systems were calculated. Since some difficulties have been mentioned regarding the spin state ordering of  $\text{Mn}^{\text{II}}\text{P}$  with respect to experimental data, the vertical transition from the  $S = 5/2$  was also investigated. The starting geometries were based on X-ray diffraction experimental structures (see Figure 1).

Starting from these structures, we removed the four meso phenyl groups and replaced them with hydrogen atoms leading to the porphin (P) ligand. The use of simple MnP species as a model for more complex architectures has been justified in previous calculations.<sup>16</sup> The resulting structures will be referred to as **a** for the Mn(II) porphyrin and **b** for the Mn(III) porphyrin MnPCL.

Structures **a** and **b** were used as starting geometries to optimize  $\text{Mn}^{\text{II}}\text{P}$  and  $\text{Mn}^{\text{III}}\text{PCL}$ , respectively (see Computational Details). As it was not possible to optimize the structures with CASPT2 methods, we chose to use reference equilibrium geometries resulting from spin-dependent optimizations done with a given GGA (i.e., PBE), with the various possible spin states. The  $\text{Mn}^{\text{II}}\text{P}$  optimized structure starting from **a** for sextuplet ( $S = 5/2$ ) and quartet ( $S = 3/2$ ) states will be noted as **a-6** and **a-4**, respectively. The  $\text{Mn}^{\text{III}}\text{PCL}$  optimized structure from **b** for triplet ( $S = 1$ ), quintet ( $S = 2$ ), and septet ( $S = 3$ ) states will be noted as **b-3**, **b-5**, and **b-7**, respectively.

## 2. COMPUTATIONAL DETAILS

All of our CASSCF and CASPT2 calculations were performed with the Molcas7.0 package,<sup>58</sup> including atomic natural orbitals (ANO-RCC) as basis sets.<sup>59–61</sup> The one-electron basis sets



**Figure 2.** The  $\pi$ -type (left) and  $\sigma$ -type (right) orbitals of the [MnP] complexes from CASSCF calculations.

employed to describe the molecular orbitals (MOs) are derived from primitive ANO-RCC (21s,15p,10d,6f,4 g,2 h), (17s,12p,5d,4f,2 g), (14s,9p,4d,3f,2 g), (14s,9p,4d,3f,2 g), and (8s,4p,3d,1f) for the manganese, chlorine, nitrogen, carbon, and hydrogen atoms, respectively. Following the atomic natural orbital contractions of Widmark, these basis sets were contracted into [7s,6p,5d,3f,2 g,1 h], [4s,3p,1d], [3s,2p,1d], [3s,2p,1d], and [2s1p]. Finally, to avoid the presence of intruder states and to provide a balanced description of open and closed shells, imaginary level and IPEA shifts of 0.20 and 0.25 au (atomic units) were used in the CASPT2 calculations, respectively. All electrons were correlated except those in the core parts. Depending on the number of d electrons, five for the Mn(II) complex and four for the Mn(III), different active spaces can be used for the [MnP] and [MnP]Cl species:

1. CAS[5,5] and CAS[4,5] are the minimal active spaces that consider nothing but the d orbitals and electrons.
2. CAS[14,13] and CAS[13,13] add to the previous one a set of bonding and antibonding porphyrin-localized ( $\pi/\pi^*$ ) orbitals (see Figure 2) within each irreducible representation.
3. CAS[15,14] and CAS[14,14] finally consider a supplementary  $\sigma$ -type orbital (see Figure 2) representing the ion–porphyrin  $\sigma$  bond. Whereas the enlargement from the five-orbital active spaces to the 13-orbital ones is rather natural in light of the extended  $\pi$  system over the porphyrin ring, the inclusion of an additional  $\sigma$ -type orbital deserves some explanation. The importance of ligand-to-metal charge transfers (LMCT), in particular along the  $\sigma$  channel, has been stressed in the ground state from previous DFT-based calculations.<sup>62,63</sup> The 14-orbital active spaces do not discriminate between the  $\sigma$  and  $\pi$  manifolds and allow one to estimate the relative importance of these LMCTs accessible from the wave function expansion.

Dynamical correlation effects were added through the CASPT2<sup>59,64</sup> method that has proven to be an impressive tool used to accurately investigate spectroscopy issues.<sup>65,66</sup> However, extended basis sets combined with rather large active spaces are necessary to reach experimental agreement.<sup>67</sup>

The DFT calculations were performed with the ADF2010 package<sup>68–70</sup> using an all-electron Slater type basis of triple- $\zeta$  quality on each atom with polarization functions. Geometry optimizations were made with TZP all-electron basis sets (one polarization function) and single points with TZ2P all-electron basis sets (two polarization functions).<sup>70</sup> All of our calculations were performed using an unrestricted formalism to describe the various spin states. The convergence criteria were fixed to  $10^{-6}$  Hartree for the energy. Several checks were made on calculations (optimizations or single points) carried out with or without symmetry ( $D_{4h}$  or  $C_{2v}$  for MnP,  $C_{4v}$  or  $C_{2v}$  for MnPCL), showing that symmetry constraints do not greatly affect the results. For

instance, optimizations lead to distance differences less than 0.02 Å and energy differences less than 0.04 eV.

### 3. MN(II) PORPHYRINS

Geometry optimizations for the sextet and quadruplet states performed with the PBE functional gave flat porphyrin systems, as already pointed out in the literature.<sup>20,27</sup> Indeed, it is known that the porphyrin ring is rather flexible, precluding a reliable investigation of minima on the potential energy surface.<sup>71,72</sup> The Mn–N distance for a-6 species was 2.07 Å (exp. 2.085 Å<sup>26</sup>) and 2.00 Å for a-4. Optimizations with BLYP, OBPE, B3LYP, and M06 gave very similar geometries with distances ranging from 2.07 to 2.085 Å for a-6 and 2.00 to 2.02 Å for a-4. The decrease of Mn–N distances between the  $S = 5/2$  and the  $S = 3/2$  structures is in line with the depopulation of the mainly  $d_{x^2-y^2}$  antibonding orbital.

In order to check their consistency and convergence, the CASPT2 results were calibrated using several active spaces on the experimental structure a (see Table 1). In agreement with experimental observations, the three active spaces CAS[5,5], CAS[14,13], and CAS[15,14] give rise to a sextet ( $S = 5/2$ ) ground state. As seen in Table 1, the spin gap between the  $S = 5/2$  (HS) and  $S = 3/2$  (IS) states is almost not affected by the active space enlargement. The comparison between CAS[14,13] and CAS[15,14] is however instructive. Indeed, the inclusion of the  $\sigma$ -type orbital does not lead to any significant modification of the wave functions' structures. It should be stressed that for both spin-states, the main configuration holds a similar weight.

Both vertical and adiabatic energy differences (see Figure 3) were computed.  $\Delta E_a^{\text{vert}}$  uses the high-spin state geometry a-6, whereas  $\Delta E_{4-6}^{\text{adia}}$  relies on the quartet and sextet optimized geometries a-4 and a-6.

Table 2 gathers the results obtained with the various functionals defined above, compared to the reference CAS[15,14]PT2 values for both types of transition energy.

The first check was to correctly reproduce the experimentally known ordering, i.e., the sextet being the ground state thus corresponding to a positive energy gap. Clearly, GGAs and meta-GGAs are unable to reproduce this ordering, yielding a negative gap. It should be mentioned nevertheless that using the OPTX exchange potential (OLYP and OPBE) results in a weak spin gap, which becomes even positive with OLYP for the vertical transition. The most satisfactory behavior is obtained when including ex. ex., in hybrids or hybrid meta-GGAs functionals. But even in that case, the expected ordering is not obtained in all cases as B3LYP\* (15%), B3LYP (20%), and PBE0 (25%) give again a wrong ordering of the gaps for the vertical transition.

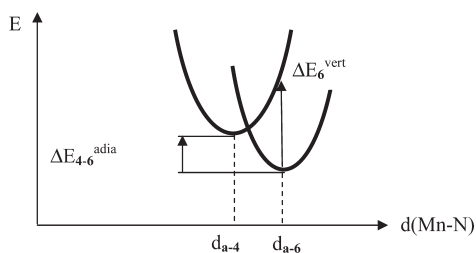
At this stage, the PBE optimized geometries might be questionable. Therefore, the geometries and corresponding adiabatic energy differences were calculated for a selection of functionals (see Table 3).

It is remarkable that OPBE and B3LYP give now qualitatively good ordering, with a sextet ground state, while BLYP and BPE still fail in reproducing the expected ordering.

From a qualitative point of view, there are clearly three major behaviors: GGAs and meta-GGAs, which are far from experimental agreement (except the OPTX-based ones), OPTX derived (OPBE and OLYP) and hybrid functionals with a low percentage of ex. ex. for which the result is ambiguous and depends on the molecular structure, and finally the hybrids with

**Table 1. CASPT2 Quartet–Sextet Vertical Spin Gap (eV) for the Experimentally Derived Structure a, Using Different Active Spaces**

active space	gap (eV)	weight of the main configuration for $S = 5/2$	weight of the main configuration for $S = 3/2$
CAS[5,5]	1.40	1.00	0.98
CAS[14,13]	1.36	0.85	0.83
CAS[15,14]	1.37	0.84	0.83

**Figure 3.** Potential energy curves for sextet and quadruplet spin states vs.  $d(\text{Mn}-\text{N})$  distance in the porphyrin and definitions of the spin gaps as estimated in the CASPT2 and DFT evaluations.

at least 50% ex. ex. and hybrid meta-GGAs which always lead to the expected sextet ground state.

We turn now to a more quantitative comparison, based on the CAS[15,14]PT2 values. Hybrids with less than 25% ex. ex., even if they give the expected sign for the vertical gap, do not perform well for numerical values. The best agreement—qualitatively and quantitatively—is obtained for hybrid functionals with a higher HF percentage, i.e., BHandHLYP and KMLYP, being the only ones to reproduce the good ordering in the vertical transition, although the numerical agreement is less favorable. Among the hybrid meta-GGAs, TPSSH gives a too low gap while M06 and M06-2X give a rather good agreement with the CAS[15,14]PT2 calculation.

We should mention that we checked that the electronic configurations of both sextet and quadruplet states obtained with both DFT and CAS approaches (the latter showing only one major configuration as above-mentioned) were the same. For the high spin state, the expected  $(d_{z^2})^1(d_{xy})^1(d_{xz},d_{yz})^2(d_{x^2-y^2})^1$  configuration was obtained. The lowest quartet state was also found to be  $(d_{z^2})^1(d_{xy})^1(d_{xz},d_{yz})^3(d_{x^2-y^2})^0$ , in agreement with the CAS[15,14] leading configuration. The spin contaminations that may reveal some mixing with higher states were checked for both spin states. It was found to be very low, as could be expected for the high spin state ( $S^2 = 8.75$  to  $8.77$  for an expected one of  $8.75$ ). For the quartet state with an expected value of  $S^2 = 3.75$ , we obtained, most of the time, weak spin contamination with  $S^2$  between 3.77 and 4 (i.e., < 10%) for GGAs, meta-GGAs, and hybrids, while three local functionals (M06-L, OPBE, and OLYP) gave a slightly larger value close to 4.2.

We should also mention that due to the differences obtained in the various B3LYP gaps (almost 0 eV for PBE optimized geometries and 0.58 eV for B3LYP optimized geometries), we also checked the consistency of the electronic configuration by calculating the gap on one geometry (B3LYP or PBE ones) restarting with the electron density obtained from the other

**Table 2. Spin Gaps  $\Delta E_6^{\text{vert}}$  and  $\Delta E_{4,6}^{\text{adia}}$  in eV Calculated by DFT Methods and CAS[15,14]PT2<sup>a</sup>**

	$\Delta E_6^{\text{vert}}$	$\Delta E_{4,6}^{\text{adia}}$
GGA		
PBE	−0.28	−0.49
OPBE	−0.16	−0.14
BLYP	−0.4	−0.59
OLYP	0.04	−0.19
BP86	−0.28	−0.55
hybrid		
B3LYP*	0.06	−0.2
B3LYP	0.23	−0.03
PBE0	0.5	−0.01
BHandHLYP	1.07	0.92
KMLYP	0.99	1.03
meta-GGA		
M06-L	−0.18	−0.62
TPSS	−0.43	−0.76
meta-hybrid		
M06	0.71	0.83
M06-2X	1.49	1.18
TPSSH	0.32	0.16
CASPT2	1.03	0.49

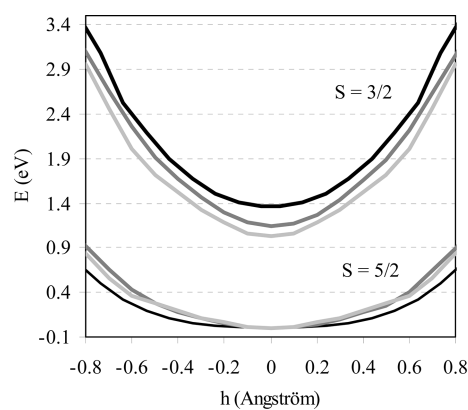
<sup>a</sup> All structures were optimized using the PBE functional.

**Table 3. Adiabatic Energy Gap between the Quadruplet and Sextet States, with Geometries Optimized for Each Functional**

functional	$\Delta E_{4,6}^{\text{adia}}$
PBE	−0.39
OPBE	0.06
BLYP	−0.43
B3LYP	0.58
M06	0.30

geometry (respectively PBE or B3LYP ones). This was accompanied by a check of the electron configurations, especially in the quartet states. The results were very similar to the ones above-mentioned (obtained without restarting densities), with a gap of 0.07 eV for the PBE geometries using the restart from B3LYP geometries and a gap of 0.55 eV for the B3LYP geometries using the restart from PBE geometries. The electronic structures were also checked to be consistent with the sextet and quartet state configurations described above.

As another check, we also explored the comparison of the potential energy curves as a function of the Mn height above the porphyrin plane, calculated for the  $[\text{Mn}^{\text{II}}\text{P}]$  complex from experimentally derived structure a. This key parameter was chosen due to ambiguous determinations through previous X-ray and theoretical studies.<sup>26,27,73</sup> As depicted in Figure 4, whatever the position of the manganese ion, the vertical quartet–sextet gap is at least 1.37 eV. Let us stress that the  $S = 5/2$  potential energy curve is rather flat, the energy variation being less than 0.1 eV for  $h \leq 0.45$  Å. This result might support the reported difficulties in  $[\text{MnP}]$  structure determinations since the



**Figure 4.** Potential energy curves with respect to the Mn ion displacement  $h$  for the [MnP] complex **a** for  $S = 5/2$  and  $S = 3/2$  states. Black: CASPT2. Dark gray: DFT/KMLYP. Light gray: DFT/M06. Zero energy reference taken as the  $h = 0$  point for  $S = 5/2$  state.

Mn(II) ion has the ability to be displaced out of the plane under weak external perturbations.<sup>25,26</sup>

At this stage, the KMLYP functional which displays the best overall agreement also gives a potential energy curve quite close to the CASPT2 result. M06 also closely reproduces the results from KMLYP, but with a smaller quartet–sextet gap, as was already observed in Table 2. Yet it gives a very satisfactory reproduction of the ground state curve.

From our comparison and previously reported ones, it is clear that the correct reproduction of the ground state of Mn<sup>II</sup>P species remains challenging for quantum chemical modeling. The inclusion of HF exchange helps to recover the high spin ground state, in line with the classical observation that high spin states are favored by including some exact exchange in the functional.<sup>28,29</sup>

#### 4. MN(III) PORPHYRINS

The spin energetics of the [MnPCl] species were investigated using again optimized structures starting from the experimental structure **b**, for each of the three spin states: triplet, quintet, and septet. The latter state is very often encountered in the Mn<sup>III</sup> porphyrin physicochemistry and corresponds to a formally Mn<sup>II</sup> ( $S = 5/2$ ) ion ferromagnetically coupled to a  $S = 1/2$  radical on the porphyrin. The agreement between the structural parameters optimized for the quintet ground state with experimental data was very satisfactory, with a mean Mn–N distance of 2.04 Å (exp. 2.015 Å),  $d(\text{Mn}–\text{Cl}) = 2.30$  Å (exp. 2.30 Å), and a Mn out of plane displacement of 0.30 Å (exp. 0.32 Å).

The geometry optimization for the  $S = 3$  state yields longer Mn–N distances (2.098 Å). Accordingly, the Kohn–Sham orbitals clearly show that the  $d_{x^2-y^2}$  is occupied—as expected for a high spin Mn(II) ion—and the spin densities give 4.5 on the Mn ion and a total of 1.06 delocalized on the four *meso* carbon atoms. This electronic structure supports the [Mn<sup>III</sup>P]<sup>+</sup> nature of the septet [Mn<sup>III</sup>P]<sup>+</sup> state in MnPCL.

In a first step, we explored the active space for CASPT2 calculations, starting from the experimentally derived structure **b**. The results are shown in Table 4. As for Mn(II) complexes, a first minimal active space (CAS[4,5]) has been used and leads to an  $S = 2$  ground state ( $^5A_2$ ) followed by two low-lying triplets  $^3B_2$  and  $^3B_1$  at 1.65 and 2.11 eV, respectively. Then, by enlarging the active space to CAS[12,13], a charge transfer state [Mn(II)P<sup>+</sup>]<sup>+</sup>

**Table 4.** CASPT2 Low-Energy Vertical Spectroscopy (eV) of the [MnPCl] Complex Calculated with Different Active Spaces from Structure **b**<sup>a</sup>

active space	$^3A_2$	$^3B_2$	$^7A_2$
CAS[4,5]	2.11	1.65	
CAS[12,13]	2.05	1.61	1.41
CAS[14,14]	1.44	1.53	1.74

<sup>a</sup>The reference energy is the quintet state  $^5A_2$  ( $C_{2v}$ ).

**Table 5.** Adiabatic Energy Gaps (eV) of the Lowest Triplet ( $S = 1$ ) and Septet ( $S = 3$ ) States, Respectively<sup>a</sup>

	$\Delta E_{3-5}^{\text{adia}}$	$\Delta E_{7-5}^{\text{adia}}$
	GGA	
PBE	0.42	1.55
OPBE	0.83	1.32
BLYP	0.37	1.50
OLYP	0.78	1.29
BP86	0.42	1.52
	hybrid	
B3LYP*	0.7	1.39
B3LYP	0.45	1.47
PBE0	1.14	1.67
BHandHLYP	1.21	1.07
KMLYP	1.28	1.11
	meta-GGA	
M06-L	1.27	1.82
TPSS	0.37	1.80
	meta-hybrid	
M06	1.45	1.54
M06-2X	1.83	1.15
TPSSH	0.55	1.72
CAS[14,14]PT2	1.51	1.10

<sup>a</sup> $\Delta E_{3-5}^{\text{adia}}$  and  $\Delta E_{7-5}^{\text{adia}}$ , with reference to the quintet ground state.

( $^7A_2$ ) can be described. It arises from the promotion of a  $\pi$  electron of the porphyrin ring into the vacant  $d_{x^2-y^2}$  orbital of Mn as already mentioned above in the DFT study. This state appears to be the first excited state at the CAS[12,13]PT2 level. In order to properly account for any charge redistributions between the Mn ion and the porphyrin ring, the 14-orbital active space has been tested, including both  $\sigma$  and  $\pi$  channels.  $^5A_2$  remains the ground state.  $^7A_2$  is shifted to much higher energies, i.e., 1.74 eV above  $^5A_2$ , while the triplets are stabilized with respect to the previous calculation. Thus, it is clear from this monitoring that the correct description of the  $\sigma$  transfer is crucial in the ordering of spin-states of Mn(III) complexes. Any successive enlargement of the active space led to no visible modification of the spectroscopy.

We have then used the CAS[14,14]PT2 results as reference to compare all DFT results obtained for each optimized geometry in each spin state. The results are summarized in Table 5.

The first conclusion is that whatever the functional, the expected quintet ground state is obtained. This is in contrast with the various differences observed with the Mn<sup>II</sup>P species. Apart from this first qualitative observation, the concern remains

about the ordering of excited spin states. In these relaxed geometries, the  $S = 3$  spin state at the CAS[14,14]PT2 level is found to be lower than the  $S = 1$  state, whereas GGAs, meta-GGAs, and some hybrids yield the opposite result. The two hybrids KMLYP and BHandHLYP give a very satisfactory quantitative agreement, which again can be related to the stabilization of high spin states due to a high ex. ex. content within these functionals. The results of the M06 and M06-2X meta-hybrid functionals are also in rather good agreement with the CASPT2 orderings.

The spin contamination has been checked for the three states ( $S^2 = 6$  for the ground state,  $S^2 = 2$  for the triplet, and  $S^2 = 12$  for the septet state). Almost no deviation is observed for the quintet state with all  $S^2$  values in the range 6.03–6.09, while deviations are below 10% for the triplet state (almost all being between 2.02 and 2.07). Finally, for the septet state, again very low deviations are obtained with all values between 12.04 and 12.16.

As for Mn(II) species, we have analyzed the Kohn–Sham orbitals of the lowest triplet state that was obtained in DFT and compared it to the major configuration given by the CAS approach. For the high spin state, the expected  $(d_{xz}, d_{yz})^2(d_{xy})^1(d_{z^2})^1(d_{x^2-y^2})^0$  configuration was obtained, resulting in a  $^5A_2$  symmetry in the  $C_{2v}$  point group as obtained in the CASPT2  $C_{2v}$  calculations. The lowest triplet state yielded the same configuration:  $(d_{xz}, d_{yz})^2(d_{xy})^2(d_{z^2})^0(d_{x^2-y^2})^0$  as the main one found in the CASPT2 calculation and corresponding to the  $^3A_2$  state in  $C_{2v}$  symmetry.

At this stage, we would like to comment on the choice of CASPT2 calculations as a reference in this particular case. Indeed, to obtain a balanced description of open and closed shells along the perturbative treatment, one has to include a so-called IPEA (ionization potential–electronic affinity) shift in the zeroth-order Hamiltonian.<sup>74</sup> The default value of 0.25 au as set by default in the current Molcas package usually gives excellent results. Nevertheless, this choice has been questioned on two occasions: (i) the case of magnetically coupled metals where it was found that the originally proposed zeroth-order Hamiltonian (corresponding to an IPEA set to 0.00 au) led to better description of the magnetic coupling<sup>75</sup> and (ii) in the evaluation of the adiabatic gap (between  $S = 0$  and  $S = 2$  spin states) for Fe(II) spin-crossover systems where it is suggested that a proper description of the gap requires an IPEA shift no less than 0.50 au.<sup>76</sup> For the former, the spectroscopy is characterized by states with identical numbers of open shells. In contrast, such number changes along the  $S = 0$  to  $S = 2$  transition. In a first step, we checked the impact of the IPEA shift on the low-energy spectroscopy of the Mn(II) species. The change in the quartet–sextet adiabatic gap appears to be less than 0.20 eV when going from the default value 0.25 au to the very high 0.75 au. In particular, no change is observed in the ordering of spin states. Thus, the standard 0.25 au value is suitable for calculations upon Mn(II) species, keeping in mind an error bar of  $\pm 0.10$  eV. On the other hand, the low-energy spectrum of Mn(III) species has to be treated more carefully. Even for a large IPEA shift value up to 1.00 au, the quintet state remains the ground state. However, the septet–quintet gap appears to be more sensitive to this parameter. Let us stress that the number of open shells reaches six for the heptet state, which involves an intramolecular electron transfer. As previously reported in the literature, the description of such a phenomenon requires larger values of the IPEA shift ( $\sim 0.5$  au). Considering the CASPT2 limitations, one may

conclude at this stage that the excited triplet and septet states are expected to lie relatively close in energy.

## 5. DISCUSSION

The present study has been intended to check the behavior of various types of functionals to reproduce qualitatively and—when possible—quantitatively the energetic ordering of the lowest spin states in Mn(II) and Mn(III) porphyrins. Our conclusion is that in Mn(III) species, all types of functionals are able to reproduce at least qualitatively and—for some—quantitatively the spin states ordering. However, the story is completely different for Mn(II) porphyrins, and we will focus the discussion on these systems.

In the following, we will first position our own results in the light of recent literature in spin states DFT benchmarking. Then, we will give some comments based first on the chemical nature of Mn(II), then on the choice of functionals that seems to result from this particular nature.

Our conclusion about the good efficiency of high exchange hybrid functionals is not completely in line with previous studies performed on the functionals to reproduce spin states ordering in 3d transition metal complexes. Indeed, most analogous comparative studies have been conducted on Fe<sup>II</sup> and Fe<sup>III</sup> complexes<sup>23,28,30,32,39,42,77–80</sup> because they are ubiquitous in active sites of metalloenzymes, within mononuclear or polynuclear clusters with magnetic coupling, and because this transition metal is very much used in spin-crossover materials. As mentioned in the Introduction, some constant conclusions emerged from these numerous studies, about the very efficient behavior of the exchange potential OPTX and about the good behavior of B3LYP\* with a decrease of ex. ex. compared to the standard B3LYP functional. Yet the efficiency of B3LYP\* for spin state orderings proved not to be universal.<sup>32,42,80,81</sup> Finally, all of these observations were nicely rationalized by some authors, on the basis of the nature of the ligand bonding.<sup>29,32,39</sup> High spin complexes are favored with ligands giving rise to more ionic bonding such as O or N donor ligands, whereas S, P, or C donor ligands are bonded with more covalent character, thus favoring a higher ligand field and lower spin ground states. Thus, hybrid functionals with a higher ex. ex. percentages are expected to perform better within ionic complexes with low covalence, whereas with complexes involving a more covalent bonding, hybrid functionals with a low percentage of ex. ex. are expected to be better.<sup>32,39,47</sup>

Thus, this rationalization in terms of chemical features can be extended to our case. Indeed, Mn(II) is very special among the 3d transition series. It is well-known by coordination chemists that this ion gives almost exclusively high spin ground state complexes<sup>82</sup> except with very strong field ligands such as CO, cyanide, or alkyl/aryl ligands. The manganese(II) ion in most of the ligand environments is known to behave as a noninteracting large sphere. The metallocene family is a typical example of this unique character, as all 3d metallocenes exhibit a low spin ground state except Mn(II). The manganocene complex (Mn<sup>II</sup>Cp<sub>2</sub>, Cp = C<sub>5</sub>H<sub>5</sub>) has been much studied using EPR and NMR solid state spectroscopies.<sup>83,84</sup> At very low temperatures, it exhibits antiferromagnetic coupling of  $S = 5/2$  units, due to a solid-state structure where the Cp ligands bridge two Mn(II) ions.<sup>84</sup> Moreover, NMR paramagnetic studies also revealed that manganocene was unique among other metallocenes because the metal–Cp ring bond was much less covalent than that of

nickelocene and cobaltocene.<sup>84</sup> Another remarkable feature is that when substituting the Cp rings with alkyl groups, the derived manganocene exhibited a low spin ground state.<sup>84</sup> This particular example illustrates several features of Mn(II) chemistry: (i) Mn(II) gives preferentially low covalent complexes and a high spin ground state even with the Cp ligand, known to favor low spin complexes, and it is necessary to substitute this ligand with alkyl groups in order to get a low spin ground state. (ii) The low covalence is correlated to the high spin ground state. (iii) DFT benchmarking on this particular case has yielded ambiguous results.<sup>28,29</sup>

The particular stability of the half filled d shell in Mn(II) due to a high exchange contribution is a key point. But Fe<sup>III</sup>—also a d<sup>5</sup> ion—does not behave the same. Yet, the latter is more positively charged, allowing ligands to get closer, thus favoring strong field environments, thus lower spin states. In fact, Fe(III) complexes exhibit various types of spin states, depending on their coordination environment.

All of these considerations give support to our results for Mn<sup>II</sup>P species, pointing to the need of increasing the ex. ex. contribution in hybrid functionals in order to properly describe the nature of Mn(II) complexes, unless particularly covalent bonds are expected. Indeed, the proper description of the exchange term and its correct balance with the correlation part is clearly very critical for Mn(II) species. In this context, the OPTX exchange constructed by fitting the HF exchange on atoms<sup>41</sup> behaves remarkably well and is the only GGA able to compete with hybrid functionals. But in order to get a quantitative agreement, inclusion of a high content of ex. ex. is desirable, and the tuning of the percentage has a direct effect on the efficiency of the hybrid, as already above-mentioned.

## 6. CONCLUSION

The present study was devoted to the evaluation of the low-energy spectroscopy of Mn-porphyrins through *ab initio* methods and to the comparison with DFT methods. For both Mn(II) and Mn(III) species, CASPT2 calculations were first conducted with different active spaces. In the case of Mn(II) compounds, the experimental spin ordering is recovered even with a minimal active space. However, the spectroscopy of Mn(III) complexes requires a larger active space to reach a correct description of the low-energy spectrum, featuring a capital part played by both  $\sigma$ - and  $\pi$ -type orbitals. This is reminiscent of the experimentally very different spectroscopies of Mn(II)P and Mn(III)P, well documented in the literature (see for instance ref 1).

These studies have allowed us to design a reliable DFT approach of Mn-porphyrin complexes. Indeed, as already mentioned in many instances, we fuel some other arguments—if necessary—that the choice of a functional is particularly crucial in transition metal complexes, where the spin state ordering may be complicated due to subtle interplay between d shell effects and the metal–ligand bonding nature. As in previously published studies, GGAs and meta-GGAs fail to reproduce the correct ordering in the Mn(II) sextet–quartet states as well as some hybrids, whereas the correct one is recovered with hybrids with a high content of ex. ex. such as KMLYP, BHANDHLYP, or to a lesser extent M06.

We must emphasize again the very special nature of the Mn(II) ion among other 3d metal ions, being a very large ion with a half filled d shell, that is known by coordination chemists to give mainly high spin species (except with very strong field

ligands).<sup>82</sup> We believe that this case is another “niche” where we cannot simply apply standard DFT methods. We must think of the chemistry behind the system in order to adapt the method and make an extensive usage of the experiences in the various theoretical approaches already published. The present study is another illustration of what has already been pointed out in some reflections published in recent years, about the pursuit of the “divine” functional.<sup>43,78,85,86</sup>

## AUTHOR INFORMATION

### Corresponding Author

\*E-mail: vincent.robort@ens-lyon.fr, vrobert@unistra.fr (V.R.); pascale.maldivi@cea.fr (P.M.).

### Present Addresses

<sup>||</sup>LITEN/DEHT/LPCEM, CEA-Grenoble, 17 rue des Martyrs, 38054 Grenoble cedex 9, France

## ACKNOWLEDGMENT

A.C. wishes to thank research funding from the European Community under the FP6 - Marie Curie Host Fellowships for Early Stage Research Training (EST) “CHEMTRONICS” Contract Number MEST-CT-2005-020513”.

## REFERENCES

- (1) Boucher, L. J. *Coord. Chem. Rev.* **1972**, *7*, 289–329.
- (2) Scheidt, W. R. *Acc. Chem. Res.* **1977**, *10*, 339–345.
- (3) Zerner, M.; Gouterma, M. *Theor. Chim. Acta* **1966**, *4*, 44.
- (4) Gunter, M. J.; Turner, P. *Coord. Chem. Rev.* **1991**, *108*, 115–161.
- (5) La Mar, G. N.; Walker, F. A. *The Porphyrins*; Academic Press: New York, 1979; Vol. 4.
- (6) Wang, C. Q.; Shalyaev, K. V.; Bonchio, M.; Carofiglio, T.; Groves, J. T. *Inorg. Chem.* **2006**, *45*, 4769–4782.
- (7) Collman, J. P.; Zhang, X. M.; Lee, V. J.; Uffelman, E. S.; Brauman, J. I. *Science* **1993**, *261*, 1404–1411.
- (8) Groves, J. T. *Proc. Natl. Acad. Sci. U.S.A.* **2003**, *100*, 3569–3574.
- (9) Beletskaya, I.; Tyurin, V. S.; Tsivadze, A. Y.; Guillard, R.; Stern, C. *Chem. Rev.* **2009**, *109*, 1659–1713.
- (10) Kulikov, O. V.; Schmidt, I.; Muresan, A. Z.; Lee, M. A. P.; Bocian, D. F.; Lindsey, J. S. *J. Porph. Phtalo.* **2007**, *11*, 699–712.
- (11) Li, C.; Fan, W. D.; Lei, B.; Zhang, D. H.; Han, S.; Tang, T.; Liu, X. L.; Liu, Z. Q.; Asano, S.; Meyyappan, M.; Han, J.; Zhou, C. W. *App. Phys. Lett.* **2004**, *84*, 1949–1951.
- (12) Li, Q. L.; Mathur, G.; Gowda, S.; Surthi, S.; Zhao, Q.; Yu, L. H.; Lindsey, J. S.; Bocian, D. F.; Misra, V. *Adv. Mater.* **2004**, *16*, 133–137.
- (13) Duclairioir, F.; Dubois, L.; Calborean, A.; Fateeva, A.; Fleury, B.; Kalaiselvan, A.; Marchon, J. C.; Maldivi, P.; Billon, M.; Bidan, G.; de Salvo, B.; Delapierre, G.; Buckley, J.; Huang, K.; Barattin, R.; Pro, T. *Int. J. Nanotechnol.* **2010**, *7*, 719–737.
- (14) Daku, L. M. L.; Castaings, A.; Marchon, J. C. *Inorg. Chem.* **2009**, *48*, 5164–5176.
- (15) Liao, M. S.; Scheiner, S. J. *Comput. Chem.* **2002**, *23*, 1391–1403.
- (16) Liao, M. S.; Scheiner, S. J. *Chem. Phys.* **2002**, *117*, 205–219.
- (17) Ghosh, A.; Vangberg, T.; Gonzalez, E.; Taylor, P. J. *J. Porph. Phtalo.* **2001**, *5*, 345–356.
- (18) Liu, C. G.; Guan, W.; Song, P.; Yan, L. K.; Su, Z. M. *Inorg. Chem.* **2009**, *48*, 6548–6554.
- (19) Scherlis, D. A.; Estrin, D. A. *Int. J. Quantum Chem.* **2002**, *87*, 158–166.
- (20) Leung, K.; Rempe, S. B.; Schultz, P. A.; Sproviero, E. M.; Batista, V. S.; Chandross, M. E.; Medforth, C. J. *J. Am. Chem. Soc.* **2006**, *128*, 3659–3668.

- (21) Baerends, E. J.; Ricciardi, G.; Rosa, A.; van Gisbergen, S. J. A. *Coord. Chem. Rev.* **2002**, *230*, 5–27.
- (22) Liao, M. S.; Watts, J. D.; Huang, M. J. *J. Comput. Chem.* **2006**, *27*, 1577–1592.
- (23) Vancoillie, S.; Zhao, H. L.; Radon, M.; Pierloot, K. *J. Chem. Theory Comput.* **2010**, *6*, 576–582.
- (24) Gonzalez, B.; Kouba, J.; Yee, S.; Reed, C. A. *J. Am. Chem. Soc.* **1975**, *97*, 3247–3249.
- (25) Kirner, J. F.; Reed, C. A.; Scheidt, W. R. *J. Am. Chem. Soc.* **1975**, *97*, 2557–2563.
- (26) Kirner, J. F.; Reed, C. A.; Scheidt, W. R. *J. Am. Chem. Soc.* **1977**, *99*, 1093–1101.
- (27) Liao, M. S.; Watts, J. D.; Huang, M. J. *Inorg. Chem.* **2005**, *44*, 1941–1949.
- (28) Salomon, O.; Reiher, M.; Hess, B. A. *J. Chem. Phys.* **2002**, *117*, 4729–4737.
- (29) Swart, M. *Inorg. Chim. Acta* **2007**, *360*, 179–189.
- (30) Swart, M.; Groenhof, A. R.; Ehlers, A. W.; Lammertsma, K. *J. Phys. Chem. A* **2004**, *108*, 5479–5483.
- (31) Liao, M. S.; Watts, J. D.; Huang, M. J. *J. Phys. Chem. A* **2007**, *111*, 5927–5935.
- (32) Pierloot, K.; Vancoillie, S. *J. Chem. Phys.* **2008**, *128*, 34104.
- (33) Ghosh, A.; Persson, B. J.; Taylor, P. R. *J. Biol. Inorg. Chem.* **2003**, *8*, 507–511.
- (34) Conradie, J.; Ghosh, A. *J. Chem. Theory Comput.* **2007**, *3*, 689–702.
- (35) Noodleman, L.; Han, W. G. *J. Biol. Inorg. Chem.* **2006**, *11*, 674–694.
- (36) Reiher, M.; Salomon, O.; Hess, B. A. *Theor. Chem. Acc.* **2001**, *107*, 48–55.
- (37) Shaik, S.; Chen, H.; Janardanan, D. *Nature Chem.* **2011**, *3*, 19–27.
- (38) Sorkin, A.; Iron, M. A.; Truhlar, D. G. *J. Chem. Theory Comput.* **2008**, *4*, 307–315.
- (39) Swart, M. *J. Chem. Theory Comput.* **2008**, *4*, 2057–2066.
- (40) Schultz, N. E.; Zhao, Y.; Truhlar, D. G. *J. Phys. Chem. A* **2005**, *109*, 4388–4403.
- (41) Handy, N. C.; Cohen, A. J. *Mol. Phys.* **2001**, *99*, 403–412.
- (42) Fouqueau, A.; Casida, M. E.; Lawson Daku, L. M.; Hauser, A.; Neese, F. *J. Chem. Phys.* **2005**, *122*, 44110.
- (43) Perdew, J. P.; Ruzsinszky, A.; Constantin, L. A.; Sun, J. W.; Csonka, G. I. *J. Chem. Theory Comput.* **2009**, *5*, 902–908.
- (44) Zhao, Y.; Truhlar, D. G. *Theor. Chem. Acc.* **2008**, *120*, 215–241.
- (45) Staroverov, V. N.; Scuseria, G. E.; Tao, J. M.; Perdew, J. P. *J. Chem. Phys.* **2003**, *119*, 12129–12137.
- (46) Perdew, J. P.; Kurth, S.; Zupan, A.; Blaha, P. *Phys. Rev. Lett.* **1999**, *82*, 2544–2547.
- (47) Ganzenmuller, G.; Berkaine, N.; Fouqueau, A.; Casida, M. E.; Reiher, M. *J. Chem. Phys.* **2005**, *122*.
- (48) Roos, B. O.; Taylor, P. R. *Chem. Phys.* **1980**, *48*, 157–173.
- (49) Becke, A. D. *Phys. Rev. A* **1988**, *38*, 3098–3100.
- (50) Perdew, J. P. *Phys. Rev. B* **1986**, *33*, 8822–8824.
- (51) Perdew, J. P.; Burke, K.; Ernzerhof, M. *Phys. Rev. Lett.* **1996**, *77*, 3865–3868.
- (52) Lee, C. T.; Yang, W. T.; Parr, R. G. *Phys. Rev. B* **1988**, *37*, 785–789.
- (53) Zhao, Y.; Truhlar, D. G. *J. Chem. Phys.* **2006**, *125*, 18.
- (54) Becke, A. D. *J. Chem. Phys.* **1993**, *98*, 1372–1377.
- (55) Perdew, J. P.; Burke, K.; Ernzerhof, M. *J. Chem. Phys.* **1996**, *105*, 9982.
- (56) Kang, J. K.; Musgrave, C. B. *J. Chem. Phys.* **2001**, *115*, 11040.
- (57) Cheng, B. S.; Fries, P. H.; Marchon, J. C.; Scheidt, W. R. *Inorg. Chem.* **1996**, *35*, 1024–1032.
- (58) Karlstrom, G.; Lindh, R.; Malmqvist, P. A.; Roos, B. O.; Ryde, U.; Veryazov, V.; Widmark, P. O.; Cossi, M.; Schimmelpfennig, B.; Neogrady, P.; Seijo, L. *Comput. Mater. Sci.* **2003**, *28*, 222–239.
- (59) Andersson, K.; Malmqvist, P. A.; Roos, B. O. *J. Chem. Phys.* **1992**, *96*, 1218–1226.
- (60) Roos, B. O.; Lindh, R.; Malmqvist, P. A.; Veryazov, V.; Widmark, P. O. *J. Phys. Chem. A* **2004**, *108*, 2851–2858.
- (61) Roos, B. O.; Lindh, R.; Malmqvist, P. A.; Veryazov, V.; Widmark, P. O. *J. Phys. Chem. A* **2005**, *109*, 6575–6579.
- (62) Cheng, R. J.; Lee, C. H.; Chao, C. W. *Chem. Commun.* **2009**, 2526–2528.
- (63) Cheng, R. J.; Wang, Y. K.; Chen, P. Y.; Han, Y. P.; Chang, C. C. *Chem. Commun.* **2005**, 1312–1314.
- (64) Andersson, K.; Malmqvist, P. A.; Roos, B. O.; Sadlej, A. J.; Wolinski, K. *J. Phys. Chem.* **1990**, *94*, 5483–5488.
- (65) Sadoc, A.; Broer, R.; De Graaf, C. *J. Chem. Phys.* **2007**, *126*, 134709.
- (66) Sadoc, A.; De Graaf, C.; Broer, R. *Phys. Rev. B* **2007**, *75*, 165116.
- (67) Kepenekian, M.; Robert, V.; Le Guennic, B.; De Graaf, C. *J. Comput. Chem.* **2009**, *30*, 2327–2333.
- (68) SCM; Theoretical Chemistry, Vrije Universiteit: Amsterdam, The Netherlands, 2010.
- (69) Fonseca Guerra, C.; Snijders, J. G.; te Velde, G.; Baerends, E. J. *Theor. Chem. Acc.* **1998**, *99*, 391–403.
- (70) te Velde, G.; Bickelhaupt, F. M.; Baerends, E. J.; Fonseca Guerra, C.; Van Gisbergen, S. J. A.; Snijders, J. G.; Ziegler, T. *J. Comput. Chem.* **2001**, *22*, 931–967.
- (71) Vangberg, T.; Ghosh, A. *J. Am. Chem. Soc.* **1999**, *121*, 12154–12160.
- (72) Paulat, F.; Praneeth, V. K. K.; Nather, C.; Lehnert, N. *Inorg. Chem.* **2006**, *45*, 2835–2856.
- (73) Kirner, J. F.; Scheidt, W. R. *Inorg. Chem.* **1975**, *14*, 2081–2086.
- (74) Ghigo, G.; Roos, B. O.; Malmqvist, P. A. *Chem. Phys. Lett.* **2004**, *396*, 142–149.
- (75) Queralt, N.; Taratiel, D.; de Graaf, C.; Caballor, R.; Cimiraaglia, R.; Angeli, C. *J. Comput. Chem.* **2008**, *29*, 994–1003.
- (76) Kepenekian, M.; Robert, V.; Le Guennic, B. *J. Chem. Phys.* **2009**, *131*.
- (77) Deeth, R. J.; Fey, N. *J. Comput. Chem.* **2004**, *25*, 1840–1848.
- (78) Ghosh, A. *J. Biol. Inorg. Chem.* **2006**, *11*, 712–724.
- (79) Conradie, J.; Ghosh, A. *J. Phys. Chem. B* **2007**, *111*, 12621–12624.
- (80) Fouqueau, A.; Mer, S.; Casida, M. E.; Daku, L. M. L.; Hauser, A.; Mineva, T.; Neese, F. *J. Chem. Phys.* **2004**, *120*, 9473–9486.
- (81) Zein, S.; Borshch, S. A.; Fleurat-Lessard, P.; Casida, M. E.; Chermette, H. *J. Chem. Phys.* **2007**, *126*, 14105.
- (82) *Advanced Inorganic Chemistry*, 6th ed.; Cotton, F. A., Wilkinson, G., Eds.; Wiley - Interscience: New York, 1999.
- (83) Hebenanz, N.; Kohler, F. H.; Muller, G.; Riede, J. *J. Am. Chem. Soc.* **1986**, *108*, 3281–3289.
- (84) Heise, H.; Kohler, F. H.; Xie, X. L. *J. Magn. Reson.* **2001**, *150*, 198–206.
- (85) Ghosh, A. *J. Biol. Inorg. Chem.* **2006**, *11*, 671–673.
- (86) Mattson, A. E. *Science* **2002**, *298*, 759–760.

# Energy-Specific Linear Response TDHF/TDDFT for Calculating High-Energy Excited States

Wenkel Liang,<sup>†</sup> Sean A. Fischer,<sup>†</sup> Michael J. Frisch,<sup>‡</sup> and Xiaosong Li<sup>\*,†</sup>

<sup>†</sup>Department of Chemistry, University of Washington, Seattle, Washington, United States 98195

<sup>‡</sup>Gaussian, Inc., 340 Quinnipiac St Bldg 40, Wallingford, Connecticut 06492, United States

**ABSTRACT:** An energy-specific TDHF/TDDFT method is introduced in this article for excited state calculations. This approach extends the conventional TDHF/TDDFT implementation to obtain excited states above a predefined energy threshold. The method introduced and developed in this work enables computationally efficient yet rigorous calculations of energy-specific spectra, e.g., X-ray absorption involving extremely high-energy transitions. All transitions are solved in the full molecular orbital space, and orthogonality to the ground state and lower-lying excited states is preserved for each high-energy excited state. Encouraging computational savings are observed in calculating the targeted energy spectrum, while the transition energies, as well as oscillator strengths, remain identical to the results from the standard implementation.

## I. INTRODUCTION

Single-reference methods such as configuration interaction singles (CIS)<sup>1</sup> and the linear-response variants of time-dependent Hartree–Fock (TDHF)<sup>2,3</sup> and time-dependent density functional theory (TDDFT)<sup>4–7</sup> are widely used for *ab initio* calculations of electronic excited states for large molecular systems because of their balance of computational efficiency and accuracy for practical applications.<sup>8–11</sup> Highly correlated methods, such as symmetry adapted cluster/configuration interaction (SAC–CI<sup>12</sup>), linear response coupled cluster (LRCC<sup>13</sup>), and equation-of-motion coupled cluster (EOM-CC<sup>14,15</sup>) and multi-reference approaches, such as multireference configuration interaction (MRCI<sup>16</sup>) and multireference perturbation theories (MRMP<sup>17</sup> and CASPT2<sup>18</sup>) are capable of providing more accurate treatments of excited states, including those with multielectron excitation character. However, these methods are generally computationally prohibitive for large molecules.

Conventional TDHF and TDDFT are subject to some non-trivial problems. For example, excitation energies for Rydberg and charge transfer states are often underestimated. The latter can be improved with the range-separated class of hybrid DFT functionals.<sup>19–22</sup> Neither TDHF nor TDDFT properly include the effects of dispersion. However, efforts have been made to include dispersion in DFT functionals, and promising results have been obtained.<sup>23–25</sup> The lack of correlation, or approximate nature of the treatment thereof, in standard implementations of TDHF and TDDFT results in the methods being unable to correctly describe excited states with multielectron excitation character.<sup>26,27</sup> In spite of these limitations, the TDHF and TDDFT methods can generally be expected to reproduce trends for one-electron valence excitations, which contribute to a majority of transitions of photochemical interest. TDDFT using hybrid density functionals, in particular, has been successful in modeling the optical absorption spectra of large molecules.<sup>28</sup> Recent studies have further extended the application of TDDFT to predict very high-energy, core–electron excitations that account for the pre-edge features in X-ray absorption spectroscopy (XAS).<sup>29–31</sup>

A simple frozen-orbital method has been proposed and found to be very effective in obtaining core orbital excitations.<sup>32–35</sup>

Although linear-response TDHF and TDDFT are among the most tractable methods for excited state calculations, they can still be computationally demanding for large molecular systems of photochemical interest. The numerical cost of solving the TDHF/TDDFT equations using iterative techniques formally scales as  $O(MN^4)$ , where  $N$  is the total number of basis functions and  $M$  is the number of excited states sought. With development of effective Krylov subspace algorithms and linear-scaling methods for direct Fock/Kohn–Sham operator builders, conventional implementations of the computational scaling of linear-response TDHF and TDDFT equations can be reduced to  $O(MN^2) - O(MN^3)$  in complexity. Detailed studies on the numerical algorithms for solving the TDHF/TDDFT equations are available in refs 36 and 37. Notably, if many states are to be obtained simultaneously, efficiency degrades considerably as a result of increased memory and I/O requirements. For large-scale systems, the computational bottleneck of orbital transformation can be avoided by using “orbital-free” approaches.<sup>38,39</sup>

In cases where a certain high-energy excited state is the subject of interest, it is possible to obtain an approximate solution of the linear-response equation of TDHF/TDDFT only in the small energy-range of interest. Along those lines, Kauczor et al.<sup>37</sup> have demonstrated that the optimal algorithm for solving the standard and damped complex response equation<sup>40,41</sup> is the preconditioned iterative subspace algorithm with symmetrized trial vectors, and the use of complex damping allows for the determination of higher excited states without knowledge of lower state solutions. Tretiak et al.<sup>36</sup> proposed using a symmetric Wilkinson shift<sup>42</sup> to acquire higher-energy excited states when solving the TDHF/TDDFT equations in an orbital-independent formulation. A response function using only a subset (e.g., core orbitals) of the molecular orbitals has also been developed to reduce the

Received: July 11, 2011

Published: September 26, 2011



cost of excited state calculations using TDHF/TDDFT. For example, the frozen-orbital method truncates the molecular orbital (MO) space and considers only the transitions between core and valence orbitals.<sup>32,33</sup> These methods are efficient in obtaining approximate energies and electronic characteristics of high-energy excited states; however, the application is limited to core orbitals with restricted radial extent that are nearly orthogonal to all other occupied orbitals and excited states. When such orthogonality is not well preserved, calculated excited states with incomplete orbital space can lead to inaccurate oscillator strengths and unphysical electron distributions.

In this article, we introduce an energy-specific TDHF/TDDFT (ES-TD) approach to selectively calculate absorption spectra above a predefined energy threshold while maintaining the orthogonality between excited states and the ground state. The algorithm introduced herein is rigorous because the solutions are exact in the full molecular orbital space. It is based on a simple yet effective idea to bracket high-energy spectra using a Davidson-like iterative algorithm<sup>43,44</sup> of the spin-unrestricted TDHF/TDDFT implementation.<sup>9</sup> Computational performance and accuracy are compared for Rydberg excited states of an alanine dimer, ligand-to-metal charge transfer transitions in Mn<sup>2+</sup>-doped ZnO semiconductor nanocrystals, and high-energy X-ray absorptions in a set of metal tetrachlorides.

## II. METHODOLOGY

In conventional linear-response TDHF/TDDFT theory, excitation energies  $\omega$  can be determined by solving the non-Hermitian eigenvalue equation, given in matrix form as<sup>6,8,9</sup>

$$\begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B} & \mathbf{A} \end{pmatrix} \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} = \omega \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix} \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} \quad (1)$$

with the matrices for TDHF

$$\begin{aligned} A_{ia,jb} &= \delta_{ij}\delta_{ab}(\varepsilon_a - \varepsilon_i) + (ia|jb) - (ib|ja) \\ B_{ia,jb} &= (ia|bj) - (ij|ba) \end{aligned} \quad (2)$$

and for TDDFT

$$\begin{aligned} A_{ia,jb} &= \delta_{ij}\delta_{ab}(\varepsilon_a - \varepsilon_i) + (ia|jb) - \alpha(ib|ja) + (ia|f_{xc}|jb) \\ B_{ia,jb} &= (ia|bj) - \alpha(ij|ba) + (ia|f_{xc}|bj) \end{aligned} \quad (3)$$

where  $\mathbf{X}$  and  $\mathbf{Y}$  are the first order electron density responses determined by solving this system of linear equations. The regular two-electron integrals are expressed in Mulliken notation. For hybrid DFT, the HF exchange integral takes on a fractional value scaled by a nonzero scaling factor  $\alpha$ , while  $\alpha = 0$  for pure DFT kernels. The response of the exchange-correlation (xc) potential term, also called the xc kernel, is given as

$$(ia|f_{xc}|jb) = \int \int \phi_i^*(r) \phi_a(r) \frac{\delta^2 E_{xc}}{\delta \rho(r) \delta \rho(r')} \phi_j^*(r') \phi_b(r') dr dr' \quad (4)$$

The  $i$  and  $j$  and the  $a$  and  $b$  indices represent occupied and virtual molecular orbitals (MOs), respectively, in the HF/Kohn–Sham ground state configuration.

For real orbitals, eq 1 can be reduced to a non-Hermitian (eq 5) or Hermitian (eq 6) eigenvalue equation with half the

dimension

$$(\mathbf{A} - \mathbf{B})(\mathbf{A} + \mathbf{B})|\mathbf{X} + \mathbf{Y}\rangle = \omega^2|\mathbf{X} + \mathbf{Y}\rangle \quad (5)$$

$$(\mathbf{A} - \mathbf{B})^{1/2}(\mathbf{A} + \mathbf{B})(\mathbf{A} - \mathbf{B})^{1/2}|\mathbf{T}\rangle = \omega^2|\mathbf{T}\rangle \quad (6)$$

$$|\mathbf{T}\rangle = (\mathbf{A} - \mathbf{B})^{-1/2}(\mathbf{X} + \mathbf{Y}) \quad (7)$$

If all solutions of these equations are sought,  $\mathbf{A}$  and  $\mathbf{B}$  include all transitions between occupied and unoccupied molecular orbitals. The size of  $\mathbf{A}$  and  $\mathbf{B}$  in the molecular orbital space is  $(N_{\text{occ}} \times N_{\text{unocc}})^2$  where  $N_{\text{occ}}$  and  $N_{\text{unocc}}$  are the numbers of occupied and unoccupied molecular orbitals. Such a full treatment has a significantly large computational cost. An efficient algorithm to solve the response equation of TDHF/TDDFT was introduced by Stratmann et al. based on the nonsymmetric Davidson diagonalization algorithms of Hirao and Nakatsuji<sup>44</sup> and Bouman et al.<sup>45</sup> for obtaining lower-lying excited states.<sup>9</sup> The idea is to solve the response function in the reduced form of eq 5 or 6. The solutions can be obtained by a symmetric diagonalization or Davidson's algorithm.<sup>43,44</sup> In this work, we introduce additional algorithms to bracket trial vectors and eigenvalues in the predefined energy range to eventually obtain high-energy excitation energies and transitions without the effort of scanning through lower-lying excited states. We choose the Stratmann method<sup>9</sup> as the basic equation solver integrated with an energy screening and bracketing idea. For simplicity, we will only present the algorithm within the framework of the Hermitian eigenvalue equation (eq 6). The non-Hermitian version of the method can be readily obtained with simple matrix transformations.

The following discussions assume that  $M$  excited states with energies greater than  $\omega_0$  are the subject of interest. The  $(\mathbf{A} + \mathbf{B})$  and  $(\mathbf{A} - \mathbf{B})$  matrices in eq 6 are first projected onto a subspace spanned by a set of zeroth order trial vectors  $\mathbf{C} = \{\mathbf{b}_1, \dots, \mathbf{b}_l\}$  where  $l > M$

$$\tilde{\mathbf{M}}^+ = \mathbf{C}^T(\mathbf{A} + \mathbf{B})\mathbf{C} \quad (8)$$

$$\tilde{\mathbf{M}}^- = \mathbf{C}^T(\mathbf{A} - \mathbf{B})\mathbf{C} \quad (9)$$

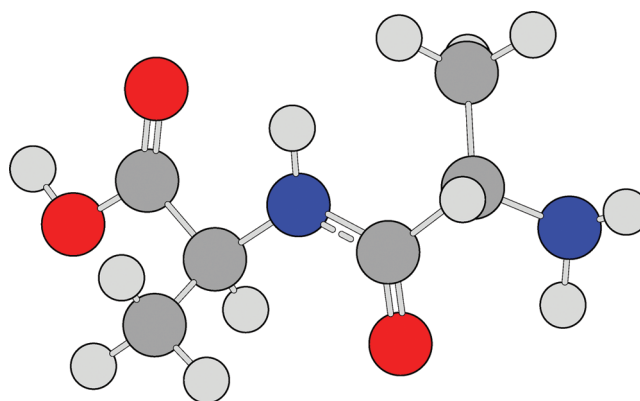
$$\tilde{\mathbf{M}} = (\tilde{\mathbf{M}}^-)^{1/2}(\tilde{\mathbf{M}}^+)(\tilde{\mathbf{M}}^-)^{1/2} \quad (10)$$

where the dimension of the resulting matrices is  $l$ . Because the number of trial vectors  $l$  is much smaller than  $N_{\text{occ}} \times N_{\text{unocc}}$ , the computational cost of directly generating the resulting matrices in eq 6 is greatly reduced. On the other hand, the initial  $l$  needs to be much larger than the requested number of excited states to include all MO transition candidates that may contribute significantly to excitations in the desired energy range. In the current implementation, we construct initial trial vectors by sampling the Koopmans' MO transitions for the requested energy range.  $l = 4M$  trial vectors are generated in the first step, corresponding to the lowest energy Koopmans' transitions with a constraint of

$$\varepsilon_a - \varepsilon_i \geq \omega_0 + \delta\omega \quad (11)$$

The  $\delta\omega$  energy shift is used to approximate corrections for the errors in Koopmans' transitions for a better selection of initial trial vectors. Note that a good selection of initial trial vectors is important for a fast convergence but generally does not affect the quality of the final results because vectors that are found to

**Table 1.** Comparison of Select Rydberg States of an Alanine Dimer Computed Using the Regular TDHF, ES-TDHF, and a Subset Space with Removal of HOMO TDHF Methods with the aug-cc-pvdz Basis Set



regular TDHF		ES-TDHF		subset TDHF	
excitation energy (eV)	oscillator strength (a.u.)	excitation energy (eV)	oscillator strength (a.u.)	excitation energy (eV)	oscillator strength (a.u.)
8.0767	0.1048	8.0767	0.1048	8.3037	0.0656
8.4052	0.0647	8.4052	0.0647	8.8951	0.0848
8.6152	0.0280	8.6152	0.0280	9.0088	0.0298
9.0680	0.0061	9.0680	0.0061	9.1274	0.0085
9.1148	0.0469	9.1148	0.0469	9.2649	0.0739

contribute to the excitations will be added into subspace **C** during later iterations until the convergence is achieved.

Once the reduced subspace is constructed, diagonalization of  $\tilde{\mathbf{M}}$  generates a set of eigenvalues  $\tilde{\omega}$  and eigenvectors  $\tilde{\mathbf{T}}$  in the reduced space. In the reduced space, qualified eigenvalues and eigenvectors are selected according to the predefined requirement for excitation energy and number of excited states:

$$\tilde{\omega}_i \geq \omega_0, i = n \dots (n + M) \quad (12)$$

where  $\tilde{\omega}_n$  is the lowest eigenvalue that is greater than the energy threshold  $\omega_0$ . We define eigenvalues and eigenvectors that satisfy eq 12 as the qualified candidates in the reduced space, denoted as  $\tilde{\omega}_M$  and  $\tilde{\mathbf{T}}_M$ . The corresponding collective transition densities,  $(\tilde{\mathbf{X}} + \tilde{\mathbf{Y}})_M$  and  $(\tilde{\mathbf{X}} - \tilde{\mathbf{Y}})_M$ , can be obtained as well. These candidates can be transformed from the reduced space to the full MO space,

$$(\mathbf{X}' + \mathbf{Y}')_M = \mathbf{C}(\tilde{\mathbf{X}} + \tilde{\mathbf{Y}})_M \quad (13)$$

$$(\mathbf{X}' - \mathbf{Y}')_M = \mathbf{C}(\tilde{\mathbf{X}} - \tilde{\mathbf{Y}})_M \quad (14)$$

$$\omega'_M = \tilde{\omega}_M \quad (15)$$

where the primed notation refers to approximate solutions in the complete MO space. In order to estimate the errors associated with the approximate solutions, residual vectors can be defined as (see ref 9 for detailed discussions)

$$\mathbf{W}_M^L = (\mathbf{A} + \mathbf{B})(\mathbf{X}' + \mathbf{Y}')_M - \omega'_M(\mathbf{X}' - \mathbf{Y}')_M \quad (16)$$

$$\mathbf{W}_M^R = (\mathbf{A} - \mathbf{B})(\mathbf{X}' - \mathbf{Y}')_M - \omega'_M(\mathbf{X}' + \mathbf{Y}')_M \quad (17)$$

When the norm of a residual vector is below a certain small threshold ( $10^{-6}$  au in this work), the associated excited state is

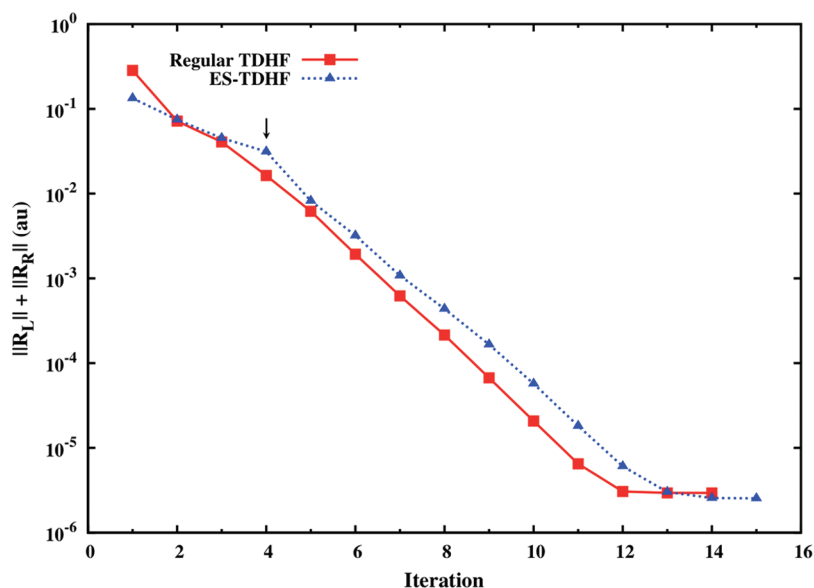
considered converged. For those unconverged excited states, a new set of vectors can be constructed following the Davidson algorithm:

$$\mathbf{Q}_M^L = (\omega'_M - \Delta\varepsilon)^{-1} \mathbf{W}_M^L \quad (18)$$

$$\mathbf{Q}_M^R = (\omega'_M - \Delta\varepsilon)^{-1} \mathbf{W}_M^R \quad (19)$$

where the  $\Delta\varepsilon$  is the orbital energy difference. These new vectors will be orthonormalized and added into the subspace **C**, and new iteration starts from eq 8 until all vectors are converged. Note that a monotonic convergence of the reduced eigenspace can be observed because MacDonald's theorem<sup>46</sup> applies when the reduced subspace is unchanged.<sup>9</sup> Usually after a few initial iterations, the subspace of interest becomes well-defined and remains the same. A monotonic convergence can then be observed.

The algorithm introduced above is a simple extension with subspace bracketing and energy spectrum selection to the Stratmann method based on the Davidson algorithm. It is worth noting that because linear transformations are always carried out of  $(\mathbf{A} + \mathbf{B})$  and  $(\mathbf{A} - \mathbf{B})$  matrices, symmetrized trial vectors are introduced in each iteration, which gives optimal efficiency for standard response equation as suggested by Kauczor et al.<sup>37</sup> This approach is different from other approximate methods<sup>34,35</sup> that make use of an incomplete MO space and neglect contributions from other minor transitions. Note that the reduced subspace is used merely for the sake of obtaining solutions that correspond to a desired energy range at low computational cost. The convergence is verified in the full MO space. The size of the reduced subspace expands during iterations to include all significant transition pairs. The final results are true eigenvalues and eigenvectors of the linear-response TDHF/TDDFT equations. The resulting



**Figure 1.** Convergence performance of the regular TDHF and the ES-TDHF methods for the first Rydberg state of alanine dimer. Residual norm is plotted against iteration number in the Davidson algorithm. The arrow indicates when the selected subspace starts to become stable and remain constant.

high-energy excited states are intrinsically orthogonal to lower lying states because eigenvectors in the full space corresponding to different eigenvalues are automatically biorthogonal. Numerical tests in the next section will show that not only can this method directly and accurately obtain high-energy excited states without scanning through lower-lying ones, the calculated excited states also maintain orthogonality to the ground state and to lower-lying states that are skipped in the calculations.

### III. BENCHMARKS AND DISCUSSION

Calculations were carried out on a Dell PowerEdge R610 Server (dual quad-core 2.4 GHz Intel Xeon with 16 GB of RAM), using the development version of the Gaussian series of programs<sup>47</sup> with the addition of energy-specific linear-response TDHF/TDDFT approach presented here. The computational time reported in this article is the absolute total CPU time. In the next two sections, we will test the ES method on the Rydberg states of an alanine dimer, charge transfer excitation in a 1.0 nm quantum dot doped with a transition metal, and X-ray absorption spectra of a series of metal tetrachlorides.

**A. Rydberg States of Alanine Dimer.** Rydberg states of a molecule are generally associated with characteristics of high excitation energies and diffusive electronic distributions. Calculations of Rydberg states usually require large basis sets with diffusive functions, and the transition vectors of these states strongly depend on many orbitals. As a result, the excited states are very sensitive to the quality of calculations. This can be considered a stringent test case for the ES method developed herein. Table 1 lists the excitation energies and oscillator strengths of select Rydberg states of an alanine dimer molecule computed at the TDHF/aug-cc-pvdz level of theory. The first excited state lies  $\sim 6.6$  eV above ground state. An energy threshold of 8 eV is used in the ES-TDHF method, which skips three lower energy valence states. The results from the first two different calculations using the regular approach and the ES-TDHF method are essentially identical. Residual norms of eqs 16 and 17 are plotted against iteration number in the Davidson algorithm in Figure 1. The

convergence performance of the two methods are very similar for the first Rydberg state, though ES-TDHF may take a few more iterations to reach final convergence due to smaller initial expansion vector space. Most importantly, both methods exhibit a monotonic convergence. Such behavior has been shown mathematically in a recent work by Kauczor et al.<sup>37</sup> The transition vectors are also in perfect agreement (errors  $< 10^{-5}$ ) in both magnitude and sign with those obtained from the regular TDHF calculation. This agreement suggests that even though lower-energy valence states were skipped in the ES-TDHF method, the resulting high-energy Rydberg states properly maintain orthogonality to the lower ones. In Table 1, we also included results from calculations using only a subset of occupied orbitals by forbidding transitions from the highest occupied orbital. Using such a frozen-orbital approach with TDHF, the calculated Rydberg type transition energies and oscillator strengths significantly deviate from the reference values. Such deviations are a result of neglecting the strong couplings between active and frozen orbitals and the inability to preserve the orthogonality to lower-lying excited states.

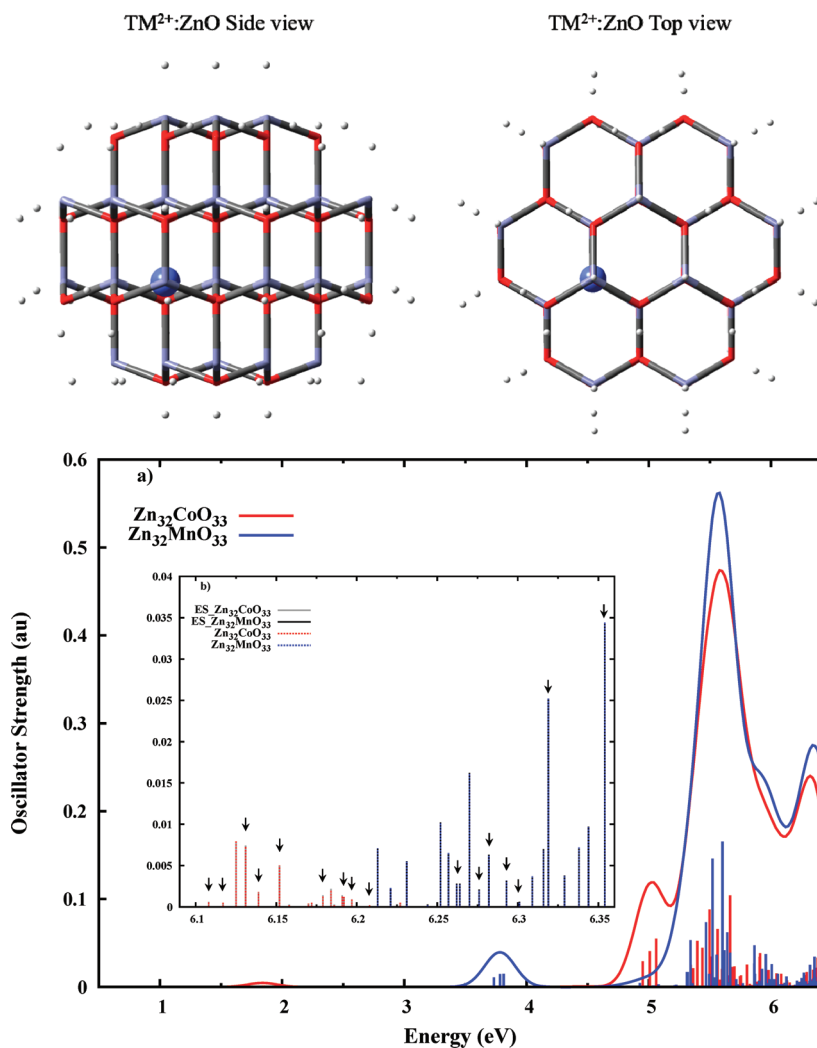
**B.  $L_{VB}MCT$  Transitions in Semiconductor Nanocrystals.** In our recent studies<sup>48,49</sup> of  $Co^{2+}$ - and  $Mn^{2+}$ -doped ZnO semiconductor nanocrystals, two types of charge transfer (CT) transitions were characterized by TDDFT: metal-to-ligand CT as the promotion of a transition metal (TM) dopant d electron to the ZnO conduction band ( $ML_{CB}CT$ ) and ligand-to-metal CT as the promotion of a ZnO valence band electron to a vacant transition metal dopant d orbital ( $L_{VB}MCT$ ). On the basis of theoretical calculation and experimental observations,<sup>49–51</sup> the  $L_{VB}MCT$  transitions are always higher in energy than the  $ML_{CB}CT$  ones in ZnO nanocrystals. The  $L_{VB}MCT$  transition is at  $\sim 6.0$  eV for the  $Zn_{32}TMO_{33}$  nanocrystal of  $\sim 1.0$  nm diameter, while the  $ML_{CB}CT$  transitions take place above  $\sim 1.75$  eV.<sup>52</sup> For detailed discussions about the characteristic transitions in a diluted magnetic semiconductor, we refer readers to our recent work using linear response TDDFT.<sup>49</sup>

Doped nanocrystal structures were constructed on the basis of the scheme described in ref 48, and the ground state electronic

**Table 2.** Comparison of Computational Costs for Obtaining  $L_{VB}MCT$  Transitions in  $Mn^{2+}$ - and  $Co^{2+}$ -Doped  $ZnO$  Nanocrystals<sup>a</sup>

nanocrystal systems	total number of AO	regular TDDFT			ES-TDDFT		
		total states (range in eV)	size of subspace	CPU time in h	total states (range in eV)	size of subspace	CPU time in h
$Zn_{32}CoO_{33}$	1015	120 (0.89–6.43)	1580	331.2 (1.0)	20 (6.11–6.23)	594	99.4 (0.30)
$Zn_{32}MnO_{33}$	1015	120 (3.73–6.45)	1526	319.0 (1.0)	20 (6.21–6.35)	576	92.5 (0.29)

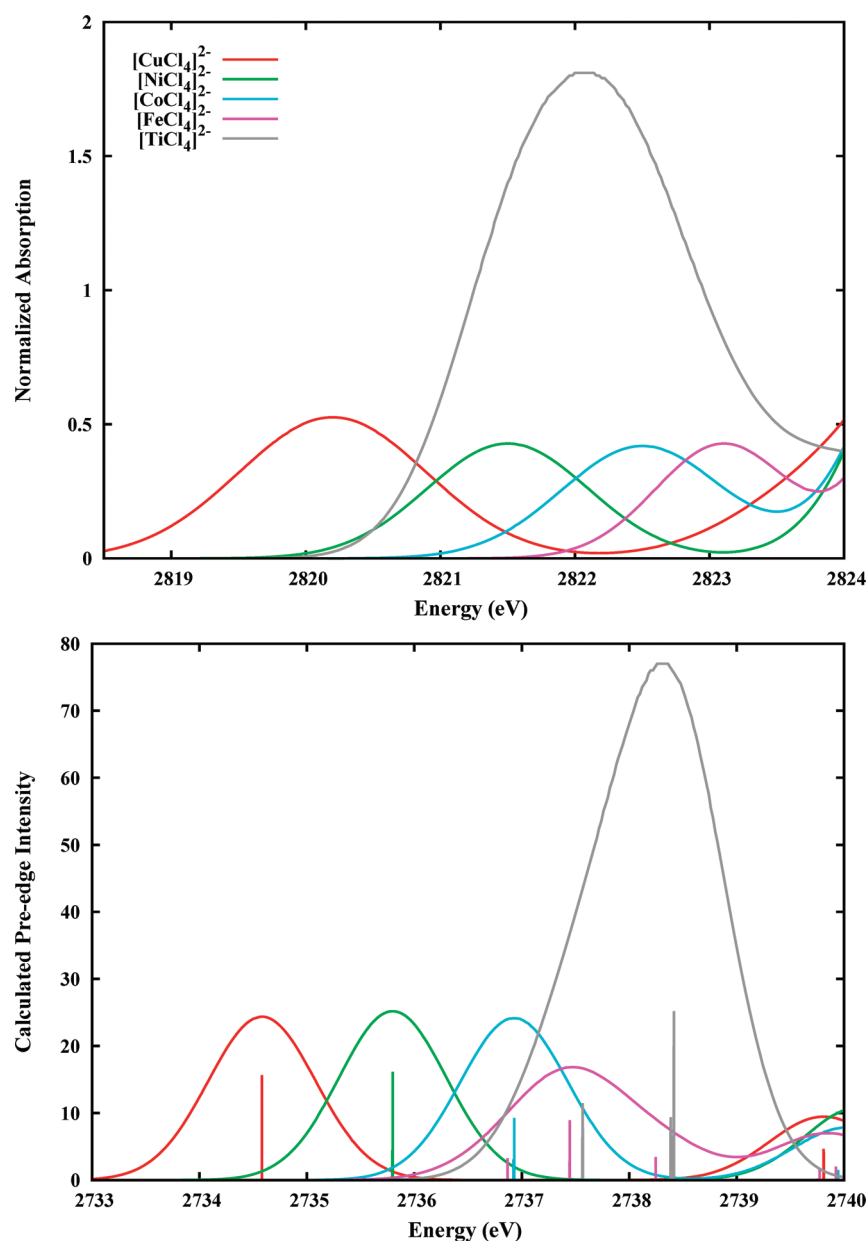
<sup>a</sup>The computational cost is evaluated using the total CPU time of the regular TDDFT method as the unit reference. Note that the wall clock time is one eighth of the total CPU time.



**Figure 2.** Top panel: Structure of the  $Zn_{32}TMO_{33}H^*_{60}$  nanocrystal. The  $TM^{2+}$  dopant ion is placed close to the center of the nanocrystal and is shown as a ball. Bottom panel: Comparison of optical transition oscillator strengths calculated for Mn- and Co-doped ZnO nanocrystals with (a) regular TDDFT method and (b) ES-TDDFT approach for higher energy  $L_{VB}MCT$  transitions. TDDFT peaks are dressed with Gaussian functions and a broadening parameter of 0.12 eV. Inset shows spectrum peaks in the high excitation energy region with  $L_{VB}MCT$  transitions identified with arrows.

structures and TDDFT spectra were obtained at the PBE1PBE<sup>53,54</sup>/LANL2DZ<sup>55–58</sup> level of theory. In order to characterize the  $L_{VB}MCT$  band in the absorption spectrum, a conventional linear response TDDFT approach would need to generate as many as 120 states to be able to reach the  $>6.0$  eV energy range. The ES-TDDFT method introduced herein can skip the lower lying excited states if a target excitation energy range is defined. Table 2 lists the computational costs of obtaining  $L_{VB}MCT$  transitions in  $Mn^{2+}$  and  $Co^{2+}$ -doped ZnO nanocrystals with both the standard

implementation and the ES-TDDFT approach. The energy threshold in ES-TDDFT is set to be 6.0 eV for the  $Zn_{32}TMO_{33}$  quantum dots. As shown in Figure 2, the modified algorithm yields excitation energies and oscillator strengths that are nearly identical to those obtained with the standard algorithm, with less than  $1 \times 10^{-4}$  eV numerical difference. The method also exhibits encouraging performance ( $\sim 30\%$  the cost of the regular calculation), even though the higher-energy transitions require more transition vectors, covering a much larger expansion vector space.



**Figure 3.** Comparison of experimental (top) Cl K-pre-edge XAS data to the calculated spectra (bottom). Experimental data are adapted from refs 67 and 68.

**C. Cl K-Edge XAS Spectra for Metal Complexes.** A series of metal tetrachlorides were constructed according to experimental X-ray structures,<sup>59–62</sup> and the Cl K-edge TDDFT XAS spectra were calculated with the BP86 functional<sup>63,64</sup> and TZVP basis set<sup>65</sup> following ref 35. XAS pre-edge features for five metal tetrachlorides ( $[\text{CuCl}_4]^{2-}$ ,  $[\text{NiCl}_4]^{2-}$ ,  $[\text{CoCl}_4]^{2-}$ ,  $[\text{FeCl}_4]^{2-}$ ,  $[\text{TiCl}_4]^0$ ) were calculated using ES-TDDFT with initial guess transitions originating from core orbitals. An energy threshold of 2700 eV was utilized to target the top edge of core electron excitations, i.e., a Cl 1s-core electron into a metal d-based MO.

For very high-energy transitions like X-ray absorptions that involve excitation of a core-electron, the spectroscopic oscillator strength needs to account for higher-order dipole interactions and takes on the form

$$f_I = f_I^{\text{ed}} + f_I^{\text{md}} + f_I^{\text{eq}} \quad (20)$$

where  $f_I^{\text{ed}}$ ,  $f_I^{\text{md}}$  and  $f_I^{\text{eq}}$  are electric dipole, magnetic dipole, and electric quadrupole oscillator strengths for the  $I$ th transition, respectively. They are given by the following expressions in atomic units<sup>35</sup>

$$\begin{aligned} f_I^{\text{ed}} &= \frac{2}{3} \omega_I |\langle \Psi_0 | \hat{r} | \Psi_I \rangle|^2 \\ f_I^{\text{md}} &= \frac{2}{3} \alpha^2 \omega_I |\langle \Psi_0 | \hat{l} + 2\hat{s} | \Psi_I \rangle|^2 \\ f_I^{\text{eq}} &= \frac{1}{20} \alpha^2 \omega_I^3 \sum_{i,j} \left| \left\langle \Psi_0 \left| \hat{r}_i \hat{r}_j - \frac{1}{3} r^2 \delta_{ij} \right| \Psi_I \right\rangle \right|^2 \end{aligned} \quad (21)$$

where  $\alpha$  is the dimensionless fine-structure constant given as  $1/137.03599$  and  $\omega_I$  is the excitation energy. Other notations

**Table 3. Calculated Electric Dipole, Magnetic Dipole and Electric Quadrupole Oscillator Strength Contributions to Major Core–Electron Excitations Shown in Figure 3**

metal complexes	excitation energy (eV)	electric dipole	magnetic dipole	electric quadrupole
[CuCl <sub>4</sub> ] <sup>2-</sup>	2734.58	<10 <sup>-3</sup>	~5 × 10 <sup>-3</sup>	15.67
	2739.75	0	0	0.21
	2739.81	<10 <sup>-3</sup>	0	4.66
[NiCl <sub>4</sub> ] <sup>2-</sup>	2735.80	<10 <sup>-3</sup>	0	16.14
	2740.17	<10 <sup>-3</sup>	0	4.32
[CoCl <sub>4</sub> ] <sup>2-</sup>	2736.93	<10 <sup>-3</sup>	0	9.23
	2739.94	<10 <sup>-3</sup>	0	1.51
	2740.08	<10 <sup>-3</sup>	0	2.55
[FeCl <sub>4</sub> ] <sup>2-</sup>	2736.87	<10 <sup>-3</sup>	0	3.27
	2737.44	<10 <sup>-3</sup>	0	8.95
	2738.25	0	~10 <sup>-3</sup>	3.46
	2739.77	<10 <sup>-3</sup>	0	1.70
	2739.92	<10 <sup>-3</sup>	0	2.01
[TiCl <sub>4</sub> ] <sup>2-</sup>	2737.57	~10 <sup>-3</sup>	<10 <sup>-3</sup>	11.99
	2738.39	~10 <sup>-3</sup>	~10 <sup>-3</sup>	9.80
	2738.42	~10 <sup>-3</sup>	0	26.48

used are  $\hat{r}$ , the position operator;  $\hat{l}$ , the angular momentum operator; and  $\hat{s}$ , the spin operator.

Figure 3 shows a comparison of the experimental Cl K pre-edge data to the calculated pre-edge spectra where the oscillator strengths account for magnetic dipole and electric quadrupole interactions in accordance with eqs 20 and 21. Oscillator strengths of electric dipole and electric quadrupole transitions are listed in Table 3. The calculated energies are underestimated (~85.6 eV on average), due to the limitations of DFT in modeling potentials near the nucleus, resulting in a Cl-1s orbital that is too high in energy relative to the valence orbitals.<sup>35,66</sup> The investigation of a systematic error in DFT such as this is beyond the scope of the present work. On the other hand, the calculated spectra are in good agreement with previously reported results from TDA type calculations.<sup>35</sup> The relative energies and intensities are also consistent with experimental observations.<sup>67,68</sup> The slight deviations from previous calculated spectra may arise from the use of an augmented basis set for the metal center and the dielectric continuum solvent that can lead to further stabilization of the valence orbitals.

#### IV. CONCLUSION

This paper presents an energy-specific TDHF/TDDFT method based on Davidson's iterative subspace algorithm. The method allows for the flexibility of only obtaining excitations above a predefined energy threshold enabling the prediction of an absorption spectrum over a specific energy range, such as extremely high-energy X-ray absorptions. Despite the fact that higher-energy transitions usually require a considerably large expansion vector space, this method shows an encouraging efficiency for large-scale systems. In about 30% of the time required for obtaining the full absorption spectra via the standard implementation of TDDFT, ES-TDDFT calculated the higher-energy ligand to metal charge transfer states for Mn<sup>2+</sup>- and Co<sup>2+</sup>-doped ZnO semiconducting nanocrystals. Most importantly, all of the calculated transitions are solved in the full MO space, ensuring that all

of the excited states maintain orthogonality to the ground state and lower-lying excited states. Because the MO space is not restricted while solving for the transitions of interest, the calculated eigenvalues and eigenvectors are true solutions to the linear-response TDHF/TDDFT equations. As a result, characteristics of excited states obtained using the ES-TDHF/TDDFT method are essentially identical to those computed using the regular TDHF/TDDFT implementations.

#### AUTHOR INFORMATION

##### Corresponding Author

\*E-mail: li@chem.washington.edu.

#### ACKNOWLEDGMENT

This work was supported by the U.S. National Science Foundation (CHE-CAREER 0844999 and CRC 0628252). Additional support from Gaussian Inc. and the University of Washington Student Technology Fund is gratefully acknowledged. Discussions with Ben Van Kuiken are greatly appreciated.

#### REFERENCES

- (1) Foresman, J. B.; Head-Gordon, M.; Pople, J. A.; Frisch, M. J. *J. Phys. Chem.* **1992**, *96*, 135.
- (2) Jorgensen, P.; Lindenberg, J. *Int. J. Quantum Chem.* **1970**, *4*, 587.
- (3) Olsen, J.; Jensen, H. J. A.; Jorgensen, P. *J. Comput. Phys.* **1988**, *74*, 265.
- (4) Runge, E.; Gross, E. K. U. *Phys. Rev. Lett.* **1984**, *52*, 997.
- (5) Gross, E. K. U.; Kohn, W. *Phys. Rev. Lett.* **1985**, *55*, 2850.
- (6) Casida, M. E. *Recent Adv. Comput. Chem.* **1995**, *1*, 155.
- (7) Casida, M. E.; Jamorski, C.; Casida, K. C.; Salahub, D. R. *J. Chem. Phys.* **1998**, *108*, 4439.
- (8) Hirata, S.; Head-Gordon, M.; Bartlett, R. J. *J. Chem. Phys.* **1999**, *111*, 10774.
- (9) Stratmann, R. E.; Scuseria, G. E.; Frisch, M. J. *J. Chem. Phys.* **1998**, *109*, 8218.
- (10) Burke, K.; Werschnik, J.; Gross, E. K. U. *J. Chem. Phys.* **2005**, *123*, 062206.
- (11) Dreuw, A.; Head-Gordon, M. *Chem. Rev.* **2005**, *105*, 4009.
- (12) Nakatsuji, H. *Chem. Phys. Lett.* **1979**, *67*, 329.
- (13) Koch, H. J. *Chem. Phys.* **1990**, *93*, 3345.
- (14) Stanton, J. F.; Bartlett, R. J. *J. Chem. Phys.* **1993**, *98*, 7029.
- (15) Krylov, A. I. *Annu. Rev. Phys. Chem.* **2008**, *59*, 433.
- (16) Dallos, M.; Lischka, H.; Shepard, R.; Yarkony, D. R.; Szalay, P. G. *J. Chem. Phys.* **2004**, *120*, 7330.
- (17) Kobayashi, Y.; Nakano, H.; Hirao, K. *Chem. Phys. Lett.* **2001**, *336*, 529.
- (18) Finley, J.; Malmqvist, P.-A.; Roos, B. O.; Serrano-Andres, L. *Chem. Phys. Lett.* **1998**, *288*, 299.
- (19) Iikura, H.; Tsuneda, T.; Yanai, T.; Hirao, K. *J. Chem. Phys.* **2001**, *115*, 3540.
- (20) Heyd, J.; Scuseria, G. E.; Ernzerhof, M. *J. Chem. Phys.* **2003**, *118*, 8207.
- (21) Tawada, Y.; Tsuneda, T.; Yanagisawa, S.; Yanai, T.; Hirao, K. *J. Chem. Phys.* **2004**, *120*, 8425.
- (22) Rohrdanz, M. A.; Martins, K. M.; Herbert, J. M. *J. Chem. Phys.* **2009**, *130*, 054112.
- (23) Dion, M.; Rydberg, H.; Schroeder, E.; Langreth, D. C.; Lundqvist, B. I. *Phys. Rev. Lett.* **2004**, *92*, 246401.
- (24) Grimme, S. *J. Comput. Chem.* **2006**, *27*, 1787.
- (25) Vydrov, O. A.; Van, V. T. *J. Chem. Phys.* **2010**, *133*, 244103.
- (26) Maitra, N. T.; Zhang, F.; Cave, R. J.; Burke, K. *J. Chem. Phys.* **2004**, *120*, 5932.
- (27) Casida, M. E. *J. Chem. Phys.* **2005**, *122*, 054111.

- (28) Jacquemin, D.; Wathelet, V.; Perpète, E. A.; Adamo, C. *J. Chem. Theory Comput.* **2009**, *5*, 2420.
- (29) Casarin, M.; Finetti, P.; Vittadini, A.; Wang, F.; Ziegler, T. *J. Phys. Chem. A* **2007**, *111*, 5270.
- (30) Kozimor, S. A.; Yang, P.; Batista, E. R.; Boland, K. S.; Burns, C. J.; Christensen, C. N.; Clark, D. L.; Conradson, S. D.; Hay, P. J.; Lezama, J. S.; Martin, R. L.; Schwarz, D. E.; Wilkerson, M. P.; Wolfsberg, L. E. *Inorg. Chem.* **2008**, *47*, 5365.
- (31) Besley, N. A.; Peach, M. J. G.; Tozer, D. J. *Phys. Chem. Chem. Phys.* **2009**, *11*, 10350.
- (32) Stener, M.; Fronzoni, G.; de, S. M. *Chem. Phys. Lett.* **2003**, 373, 115.
- (33) Besley, N. A.; Noble, A. *J. Phys. Chem. C* **2007**, *111*, 3333.
- (34) DeBeer, G. S.; Petrenko, T.; Neese, F. *J. Phys. Chem. A* **2008**, *112*, 12936.
- (35) DeBeer, G. S.; Petrenko, T.; Neese, F. *Inorg. Chim. Acta* **2008**, *361*, 965.
- (36) Tretiak, S.; Isborn, C. M.; Niklasson, A. M. N.; Challacombe, M. *J. Chem. Phys.* **2009**, *130*, 054111.
- (37) Kauczor, J.; Jorgensen, P.; Norman, P. *J. Chem. Theory Comput.* **2011**, *7*, 1610.
- (38) Niklasson, A. M. N.; Challacombe, M. *Phys. Rev. Lett.* **2004**, *92*, 193001.
- (39) Coriani, S.; Hoest, S.; Jansik, B.; Thøgersen, L.; Olsen, J.; Jorgensen, P.; Reine, S.; Pawłowski, F.; Helgaker, T.; Salek, P. *J. Chem. Phys.* **2007**, *126*, 154108.
- (40) Norman, P.; Bishop, D. M.; Jensen, H. J. A.; Oddershede, J. *J. Chem. Phys.* **2001**, *115*, 10323.
- (41) Norman, P.; Bishop, D. M.; Jensen, H. J. A.; Oddershede, J. *J. Chem. Phys.* **2005**, *123*, 194103.
- (42) Wilkinson, J. H. *The Algebraic Eigenvalue Problem*; Clarendon Press: Oxford, U. K., 1965; pp 582.
- (43) Davidson, E. R. *J. Comput. Phys.* **1975**, *17*, 87.
- (44) Hirao, K.; Nakatsuji, H. *J. Comput. Phys.* **1982**, *45*, 246.
- (45) Bouman, T. D.; Hansen, A. E.; Voigt, B.; Rettrup, S. *Int. J. Quantum Chem.* **1983**, *23*, 595.
- (46) MacDonald, J. K. L. *Phys. Rev.* **1933**, *43*, 830.
- (47) Frisch, M. J.; Trucks, G. W.; Schlegel, H. B.; Scuseria, G. E.; Robb, M. A.; Cheeseman, J. R.; Scalmani, G.; Barone, V.; Mennucci, B.; Petersson, G. A.; Nakatsuji, H.; Caricato, M.; Li, X.; Hratchian, H. P.; Izmaylov, A. F.; Bloino, J.; Zheng, G.; Sonnenberg, J. L.; Liang, W.; Hada, M.; Ehara, M.; Toyota, K.; Fukuda, R.; Hasegawa, J.; Ishida, M.; Nakajima, T.; Honda, Y.; Kitao, O.; Nakai, H.; Vreven, T.; Montgomery, J. A., Jr.; Peralta, J. E.; Ogliaro, F.; Bearpark, M.; Heyd, J. J.; Brothers, E.; Kudin, K. N.; Staroverov, V. N.; Keith, T.; Kobayashi, R.; Normand, J.; Raghavachari, K.; Rendell, A.; Burant, J. C.; Iyengar, S. S.; Tomasi, J.; Cossi, M.; Rega, N.; Millam, J. M.; Klene, M.; Knox, J. E.; Cross, J. B.; Bakken, V.; Adamo, C.; Jaramillo, J.; Gomperts, R.; Stratmann, R. E.; Yazyev, O.; Austin, A. J.; Cammi, R.; Pomelli, C.; Ochterski, J. W.; Martin, R. L.; Morokuma, K.; Zakrzewski, V. G.; Voth, G. A.; Salvador, P.; Dannenberg, J. J.; Dapprich, S.; Parandekar, P. V.; Mayhall, N. J.; Daniels, A. D.; Farkas, O.; Foresman, J. B.; Ortiz, J. V.; Cioslowski, J.; Fox, D. J. *Gaussian Development Version H.12+*; Gaussian, Inc.: Wallingford, CT, 2011.
- (48) Badaeva, E.; Feng, Y.; Gamelin, D. R.; Li, X. *New J. Phys.* **2008**, *10*.
- (49) Badaeva, E.; Isborn, C. M.; Feng, Y.; Ochsenbein, S. T.; Gamelin, D. R.; Li, X. *J. Phys. Chem. C* **2009**, *113*, 8710.
- (50) Norberg, N. S.; Kittilstved, K. R.; Amonette, J. E.; Kukkadapu, R. K.; Schwartz, D. A.; Gamelin, D. R. *J. Am. Chem. Soc.* **2004**, *126*, 9387.
- (51) Kittilstved, K. R.; Liu, W. K.; Gamelin, D. R. *Nat. Mater.* **2006**, *5*, 291.
- (52) Liu, W. K.; Mackay, S. G.; Gamelin, D. R. *J. Phys. Chem. B* **2005**, *109*, 14486.
- (53) Perdew, J. P.; Burke, K.; Ernzerhof, M. *Phys. Rev. Lett.* **1996**, *77*, 3865.
- (54) Perdew, J. P.; Burke, K.; Ernzerhof, M. *Phys. Rev. Lett.* **1997**, *78*, 1396.
- (55) Dunning, T. H., Jr.; Hay, P. J. *Mod. Theor. Chem.* **1977**, *3*, 1.
- (56) Hay, P. J.; Wadt, W. R. *J. Chem. Phys.* **1985**, *82*, 299.
- (57) Hay, P. J.; Wadt, W. R. *J. Chem. Phys.* **1985**, *82*, 270.
- (58) Wadt, W. R.; Hay, P. J. *J. Chem. Phys.* **1985**, *82*, 284.
- (59) Pauling, P. *Inorg. Chem.* **1966**, *5*, 1498.
- (60) McGinnety, J. A. *J. Am. Chem. Soc.* **1972**, *94*, 8406.
- (61) Lauher, J. W.; Ibers, J. A. *Inorg. Chem.* **1975**, *14*, 348.
- (62) Dawson, A.; Parkin, A.; Parsons, S.; Pulham, C. R.; Young, A. L. C. *Acta Crystallogr., Sect. E: Struct. Rep. Online* **2002**, *E58*, i95.
- (63) Becke, A. D. *Phys. Rev. A* **1988**, *38*, 3098.
- (64) Perdew, J. P.; Burke, K.; Wang, Y. *Phys. Rev. B: Condens. Matter* **1996**, *54*, 16533.
- (65) Schaefer, A.; Huber, C.; Ahlrichs, R. *J. Chem. Phys.* **1994**, *100*, 5829.
- (66) Leeuwen, R. v.; Gritsenko, O. V.; Baerends, E. J. *Top. Curr. Chem.* **1996**, *180*, 107.
- (67) Shadle, S. E.; Hedman, B.; Hodgson, K. O.; Solomon, E. I. *J. Am. Chem. Soc.* **1995**, *117*, 2259.
- (68) DeBeer, G. S.; Brant, P.; Solomon, E. I. *J. Am. Chem. Soc.* **2005**, *127*, 667.

# Two-Dimensional Scan of the Performance of Generalized Gradient Approximations with Perdew–Burke–Ernzerhof-Like Enhancement Factor

E. Fabiano,<sup>\*,†</sup> Lucian A. Constantin,<sup>‡</sup> and F. Della Sala<sup>†,‡</sup>

<sup>†</sup>National Nanotechnology Laboratory (NNL), Istituto Nanoscienze–CNR, Via per Arnesano 16, I-73100 Lecce, Italy

<sup>‡</sup>Center for Biomolecular Nanotechnologies @UNILE, Istituto Italiano di Tecnologia (IIT), Via Barsanti, I-73010 Arnesano, Italy

**S** Supporting Information

**ABSTRACT:** We assess the performance of the whole class of functionals defined by the Perdew–Burke–Ernzerhof (PBE) exchange–correlation enhancement factor, by performing a two-dimensional scan of the  $\mu$  and  $\kappa$  parameters (keeping  $\beta$  fixed by the recovery of the local density approximation linear response). We consider molecular (atomization energies, bond lengths, and vibrational frequencies), intermolecular (hydrogen-bond and dipole interactions), and solid-state (lattice constant and cohesive energies) properties. We find, for the energetical properties, a whole family of functionals (with  $\mu$  and  $\kappa$  interrelated) giving very similar results and the best accuracy. Overall, we find that the original PBE and the recently proposed APBE functional [*Phys. Rev. Lett.* **2011**, *106*, 186406], based on the asymptotic expansion of the semiclassical neutral atom, give the highest global accuracy, with a definite superior performance of the latter for all of the molecular properties.

## 1. INTRODUCTION

Ground-state density functional theory<sup>1,2</sup> (DFT) in the Kohn–Sham<sup>3</sup> (KS) self-consistent formalism is nowadays one of the most popular computational methods in electronic calculations of quantum chemistry and solid-state physics. The central quantity in DFT is the exchange–correlation (XC) functional, which collects all of the “unknown” terms of the electron–electron interaction. Over the years, many approximations have been developed for the XC functional, which form the so-called “Jacob’s ladder” of DFT.<sup>4</sup>

The ladder is grounded on the Hartree approximation (i.e., no XC contribution) and has at the first rung the local spin-density approximation<sup>3</sup> (LSDA), which only contains as ingredients the electron spin-densities  $\rho_{\uparrow}(\mathbf{r})$  and  $\rho_{\downarrow}(\mathbf{r})$ . The second rung of the ladder is formed by those functionals depending on the gradient of the electron spin-densities ( $\nabla\rho_{\uparrow}(\mathbf{r})$ ,  $\nabla\rho_{\downarrow}(\mathbf{r})$ ) as well as on the electron spin-densities themselves. The generalized gradient approximations<sup>5</sup> (GGAs) and the second-order gradient expansions, derived from small perturbations of the uniform electron gas (GE2)<sup>6</sup> or from the semiclassical theory of neutral atoms (MGE2),<sup>7–12</sup> are the most important representatives of this class. On the third rung of Jacob’s ladder, formed by meta-GGAs, the additional ingredients of the positive KS kinetic energy spin-densities ( $\tau_{\uparrow}(\mathbf{r})$ ,  $\tau_{\downarrow}(\mathbf{r})$ ) or the Laplacian of the spin-densities ( $\nabla^2\rho_{\uparrow}(\mathbf{r})$ ,  $\nabla^2\rho_{\downarrow}(\mathbf{r})$ ) are considered, in order to satisfy more exact constraints of the XC energy and potential. On the fourth rung, the dependence on occupied KS orbitals is taken into account, in search of an improved description of the exchange energy or of a fully nonlocal correlation energy compatible with exact exchange.<sup>13</sup> Hybrid<sup>14</sup> and orbital-dependent<sup>15–17</sup> functionals, as well as hyper-GGAs,<sup>13</sup> belong to this rung. Finally, on the fifth rung of the ladder, we find functionals including an explicit

dependence on virtual KS orbitals,<sup>18–23</sup> which allow one to describe exactly nonlocal parts of the correlation energy density.

Despite the remarkable accuracy demonstrated by the functionals belonging to the fourth and fifth rungs of Jacob’s ladder in different test studies, the majority of practical DFT applications are based on GGAs and hybrid functionals, which provide the best compromise between accuracy and computational effort. In particular, GGA functionals provide an efficient tool for the study of large systems (e.g., for biology and solid-state physics) and still outperform hybrid functionals for organometallic and transition metal complexes.<sup>24–27</sup> In addition, they attract basic theoretical interest, because they constitute the basis on which meta-GGA, hyper-GGA, and hybrid functionals are constructed.

Among different GGA functionals, the generalized gradient approximation proposed in 1986 by Perdew, Burke, and Ernzerhof (PBE)<sup>28</sup> has gained large popularity in both quantum chemistry and condensed-matter physics, due to its simplicity and good performance in a broad range of applications. The PBE functional contains no empirical parameters and fulfills a number of exact constraints for the XC energy. The correlation part of the PBE functional was constructed from a sharp cutoff of the GE2 correlation hole (in the high density limit)<sup>29</sup> and is defined as

$$E_c^{\text{PBE}}[\rho_{\uparrow}, \rho_{\downarrow}] = \int \rho[\varepsilon_c^{\text{unif}}(r_s, \zeta) + H(r_s, \zeta, t)] \, \text{d}\mathbf{r} \quad (1)$$

where

$$H(r_s, \zeta, t) = \gamma\phi^3 \ln\left(1 + \frac{\beta}{\gamma} \frac{t^2 + At^4}{1 + At^2 + A^2t^4}\right) \quad (2)$$

Received: July 26, 2011

Published: September 20, 2011



with  $\rho = \rho_{\uparrow} + \rho_{\downarrow}$  being the total electron density,  $r_s = [(4\pi/3)\rho]^{1/3}$  being the local Seitz radius,  $\zeta = (\rho_{\uparrow} - \rho_{\downarrow})/\rho$  being the relative spin polarization,  $\phi = ((1 + \zeta)^{2/3} + (1 - \zeta)^{2/3})/2$  being a spin scaling factor,  $\epsilon_c^{\text{unif}}(r_s, \zeta)$  being the correlation energy per particle of the uniform electron gas,  $A$  being a function of  $\epsilon_c^{\text{unif}}$  and  $\phi$ , and  $t = |\nabla\rho|/(2\phi k_F \rho)$  being the correlation density gradient that measures the density variations over a Thomas–Fermi screened wave-number  $k_s = (4k_F/\pi)^{1/2}$ , where  $k_F = (3\pi^2\rho)^{1/3}$  is the Fermi wave vector. The parameter  $\gamma = (1 - \ln 2)/\pi^2 \approx 0.031091$  is fixed by uniform scaling to the high-density limit of the (spin-unpolarized) correlation energy, and the parameter  $\beta = \beta^{\text{PBE}} = 0.066725$  is the second-order gradient expansion coefficient of the correlation energy in the high-density limit.

The exchange part of the PBE functional has as the enhancement factor a simple Padé-polynomial formula originally proposed by Becke:<sup>30</sup>

$$F_x^{\text{PBE}}(s) = 1 + \kappa - \frac{\kappa}{1 + \frac{\mu}{s^2}} \quad (3)$$

where  $s = |\nabla\rho|/(2k_F\rho)$  is the reduced gradient. The exchange energy for a spin-unpolarized system is then

$$E_x^{\text{PBE}}[\rho] = \int \rho \epsilon_x^{\text{unif}}(\rho) F_x^{\text{PBE}}(s) \, d\mathbf{r} \quad (4)$$

where  $\epsilon_x^{\text{unif}}(\rho)$  is the exchange energy per particle of the uniform electron gas, while for any spin-polarized system

$$E_x[\rho_{\uparrow}, \rho_{\downarrow}] = \frac{E_x[2\rho_{\uparrow}] + E_x[2\rho_{\downarrow}]}{2} \quad (5)$$

from the spin-scaling relation of the exchange energy.<sup>31</sup>

The exchange enhancement factor in eq 3 is very simple and satisfies two important limits: for small  $s$ , we have  $F_x^{\text{PBE}}(s) \approx 1 + \mu s^2$ ; while for large  $s$ ,  $F_x^{\text{PBE}}(s) \rightarrow 1 + \kappa$ . The parameter  $\kappa = \kappa^{\text{PBE}} = 0.804$  is fixed by the Lieb–Oxford bound for the exchange energy, and the parameter  $\mu$  is fixed to satisfy the correct linear response of the spin-unpolarized uniform electron gas, i.e.

$$\mu = \beta \frac{\pi^2}{3} \quad (6)$$

which leads to  $\mu = \mu^{\text{PBE}} = 0.21951$ .

Since its introduction, many variations of the original PBE functional have been presented.<sup>12,32–45</sup> Some of them keep the same functional form for the exchange and correlation but employ different parameters.<sup>12,32,35,37,38,41,43</sup> These functionals can be represented by a triplet of parameters  $(\mu; \beta; \kappa)$ , see also ref 46. Among them, we recall:

- (i) revPBE,<sup>32</sup> an empirical functional constructed for molecules with

$$(\mu = \mu^{\text{PBE}}; \beta = \beta^{\text{PBE}}; \kappa = 1.245)$$

where  $\kappa$  was fitted to atoms

- (ii) PBEsol,<sup>37,38</sup> a functional for solids and surfaces, with

$$\left( \mu = \mu^{\text{GE2}} = \frac{10}{81}; \beta = 0.046; \kappa = \kappa^{\text{PBE}} \right)$$

where  $\mu^{\text{GE2}}$  is the exact second-order gradient expansion coefficient of the exchange energy and  $\beta$  was fitted to jellium surfaces

- (iii) APBE,<sup>12</sup> a nonempirical functional accurate for molecular systems, constructed from the semiclassical theory

of neutral atoms, with

$$\left( \mu = \mu^{\text{MGE2}} = 0.26; \beta = \frac{3\mu^{\text{MGE2}}}{\pi^2} = 0.079; \kappa = \kappa^{\text{PBE}} \right)$$

where  $\mu^{\text{MGE2}}$  is the coefficient of the modified second-order gradient expansion and  $\beta$  was chosen to recover the LSDA linear response. We note that the construction of the APBE functional shares some similarities with that of PBE(Jr,Gx).<sup>41</sup> However, the latter uses  $\mu = \mu^{\text{GE2}} = 10/81$ , recovering GE2 and not MGE2. It thus behaves similarly to PBEsol and rather differently than APBE. This fact highlights the importance of the modified second-order gradient expansion for the exchange in the construction of APBE.

Other functionals modify the functional form of the exchange enhancement factor,<sup>33,34,36,39,40,42,44,45</sup> often producing a faster increase of  $F_x^{\text{PBE}}(s)$  with  $s$ . Of particular relevance we recall the following: Wu–Cohen (WC)<sup>36</sup> and the second-order GGA<sup>40</sup> (developed for better solid-state properties), PBEint<sup>44</sup> and RPBE<sup>35</sup> (developed for hybrid interfaces), and the regularized gradient expansion.<sup>42</sup> These functionals can also be written as a parameter triplet provided that  $\mu$  is expressed as a function of  $s$ . For the PBEint functional, e.g., we have

$$\left( \mu = \mu(s) = \mu^{\text{GE2}} + \frac{(\mu^{\text{PBE}} - \mu^{\text{GE2}})\alpha s^2}{1 + \alpha s^2}; \beta = 0.052; \kappa = \kappa^{\text{PBE}} \right)$$

where  $\alpha = (\mu^{\text{GE2}})^2/(\kappa(\mu^{\text{PBE}} - \mu^{\text{GE2}})) = 0.197$  is determined by the requirement of a smooth functional derivative. Thus,  $\mu(s)$  interpolates between the GE2 and PBE coefficients, whereas  $\beta$  is fitted to jellium surfaces.

By changing the parameters in the PBE functional form, important exact constraints for solids, surfaces, or molecular systems can be recovered, and improved accuracy can be achieved for special classes of problems. However, it is not possible to satisfy all of the different constraints at once. In fact, no GGA functional can be accurate for both solid-state properties and atoms.<sup>10</sup>

In recent years, many different investigations about the performance of the PBE-like functionals have been presented.<sup>46,47</sup> Some studies focused on the relevance of the  $\kappa$  parameter,<sup>40,48</sup> the  $\mu$  parameter,<sup>46</sup> or both.<sup>12,37,38,41</sup> However, so far, only few points in the three-dimensional  $(\mu; \beta; \kappa)$  space have been investigated and mainly for few selected properties.<sup>46</sup>

In this paper, we aim at critically assessing this issue and explore the dependence of the performance of a whole family of PBE-like functionals on the values of the  $\mu$  and  $\kappa$  parameters. For the sake of clarity, we restricted our attention to those PBE-like functionals that satisfy the LSDA linear response. This choice is also motivated by the fact that this is the only known exact constraint for the correlation energy of importance for molecular or solid-state systems. The second-order gradient expansion for the correlation has been in fact demonstrated to be of minor importance for real systems.<sup>49</sup> Therefore, we performed a two-dimensional scan of PBE-functionals of the type

$$[\mu, \kappa] = \left( \mu; \beta = \frac{3\mu}{\pi^2}; \kappa \right)$$

Previous studies<sup>41</sup> and preliminary test calculations (see the Supporting Information) indicate that different choices of  $\beta$  do not modify the results significantly. Nevertheless, we

cannot exclude more important effects for particular properties and/or systems.

We consider many different tests, namely, atomization energies and bond lengths of organic molecules, transition metals, and metal complexes; harmonic vibrational frequencies of organic molecules; binding energies of hydrogen-bonded and dipole-interacting molecular systems; equilibrium lattice constants; and cohesive energies of solids. For each test, the results are presented as contour graphs showing the accuracy of the PBE-like functionals for different combinations of the  $(\mu, \kappa)$  parameters, and the performance of standard PBE-like functionals (i.e., PBE, APBE, revPBE) is analyzed. In addition, global errors for different classes of problems are considered. We find that the nonempirical APBE and PBE functionals are the most representative functionals for a broad palet of properties and systems and yield the higher accuracy over a large number of tests. In addition, APBE outperforms PBE for molecular properties.

We note finally that, of course, for practical reasons, our selection of tests is necessarily limited. For example, it does not consider quantities such as transition barriers, isomerization energies, or reaction energies, just to mention a few relevant to the ground-state energetics. At the same time, for computational reasons, it is essentially restricted to considering small molecules, while tests including large organic molecules<sup>43,50</sup> or extended metal systems are not considered (these would require in addition the consideration of further theoretical issues as for example dispersion corrections<sup>51</sup>). This facts suggest that caution must be used in drawing conclusions from the present work, as from any other similar broad-range assessment work, because of a possible bias introduced by the selection of specific test sets (see in this respect the discussion in ref 52). Nevertheless, we believe that the present selection of tests is representative of a fairly large class of the most fundamental and important problems in quantum chemistry and solid-state physics and can thus provide useful insight into the performance of the family of PBE-like functionals.

## 2. COMPUTATIONAL DETAILS

Several properties of molecules and bulk solids were investigated by employing a PBE-like functional (eqs 1–6) with  $\mu$  and  $\kappa$  values in the intervals  $\mu \in [0.1, 0.3]$  and  $\kappa \in [0.5, 1.5]$  and the  $\beta$  parameter fixed by the relation  $\beta = 3\mu/\pi^2$  in order to preserve the accurate LSDA linear response. The parameter  $\mu$  ( $\kappa$ ) was varied in steps of 0.01 (0.1). In total, we tested 231 different PBE-like functionals.

We considered the following properties and test sets:

AE6: Atomization energies were computed for SiH<sub>4</sub>, SiO, S<sub>2</sub>, CH<sub>4</sub>, C<sub>2</sub>O<sub>2</sub>H<sub>2</sub>, and C<sub>4</sub>H<sub>8</sub>; accurate reference data were taken from ref 53.

TMAE4: Atomization energies were computed for the Cr<sub>2</sub>, Cu<sub>2</sub>, V<sub>2</sub>, and Ag<sub>2</sub> transition metal complexes; reference data were taken from ref 25.

MCAE6: Atomization energies were calculated for the AgH, BeO, FeS, LiCl, MgO, and VS metal complexes; reference data were obtained from ref 26.

HBL9: Optimization of bond lengths involving at least one hydrogen atom are provided. The following molecules were considered: H<sub>2</sub>, CH<sub>4</sub>, NH<sub>3</sub>, H<sub>2</sub>O, HF, C<sub>2</sub>H<sub>2</sub> (C–H bond), HCN (C–H bond), H<sub>2</sub>CO (C–H bond), and OH; reference values were taken from ref 54.

NHBL10: Bond lengths were optimized for CO, N<sub>2</sub>, F<sub>2</sub>, C<sub>2</sub>H<sub>2</sub> (C–C bond), HCN (C–N bond), H<sub>2</sub>CO (C–O bond), CO<sub>2</sub>, N<sub>2</sub>O, and Cl<sub>2</sub>; reference values were taken from ref 54. TMBL4: Bond lengths were computed for the Cr<sub>2</sub>, Cu<sub>2</sub>, V<sub>2</sub>, and Ag<sub>2</sub> transition metal complexes; reference data were taken from ref 25.

MCBL6: Bond lengths were optimized for the AgH, BeO, FeS, LiCl, MgO, and VS metal complexes; reference data were obtained from ref 26.

F38: Harmonic vibrational frequencies were calculated for H<sub>2</sub>, CH<sub>4</sub>, NH<sub>3</sub>, H<sub>2</sub>O, HF, CO, N<sub>2</sub>, F<sub>2</sub>, C<sub>2</sub>H<sub>2</sub>, HCN, H<sub>2</sub>CO, CO<sub>2</sub>, N<sub>2</sub>O, Cl<sub>2</sub>, and OH; reference data were obtained from ref 55.

HB6/04: The binding energies of hydrogen-bond interacting systems were calculated for (H<sub>2</sub>O)<sub>2</sub>, (HCONH<sub>2</sub>)<sub>2</sub>, (HCOOH)<sub>2</sub>, (HF)<sub>2</sub>, (NH<sub>3</sub>)<sub>2</sub>, and NH<sub>3</sub>–H<sub>2</sub>O; reference data were taken from ref 56.

DI6/04: The binding energies of dipole-interacting systems were calculated for CH<sub>3</sub>Cl–HCl, CH<sub>3</sub>SH–HCl, CH<sub>3</sub>SH–NCH, (H<sub>2</sub>S)<sub>2</sub>, (HCl)<sub>2</sub>, and H<sub>2</sub>S–HCl; reference data were taken from ref 56.

SOLIDS: Equilibrium lattice constants and cohesive energies were computed for bulk Na (simple metal), Ag, Cu (transition metals), Si, GaAs (semiconductors), and NaCl (ionic solid); reference values were taken from refs 47 and 57. This small test of six solids can reproduce the mean absolute errors of the functionals for larger set of solids well.<sup>44</sup> For example, for 60 solids,<sup>45</sup> the PBE lattice constant mean absolute error is 0.054 Å, whereas our small test of solids gives a PBE error of 0.0597 Å.

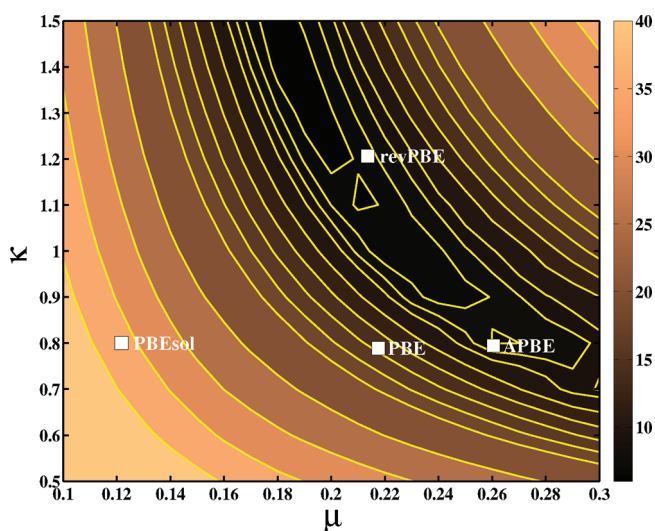
In all calculations, fully relaxed geometries were considered. The simulations of molecular systems were performed with the TURBOMOLE program package,<sup>58</sup> using a def2-TZVP<sup>59,60</sup> basis set. The simulations of solid-state properties were performed employing the FHI-AIMS program<sup>61,62</sup> using the light basis set and a 18 × 18 × 18 *k*-point grid. Scalar relativistic effects were included, where needed, through the zeroth-order relativistic approximation (ZORA).<sup>63</sup>

## 3. RESULTS

In this section, we report the performance of PBE-like functionals using different values of the  $\mu$  and  $\kappa$  parameters. For each test, the results are reported as a two-dimensional plot showing the mean absolute error (MAE) as a function of  $\mu$  and  $\kappa$ . In the figures also the combinations of  $(\mu, \kappa)$  corresponding to PBE, APBE, revPBE, and mPBESol are indicated for reference. Here, mPBESol indicates a functional having the same  $(\mu, \kappa)$  values as the original PBESol<sup>37,38</sup> ( $\mu = 10/81$ ,  $\kappa = 0.804$ ), but a different value of  $\beta$  (0.037) imposed by the constraint of the LSDA response satisfaction. Actually, this functional was already considered in ref 41, where it was indicated as PBE(J<sub>r</sub>,G<sub>x</sub>), yielding a very similar performance with the original PBESol.

In the following discussion, we will also compare the results of the PBE-like functional with other common ones, i.e., BLYP,<sup>64,65</sup> OLYP,<sup>65–67</sup> PBEint,<sup>44</sup> TPSS meta-GGA,<sup>68</sup> and the global hybrid PBE0.<sup>69</sup> These results are reported in Table S1, in the Supporting Information.

To discuss the performances of different functionals, we introduce an exchange (X) nonlocality measure  $\Lambda$  for PBE-like functionals. The true nonlocality of a GGA functional is given by the XC enhancement factor  $F_{XC}(r_s, \zeta_s)$ .<sup>70</sup> However, this function



**Figure 1.** Mean absolute error (kcal/mol) for the atomization energy of molecules of the AE6 set as a function of  $\mu$  and  $\kappa$ . The positions corresponding to PBE, APBE, revPBE, and mPBESol are denoted by white boxes.

of three variables is too complicated to yield a simple measure, as the one required in this work. For this reason, we consider simply  $F_{XC}(r_s = 0, \zeta = 0, s) = F_X(s)$ , i.e., the leading exchange part. For a given functional, the X-nonlocality is thus defined as

$$I_{\text{func}} \equiv \int_0^{s_{\text{max}}} (F_{\text{func}}(s) - 1) ds \quad (7)$$

where  $F$  is the enhancement factor of the functional and  $s_{\text{max}}$  is the maximum value of the reduced gradient  $s$  that contributes to the integration of the exchange energy (eq 4). Here, we used  $s_{\text{max}} = 6$ ; however, our final result will turn out to be independent of the value of  $s_{\text{max}}$ , provided that it is large enough. Using the PBE-like exchange enhancement factor form (eq 3), performing the integration, and after some algebra, we find

$$I_{\text{func}} = \frac{\kappa_{\text{func}}^2}{\sqrt{\mu_{\text{func}} \kappa_{\text{func}}}} G\left(\frac{\mu_{\text{func}} s_{\text{max}}}{\sqrt{\mu_{\text{func}} \kappa_{\text{func}}}}\right) \quad (8)$$

with  $G(x) = x - \arctan(x)$ . In the range of  $\mu$  and  $\kappa$  values considered in this work, we can approximate the function  $G$  as

$$G(x) \approx \frac{x^2}{c s_{\text{max}}} \quad (9)$$

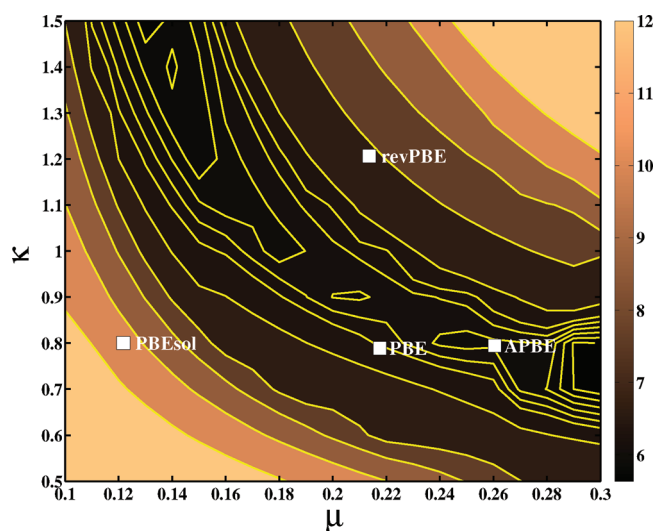
where  $c = 1.05$  is a fitting parameter. The X-nonlocality of the functional can thus be written as

$$I_{\text{func}} = \frac{\sqrt{\mu_{\text{func}} \kappa_{\text{func}} s_{\text{max}}}}{c} \quad (10)$$

The X-nonlocality measure can be finally defined as the X-nonlocality of the given functional relative to PBE, i.e.

$$\Lambda_{\text{func}} = \frac{I_{\text{func}}}{I_{\text{PBE}}} = \frac{\sqrt{\mu_{\text{func}} \kappa_{\text{func}}}}{\sqrt{\mu_{\text{PBE}} \kappa_{\text{PBE}}}} \quad (11)$$

The values of the X-nonlocality measure for the different PBE-like functionals considered in this paper are reported in Figure S1, in the Supporting Information. For standard PBE-like



**Figure 2.** Mean absolute error (kcal/mol) for the atomization energy of the TMAE4 set, as a function of  $\mu$  and  $\kappa$ . The positions corresponding to PBE, APBE, revPBE, and mPBESol are denoted by white boxes.

functionals, we have 0.75, 1.00, 1.09, and 1.24 for mPBESol, PBE, APBE, and revPBE, respectively.

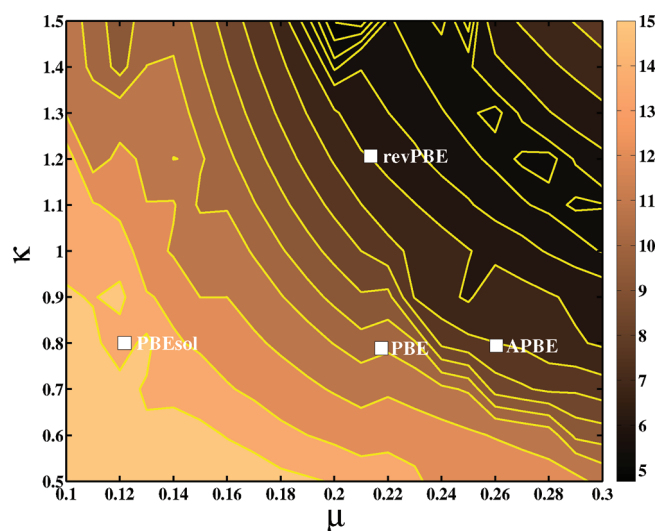
**3.1. Molecules.** **3.1.1. Atomization Energies.** In Figure 1, we report the MAE for the AE6 test as a function of  $\mu$  and  $\kappa$ . The best results, with MAEs below 10 kcal/mol, are found for combinations of the two parameters including either high values of  $\kappa$  and intermediate values of  $\mu$  or intermediate values of  $\kappa$  and high values of  $\mu$ . Both the APBE (MAE 7.9 kcal/mol) and the revAPBE (MAE 8.84 kcal/mol) functionals belong to this region and indeed yield for this test a performance close to the best GGAs (e.g., OLYP<sup>65–67</sup> has a MAE of 4.3 kcal/mol) and to the TPSS<sup>68</sup> meta-GGA (MAE 5.4 kcal/mol) and hybrid PBE0<sup>69</sup> functional (MAE 5.4 kcal/mol).

The PBE functional lays just outside the region of minimum MAEs and gives a mean absolute error of 14.5 kcal/mol. A very poor performance for the AE6 test is found, as expected, for the mPBESol functional, which strongly overestimates atomization energies, because of its reduced X-nonlocality ( $\Lambda_{\text{mPBESol}} = 0.75$ ).

Interestingly, the PBE performance can be improved, increasing the X-nonlocality by either increasing the value of the  $\mu$  parameter, yielding APBE ( $\Lambda_{\text{APBE}} = 1.09$ ), or increasing the value of the  $\kappa$  parameter, yielding revPBE ( $\Lambda_{\text{revPBE}} = 1.24$ ). However, a too pronounced X-nonlocality (functionals in the top-right region of Figure 1) makes the result worse (significant underestimation). This appears also from the comparison of APE and revPBE results. Indeed, the values of  $\mu$  and  $\kappa$  defining the APBE functional correspond approximately to a minimum for the MAE of the AE6 test.<sup>12</sup>

In Figure 2, we show the MAE for the TMAE4 test as a function of  $\mu$  and  $\kappa$ . A similar trend is observed as in the case of the AE6 test, but the region with lower MAE is shifted toward lower  $\mu$  and  $\kappa$ . Thus, a slightly lower X-nonlocality of the functionals is required: in fact, the smaller errors (about 6–7 kcal/mol) are obtained for functionals with  $\Lambda \sim 1.04–1.1$ . This finding is not surprising, as the TMAE4 test considers large atoms. In fact, even a slightly lower nonlocality might be expected to be needed if larger systems are considered.<sup>71</sup>

Both the PBE (MAE 6.3 kcal/mol) and APBE (MAE 6.1 kcal/mol) functionals perform very well and better than the



**Figure 3.** Mean absolute error (kcal/mol) for the atomization energy of the MCAE6 test set, as a function of  $\mu$  and  $\kappa$ . The positions corresponding to PBE, APBE, revPBE, and mPBESol are denoted by white boxes.

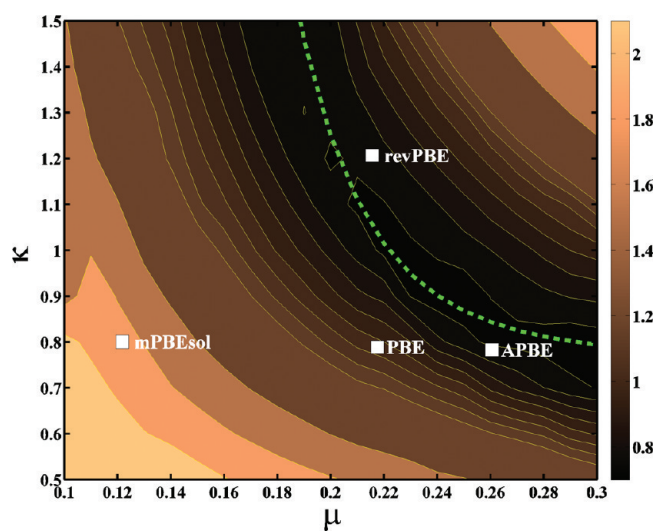
TPSS meta-GGA functional (MAE 6.6 kcal/mol). The revPBE functional yields a MAE that is about 1.5 kcal/mol higher than PBE, mainly because of its excessive X-nonlocality, given by a too high value of the  $\kappa$  parameter in combination with  $\mu = 0.2195$ . On the other hand, the mPBESol functional yields again a rather poor result (overestimated atomization energies), because of its too low degree of X-nonlocality.

In Figure 3, the MAE for the MCAE6 test is shown. Unlike the previous two tests, in this case, a higher level of X-nonlocality is necessary to appropriately describe the metal complexes' atomization energies. The region of the minima is in fact significantly moved toward the top-right corner of the plot, corresponding to PBE-like functionals with  $\Lambda \sim 1.35$ . None of the commonly used PBE-like functionals possesses such a high level of X-nonlocality, and therefore none of them yields results close to the best possible performance. The smallest MAE is found with the revPBE functional (MAE 6.3 kcal/mol), which improves with respect to PBE because of the higher value of  $\kappa$ , and it is close to the best GGA, i.e., OLYP with a MAE of 5.4 kcal/mol. Rather good results are also found for the APBE functional (MAE 7.4 kcal/mol), which performs better than TPSS (MAE 7.7 kcal/mol) and PBE0 (MAE 10.3 kcal/mol). The PBE functional gives instead a MAE of 10.5 kcal/mol and turns out to be unable to provide a completely reliable description of these systems. Finally, mPBESol displays very poor performance (overestimating atomization energies).

To have a global assessment of all atomization energies, in Figure 4, we report the global mean absolute error for the atomization energies of the AE6, TMAE4, and MCAE6 test sets, normalized to the PBE value (taken as a reference value), i.e.

$$\text{GMAE}(\mu, \kappa) = \frac{1}{3} \sum_i \frac{\text{MAE}_i(\mu, \kappa)}{\text{MAE}_i(\text{PBE})} \quad (12)$$

where  $i$  runs over AE6, TMAE4, and MCAE6. An inspection of the figure shows that for the GMAE there exists an almost continuous distribution of minima, all with very similar GMAEs (about 0.7). The locus of these minima can be



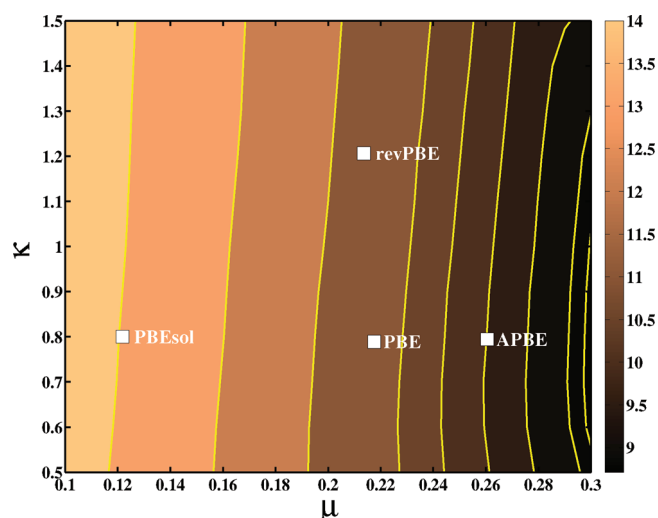
**Figure 4.** Global mean absolute error for the atomization energies of AE6, TMAE4, and MCAE6, normalized to the PBE value. The positions corresponding to PBE, APBE, revPBE, and mPBESol are denoted by white boxes. The family of functionals defined by eq 13 is shown with a green line.

described well by the relation

$$\kappa = a + \frac{b}{\mu^\gamma} \quad (13)$$

with  $a = 0.7457$ ,  $b = 6.532 \times 10^{-6}$ , and  $\gamma = 6.968$ . Equation 13 defines a family of PBE-like functionals optimized for the atomization energies of molecular systems and is shown in Figure 4 as a green dashed-line. We note that eq 13 was obtained considering test sets including only small molecules; therefore, it may be expected to reflect a slight preference for moderately high levels of nonlocality. Indeed, the need for a slightly reduced nonlocality in PBE-like functionals was already evidenced in the case of gold nanostructures of increasing size.<sup>71</sup> Moreover, eq 13 is only a simple empirical fit to the data of Figure 4; therefore, it cannot be employed to obtain accurate numerical results (note also that eq 13, because of its form, is prone to numerical noise). Nevertheless, we can use eq 13 to discuss some important results:

- (i) For accurate atomization energies,  $\kappa$  displays a lower bound ( $\kappa \geq a = 0.7457$ ) close to the nonempirical  $\kappa^{\text{PBE}}$ . Thus, we can extrapolate that a PBE-like functional with  $\mu \rightarrow \infty$  and  $\kappa = a = 0.7457$  will give extremely high total energies of atoms and molecules but still will be accurate for atomization energies. This shows that the atomization energies are dominated by the valence regions where the reduced gradient is relatively big ( $s \geq 2$ ).
- (ii) For  $\mu = \mu^{\text{GE2}} = 10/81$ , we can extrapolate  $\kappa \approx 14.7$  (eq 13 might not be very accurate in this region; thus, the following discussion is only qualitative). This very large value of  $\kappa$  violates largely the Lieb–Oxford bound<sup>48</sup> and has little physical meaning, showing that accurate atomization energies cannot be recovered by any reasonable PBE-like functional when  $\mu = \mu^{\text{GE2}}$  (they can be however obtained by relaxing slightly the PBE form, as in the PBEint functional<sup>44</sup>). Additionally, we note that a functional with  $\mu = \mu^{\text{GE2}}$  and  $\kappa \approx 14.7$  would yield very overestimated total energies and also extremely poor



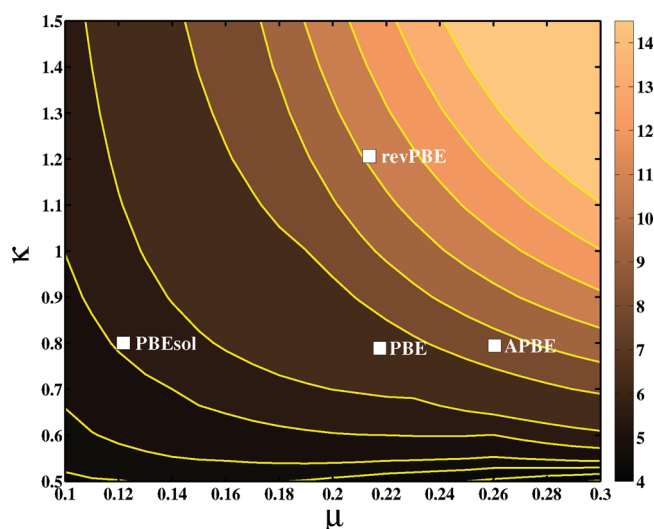
**Figure 5.** Mean absolute error (mÅ) for the equilibrium bond lengths of the HBL9 test set, as a function of  $\mu$  and  $\kappa$ . The positions corresponding to PBE, APBE, revPBE, and mPBESol are denoted by white boxes.

results for solid-state systems, worsening much over mPBESol because of its very high X-nonlocality measure ( $\Lambda = 3.2$ ).

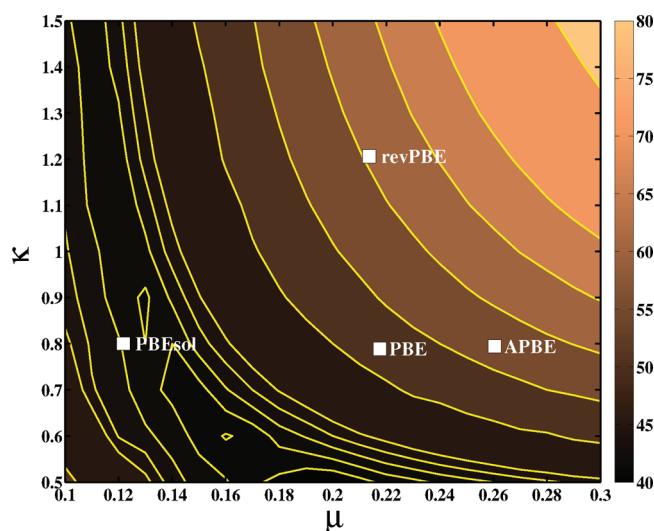
- (iii) Both APBE and revPBE are practically members of the PBE-like family defined by eq 13. APBE is however closer to the curve than revPBE. In fact, the GMAE for atomization energies (see eq 12) is 0.74 and 0.82 for APBE and revPBE, respectively.
- (iv) The functionals defined by eq 13, although yielding very similar global MAEs for the atomization energies, possess a very different X-nonlocality measure. Using eqs 11 and 13, we find in fact  $\Lambda_{(\mu,\kappa)} = (a\mu + b/\mu^{\gamma-1})^{1/2} / (\mu_{\text{PBE}}\kappa_{\text{PBE}})^{1/2}$ , which is very high for small values of  $\mu$  and close to 1.1 for  $0.22 \leq \mu \leq 0.3$  (it has a minimum of 1.094 at  $\mu = 0.243$ ). This implies that only the functionals with a relatively high value of the  $\mu$  parameter, and correspondingly  $\kappa \sim 0.8$ , e.g., APBE, can be expected to work well for atomization energies as well as for problems that require a rather reduced level of X-nonlocality such as bond lengths and solid-state properties (see later).

**3.1.2. Bond Lengths.** To perform an assessment for the bond lengths of organic molecules, we divided the systems of the MGBL19<sup>54</sup> into two groups: The HBL9 test set contains bonds that involve at least one hydrogen atom; the NHBL10 test set instead contains only bonds which do not involve hydrogen. The two sets in fact turn out to have completely different behaviors and need to be analyzed separately (see Figure S2 in the Supporting Information).

In Figure 5, we report the results of the HBL9 test for the bond lengths of several organic molecules containing hydrogen. In this case, unlike for the atomization energies, the results of the test do not appear to be directly related to the X-nonlocality measure of the functionals, and the figure does not show the characteristic hyperbolic pattern of the previous cases. The results are almost independent from the value of the  $\kappa$  parameter and only vary with  $\mu$ . In more detail (see Figure S2 in Supporting Information), using higher values of  $\mu$  leads to a reduction of the bond lengths, especially for the H–H bond. Thus, because all of the bond



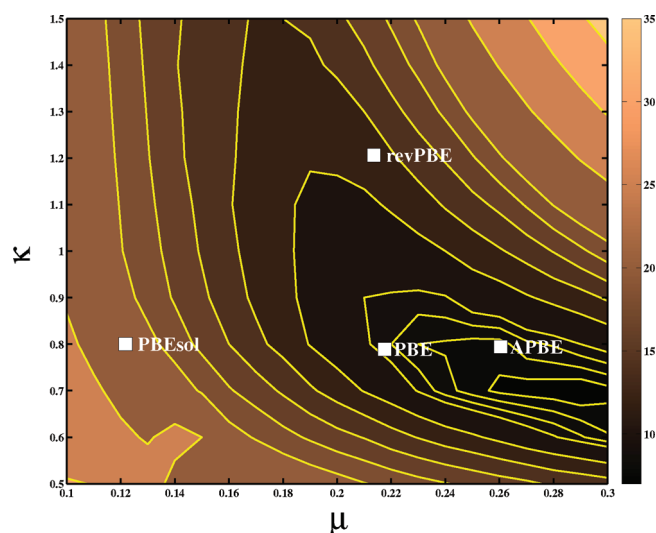
**Figure 6.** Mean absolute error (mÅ) for the equilibrium bond lengths of the NHBL10 test set, as a function of  $\mu$  and  $\kappa$ . The positions corresponding to PBE, APBE, revPBE, and mPBESol are denoted by white boxes.



**Figure 7.** Mean absolute error (mÅ) for the equilibrium bond lengths of the TMBL4 test set, as a function of  $\mu$  and  $\kappa$ . The positions corresponding to PBE, APBE, revPBE, and mPBESol are denoted by white boxes.

lengths in this test are generally overestimated, a better agreement with the reference values is found at high  $\mu$ . As a result, the smallest MAE is found, for standard PBE-like functionals, at the APBE level, with 9.6 mÅ.

In Figure 6, we report the results of the NHBL10 test for the bond lengths of several organic molecules, excluding bonds with hydrogen atoms. The plot, in contrast to Figure 5, shows the characteristic hyperbolic pattern already observed for the atomization energies and indicates that the best performance is achieved by functionals having a rather small X-nonlocality measure. Among the standard PBE-like functionals, in fact, mPBESol yields the smallest MAE (5.6 mÅ), while the worst results are obtained by revPBE (MAE 11.1 mÅ). The bond distances in the NHBL10 test set are all increased when the



**Figure 8.** Mean absolute error (mÅ) for the equilibrium bond lengths of the MCBL6 test set, as a function of  $\mu$  and  $\kappa$ . The positions corresponding to PBE, APBE, revPBE, and mPBESol are denoted by white boxes.

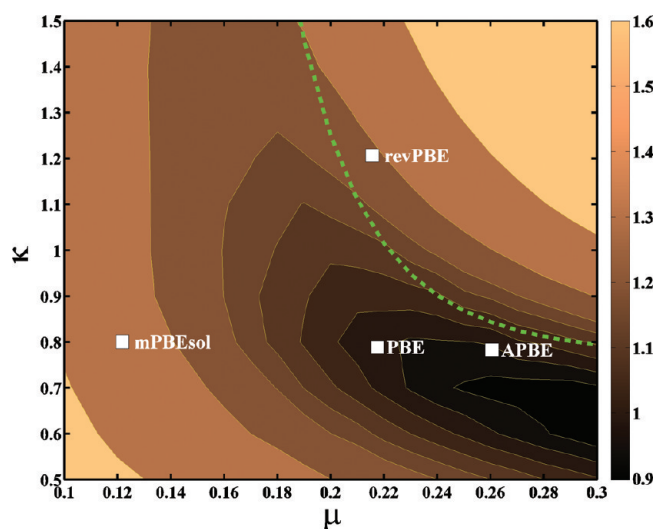
X-nonlocality of the functional is increased (see Figure S2 in the Supporting Information). Since for small values of the X-nonlocality measure most bonds are very well described, this causes a worsening of the accuracy for high values of the X-nonlocality.

Because of the different behavior of the two tests with respect to the variations of  $\mu$  and  $\kappa$ , good results cannot be achieved at the same time for both sets. The best performance for all organic molecules (i.e., the whole MGBL19 test) is obtained overall for high values of  $\mu$  and very low values of  $\kappa$  (see Figure S3 in the Supporting Information), because this combination best balances the two opposing trends. In this case, MAEs of about 6.5 mÅ are obtained, which compare well with the results of TPSS (MAE 6.9 mÅ) and PBE0 (MAE 6.3 mÅ) calculations. PBE, mPBESol, and APBE give all the same accuracy, while revPBE works quite badly (MAE 11.4 mÅ).

Considering together atomization energy and bond length for organic molecular systems, Figure S4 in the Supporting Information shows that APBE is the best choice ( $\text{MAE}_{\text{APBE}}/\text{MAE}_{\text{PBE}} = 0.78$ ). APBE not only outperforms both PBE and revPBE but it has the maximum accuracy among all of the PBE-like functionals considered.

The MAEs of the bond lengths of the TMBL4 test set as functions of  $\mu$  and  $\kappa$  are shown in Figure 7. The best results are found when small values of the  $\mu$  parameter are considered and, despite the fact that when this condition is satisfied, the MAE is rather independent of the value of the  $\kappa$  parameter, in general a low X-nonlocality is needed ( $\Lambda < 1$ ). Thus, the use of a small  $\kappa$  ( $\sim 0.6$ ) and medium-small  $\mu$  ( $\sim 0.16$ ) yields the smallest error for bond lengths with a MAE of 40.9 mÅ. This error compares favorably with that obtained using the TPSS meta-GGA functional (MAE 42.6 mÅ). Very large errors are found on the other hand for functionals exploiting a high level of X-nonlocality.

According to this analysis, among the standard PBE-like functionals, the mPBESol functional yields the best performance for the description of bond lengths of transition metal dimers (MAE 43.1 mÅ). Note that mPBESol was instead very bad for the atomization energies of these systems (TMAE4; Figure 2). Larger errors are found in the order PBE (MAE 52.8 mÅ) and APBE (MAE 57.3 mÅ), because of the increasing X-nonlocality of the functionals. The revPBE functional yields finally the

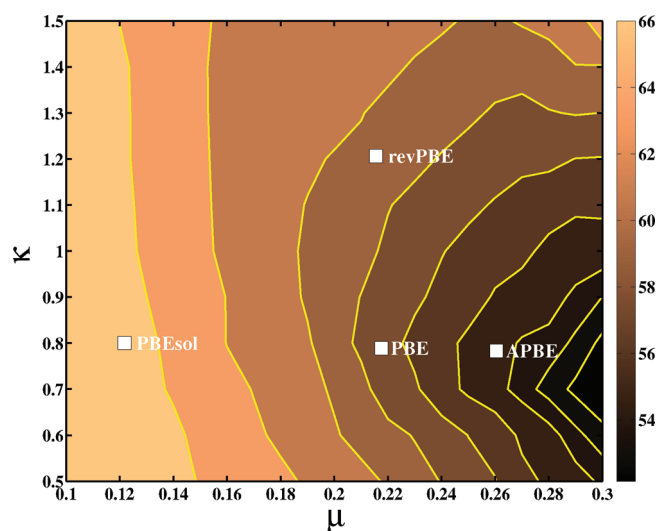


**Figure 9.** Global mean absolute error for the bond lengths of MGBL19, TMBL4, and MCBL6, normalized to the PBE value. The positions corresponding to PBE, APBE, revPBE, and mPBESol are denoted by white boxes. The family of functionals defined by eq 13 is shown with a green line.

poorest performance with a MAE of 62.2 mÅ. While mPBESol is the best for TMBL4, it was the worst for TMAE4. Considering together atomization and bond length for transition metal systems, Figure S5 in the Supporting Information shows that PBE and APBE yield high and comparable accuracy, outperforming both mPBESol and revPBE. In this case, however, Figure S5 shows that the best functional should have  $\mu$  as in mPBESol but a very high  $\kappa$ .

In Figure 8, we report the results of the MCBL6 test on the bond lengths of six metal complexes. The smallest mean average errors are found for combinations of  $\mu$  and  $\kappa$ , giving a medium value of the X-nonlocality measure  $\Lambda \sim 1$ , i.e., for relatively high values of the  $\mu$  parameter and  $\kappa \sim 0.7/0.8$ . The APBE and PBE functionals thus perform well with a MAE of 8.3 and 9.2 mÅ, respectively. For comparison, MAEs of 8.1 and 14.2 mÅ are found at the TPSS and PBE0 levels, respectively. The revPBE (MAE 14.5 mÅ) and mPBESol (MAE 22.4 mÅ) functionals yield instead significantly worse results because of the too large/small X-nonlocality of revPBE/mPBESol. Considering together atomization energies and bond lengths for metal complexes, Figure S6 in the Supporting Information shows that APBE has the maximum accuracy ( $\text{MAE}_{\text{APBE}}/\text{MAE}_{\text{PBE}} = 0.80$ ) among all of the PBE-like functionals considered.

Figures 5–8 show overall that for bond lengths, a more important role is played by the value of the  $\mu$  parameter, which must be high for HBL10 and MCBL6 and small for TMBL4, while the value of the  $\kappa$  parameter is less important. Moreover, a moderate/small level of X-nonlocality is requested to obtain accurate results. The importance of  $\mu$  traces back to the fact that for bond lengths a fundamental role is played by the first derivative of the XC potential with respect to the nuclear positions, which is related in the present context to the derivative of the exchange enhancement factor with respect to  $s^2$ . This latter term, once the X-nonlocality measure of the functional is fixed ( $\mu/\kappa \approx \text{const.}$ ), depends in first approximation only on  $\mu$ , which is then mainly determining the performance of the functionals for the problem.



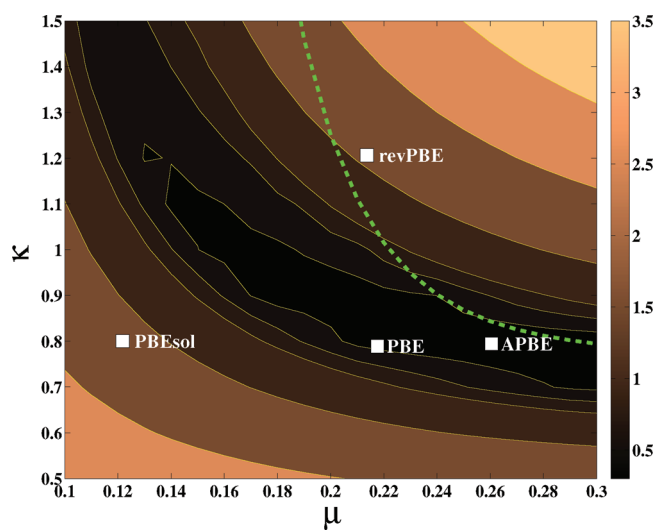
**Figure 10.** Mean absolute error ( $\text{cm}^{-1}$ ) for the vibrational frequencies of molecules of the F38 test set, as a function of  $\mu$  and  $\kappa$ . The positions corresponding to PBE, APBE, revPBE, and mPBESol are denoted by white boxes. The family of functionals defined by eq 13 is shown with a green line.

The overall performance of different PBE-like functionals for the computation of equilibrium bond lengths is summarized in Figure 9, where we report the global MAE for the bond lengths tests (see eq 12). Because different values of  $\mu$  are required by different tests, in this case, we do not find a full family of PBE-like functionals with the same and high accuracy. Instead, the best performance corresponds to a well-defined region at site  $0.25 < \mu < 0.3$  and  $0.6 < \kappa < 0.7$ . Therefore, among the commonly used PBE-like functionals, the best results are given by APBE with a GMAE (with respect to PBE) of 0.99. Note, however, that the global performance for bond lengths results mainly from an error balancing between the values obtained for the organic molecules and the transition-metal dimers.

As previously discussed, the APBE functional provides also the best compromise for the simultaneous accurate calculations of bond lengths and atomization energies (see Figures S4–S6 in the Supporting Information). Finally, APBE is also the best functional considering together atomization energies and bond lengths of all molecular systems, with a GMAE (with respect to PBE) of 0.82 (see Figure S7 in the Supporting Information).

**3.1.3. Vibrational Frequencies.** For molecular systems, we also consider harmonic vibrational frequencies (the F38 test). First of all, we note that Figure 10 does not show a hyperbolic pattern, but a strong dependence of the functionals' performance from the value of the  $\mu$  parameter and a minor role of the  $\kappa$  value. This finding can be explained, in analogy with the case of bond lengths where the first derivative of the exchange enhancement factor was important, by the importance of the second derivative of the enhancement factor for harmonic vibrations. This term in fact is, in first approximation, linearly dependent on the  $\mu$  parameter and independent from the  $\kappa$  parameter (once the X-nonlocality of the functional is fixed, in this case, to a moderate value  $\Lambda \sim 1.1$ ).

Moreover, we observe that all of the functionals belonging to the PBE family yield similar results, with maximum differences on the order of  $10 \text{ cm}^{-1}$ , irrespective of the values used for  $\mu$  and  $\kappa$ . None of the combinations of  $\mu$  and  $\kappa$  considered in this work proved to be able to yield very accurate results.



**Figure 11.** Mean absolute error (kcal/mol) for the binding energies of the HB6/04 benchmark test of hydrogen-bond interacting systems, as a function of  $\mu$  and  $\kappa$ . The positions corresponding to PBE, APBE, revPBE, and mPBESol are denoted by white boxes.

In fact, the minimum MAE in Figure 10 is  $52 \text{ cm}^{-1}$ , while the best GGA and hybrid methods give much smaller errors (e.g., OLYP MAE is  $40.1 \text{ cm}^{-1}$ , B3LYP MAE is  $33 \text{ cm}^{-1}$ <sup>55</sup>). However, ref 55 shows that F38 is well described only by the nonlocal rungs of Jacob's ladder (e.g., double-hybrids with a MAE of  $18 \text{ cm}^{-1}$ ), whereas the semilocal rungs (including meta-GGAs) give in general modest accuracy, so they cannot be used in spectroscopic studies.

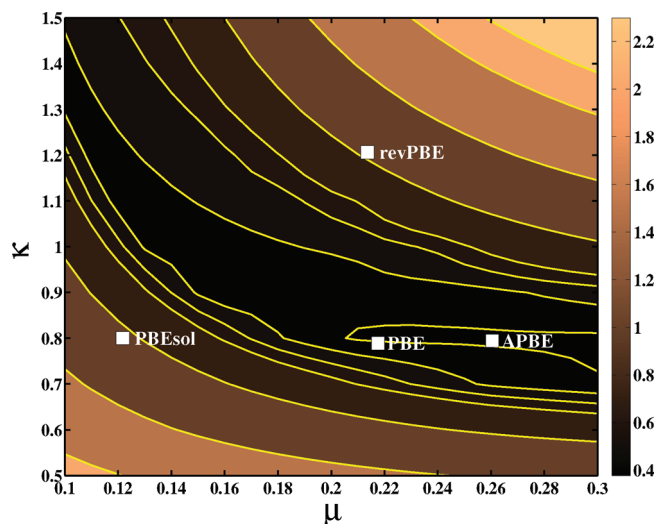
PBE-like functionals with high values of the  $\mu$  parameter and  $\kappa \sim 0.7$  display the best performance, while poor results are found for small values of  $\mu$ , for any  $\kappa$  value. As a consequence, among the standard PBE-like functionals, the smallest MAE is obtained with the APBE functional ( $55.0 \text{ cm}^{-1}$ ), and the worst result is achieved by the mPBESol functional ( $67.4 \text{ cm}^{-1}$ ).

**3.1.4. Nonbonded Interaction.** Despite the fact that GGA functionals cannot correctly describe nonbonded interactions due to the missing/incorrect dispersion forces,<sup>72</sup> very good performances were obtained by the PBE functional for hydrogen-bond and dipole–dipole interaction systems.<sup>56</sup>

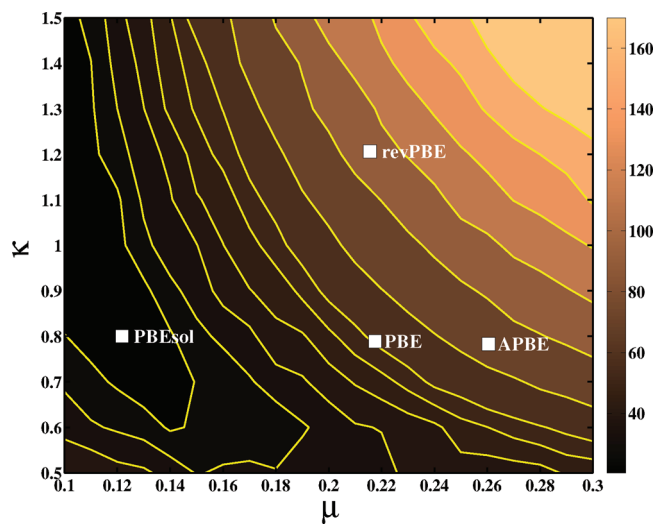
In Figure 11, we report the plot of the MAE of HB6/04 binding energies (kcal/mol) as a function of  $\mu$  and  $\kappa$ . In this case, the plot resembles the one for atomization energies. In fact, both of them are total energy differences between interacting and noninteracting subsystems.

The best results are obtained from functionals displaying a medium X-nonlocality measure ( $\Lambda \sim 1/1.1$ ), while a too high/too low X-nonlocality leads to poor results, corresponding to a general underestimation/overestimation of the interaction energy. The absolute minimum in our plot is found for the APBE functional ( $\mu = 0.26$ ,  $\kappa = 0.8$ ) with a MAE of 0.32 kcal/mol. Note that this is, to our knowledge, the best performance in the literature for the hydrogen-bond problem,<sup>56</sup> twice as good as the best meta-GGA (TPSS MAE is 0.60 kcal/mol) and slightly better than the best hybrid functionals (PBE0 MAE 0.42 kcal/mol). Good results are obtained also from the PBE functional (MAE 0.38 kcal/mol), while poor results are found from revPBE (MAE 1.90 kcal/mol) and mPBESol (MAE 1.86 kcal/mol) calculations.

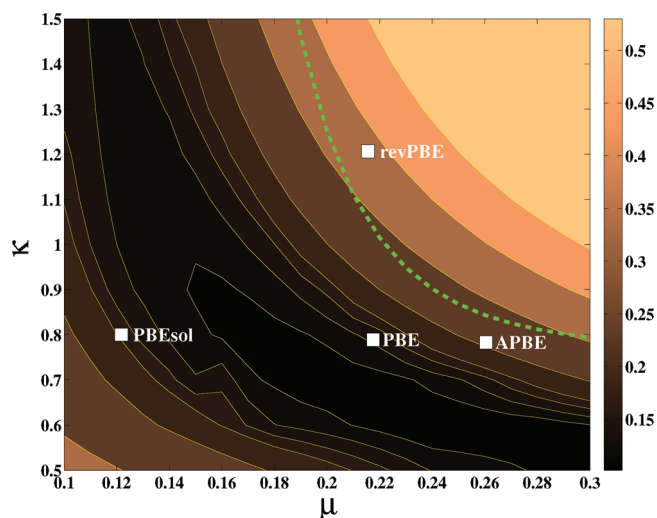
In Figure 12, we report the MAE of DI6/04 binding energies (kcal/mol) as a function of  $\mu$  and  $\kappa$ . A similar behavior as for the



**Figure 12.** Mean absolute error (kcal/mol) for the binding energies of the DI6/04 benchmark test of dipole interacting systems, as a function of  $\mu$  and  $\kappa$ . The positions corresponding to PBE, APBE, revPBE, and mPBESol are denoted by white boxes.



**Figure 14.** Mean absolute error (mÅ) for the lattice constants of six solids, as a function of  $\mu$  and  $\kappa$ . The positions corresponding to PBE, APBE, revPBE, and mPBESol are denoted by white boxes.



**Figure 13.** Mean absolute error (eV/atom) for the cohesive energies of six solids, as a function of  $\mu$  and  $\kappa$ . The positions corresponding to PBE, APBE, revPBE, and mPBESol are denoted by white boxes. The family of functionals defined by eq 13 is shown with a green line.

HB6/04 test is obtained, with functionals characterized by a medium X-nonlocality measure ( $\lambda \sim 1$ ) performing best, while functionals with a high/low degree of X-nonlocality tend to underestimate/overestimate the interaction energy.

The PBE and APBE functionals yield the same MAE (0.38 kcal/mol), while larger errors are found for revPBE (1.22 kcal/mol) and mPBESol (1.09 kcal/mol). The results obtained from the PBE and APBE functionals are among the best achievable at the GGA level,<sup>56</sup> comparable with the best meta-GGA functionals and slightly worse than the best hybrid approaches.<sup>56</sup>

**3.2. Solid State.** In this section, we report briefly on the performance of PBE-like functionals for the description of equilibrium properties of solid-state systems. In particular, we focus on lattice constants and cohesive energies.

**3.2.1. Cohesive Energies.** In Figure 13, we report the MAE of the cohesive energies of six solids as a function of  $\mu$  and  $\kappa$ . The plot resembles roughly that of Figure 2, where the performances for the atomization energies of transition metal dimers are reported. The cohesive energy of a bulk solid (of the chemical element Y) seems in fact to be the upper bound of the atomization energy of any neutral cluster of the same chemical element.<sup>71,73</sup> However, in the bulk, the density is more slowly varying than in molecular systems, and thus the set of functionals given by eq 13 becomes too nonlocal (see Figure 13). The smaller errors are obtained for functionals with a relatively small X-nonlocality measure ( $\Lambda \sim 0.9$ ). This value of the X-nonlocality is in fact needed to provide balance between the description of the bulk solid, which is well described by rather local functionals, and the description of isolated atoms, which require a larger X-nonlocality in the functional.

Good results are obtained from the PBE functional with a MAE of 0.15 eV/atom. The mPBESol functional instead yields a MAE of 0.21 eV/atom, because of its poor performance for the atomic energies. On the other hand, high errors are also obtained from APBE (MAE 0.26 eV/atom) and revPBE (MAE 0.41 eV/atom), because of the too high X-nonlocality included in these functionals. We note finally that both PBE and APBE performance can be improved reducing kappa, as found in ref 40.

**3.2.2. Lattice constants.** In Figure 14, we consider the ability of PBE-like functionals with different  $\mu$  and  $\kappa$  values to describe the lattice constant of different solids. The best performance is obtained for the functionals characterized by small values of  $\mu$  and in general by a low X-nonlocality. In fact, for higher values of  $\mu$ , reasonably small MAEs are found in conjunction with very small values of  $\kappa$ , while large errors are obtained when both  $\mu$  and  $\kappa$  are large. These findings resemble the results of the bond lengths of transition metal dimers (Figure 7).

The mPBESol functional is the best standard PBE-like functional for this problem with a MAE of 19 mÅ, while the revPBE functional yields very poor results (MAE 105 mÅ). Large errors (overestimated bond-lengths) are also found at the PBE (MAE 60 mÅ) and APBE (MAE 79 mÅ) levels because of the too high value of  $\mu$  and X-nonlocality in these two functionals. We note



that mPBESol represents almost a global minimum on the plot of Figure 14 and can thus be hardly improved for this property within the PBE GGA form.

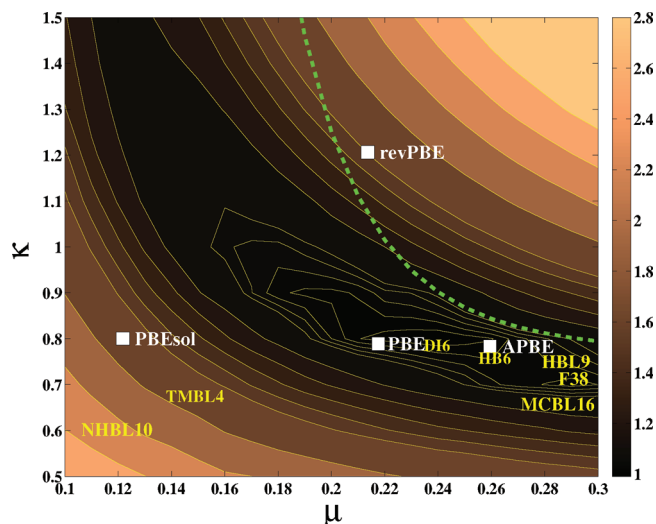
Many studies have been carried out for the construction of accurate GGAs for solids. The AM05 GGA<sup>74</sup> was constructed using the Airy gas, the uniform electron gas, and the jellium surfaces as reference systems. The PBEsol GGA, instead, recovers the second-order gradient expansion of the exchange energy, which is the right, exact constraint for solids, as also shown in Figure 14. Several other GGAs for solids have been proposed in refs 36, 40, 41, 46, and 75, and all of them have a reduced X-nonlocality, showing small gradient corrections to LSDA for small values of  $s$ . We recall that LSDA is remarkably accurate for solid-state physics. In particular, when  $\mu$  is fixed to  $\mu = \mu_{\text{PBE}} = 0.2195$ , the following values of  $\kappa$  are needed for accurate lattice constants:  $\kappa \approx 0.5$  for 4d transition metals,  $\kappa \approx 0.3$  for 5d transition metals,  $\kappa \sim \kappa^{\text{PBE}} = 0.804$  for 3d metals.<sup>75</sup>

#### 4. SUMMARY AND GLOBAL RESULTS

In the previous section, we conducted a survey of the PBE-like functionals with different  $(\mu, \kappa)$  parameters (fixing  $\beta = 3\mu/\pi^2$ ), for a broad set of tests and properties of molecular and solid-state systems. We have shown that different requirements are necessary in order to have PBE-like functionals accurate for different energetic or structural properties of molecules or solids. In this section, we summarize global results for molecular and solid-state systems. In order to be able to compare the performance of different functionals for different problems, we will normalize all MAEs to the PBE value; i.e., we will consider  $\text{MAE}/\text{MAE}(\text{PBE})$  for each property and functional and use these to compute a global MAE, denoted  $\text{MAE}_{\text{PBE}}$ .

For molecular properties, including atomization energies, bond lengths, harmonic vibrational frequencies, and noncovalent interaction energies, we found that PBE-like functionals displaying a medium X-nonlocality ( $\Lambda \sim 1/1.1$ ) yield the best overall performance thanks to the right balance between situations where a relatively small X-nonlocality is favored (TMAE4, bond lengths) and problems where a higher X-nonlocality is needed (AE6, MCAE6). The APBE functional is thus the best one when a global average is considered, with a  $\text{MAE}_{\text{PBE}}$  of 0.90 (see Table S1 in the Supporting Information), showing performance superior to that of PBE ( $\text{MAE}_{\text{PBE}} = 1.00$ ), TPSS ( $\text{MAE}_{\text{PBE}} = 0.91$ ), and the hybrid functionals (e.g.,  $\text{MAE}_{\text{PBE}}$  of PBE0 = 1.32). The performance of the latter is very poor because it largely fails for transition metal dimers. This finding confirms the importance of the semiclassical neutral atom as the reference system used in the construction of APBE, for molecules.<sup>12</sup> The revPBE functional can give accurate atomization energies, although it only outperforms APBE for MCAE6 but is not accurate for bond lengths and noncovalent interactions. Thus, it gives a total  $\text{MAE}_{\text{PBE}}$  of 1.7, showing severe limitations for broad applicability in molecular calculations. Finally, the mPBESol functional yields a  $\text{MAE}_{\text{PBE}}$  of 1.97, demonstrating its limits for the description of molecular properties.

For solid-state properties, a lower level of X-nonlocality is required, and the best overall performance is obtained with functionals having  $\Lambda \sim 0.85$ . This value constitutes a balance between the requirements of the cohesive-energy problem (medium X-nonlocality) and the lattice-constant determination (low X-nonlocality). Among the standard PBE-like functionals, the best overall performance is obtained with the mPBESol



**Figure 15.** Global mean absolute error for all properties, normalized to the PBE value. The positions corresponding to PBE, APBE, revPBE, and mPBESol are denoted by white boxes. The family of functionals defined by eq 13 is shown with a green line.

functional, having a  $\text{MAE}_{\text{PBE}}$  of 0.85 (see Table S2 in the Supporting Information). This originates mainly from its excellent performance for lattice constants, while not so accurate results are obtained for cohesive energies at the mPBESol level. The opposite occurs for PBE, which shows the smallest MAE for the cohesive energies but an error 3 times larger than that of mPBESol for lattice constants. In fact, a very good overall performance is obtained using the PBEint functional ( $\text{MAE}_{\text{PBE}} = 0.76$ ), which can correctly describe the slowly varying density regime, relevant for lattice constants and partly for cohesive energies, and the rapidly varying density limit, which is essential for the description of atomic energies used to evaluate the cohesive energies. Similar results for the PBEint functional were already found concerning the energy and structural properties of metal clusters.<sup>71</sup>

To conclude, we report in Figure 15, the global mean absolute error of all properties and test sets, normalized to the PBE value. This plot shows why PBE has been the workhorse of electronic calculations for more than a decade: this nonempirical functional shows in fact almost the best average accuracy for a large number of properties of different systems, resulting in a good choice in almost any electronic-structure problem. This finding supports the idea behind the construction of the PBE functional, which is based on a wise selection of the most important exact constraints of the exchange-correlation energy for both molecules and solids. The same global average performance is also found for the APBE functional, which has the same global MAE as PBE. Larger values of  $\text{MAE}_{\text{PBE}}$  are instead found for revPBE (1.74) and mPBESol (1.85), which do not show a broad applicability but must be instead considered specialized functionals.

We note also that, among the standard PBE-like functionals, APBE is the one which is closer, in the  $(\mu, \kappa)$  space, to the largest number of minima for different problems considered in this work (the name of each test in Figure 15 indicates approximately the position of the corresponding minimum MAE). This means that, within these functionals, it is the one that yields the best MAE for the largest number of the properties. On the other hand, APBE

has a lower accuracy for NHBL10 and TMBL4, which require a low level of X-nonlocality.

In conclusion, the APBE functional proved to be very accurate for molecular properties, competing with meta-GGA and hybrid approaches. This result confirms the importance of the recent work on semiclassical theory,<sup>7–11</sup> which brought new frontiers in density functional theory. It supports especially the role of the modified second-order gradient expansion (MGE2),<sup>9,11</sup> built from the semiclassical neutral atom theory, which can be used at the GGA level as a powerful tool for the development of XC and kinetic energy functionals.<sup>12</sup>

The results presented in this work are of great importance for the assessment of PBE-like GGA XC functionals and to understand the merits and limitations of the presently available PBE-like approximations. We showed in fact the existence of an interrelation between the values of the  $\mu$  and  $\kappa$  parameters, which must balance each other for the best performance, and the importance of considering properly the X-nonlocality measure of the functionals. The former property was recently also evidenced for kinetic energy functionals with the PBE-like form.<sup>76</sup> Thus, the present results can serve as a guide for the development and optimization of density functionals, in search of approximations having higher accuracy and broader applicability. However, it appears from the present study that there is little room for improvement within the PBE functional form, and new developments must be based on more flexible GGA expressions or based on higher rungs of the DFT Jacob's ladder. We recall that the nonempirical meta-GGAs (TPSS,<sup>68</sup> revTPSS,<sup>49,77</sup> and JS<sup>78</sup>) as well the hyper-GGA<sup>13</sup> have all been constructed using the PBE functional form. Moreover, optimization of the  $\mu$  parameter in the TPSS functional form<sup>79</sup> (that is responsible for the behavior of the meta-GGA at large  $s$ ) revealed that the use of  $\mu^{\text{APBE}} = 0.26$  improves over the original TPSS for the atomization energies of molecules, the molecular enthalpies of formation, and the barrier heights without worsening the XC jellium surface energies.<sup>79</sup> Thus, further work needs to be done for implementing the APBE ideas in meta- and hyper-GGAs.

## ■ ASSOCIATED CONTENT

**S Supporting Information.** Two-dimensional plot of the X-nonlocality measure as a function of  $\mu$  and  $\kappa$ , tests of atomization and interaction energies with  $\beta$  parameter fixed to 0.06672, performance for individual bond lengths of organic molecules, global performance for bond lengths of organic molecules, global performance for molecular systems, global performance for organic molecules, global performance for transition-metal dimers, global performance for organic-metal complexes, global performance for solid-state systems, and results for selected functionals. This information is available free of charge via the Internet at <http://pubs.acs.org/>.

## ■ AUTHOR INFORMATION

### Corresponding Author

\*E-mail: [eduardo.fabiano@nano.cnr.it](mailto:eduardo.fabiano@nano.cnr.it).

## ■ ACKNOWLEDGMENT

We thank TURBOMOLE GmbH for providing us with the TURBOMOLE program package and M. Margarito for technical support. This work was funded by the ERC Starting Grant FP7 Project DEDOM, Grant Agreement No. 207441.

## ■ REFERENCES

- (1) Hohenberg, P.; Kohn, W. *Phys. Rev.* **1964**, *136*, B864.
- (2) Parr, R. G.; Yang, W. *Density-Functional Theory of Atoms and Molecules*; Oxford University Press: Oxford, 1989; pp 1–331.
- (3) Kohn, W.; Sham, L. *Phys. Rev.* **1965**, *140*, A1133.
- (4) Perdew, J. P.; Schmidt, K. In *Density Functional Theory and Its Application to Materials*; Van Doren, V. E., Van Alsenoy, K., Geerlings, P., Eds.; American Institute of Physics: Melville, 2001; pp 1–207.
- (5) Langreth, J. P., D. C.; Perdew *Phys. Rev. B* **1980**, *21*, 5469.
- (6) Antoniewicz, P. R.; Kleinman, L. *Phys. Rev. B* **1985**, *31*, 6779–6781.
- (7) Elliot, P.; Lee, D.; Cangi, A.; Burke, K. *Phys. Rev. Lett.* **2008**, *100*, 256406.
- (8) Cangi, A.; Lee, D.; Elliot, P.; Burke, K. *Phys. Rev. Lett.* **2010**, *81*, 235128.
- (9) Elliot, P.; Burke, K. *Can. J. Chem.* **2009**, *87*, 1485.
- (10) Perdew, J. P.; Constantin, L. A.; Sagvolden, E.; Burke, K. *Phys. Rev. Lett.* **2006**, *97*, 223002.
- (11) Lee, D.; Constantin, L. A.; Perdew, J. P.; Burke, K. *J. Chem. Phys.* **2009**, *130*, 034107.
- (12) Constantin, L. A.; Fabiano, E.; Laricchia, S.; Della Sala, F. *Phys. Rev. Lett.* **2011**, *106*, 186406.
- (13) Perdew, J. P.; Staroverov, V. N.; Tao, J.; Scuseria, G. E. *Phys. Rev. A* **2008**, *78*, 052513.
- (14) Becke, A. D. *J. Chem. Phys.* **1993**, *98*, 5648.
- (15) Kümmel, S.; Kronik, L. *Rev. Mod. Phys.* **2008**, *80*, 3.
- (16) Gritsenko, O. V.; Baerends, E. J. *Phys. Rev. A* **2001**, *64*, 042506.
- (17) Della Sala, F.; Görling, A. *J. Chem. Phys.* **2001**, *115*, 5718.
- (18) Grabowski, L.; Hirata, S.; Ivanov, S.; Bartlett, R. J. *J. Chem. Phys.* **2002**, *116*, 4415–4425.
- (19) Furche, F.; Voorhis, T. V. *J. Chem. Phys.* **2005**, *122*, 164106.
- (20) Fabiano, E.; Della Sala, F. *J. Chem. Phys.* **2007**, *126*, 214102.
- (21) Schimka, L.; Harl, J.; Stroppa, A.; Grüneis, A.; Marsman, M.; Mittendorfer, F.; Kresse, G. *Nat. Mater.* **2010**, *9*, 741.
- (22) Ruzsinszky, A.; Perdew, J. P.; Csonka, G. I. *J. Chem. Theory Comput.* **2010**, *6*, 127–134.
- (23) Ren, X.; Tkatchenko, A.; Rinke, P.; Scheffler, M. *Phys. Rev. Lett.* **2011**, *106*, 153003.
- (24) Koch, W.; Holthausen, M. C. *A Chemist's Guide to Density Functional Theory*; Wiley-VCH: New York, 2001; pp 1–293.
- (25) Schultz, N. E.; Zhao, Y.; Truhlar, D. G. *J. Phys. Chem. A* **2005**, *109*, 4388–4403 PMID: 16833770.
- (26) Schultz, N. E.; Zhao, Y.; Truhlar, D. G. *J. Phys. Chem. A* **2005**, *109*, 11127–11143.
- (27) Furche, F.; Perdew, J. P. *J. Chem. Phys.* **2006**, *124*, 044103.
- (28) Perdew, J. P.; Burke, K.; Ernzerhof, M. *Phys. Rev. Lett.* **1996**, *77*, 3865.
- (29) Perdew, J. P.; Burke, K.; Wang, Y. *Phys. Rev. B* **1996**, *54*, 16533–16539.
- (30) Becke, A. D. *J. Chem. Phys.* **1986**, *84*, 4524–4529.
- (31) Oliver, G. L.; Perdew, J. P. *Phys. Rev. A* **1979**, *20*, 397–403.
- (32) Zhang, Y.; Yang, W. *Phys. Rev. Lett.* **1998**, *80*, 890.
- (33) Hammer, B.; Hansen, L. B.; Nørskov, J. K. *Phys. Rev. B* **1999**, *59*, 7413–7421.
- (34) Adamo, C.; Barone, V. *J. Chem. Phys.* **2002**, *116*, 5933–5940.
- (35) Xu, X.; W. A., G., III *J. Chem. Phys.* **2004**, *121*, 4068–4082.
- (36) Wu, Z.; Cohen, R. E. *Phys. Rev. B* **2006**, *73*, 235116.
- (37) Perdew, J. P.; Ruzsinszky, A.; Csonka, G. I.; Vydrov, O. A.; Scuseria, G. E.; Constantin, L. A.; Zhou, X.; Burke, K. *Phys. Rev. Lett.* **2008**, *100*, 136406.
- (38) Perdew, J. P.; Ruzsinszky, A.; Csonka, G. I.; Vydrov, O. A.; Scuseria, G. E.; Constantin, L. A.; Zhou, X.; Burke, K. *Phys. Rev. Lett.* **2009**, *102*, 039902.
- (39) Madsen, G. K. H. *Phys. Rev. B* **2007**, *75*, 195108.
- (40) Zhao, Y.; Truhlar, D. G. *J. Chem. Phys.* **2008**, *128*, 184109.
- (41) Pedroza, L. S.; da Silva, A. J. R.; Capelle, K. *Phys. Rev. B* **2009**, *79*, 201106.

- (42) Ruzsinszky, A.; Csonka, G. I.; Scuseria, G. E. *J. Chem. Theory Comput.* **2009**, *5*, 763–769.
- (43) Goerigk, L.; Grimme, S. *J. Chem. Theory Comput.* **2010**, *6*, 107–126.
- (44) Fabiano, E.; Constantin, L. A.; Della Sala, F. *Phys. Rev. B* **2010**, *82*, 113104.
- (45) Haas, P.; Tran, F.; Blaha, P.; Schwarz, K. *Phys. Rev. B* **2011**, *83*, 205117.
- (46) Haas, P.; Tran, F.; Blaha, P.; Pedroza, L. S.; da Silva, A. J. R.; Odashima, M. M.; Capelle, K. *Phys. Rev. B* **2010**, *81*, 125136.
- (47) Csonka, G. I.; Perdew, J. P.; Ruzsinszky, A.; Philippen, P. H. T.; Lebègue, S.; Paier, J.; Vydrov, O. A.; Ángyán, J. G. *Phys. Rev. B* **2009**, *79*, 155107.
- (48) Odashima, M. M.; Capelle, K. *J. Chem. Phys.* **2007**, *127*, 054106.
- (49) Perdew, J. P.; Ruzsinszky, A.; Csonka, G. I.; Constantin, L. A.; Sun, J. *Phys. Rev. Lett.* **2009**, *103*, 026403.
- (50) Goerigk, L.; Grimme, S. *J. Chem. Theory Comput.* **2011**, *7*, 291–309.
- (51) Grimme, S. *Wiley Interdisciplinary Reviews: Computational Molecular Science* **2011**, *1*, 211–228.
- (52) Korth, M.; Grimme, S. *J. Chem. Theory Comput.* **2009**, *5*, 993–1003.
- (53) Lynch, B. J.; Truhlar, D. G. *J. Phys. Chem. A* **2003**, *107*, 8996–8999.
- (54) Zhao, Y.; Truhlar, D. G. *J. Chem. Phys.* **2006**, *125*, 194101.
- (55) Biczysko, M.; Panek, P.; Scalmani, G.; Bloino, J.; Barone, V. *J. Chem. Theory Comput.* **2010**, *6*, 2115–2125.
- (56) Zhao, Y.; Truhlar, D. G. *J. Chem. Theory Comput.* **2005**, *1*, 415–432.
- (57) Haas, P.; Tran, F.; Blaha, P. *Phys. Rev. B* **2009**, *79*, 085104.
- (58) TURBOMOLE V6.2, 2009, a development of University of Karlsruhe and Forschungszentrum Karlsruhe GmbH, 1989–2007, TURBOMOLE GmbH, since 2007; available from <http://www.turbomole.com> (accessed September 2011).
- (59) Weigend, F.; Furche, F.; Ahlrichs, R. *J. Chem. Phys.* **2003**, *119*, 12753.
- (60) Weigend, F.; Ahlrichs, R. *Phys. Chem. Chem. Phys.* **2005**, *7*, 3297.
- (61) Blum, V.; Gehrke, R.; Hanke, F.; Havu, P.; Havu, V. H.; Ren, X.; Reuter, K.; Scheffler, M. *Comput. Phys. Commun.* **2009**, *180*, 2175–2196.
- (62) Havu, V.; Blum, V.; Havu, P.; Scheffler, M. *J. Comput. Phys.* **2009**, *228*, 8367–8379.
- (63) Faas, S.; Snijders, J. G.; van Lenthe, J. H.; van Lenthe, E.; Baerends, E. J. *Chem. Phys. Lett.* **1995**, *246*, 632–640.
- (64) Becke, A. D. *Phys. Rev. A* **1988**, *38*, 3098.
- (65) Lee, C.; Yang, W.; Parr, R. G. *Phys. Rev. B* **1988**, *37*, 785.
- (66) Handy, A. J.; Cohen, N. C. *Mol. Phys.* **2001**, *99*, 403.
- (67) Hoe, W.-M.; Cohen, A. J.; Handy, N. C. *Chem. Phys. Lett.* **2001**, *341*, 319–328.
- (68) Tao, J.; Perdew, J. P.; Staroverov, V. N.; Scuseria, G. E. *Phys. Rev. Lett.* **2003**, *91*, 146401.
- (69) Adamo, C.; Barone, V. *J. Chem. Phys.* **1999**, *110*, 6158.
- (70) Perdew, J. P.; Ernzerhof, M.; Zupan, A.; Burke, K. *J. Chem. Phys.* **1998**, *108*, 1522–1531.
- (71) Fabiano, E.; Constantin, L. A.; Della Sala, F. *J. Chem. Phys.* **2011**, *134*, 194112.
- (72) Tao, J.; Perdew, J. P.; Ruzsinszky, A. *Phys. Rev. B* **2010**, *81*, 233102.
- (73) Baletto, F.; Ferrando, R. *Rev. Mod. Phys.* **2005**, *77*, 371–423.
- (74) Armiento, R.; Mattsson, A. E. *Phys. Rev. B* **2005**, *72*, 085108.
- (75) Peltzer y Blancá, E. L.; Rodríguez, C. O.; Shitu, J.; Novikov, D. L. *Journal of Physics: Condensed Matter* **2001**, *13*, 9463.
- (76) Laricchia, S.; Fabiano, E.; Constantin, L. A.; Della Sala, F. *J. Chem. Theory Comput.* **2011**, *7*, 2439.
- (77) Perdew, J. P.; Ruzsinszky, A.; Csonka, G. I.; Constantin, L. A.; Sun, J. *Phys. Rev. Lett.* **2011**, *106*, 179902.
- (78) Constantin, L. A.; Chiodo, L.; Fabiano, E.; Bodrenko, I.; Della Sala, F. *Phys. Rev. B* **2011**, *84*, 045126.
- (79) Perdew, J. P.; Ruzsinszky, A.; Tao, J.; Csonka, G. I.; Scuseria, G. E. *Phys. Rev. A* **2007**, *76*, 042406.

# Determination of Local Spins by Means of a Spin-Free Treatment

Diego R. Alcoba,<sup>†</sup> Alicia Torre,<sup>‡</sup> Luis Lain,<sup>\*,‡</sup> and Roberto C. Boichicchio<sup>†</sup>

<sup>†</sup>Departamento de Física, Facultad de Ciencias Exactas y Naturales, Universidad de Buenos Aires and Instituto de Física de Buenos Aires, Consejo Nacional de Investigaciones Científicas y Técnicas, Ciudad Universitaria, 1428, Buenos Aires, Argentina

<sup>‡</sup>Departamento de Química Física, Facultad de Ciencia y Tecnología, Universidad del País Vasco, Apdo. 644 E-48080 Bilbao, Spain

**ABSTRACT:** This work describes a Mulliken-type partitioning of the expectation value of the spin-squared operator  $\langle \hat{S}^2 \rangle$  corresponding to an  $N$ -electron system. Our algorithms, which are based on a spin-free formulation, predict appropriate spins for the molecular fragments (at equilibrium geometries and at dissociation limits) and can be applied to any spin symmetry. Numerical determinations performed in selected closed- and open-shell systems at correlated level are reported. A comparison between these results and their counterpart ones arising from other alternative approaches is analyzed in detail.

## 1. INTRODUCTION

The study of procedures to decompose the expectation value of the spin-squared operator  $\langle \hat{S}^2 \rangle$  corresponding to an  $N$ -electron system into one- and two-center terms (local spins) has attracted attention of a considerable number of authors in the last years. This interest arises from the ability of the local spins to determine the spin state of an atom or group of atoms in a molecule, radical, cluster, etc., as well as to describe magnetic interactions between the atoms which compose the system. In fact, spin–spin coupling constants can be calculated by means of two-center local spins within the well-known Heisenberg Hamiltonian model. The partitioning of the  $\langle \hat{S}^2 \rangle$  quantity has been performed using several approaches. One of them utilizes the technique of local projection operators, in which the total spin-squared operator  $\hat{S}^2$  is decomposed into one- and two-center operators associated with the nuclei of the system; then in a subsequent step the expectation values of these operators are evaluated for different approximations of the wave function.<sup>1–7</sup> Alternatively, the partitioning of the expectation value  $\langle \hat{S}^2 \rangle$  has also been performed in a direct way.<sup>8–13</sup> Within the framework of this last procedure the quantity  $\langle \hat{S}^2 \rangle$ , expressed in terms of elements of reduced density matrices and related quantities, is partitioned in the Hilbert space of the atomic basis set according to a Mulliken-type population analysis. More recently, this technique of partitioning has also been extended to the three-dimensional physical space and its results compared with those arising from the Hilbert space.<sup>14</sup>

This work deals with the partitioning of the  $\langle \hat{S}^2 \rangle$  quantity in the Hilbert space. Determinations of local spins in that space at the level of single Slater determinant wave functions and higher correlation levels have been described in refs 8, 9, 11, and 13. These reported results are satisfactory from a chemical point of view since they show appropriate spins for the fragments at the dissociation limit and zero local spin values for closed-shell systems described at the restricted Hartree–Fock level. However, at correlated level, the algorithms used to get these results depend on the spin blocks of the second-order reduced density matrix, which, in practice, are not available in most standard codes in quantum chemistry. Besides, these matrix elements depend on the substate  $S_z$  corresponding to a determined spin  $S$

for nonsinglet states. Consequently, the values of the terms derived from that  $\langle \hat{S}^2 \rangle$  partitioning are  $S_z$  dependent. Obviously, the partitioning of a quantity into several components is usually not unique. Hence, it is important to consider other possibilities which can also produce physically reasonable results in those limit cases, provided they present additional theoretical and practical advantages. The aim of this work is to overcome the mentioned drawbacks, reporting an algorithm in terms of spin-free tools, so that the local spins of a system can be calculated for any state of any spin symmetry, fulfilling the physical requirement of uniqueness for the spin multiplet components (in absence of magnetic fields). Our algorithm is based on the use of the one-electron effectively unpaired electron density matrix<sup>15–17</sup> and the two-electron spin-free cumulant matrix of the spin-free second-order reduced density matrix;<sup>18,19</sup> both matrices are directly calculable from the spin-free first- and second-order reduced density matrices, which can be obtained from standard codes.

The organization of this work is as follows. The second section describes a straightforward derivation of the formulas used in refs 11,13 to evaluate one- and two-center local spins at correlated level. In this way, we point out their  $S_z$  dependence and the difficulties to access to elements of the cumulant matrix in the spin–orbital representation, in standard codes, mainly for nonsinglet states. In the third section, we propose an alternative algorithm which only utilizes matrix elements of spin-free quantities. In the fourth section, we describe the results obtained from both  $S_z$ -dependent and  $S_z$ -independent algorithms for some selected closed- and open-shell systems, as well as their corresponding discussion. A study of the dependence of the results on the degree of correlation used is also included in this section. Finally, in the last section we summarize the concluding remarks of this work.

## 2. PARTITIONING OF $\langle \hat{S}^2 \rangle$ AT CORRELATED LEVEL

A finite basis set of orthonormal orbitals will be denoted by  $\{i, j, k, l, \dots\}$ ; in this basis set  ${}^1D_j^i$  and  ${}^2D_{jl}^{ik}$  will stand for the

**Received:** July 29, 2011

**Published:** September 08, 2011

spin-free matrix elements corresponding to the first- and second-order reduced density matrices of an  $N$ -electron system in a state  $\Psi$ , respectively. The trace of the first-order reduced density matrix is normalized to  $\text{tr}(^1D) = N$  and that of the second-order one may be normalized to

$$\text{tr}(^2D) = \binom{N}{2}$$

or to  $\text{tr}(^2D) = N(N-1)$ ;<sup>11,13</sup> in this work we will use the former procedure. The expectation value of the spin-squared operator  $\hat{S}^2$ ,  $\langle \hat{S}^2 \rangle = \langle \Psi | \hat{S}^2 | \Psi \rangle$ , can be expressed as follows:<sup>20,21</sup>

$$\langle \hat{S}^2 \rangle = N - \frac{N^2}{4} - \sum_{i,k} {}^2D_{ik}^{ik} \quad (1)$$

Likewise, taking into account the values of those traces, eq 1 can be written as follows:

$$\langle \hat{S}^2 \rangle = \frac{3}{4} \sum_i {}^1D_i^i - \frac{1}{2} \sum_{i,k} {}^2D_{ik}^{ik} - \sum_{i,k} {}^2D_{ki}^{ik} \quad (2)$$

The decomposition of these spin-free matrix elements according to their spin orbitals, that is,  ${}^1D_i^i = ({}^1D_{i\alpha}^{\alpha\alpha} + {}^1D_{i\beta}^{\beta\beta})$  and  ${}^2D_{ij}^{ik} = ({}^2D_{i\alpha k\alpha}^{\alpha\alpha k\alpha} + {}^2D_{j\alpha i\beta}^{\beta\beta k\alpha} + {}^2D_{i\beta k\alpha}^{\alpha\alpha k\beta} + {}^2D_{j\beta i\alpha}^{\beta\beta k\beta})$ , and an appropriate permutation of spin orbitals, based on the anticommutation rules of fermion operators, leads to the following:

$$\begin{aligned} \langle \hat{S}^2 \rangle &= \frac{3}{4} \sum_i ({}^1D_{i\alpha}^{\alpha\alpha} + {}^1D_{i\beta}^{\beta\beta}) - \sum_{i,k} {}^2D_{i\beta k\alpha}^{\alpha\alpha k\alpha} \\ &\quad - \frac{1}{2} \sum_{i,k} {}^2D_{k\alpha i\alpha}^{\alpha\alpha k\alpha} - \frac{1}{2} \sum_{i,k} {}^2D_{k\beta i\beta}^{\beta\beta k\beta} - 2 \sum_{i,k} {}^2D_{k\alpha i\beta}^{\alpha\alpha k\beta} \end{aligned} \quad (3)$$

A trivial but tedious algebra, which consists in relating the second-order reduced density matrix elements in the spin-orbital representation, with the corresponding elements of its cumulant matrix  $\Gamma_{j\sigma l\sigma'}^{i\alpha k\alpha}$  ( $\sigma, \sigma' = \alpha, \beta$ ), that is as follows:<sup>22</sup>

$${}^2D_{j\sigma l\sigma'}^{i\alpha k\alpha} = \frac{1}{2} {}^1D_{j\sigma}^{\sigma\sigma} {}^1D_{l\sigma'}^{k\sigma'} - \frac{1}{2} {}^1D_{l\sigma'}^{\sigma\sigma} {}^1D_{j\sigma}^{k\sigma'} + \frac{1}{2} \Gamma_{j\sigma l\sigma'}^{i\alpha k\alpha} \quad (4)$$

provides to express eq 3 as follows:

$$\begin{aligned} \langle \hat{S}^2 \rangle &= \frac{1}{2} \sum_{i,k} (P^s)_k^i (P^s)_i^k + \frac{1}{4} \sum_{i,k} (P^s)_i^i (P^s)_k^k \\ &\quad + \frac{3}{4} \sum_i \left[ {}^1D_{i\alpha}^{\alpha\alpha} - \sum_k {}^1D_{k\alpha}^{\alpha\alpha} {}^1D_{i\alpha}^{k\alpha} + {}^1D_{i\beta}^{\beta\beta} - \sum_k {}^1D_{k\beta}^{\beta\beta} {}^1D_{i\beta}^{k\beta} \right] \\ &\quad - \sum_{i,k} \left[ 2 {}^2D_{i\beta k\alpha}^{\alpha\alpha k\alpha} - \frac{1}{2} {}^1D_{i\beta}^{\beta\beta} {}^1D_{k\alpha}^{k\alpha} \right] \\ &\quad - \frac{1}{2} \sum_{i,k} \left[ 2 {}^2D_{k\alpha i\alpha}^{\alpha\alpha k\alpha} - \frac{1}{2} {}^1D_{k\alpha}^{\alpha\alpha} {}^1D_{i\alpha}^{k\alpha} + \frac{1}{2} {}^1D_{i\alpha}^{\alpha\alpha} {}^1D_{k\alpha}^{k\alpha} \right] \\ &\quad - \frac{1}{2} \sum_{i,k} \left[ 2 {}^2D_{k\beta i\beta}^{\beta\beta k\beta} - \frac{1}{2} {}^1D_{k\beta}^{\beta\beta} {}^1D_{i\beta}^{k\beta} + \frac{1}{2} {}^1D_{i\beta}^{\beta\beta} {}^1D_{k\beta}^{k\beta} \right] \\ &\quad - 2 \sum_{i,k} \left[ 2 {}^2D_{k\alpha i\beta}^{\alpha\alpha k\beta} - \frac{1}{2} {}^1D_{k\alpha}^{\alpha\alpha} {}^1D_{i\beta}^{k\beta} \right] \end{aligned} \quad (5)$$

In eq 5,  $(P^s)_j^i = {}^1D_{j\alpha}^{\alpha\alpha} - {}^1D_{j\beta}^{\beta\beta}$  are the elements of the spin-density matrix and the second-order reduced density matrix has been

normalized by the value of its trace.

$$\text{tr}(^2D) = \binom{N}{2}$$

The derivation of this equation has required to add and to subtract terms in order to express the  $\langle \hat{S}^2 \rangle$  quantity by means of the spin-orbital components of the cumulant matrix of the second-order reduced density matrix, which are the last four brackets (see eq 4).

Formula 5, which is expressed in an orthogonal basis set, is equivalent to those reported in refs,11 and 13 expressed in nonorthogonal atomic basis sets; a simple basis transformation allows one to pass from this formula to the others. The partitioning of the quantity  $\langle \hat{S}^2 \rangle$  according to formula 5 transformed to the atomic basis set, requires to know the values of the elements of the second-order reduced density matrix in the spin-orbital representation ( ${}^2D_{j\sigma l\sigma'}^{i\alpha k\alpha}$ ,  $\sigma, \sigma' = \alpha, \beta$ ), which usually are not provided by the execution of most standard codes. Apart from this shortcoming, another aspect to take into account is that, as is well-known, those matrix elements depend on the  $S_z$  substate of the state  $\Psi$  and consequently the local spin results for nonsinglet states turn out to be  $S_z$  dependent. Thus, the requirement of uniqueness for the spin multiplet components is not fulfilled by this partitioning. This aspect has been numerically tested in the lowest triplet state of the system  $\text{HeH}^+$  (see Appendix A). As has been mentioned in the Introduction, the purpose of this work is to set up a spin-free  $S_z$ -independent algorithm that avoids these drawbacks. In the next section, we report that algorithm.

### 3. SPIN-FREE TREATMENT PROPOSAL

We will express the elements of the spin-free second-order reduced density matrix as follows:<sup>18,19</sup>

$${}^2D_{jl}^{ik} = \frac{1}{2} {}^1D_j^i {}^1D_l^k - \frac{1}{4} {}^1D_i^i {}^1D_j^k + \frac{1}{2} \Lambda_{jl}^{ik} \quad (6)$$

in which  $\Lambda_{jl}^{ik}$  stands for the elements of the spin-free cumulant matrix of that second-order reduced density matrix. These matrix elements are related with those of the cumulant matrix ones ( $\Gamma_{j\sigma l\sigma'}^{i\alpha k\alpha}$ ) by the following:<sup>23</sup>

$$\Lambda_{jl}^{ik} = -\frac{1}{2} (P^s)_i^i (P^s)_j^k + \sum_{\sigma, \sigma'} \Gamma_{j\sigma l\sigma'}^{i\alpha k\alpha} \quad (7)$$

However, we will regard the effectively unpaired electron density matrix  $u$ , initially defined by Takatsuka et al.<sup>15</sup> as follows:

$$u_j^i = 2 {}^1D_j^i - \sum_k {}^1D_k^i {}^1D_j^k \quad (8)$$

The mathematical features of this matrix have been widely studied<sup>16,17,23–25</sup> and utilized in a great variety of population analysis studies.<sup>26–32</sup> Although other formulations of the matrix  $u$  have been proposed,<sup>33–35</sup> in this work we will use that formulated by eq 8 whose relation with the  $\Lambda$  matrix turns out to be following:<sup>17</sup>

$$u_j^i = -2 \sum_k \Lambda_{jk}^{ik} \quad (9)$$

The substitution of the elements  ${}^2D_{ik}^{ik}$  and  ${}^2D_{ki}^{ik}$  according to eq 6 and the use of eqs 8 and 9 provide to express the quantity

$\langle \hat{S}^2 \rangle$  in formula 2 as follows:

$$\langle \hat{S}^2 \rangle = \frac{1}{2} \sum_i u_i^i - \frac{1}{2} \sum_{i,k} \Lambda_{ki}^{ik} \quad (10)$$

However, in order to partition the  $\langle \hat{S}^2 \rangle$  quantity into one-center and two-center terms it is more useful to express that equation in the basis set of the atomic orbitals  $\{\mu, \nu, \lambda, \gamma, \dots\}$

$$\langle \hat{S}^2 \rangle = \frac{1}{2} \sum_{\mu} (uS)_{\mu}^{\mu} - \frac{1}{2} \sum_{\mu, \nu, \lambda, \gamma} (S)_{\lambda}^{\mu} \Lambda_{\gamma\mu}^{\lambda\nu} (S)_{\nu}^{\gamma} \quad (11)$$

where  $(S)_{\nu}^{\mu} = \langle \mu | \nu \rangle$  are the elements of the overlap matrix of the atomic orbitals.

The decomposition of the expectation value  $\langle \hat{S}^2 \rangle$  in the Hilbert space of atomic orbitals into one-center terms  $\langle \hat{S}^2 \rangle_A$  and two-center terms  $\langle \hat{S}^2 \rangle_{AB}$ :

$$\langle \hat{S}^2 \rangle = \sum_A \langle \hat{S}^2 \rangle_A + \sum_{A \neq B} \langle \hat{S}^2 \rangle_{AB} \quad (12)$$

is performed assigning every atomic function  $\mu$  to one nucleus  $A$ . Although the matrix  $\Lambda_{\gamma\mu}^{\lambda\nu}$  possesses four indices, in this work we have limited to determine only one- and two-center local spins (excluding the three- and four-center contributions) so that two of those indices have been mathematically removed by means of a sum over them. Hence, the expressions for the  $\langle \hat{S}^2 \rangle_A$  and  $\langle \hat{S}^2 \rangle_{AB}$  quantities according to eq 11 are as follows

$$\langle \hat{S}^2 \rangle_A = \frac{1}{2} \sum_{\mu \in A} (uS)_{\mu}^{\mu} - \frac{1}{2} \sum_{\mu \in A, \nu \in A} \sum_{\lambda, \gamma} (S)_{\lambda}^{\mu} \Lambda_{\gamma\mu}^{\lambda\nu} (S)_{\nu}^{\gamma} \quad (13)$$

and,

$$\langle \hat{S}^2 \rangle_{AB} = -\frac{1}{2} \sum_{\mu \in A, \nu \in B} \sum_{\lambda, \gamma} (S)_{\lambda}^{\mu} \Lambda_{\gamma\mu}^{\lambda\nu} (S)_{\nu}^{\gamma} \quad (14)$$

Formulas 13 and 14 provide the means to carry out numerical determinations of one-center and two-center local spins, respectively. Because the matrix elements  $u_{\mu}^{\mu}$  and  $\Lambda_{\gamma\mu}^{\lambda\nu}$  are spin free,  $S_z$ -independent quantities, the values of  $\langle \hat{S}^2 \rangle_A$  and  $\langle \hat{S}^2 \rangle_{AB}$  arisen from these equations are independent of the quantum number  $S_z$ . Consequently, local spin evaluations can be obtained for any state of any spin symmetry, fulfilling the conditions of invariance for all the components of a multiplet state. In practice, the matrices  $\Lambda$  and  $u$  are calculated by means of eqs 6 and 8, respectively. Hence, from a computational point of view the unique required matrices to implement local spin determinations are the overlap matrix  $S$  and the first- and second-order reduced density matrices  ${}^1D$  and  ${}^2D$ , all of them in the spin-free formulation, which are usually drawn from standard codes. In the next section, we report results of local spins arising from this procedure which are compared with those obtained in other treatments.

#### 4. NUMERICAL DETERMINATIONS AND DISCUSSION

The elements of the overlap matrices and those of the spin-free first- and second-order reduced density matrices have been obtained from a modified version of the PSI 3.3 package.<sup>36</sup> In a subsequent step, we have used our own codes to evaluate local spins using eqs 13 and 14 within the spin-free treatment. We have also performed determinations of local spins by means of eq 5, expressed in the atomic basis sets, for singlet states with a unique substate  $S_z = 0$ , for which the spin blocks of the second-order

reduced density matrix can be calculated from its spin-free matrix elements.<sup>37</sup> As has been mentioned in section 2, this procedure was reported in refs 11 and 13. Table 1 gathers the results arising from both algorithms (denoted by *spin-free* and *with spin* in that Table) for singlet states, in order to carry out an appropriate comparison between them. Likewise, in Tables 2 and 3, we report results of systems in doublet and triplet spin symmetries respectively, obtained from eqs 13 and 14. The computational details are shown in these Tables, i.e., the basis sets and the experimental geometries used<sup>38–42</sup> as well as the correlation levels utilized, full configuration interaction (FCI), configuration interaction with single and double excitations (CISD), etc.

A survey of the results for singlet states reported in Table 1 shows that at equilibrium distances, the one- and two-center local spins absolute values are a little lower in the spin-free treatment than in that denominated with spin. However, these series of values become almost coincident at distances near the dissociation limits of these molecules. As can be observed in that Table, in both treatments the systems  $H_2^a$ ,  $Li_2^a$ ,  $Be_2^a$ , and  $C_2H_4^a$  exhibit values for the one-center local spins which are very close to those corresponding to the dissociated fragments. In the case of the ethylene molecule, the values reported for the system denoted by  $C_2H_4^a$  refer to its dissociation into two triplet methylene groups by stretching the bond distance C–C. However, for singlets at equilibrium distances, it seems reasonable to expect that the distribution of the  $\langle \hat{S}^2 \rangle$  quantity along the whole molecule presents not too high local spin values and consequently, from a genuine chemical point of view, the lower values found in the spin-free treatment can be regarded as a favorable tendency. This behavior is followed by all systems included in Table 1, the light ones ( $H_2$ ,  $Li_2$  and  $Be_2$ ), the hydrides of the second row ( $HF$ ,  $H_2O$  and  $NH_3$ ) and the hydrocarbons ( $CH_4$ ,  $C_2H_6$ ,  $C_2H_4$  and  $C_2H_2$ ), at the reported correlation levels. Another aspect to highlight is that both treatments predict identical signs for counterpart values of the one- and two-center local spins. As has been pointed out in refs 11 and 13 an adequate partitioning of the  $\langle \hat{S}^2 \rangle$  quantity requires that the atomic spins for atoms at large distances reproduce the spins of the free atomic fragments, as well as to predict zero spins for systems described by closed-shell restricted Hartree–Fock (RHF) wave functions, which would correspond to a pure covalent description. The spin-free algorithm that we have described in section 3 fulfills both requirements, i.e., it leads to suitable spin values at the dissociation limits and provides values  $\langle \hat{S}^2 \rangle_A = 0$  and  $\langle \hat{S}^2 \rangle_{AB} = 0$  for RHF wave functions (see eqs 13 and 14) because all the elements of the matrices  $u$  and  $\Lambda$  are zero for that type of wave functions.<sup>16,19</sup> Moreover, in the unrestricted Hartree–Fock (UHF) case the elements of those matrices are nonzero and a simple algebra shows that eqs 5 and 10 are transformed to an identical expression.

Tables 2 and 3 show results of local spin evaluations within the spin-free treatment for species (molecules and radicals) doublets and triplets at the experimental equilibrium distances, except for systems  $H_2^a$  and  $Li_2^a$  (in Table 3) which refer to the lowest triplet states at distances near the dissociation limit. These results have also been obtained from eqs 13 and 14 since they are valid for any quantum number  $S$ . As can be seen in Table 2, the radicals hydroxyl, cyano and amino present high values of the one-center contribution  $\langle \hat{S}^2 \rangle_A$  in the atoms oxygen, carbon, and nitrogen, respectively, indicating that the unpaired electron which originates the doublet spin symmetry is located on those atoms. The NO molecule shows a distribution of the spin cloud between the nitrogen and oxygen atoms although the value  $\langle \hat{S}^2 \rangle_N$  is

**Table 1. Local Spins of One- And Two-Centers ( $\langle \hat{S}^2 \rangle_A$  and  $\langle \hat{S}^2 \rangle_{AB}$ ) Arising from the Treatments Spin-Free (eqs 13 and 14) and with Spin (eq 5 in the Atomic Basis Set) for Singlet Systems in the Ground State at Experimental Equilibrium Distances ( $^a$  near Dissociation Limits)**

system	atom/ bond	spin-free		with spin		basis set/ method		
		$\langle \hat{S}^2 \rangle_A$	$\langle \hat{S}^2 \rangle_{AB}$	$\langle \hat{S}^2 \rangle_A$	$\langle \hat{S}^2 \rangle_{AB}$			
H <sub>2</sub>	H	0.100		0.116		6-31G/FCI		
	HH		-0.100		-0.116			
H <sub>2</sub> <sup>a</sup>	H	0.743		0.744		6-31G/FCI		
	HH		-0.743		-0.744			
Li <sub>2</sub>	Li	0.204		0.210		STO-3G/FCI		
	LiLi		-0.204		-0.210			
Li <sub>2</sub> <sup>a</sup>	Li	0.750		0.750		STO-3G/FCI		
	LiLi		-0.750		-0.750			
Be <sub>2</sub>	Be	0.125		0.127		STO-3G/FCI		
	BeBe		-0.125		-0.127			
Be <sub>2</sub> <sup>a</sup>	Be	0.000		0.000		STO-3G/FCI		
	BeBe		0.000		0.000			
C <sub>2</sub> H <sub>4</sub>	C	0.477		0.544		6-31G/CISD		
	H	0.058		0.067				
	CC		-0.365		-0.411			
	CH		-0.098		-0.108			
	C...H		0.042		0.041			
	HH		0.008		0.008			
	H...H		-0.005		-0.003			
	H...H		-0.006		-0.005			
	C <sub>2</sub> H <sub>4</sub> <sup>a</sup>	C	1.884		1.875			6-31G/CISD
		H	0.020		0.015			
CC			-1.850		-1.847			
C...H			-0.038		-0.038			
CH			0.019		0.023			
HH			0.000		0.000			
H...H			0.000		0.000			
H...H			0.000		0.000			
HF	F	0.050		0.059		6-31G/CISD		
	H	0.050		0.059				
	FH		-0.050		-0.059			
H <sub>2</sub> O	O	0.121		0.141		6-31G/CISD		
	H	0.055		0.064				
	OH		-0.060		-0.070			
NH <sub>3</sub>	N	0.226		0.257		6-31G/CISD		
	H	0.059		0.068				
	NH		-0.075		-0.086			
CH <sub>4</sub>	CH		0.005		0.006	6-31G/CISD		
	HH		0.008		0.009			
	C	0.428		0.470				
	H	0.066		0.077				
C <sub>2</sub> H <sub>6</sub> (staggered)	CH		-0.107		-0.118	6-31G/CISD		
	HH		0.014		0.014			
	C	0.361		0.400				
	H	0.059		0.068				
	CC		-0.130		-0.136			
	CH		-0.094		-0.104			
	C...H		0.017		0.016			

**Table 1. Continued**

system	atom/ bond	spin-free		with spin		basis set/ method
		$\langle \hat{S}^2 \rangle_A$	$\langle \hat{S}^2 \rangle_{AB}$	$\langle \hat{S}^2 \rangle_A$	$\langle \hat{S}^2 \rangle_{AB}$	
C <sub>2</sub> H <sub>2</sub>	H...H		-0.003		-0.003	6-31G/CISD
	HH		0.012		0.012	
	C	0.603		0.699		
	H	0.049		0.055		
	CC		-0.558		-0.646	
	C...H		-0.036		-0.037	
	CH		-0.081		-0.090	
	HH		-0.003		-0.002	

considerably higher than the  $\langle \hat{S}^2 \rangle_O$  one, which agrees with the well-known structural features of that molecule. The doublet radicals CH and CH<sub>3</sub> (in Table 2) and the triplet one CH<sub>2</sub> (in Table 3) also show a clear localization of the unpaired electrons on the carbon atom. In the series of radicals ethyl, vinyl, and ethynyl the main localization of unpaired electrons appears on the carbon atom linked to less hydrogen atoms, denominated as C<sup>(2)</sup> in Table 2, which is markedly larger than on the other carbon atom C<sup>(1)</sup>. The application of this methodology to the study of the allyl radical leads to show that the carbon atoms C<sup>(1)</sup> and C<sup>(3)</sup> (Table 2) are equivalent, presenting a higher  $\langle \hat{S}^2 \rangle_A$  value in these atoms than in the C<sup>(2)</sup> one. This behavior is well-known in this species and explained in terms of the resonance of double bond between the atoms C<sup>(1)</sup> and C<sup>(2)</sup> and between the atoms C<sup>(2)</sup> and C<sup>(3)</sup>. In this radical, the two-center local spin  $\langle \hat{S}^2 \rangle_{C(1)C(3)}$  turns out to be 0.247 which is a positive and non-negligible value; this fact can be interpreted in terms of delocalization of the unpaired electron. The presence of this feature, that is, positive non-negligible values for the two-center local spins, has also been found in the triplet homonuclear diatomic molecules O<sub>2</sub> and C<sub>2</sub> as well as in the B<sub>2</sub>H<sub>2</sub> molecule (Table 3). The explanation of these values must be done again in terms of delocalization of unpaired electrons. The rest of the triplet systems described in Table 3, NF, NH, and C<sub>2</sub>H<sub>4</sub> exhibit local spin features in agreement with those presented in Table 2 and consequently deserve similar comments. For the two triplets reported in Table 3 at the dissociation limit, our treatment describes values of one-center contribution  $\langle \hat{S}^2 \rangle_A = 0.750$  for both H<sub>2</sub><sup>a</sup> and Li<sub>2</sub><sup>a</sup> molecules. These results coincide with those reported for these systems in Table 1 (singlet states), leading again to right  $\langle \hat{S}^2 \rangle$  values for the dissociated atomic fragments. However, the values of two-center contributions  $\langle \hat{S}^2 \rangle_{AB} = 0.250$  found for both molecules allow a right  $\langle \hat{S}^2 \rangle = 2$  for the global triplet states described.

In Table 4, we report numerical values of local spins in order to check the dependence on the electronic correlation of this methodology. In that Table, we describe results arising from both treatments (spin-free and with spin) for the singlet systems HF, NH<sub>3</sub> and C<sub>2</sub>H<sub>6</sub>, using the configuration interaction (CI) technique at several levels; single and double excitations (CISD); single, double, and triple excitations (CISDT) and single, double, triple, and quadruple excitations (CISDTQ). In the case of the ethane molecule, we have kept frozen (without excitation) 26 of the 30 orbitals forming the 6-31G basis set (the 7 lowest occupied molecular orbitals and the 19 highest unoccupied molecular orbitals) in the procedure denoted frozen I, and 18 orbitals (the 3 lowest occupied molecular orbitals and the 15 highest

**Table 2. Local Spins of One- And Two-Centers ( $\langle \hat{S}^2 \rangle_A$  and  $\langle \hat{S}^2 \rangle_{AB}$ ) Arising from the Spin-Free Treatment (eqs 13 and 14) for Doublet Systems at Experimental Equilibrium Distances**

system	atom/ bond	spin-free		basis set/ method	
		$\langle \hat{S}^2 \rangle_A$	$\langle \hat{S}^2 \rangle_{AB}$		
OH	O	0.874		6-31G/CISD	
	H	0.063			
	OH		-0.094		
NO	N	0.670		6-31G/CISD	
	O	0.319			
	NO		-0.119		
CN	C	1.141		6-31G/CISD	
	N	0.452			
	CN		-0.421		
NH <sub>2</sub>	N	1.060		6-31G/CISD	
	H	0.069			
	NH		-0.117		
	HH		0.009		
CH	C	0.912		6-31G/CISD	
	H	0.081			
	CH		-0.122		
CH <sub>3</sub>	C	1.368		6-31G/CISD	
	H	0.069			
	CH		-0.151		
	HH		0.013		
C <sup>(1)</sup> H <sub>3</sub> -C <sup>(2)</sup> H <sub>2</sub>	C <sup>(1)</sup>	0.368		6-31G/CISD	
	C <sup>(2)</sup>	1.226			
	H <sub>(CH<sub>3</sub>)</sub>	0.076			
	H <sub>(CH<sub>2</sub>)</sub>	0.063			
	CC		-0.194		
	CH <sub>(CH<sub>3</sub>)</sub>		-0.097		
	CH <sub>(CH<sub>2</sub>)</sub>		-0.128		
	HH <sub>(CH<sub>3</sub>)</sub>		0.013		
	HH <sub>(CH<sub>2</sub>)</sub>		0.010		
	C <sup>(1)</sup> H <sub>2</sub> =C <sup>(2)</sup> H	C <sup>(1)</sup>	0.503		
C <sup>(2)</sup>	1.251				
H <sub>(CH<sub>2</sub>)</sub>	0.068				
H <sub>(CH)</sub>	0.067				
CC		-0.469			
CH <sub>(CH<sub>2</sub>)</sub>		-0.104			
CH <sub>(CH)</sub>		-0.100			
HH		0.011			
C <sup>(1)</sup> H≡C <sup>(2)</sup>	C <sup>(1)</sup>	0.639		6-31G/CISD	
C <sup>(2)</sup>	1.425				
H	0.054				
CC		-0.646			
CH		-0.090			
C <sup>(1)</sup> H <sub>2</sub> =C <sup>(2)</sup> H-C <sup>(3)</sup> H <sub>2</sub>	C <sup>(1)</sup>	0.753			6-31G/CISD
C <sup>(2)</sup>	0.392				
C <sup>(3)</sup>	0.753				
H <sub>(CH<sub>2</sub>)</sub>	0.051				
H <sub>(CH)</sub>	0.051				
C <sup>(1)</sup> C <sup>(2)</sup>		-0.296			

**Table 2. Continued**

system	atom/ bond	spin-free		basis set/ method
		$\langle \hat{S}^2 \rangle_A$	$\langle \hat{S}^2 \rangle_{AB}$	
	C <sup>(2)</sup> C <sup>(3)</sup>		-0.296	
	C <sup>(1)</sup> C <sup>(3)</sup>		0.247	
	CH <sub>(CH<sub>2</sub>)</sub>		-0.094	
	CH <sub>(CH)</sub>		-0.083	
	HH		0.008	

**Table 3. Local Spins of One- And Two-Centers ( $\langle \hat{S}^2 \rangle_A$  and  $\langle \hat{S}^2 \rangle_{AB}$ ) Arising from the Spin-Free Treatment (eqs 13 and 14) for Triplet Systems at Experimental Equilibrium Distances (<sup>a</sup> near Dissociation Limits)**

system	atom/ bond	spin-free		basis set/ method
		$\langle \hat{S}^2 \rangle_A$	$\langle \hat{S}^2 \rangle_{AB}$	
H <sub>2</sub> <sup>a</sup>	H	0.750		6-31G/FCI
	HH		0.250	
Li <sub>2</sub> <sup>a</sup>	Li	0.750		STO-3G/FCI
	LiLi		0.250	
O <sub>2</sub>	O	0.760		6-31G/CISD
	OO		0.240	
NF	N	1.879		6-31G/CISD
	F	0.141		
NH	NF		-0.010	6-31G/CISD
	N	2.277		
	H	0.082		
C <sub>2</sub>	NH		-0.180	6-31G/CISD
	C	0.761		
	CC		0.239	
B <sub>2</sub> H <sub>2</sub>	B	0.979		6-31G/CISD
	H	0.082		
	BB		0.245	
CH <sub>2</sub>	BH		-0.123	6-31G/CISD
	B...H		-0.030	
	HH		-0.001	
	C	2.604		
	H	0.083		
C <sub>2</sub> H <sub>4</sub> (triplet)	CH		-0.196	6-31G/CISD
	HH		0.007	
	C	1.352		
	H	0.064		
	CC		0.072	
C <sup>(1)</sup> H≡C <sup>(2)</sup>	CH		-0.146	6-31G/CISD
	C...H		0.005	
	HH		0.012	
	H...H		-0.002	
	H...H		-0.003	

unoccupied molecular orbitals) in the procedure denominated frozen II. As has been pointed out above, in absence of correlation that is, for RHF wave functions, formulas 5, 13 and 14 predict zero values for local spins of one- and two-center which is



**Table 4. Local Spins of One- And Two-Centers ( $\langle \hat{S}^2 \rangle_A$  and  $\langle \hat{S}^2 \rangle_{AB}$ ) Arising from the Treatments Spin-Free (eqs 13 and 14) and with Spin (eq 5 in the Atomic Basis Set) at Several Correlation Levels, With the 6-31G Basis Sets, At Experimental Equilibrium Distances**

System	correlation level	atom/ bond	spin-free		with spin	
			$\langle \hat{S}^2 \rangle_A$	$\langle \hat{S}^2 \rangle_{AB}$	$\langle \hat{S}^2 \rangle_A$	$\langle \hat{S}^2 \rangle_{AB}$
HF	CISD	F	0.050		0.059	
		H	0.050		0.059	
		HF		-0.050		-0.059
	CISDT	F	0.052		0.060	
		H	0.052		0.060	
		HF		-0.052		-0.060
	CISDTQ	F	0.056		0.064	
		H	0.056		0.064	
		HF		-0.056		-0.064
NH <sub>3</sub>	CISD	N	0.226		0.257	
		H	0.059		0.068	
		NH		-0.075		-0.086
	CISDT	HH		0.008		0.009
		N	0.234		0.264	
		H	0.061		0.071	
	CISD	NH		-0.078		-0.088
		HH		0.008		0.009
		CC		-0.130		-0.136
C <sub>2</sub> H <sub>6</sub> (staggered)	CISD (frozen I)	C	<0.001		<0.001	
		CC	<0.001		<0.001	
	CISD (frozen II)	C	0.145		0.174	
		CC	<0.001		0.007	
	CISD	C	0.361		0.400	
		CC		-0.130		-0.136

a suitable chemical requirement.<sup>11,13</sup> However, as can be observed in Table 4, the presence of correlation increases the absolute values of local spins in both procedures, although this effect turns out to be slightly less marked in the results arising from the spin-free treatment.

## 5. CONCLUSION

In this work, we have described a simple and direct partitioning of the expectation value  $\langle \hat{S}^2 \rangle$  corresponding to an  $N$ -electron system into one- and two-center terms, according to a Mulliken scheme. Our treatment, which utilizes spin-free quantities, can be applied to states of any spin symmetry  $S$  and is valid for both independent and correlated particle models of wave functions. This procedure constitutes an improvement on the previously reported treatments since it is independent of the  $S_z$  substate, fulfilling the physical requirement of uniqueness for the components of the spin multiplet. Another achievement of our approach, from a computational point of view, is that it avoids the use of the spin blocks of the second-order reduced density matrix, which are not usually available in most standard codes. The results arising from several singlet state systems show lower local spin values compared with those from other reported methods (with spin) at the used correlation levels, although these differences are not too large. We have applied our treatment to selected closed- and open-shell systems and the obtained results are chemically meaningful in all studied cases. They show a

correct behavior in limit situations; adequate atomic spin values at the dissociation limits, in agreement with those of the respective free atoms, and zero values for all one-center and two-center local spins for closed-shell RHF wave functions. Our results and those of other treatments show dependence on the electronic correlation level, although that dependence is slightly lower in our proposal than in the  $S_z$ -dependent methods.

## APPENDIX A

**Table 5. Local Spins of One- And Two-Centers ( $\langle \hat{S}^2 \rangle_A$  and  $\langle \hat{S}^2 \rangle_{AB}$ ) Arising from eq 5 (in the Atomic Basis Sets) for the Lowest Triplet HeH<sup>+</sup> System**

system	atom/bond	$\langle \hat{S}^2 \rangle_A$	$\langle \hat{S}^2 \rangle_{AB}$	basis set/method
HeH <sup>+</sup> <sup>a</sup> ( $S_z = 0$ )	He	0.700		aug-cc-pVDZ/FCI
	H	0.870		
	HeH		0.215	
HeH <sup>+</sup> <sup>a</sup> ( $ S_z  = 1$ )	He	0.815		aug-cc-pVDZ/FCI
	H	0.985		
	HeH		0.100	
HeH <sup>+</sup> <sup>b</sup> ( $S_z = 0$ )	He	0.828		aug-cc-pVDZ/FCI
	H	0.696		
	HeH		0.238	
HeH <sup>+</sup> <sup>b</sup> ( $ S_z  = 1$ )	He	0.861		aug-cc-pVDZ/FCI
	H	0.729		
	HeH		0.205	
HeH <sup>+</sup> <sup>c</sup> ( $S_z = 0$ )	He	0.763		aug-cc-pVDZ/FCI
	H	0.734		
	HeH		0.251	
HeH <sup>+</sup> <sup>c</sup> ( $ S_z  = 1$ )	He	0.758		aug-cc-pVDZ/FCI
	H	0.729		
	HeH		0.256	

<sup>a</sup>R(HeH) = 1.117 au. <sup>b</sup>R(HeH) = 2.235 au. <sup>c</sup>R(HeH) = 4.47 au.

## AUTHOR INFORMATION

### Corresponding Author

\*E-mail: qfplapel@lg.ehu.es.

## ACKNOWLEDGMENT

This report has been financially supported by Projects X017 (Universidad de Buenos Aires), PIP No. 11220090100061 (Consejo Nacional de Investigaciones Científicas y Técnicas, República Argentina), Grant No. CTQ2009-07459/BQU (the Spanish Ministry of Science and Innovation), and Grant No. GIU09/43 (Universidad del País Vasco). We thank the Universidad del País Vasco for allocation of computational resources.

## REFERENCES

- Clark, A. E.; Davidson, E. R. *J. Chem. Phys.* **2001**, *115*, 7382.
- Davidson, E. R.; Clark, A. E. *Mol. Phys.* **2002**, *100*, 373.
- Clark, A. E.; Davidson, E. R. *J. Phys. Chem. A* **2002**, *106*, 6890.
- Herrmann, C.; Reiher, M.; Hess, B. A. *J. Chem. Phys.* **2005**, *122*, 034102.
- Reiher, M. *Faraday Discuss.* **2007**, *135*, 97.
- Davidson, E. R.; Clark, A. E. *Phys. Chem. Chem. Phys.* **2007**, *9*, 1881.
- Luzanov, A. V.; Prezhdo, O. V. *Mol. Phys.* **2007**, *105*, 2879.

- (8) Mayer, I. *Chem. Phys. Lett.* **2007**, *440*, 357.
- (9) Podewitz, M.; Herrmann, C.; Malassa, A.; Westerhausen, M.; Reiher, M. *Chem. Phys. Lett.* **2008**, *451*, 301.
- (10) Alcoba, D. R.; Lain, L.; Torre, A.; Bochicchio, R. C. *Chem. Phys. Lett.* **2009**, *470*, 136.
- (11) Mayer, I. *Chem. Phys. Lett.* **2009**, *478*, 323.
- (12) Torre, A.; Alcoba, D. R.; Lain, L.; Bochicchio, R. C. *J. Phys. Chem. A* **2010**, *114*, 2344.
- (13) Mayer, I.; Matito, E. *Phys. Chem. Chem. Phys.* **2010**, *12*, 11308.
- (14) Alcoba, D. R.; Torre, A.; Lain, L.; Bochicchio, R. C. *Chem. Phys. Lett.* **2011**, *504*, 236.
- (15) Takatsuka, K.; Fueno, T.; Yamaguchi, K. *Theor. Chim. Acta* **1978**, *48*, 175.
- (16) Staroverov, V. N.; Davidson, E. R. *Chem. Phys. Lett.* **2000**, *330*, 161.
- (17) Lain, L.; Torre, A.; Bochicchio, R. C.; Ponec, R. *Chem. Phys. Lett.* **2001**, *346*, 283.
- (18) Kutzelnigg, W.; Mukherjee, D. J. *Chem. Phys.* **2002**, *116*, 4787.
- (19) Lain, L.; Torre, A.; Bochicchio, R. *J. Chem. Phys.* **2002**, *117*, 5497.
- (20) Kutzelnigg, W. *Z. Naturforsch.* **1963**, *18a*, 1058.
- (21) Torre, A.; Lain, L. *THEOCHEM* **1998**, *426*, 25.
- (22) Kutzelnigg, W.; Mukherjee, D. J. *Chem. Phys.* **1999**, *110*, 2800.
- (23) Lain, L.; Torre, A.; Alcoba, D. R.; Bochicchio, R. C. *Theor. Chem. Acc.* **2011**, *128*, 405.
- (24) Lain, L.; Torre, A.; Alcoba, D. R.; Bochicchio, R. C. *Chem. Phys. Lett.* **2009**, *476*, 101.
- (25) Alcoba, D. R.; Bochicchio, R. C.; Lain, L.; Torre, A. *J. Chem. Phys.* **2010**, *133*, 144104.
- (26) Torre, A.; Lain, L.; Bochicchio, R. *J. Phys. Chem. A* **2003**, *107*, 127.
- (27) Bochicchio, R.; Lain, L.; Torre, A. *J. Comput. Chem.* **2003**, *24*, 1902.
- (28) Lain, L.; Torre, A.; Bochicchio, R. *J. Phys. Chem. A* **2004**, *108*, 4132.
- (29) Torre, A.; Alcoba, D. R.; Lain, L.; Bochicchio, R. C. *J. Phys. Chem. A* **2005**, *109*, 6587.
- (30) Lobayan, R. M.; Bochicchio, R. C.; Lain, L.; Torre, A. *J. Chem. Phys.* **2005**, *123*, 144116.
- (31) Lobayan, R. M.; Bochicchio, R. C.; Lain, L.; Torre, A. *J. Phys. Chem. A* **2007**, *111*, 3166.
- (32) Lobayan, R. M.; Alcoba, D. R.; Bochicchio, R. C.; Torre, A.; Lain, L. *J. Phys. Chem. A* **2010**, *114*, 1200.
- (33) Head-Gordon, M. *Chem. Phys. Lett.* **2003**, *372*, 508.
- (34) Alcoba, D. R.; Bochicchio, R. C.; Lain, L.; Torre, A. *Chem. Phys. Lett.* **2006**, *429*, 286.
- (35) Karafiloglou, P. *J. Chem. Phys.* **2009**, *130*, 164103.
- (36) Crawford, T. D.; Sherrill, C. D.; Valeev, E. F.; Fermann, J. T.; King, R. A.; Leininger, M. L.; Brown, S. T.; Janssen, C. L.; Seidl, E. T.; Kenny, J. P.; Allen, W. D. *J. Comput. Chem.* **2007**, *28*, 1610.
- (37) Shamasundar, K. R. *J. Chem. Phys.* **2009**, *131*, 174109.
- (38) Lide, D. R. (Ed.), *Handbook of Chemistry and Physics*, 89th ed.; CRC: Boca Raton FL, 2008.
- (39) Johnson III R. D. (Ed.) *NIST Computational Chemistry Comparison and Benchmark Database, NIST Standard Reference Database No. 101*; National Institute of Standard and Technology, 2006. <http://srdata.nist.gov/cccbdb>.
- (40) Treboux, G.; Barthelat, J. C. *J. Am. Chem. Soc.* **1993**, *115*, 4870.
- (41) Foresman, J. B.; Frisch, A. *Exploring Chemistry with Electronic Structure Methods*; Gaussian Inc.: Pittsburgh, 1996; p 214.
- (42) Jemmis, E. D.; Pathak, B.; King, R. B.; Schaefer, H. F., III *Chem. Commun.* **2006**, *20*, 2164, and references therein.

# Comprehensive Benchmarking of a Density-Dependent Dispersion Correction

Stephan N. Steinmann and Clemence Corminboeuf\*

Laboratory for Computational Molecular Design, Institut des Sciences et Ingénierie Chimiques, Ecole Polytechnique Fédérale de Lausanne, CH-1015 Lausanne, Switzerland

**S** Supporting Information

**ABSTRACT:** Standard density functional approximations cannot accurately describe interactions between nonoverlapping densities. A simple remedy consists in correcting for the missing interactions *a posteriori*, adding an attractive energy term summed over all atom pairs. The density-dependent energy correction, dDsC, presented herein, is constructed from dispersion coefficients computed on the basis of a generalized gradient approximation to Becke and Johnson's exchange-hole dipole moment formalism. dDsC also relies on an extended Tang and Toennies damping function accounting for charge-overlap effects. The comprehensive benchmarking on 341 diverse reaction energies divided into 18 illustrative test sets validates the robust performance and general accuracy of dDsC for describing various intra- and intermolecular interactions. With a total MAD of 1.3 kcal mol<sup>-1</sup>, B97-dDsC slightly improves the results of M06-2X and B2PLYP-D3 (MAD = 1.4 kcal mol<sup>-1</sup> for both) at a lower computational cost. The density dependence of both the dispersion coefficients and the damping function makes the approach especially valuable for modeling redox reactions and charged species in general.

## INTRODUCTION

Many chemical phenomena are dominated by weak interactions, as exemplified by the highly ordered structures of biomolecules (stacking of DNA,<sup>1</sup> protein folding<sup>2</sup>) and supramolecular assemblies,<sup>3</sup> crystals arrangements of organic<sup>4</sup> and inorganic materials,<sup>5</sup> or catalysis intermediates (see, e.g., ref 6). Because of the incomparable balance of accuracy and computational cost, Kohn–Sham density functional theory<sup>7</sup> has emerged as the most widely applied methodology for investigating electronic structures and geometries of extended molecular systems. Despite this success, standard semilocal approximations do not properly describe attractive dispersion interactions that decay with  $R^{-6}$  at large intermolecular distances.<sup>8–11</sup> Even at the medium range, most semilocal density functionals fail to give an accurate description of weak interactions such as those dominating alkane isomerization energies and Pople's isodesmic bond separation equations (BSEs).<sup>12–17</sup>

Near the energy minimum, dispersion-corrected atom-centered potentials (DCAPs)<sup>18–22</sup> or carefully fitted density functionals<sup>23–28</sup> (M06-2X<sup>27</sup> is certainly the most successful functional originating from this approach) give satisfactory results. Nevertheless, both approaches intrinsically lack the ability to recover the correct long-range  $\sim R^{-6}$  attractive form. The simplest conceptual remedy,<sup>29–33</sup> first popularized by Grimme (motivated by HF-D)<sup>34–38</sup> under the DFT-D acronym,<sup>33,39,40</sup> is to correct for the missing interaction *a posteriori* by adding an attractive energy term summed over all of the atom pairs in the system. The quest for the optimal parametrization is, however, still an active field of research.<sup>40–57</sup> Recent DFT-D (e.g., D2<sup>39</sup> and D3<sup>40</sup>) gives an accurate description for intermolecular interactions, but

the proper treatment of weak intramolecular interactions is trickier.<sup>14,40,58–60</sup> Over the past three years, our group has pioneered the design of corrections which give a balanced description of both inter- and intramolecular weak interactions.<sup>43,50,57,61,62</sup> Our most recent scheme combines dispersion coefficients ( $C$ ) computed on the basis of an approximation to Becke and Johnson's<sup>63–69</sup> exchange-hole-dipole moment (XDM) formalism depending on the reduced density gradient ( $s$ )<sup>70</sup> and a genuine density dependent damping factor.<sup>57</sup> The resulting density dependent correction, called dDsC, promises substantial advantages over standard DFT computations for a broad range of applications. Following a careful validation of the dDsC scheme, we here introduce a few improvements to our original density dependent damping factor<sup>57,70</sup> and provide a comprehensive benchmarking of the density-dependent dispersion correction scheme. dDsC is tested on 18 diverse test sets featuring both intra- and intermolecular weak interaction energies together with a series of illustrative density functionals, i.e., BP86,<sup>71–73</sup> BLYP,<sup>71,74</sup> B3LYP,<sup>71,74–76</sup> PBE,<sup>77</sup> B97<sup>78</sup> and the long-range corrected LC- $\omega$ PBELYP.<sup>74,79–81</sup> Results for other schemes designed to better describe weak interactions are discussed as well: the local response dispersion (LRD) correction combined with LC-BOP,<sup>53,54</sup> two fully nonlocal density functionals, VV10<sup>82</sup> and vdW-DF-10,<sup>83</sup> the double hybrid functional B2PLYP-D3<sup>40,84</sup> and M06-2X.<sup>27</sup> The benchmark is completed by a short assessment of the dDsC schemes on geometries.

**Received:** August 26, 2011

**Published:** October 05, 2011

## THEORY

The basic form of our correction is the Tang and Toennies (TT) damping function<sup>85</sup>

$$E_{\text{disp}} = - \sum_{i=2}^{N_{\text{at}}} \sum_{j=1}^{i-1} \sum_{n=3}^{n=5} f_{2n}(bR_{ij}) \frac{C_{2n}^{ij}}{R_{ij}^{2n}} \quad (1)$$

where  $N_{\text{at}}$  is the number of atoms in the system and  $b$  is the TT-damping factor (*vide infra*). The correction is called dDsC if only the first term is retained in the multipole expansion ( $n = 3$ , corresponding to  $C_6$ ), and dDsC10 otherwise (up to  $n = 5$ , i.e., up to  $C_{10}$ ).  $f_{2n}(bR_{ij})$  are the “universal damping functions”<sup>85</sup> that are specific to each dispersion coefficient and that serve to attenuate the correction at short internuclear distances (to account for overlapping densities).

$$f_{2n}(x) = 1 - \exp(-x) \sum_{k=0}^{2n} \frac{x^k}{k!} \quad (2)$$

This section describes the determination of the damping factor  $b$  in eq 1. The dispersion coefficients themselves are obtained as described previously<sup>70</sup> and rely on a classical Hirshfeld dominant partitioning of the electron density among the atomic centers.

Classical Hirshfeld weightings are defined as<sup>86</sup>

$$w_i(r) = \frac{\rho_i^{\text{at}}(r)}{\sum_n \rho_n^{\text{at}}(r)} \quad (3)$$

where  $\rho_i^{\text{at}}$  is the sphericalized free (neutral) atomic density of atom  $i$ , weighted by the superposition of all  $\rho_i^{\text{at}}$  with all atoms  $n$  positioned as in the real molecule. The classical Hirshfeld dominant partitioning  $w_i^{\text{D}}$  is obtained by assigning each point exclusively to the atom which has the highest weight at that particular grid point. Such a partitioning is more appealing than the classical Hirshfeld populations, as it avoids overlapping atomic regions that conflict with the multipole expansion that is at the origin of the atom-pair-wise London dispersion correction.<sup>87</sup>

A key component of dDsC is the damping factor  $b$ . We showed previously<sup>50,57</sup> that the performance of the TT-damping function is improved by the introduction of a second damping function, which prevents the corrections at regions of strong density overlap (i.e., covalent distances) that are better described by density functionals.<sup>61</sup> Akin to our previous work,<sup>57</sup>  $b_{ij,\text{asym}}$ , the asymptotic value of  $b$ , accounts for the short-range effect through a multiplicative function

$$b(x) = F(x) b_{ij,\text{asym}} \quad (4)$$

$x$  and  $F(x)$  are, respectively, the damping argument and function for  $b_{ij,\text{asym}}$ , the TT-damping factor associated with two separated atoms.  $b_{ij,\text{asym}}$  is computed according to the combination rule<sup>88,89</sup>

$$b_{ij,\text{asym}} = 2 \frac{b_{ii,\text{asym}} b_{jj,\text{asym}}}{b_{ii,\text{asym}} + b_{jj,\text{asym}}} \quad (5)$$

$b_{ii,\text{asym}}$  is generally estimated from the square root of (atomic) ionization energies.<sup>90–94</sup> However, the ionization energy does not correlate well with the size of an atom that is a determinant characteristic for the damping of a dispersion term.<sup>31,39,49,95</sup> We instead propose to compute  $b_{ii,\text{asym}}$  on the basis of effective atomic polarizabilities. Note that polarizabilities as a measure of the “size” are extensively used in the closely related context of

Thole’s interacting dipole moments.<sup>96</sup> After introduction of the parameter  $b_0$ , which dictates the strength of the correction in the medium range, one obtains

$$b_{ii,\text{asym}} = b_0 \sqrt[3]{\frac{1}{\alpha_i}} = b_0 \sqrt[3]{\frac{1}{\alpha_{i,\text{free}}}} \sqrt[3]{\frac{V_{i,\text{free}}}{V_{i,\text{AIM}}}} \quad (6)$$

In the above definition,  $b_0$  includes the conversion factor from  $\text{\AA}^3$  to atomic units for  $\alpha_i$ .

The effective atom in molecule polarizabilities are estimated from scaled free atomic polarizabilities<sup>97,98</sup>

$$\begin{aligned} \alpha_i &= \frac{\langle r^3 \rangle_i}{\langle r^3 \rangle_{i,\text{free}}} \alpha_{i,\text{free}} = \frac{\int r^3 w_i^{\text{D}}(r) \rho(r) d^3r}{\int r^3 \rho_{i,\text{free}}(r) d^3r} \alpha_{i,\text{free}} \\ &= \frac{V_{i,\text{AIM}}}{V_{i,\text{free}}} \alpha_{i,\text{free}} \end{aligned} \quad (7)$$

A density cutoff of 0.002 au is applied to improve the consistency of atomic volumes between atoms at the surface and in the interior of a molecule.<sup>70,99</sup>

The  $b_{ii,\text{asym}}$  dependency on atomic polarizabilities (instead of atomic ionization energies) mostly benefits the treatment of highly polarizable atoms as shown later (e.g., neutral alkali-metal cluster like  $K_8$  of the ALK6 test set). A similar relationship could also be an advantage in force fields specifically designed to predict crystal structures. In such force fields, atomic polarizabilities have already been introduced, but  $b_{ii,\text{asym}}$  is usually determined from the molecular ionization energy with no dependency on the specific atom pair.<sup>92–94</sup> Along with the modified  $b_{ii,\text{asym}}$ , the secondary damping function is modified slightly and represented by a (steeper) exponential decay (see ref 100 for more discussion) rather than by the previously used arctan function

$$F(x) = \frac{2}{e^{a_0 x} + 1} \quad (8)$$

where the fitted parameter  $a_0$  adjusts the short-range behavior of the correction.

The last element of the correction is the damping argument  $x$

$$x = \left( 2q_{ij} + \frac{\text{abs}((Z_i - N_i^{\text{D}})(Z_j - N_j^{\text{D}}))}{r_{ij}} \right) \frac{N_i^{\text{D}} + N_j^{\text{D}}}{N_i^{\text{D}} N_j^{\text{D}}} \quad (9)$$

where  $Z_i$  and  $N_i^{\text{D}}$  are the nuclear charge and Hirshfeld dominant population of atom  $i$ , respectively.  $2q_{ij} = q_{ij} + q_{ji}$  is a covalent bond index<sup>101</sup> based on the overlap of classical Hirshfeld populations  $w_i(r)$   $q_{ij} = \int w_i(r) w_j(r) \rho(r) dr$ , and the fractional term in the parentheses is a distance-dependent ionic bond index<sup>102</sup> taken as an absolute value. Classical Hirshfeld dominant charges in the damping function resolve the inconvenience of classical Hirshfeld charges that are generally too small.<sup>57,103,104</sup> The multiplicative factor,  $(N_i^{\text{D}} + N_j^{\text{D}})/(N_i^{\text{D}} N_j^{\text{D}})$ , serves to attenuate the damping of  $b_{ii,\text{asym}}$  for heavier atoms (containing more electrons). Note that the damping function  $F(x)$  has the adequate form (i.e.,  $F(0) = 1$  and  $F(\infty) = 0$ ), given that  $x$  is large when atoms are close to each other and goes to zero with increasing distance  $r_{ij}$ . In the present form, approximated dDsC gradients are available: All derivatives of the (density dependent) parameters (the damping parameter  $b$  and the dispersion coefficients) are set

to zero, or in other words, kept fixed at their values corresponding to the energy of the geometry for which the gradient is being computed. The approximation is expected to introduce only small errors, similar to those engendered by the use of a smaller basis set for geometry optimization, followed by energy refinement with a larger basis set. Exact gradients are computationally more expensive (although simpler than those derived for the original Becke–Roussel exchange hole in ref 105) given that the contributions to the Fock matrix are needed at each SCF cycle.

To summarize, the presented dDsC correction employs electronic structure information to determine dispersion coefficients and two fitted, functional dependent, damping parameters that are the strength of the TT-damping ( $b_0$ ) and the steepness factor ( $a_0$ ).

## DETERMINATION OF THE ADJUSTABLE PARAMETERS

In line with our former work,<sup>43,50,57</sup> the chosen fitting procedure ensures a successful treatment of both weak intra- (medium-range) and inter- (long-range) molecular interactions. The two parameters ( $a_0$  and  $b_0$ ) are fitted for each functional so as to minimize the mean absolute deviation (MAD) over a representative set of 48 reactions, assessing inter- and intramolecular interactions. The detailed list of reactions in the training set is given in the Supporting Information, but in summary, 3–6 entries are taken from the following test sets (*vide infra*): BSR36, RSE43, ISO34, NBPRC, WATER27, ACONF, CYCONF, SCONF, HEAVY28, and S22. The best fit parameters are given in the Supporting Information for dDsC and dDsC10.

## TEST SETS

Eighteen test sets, corresponding to 341 reaction energies, were selected out of the 30 test sets from the GMTKN30 (database for general main group thermochemistry, kinetics, and noncovalent interactions) database<sup>106,107</sup> from where the geometries and reference values were taken. The sets are divided into three categories:

- Intramolecular interactions (5 sets, 85 reactions): ISOL22 (isomerization energies of large organic molecules),<sup>108</sup> DARC (Diels–Alder reactions energies),<sup>109</sup> BSR36 (bond separation reactions of alkanes),<sup>43,110</sup> IDISP (intramolecular dispersion interactions),<sup>14,106,111</sup> and AL2X (dimerization energies of  $\text{AlX}_3$  and  $\text{AlHX}_2$  compounds,  $X = \text{F}, \text{Cl}, \text{Br},$  and  $\text{Me}$ ).<sup>109</sup>
- Intermolecular interactions and conformational energies (7 sets, 108 reactions): S22 (binding energies of noncovalently bound dimers),<sup>112–114</sup> ADIM6 (interaction energies of  $n$ -alkane dimers),<sup>40</sup> HEAVY28 (noncovalent interaction energies between heavy element hydrides),<sup>40</sup> ACONF (relative energies of alkane conformers),<sup>115</sup> SCONF (relative energies of sugar conformers),<sup>116,117</sup> PCONF (relative energies of PHE-GLY-GLY),<sup>118</sup> and CYCONF (relative energies of cysteine conformers).<sup>119</sup>
- Mixed category of reaction energies (6 sets, 148 reactions): ALK6 (fragmentation and dissociation reactions of alkaline metal clusters and alkaline–cation benzene complexes),<sup>40</sup> BHPERI (barrier heights of pericyclic reactions),<sup>120–123</sup> RSE43 (radical stabilization energies),<sup>124</sup> NBPRC (oligomerizations and  $\text{H}_2$  fragmentations of  $\text{NH}_3/\text{BH}_3$  systems and  $\text{H}_2$  activation reactions with  $\text{PH}_3/\text{BH}_3$ ),<sup>116,125</sup> ISO34 (isomerization energies of small and medium-sized organic molecules),<sup>125,126</sup> and WATER27 (binding energies of water,  $\text{H}^+(\text{H}_2\text{O})_n$  and  $\text{OH}^-(\text{H}_2\text{O})_n$  clusters).<sup>127</sup>

**Table 1.** Mean Absolute Deviations for All Methods Tested, For All Test Sets (Overall), and the Three Individual Subcategories, i.e., Intramolecular Interactions (Intra), Intermolecular Interactions and Relative Conformational Energies (Inter+Conf), and the Mixed Test Sets (Mix)<sup>a</sup>

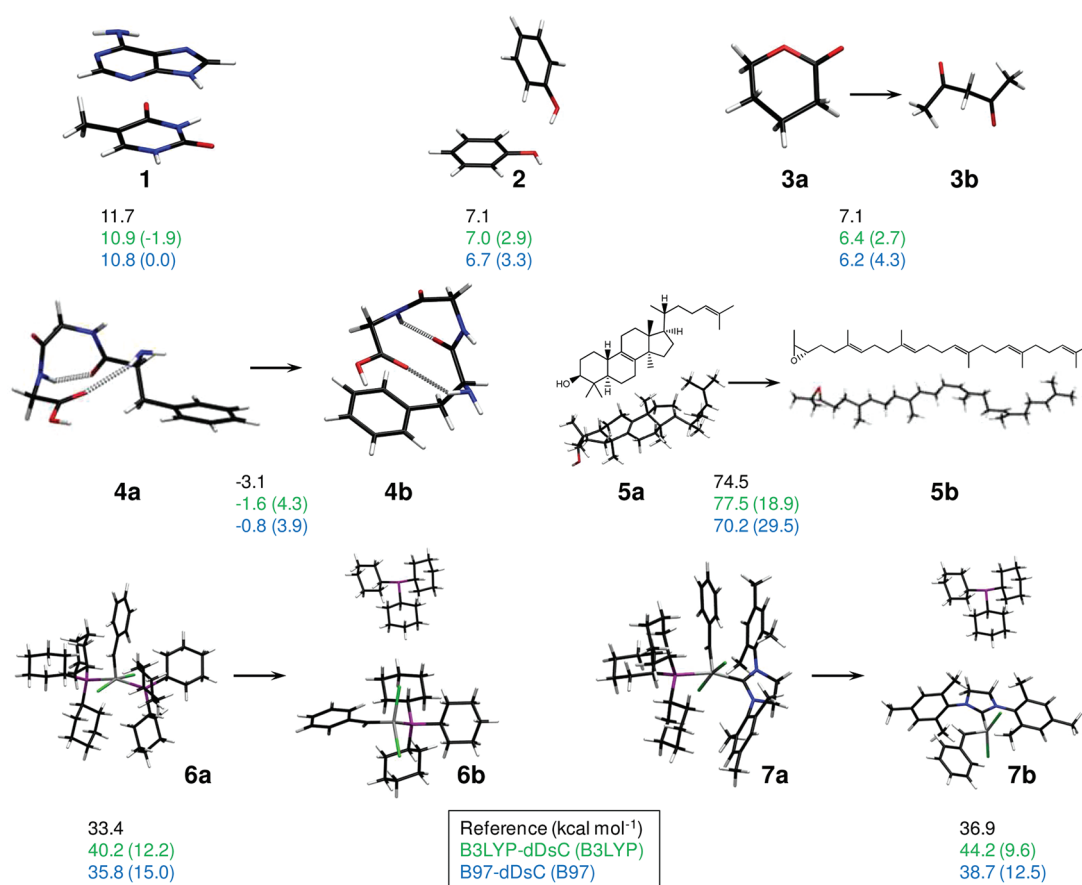
	Overall	Intra	Inter+Conf	Mix
HF	9.05	12.62	3.10	11.34
BLYP	6.85	14.38	2.53	5.67
revPBE	6.26	11.28	2.70	5.97
B3LYP	5.70	12.22	2.20	4.50
TPSSM	4.84	10.47	1.98	3.68
vdW-DF-10	4.80	11.13	0.61	4.00
BP86	4.54	9.07	2.14	3.68
B97	4.47	9.56	1.83	3.48
BHHLYP	4.40	9.10	1.77	3.63
HF-dDsC	3.74 (3.57)	5.82 (4.87)	1.25 (1.40)	4.37 (4.41)
LC- $\omega$ PBE	3.49	6.24	1.48	3.38
PBE	3.49	7.39	1.39	2.77
LC- $\omega$ PBELYP	3.35	6.14	1.26	3.26
VV10	3.34	5.50	0.43	4.22
LC-BOP	3.32	5.36	1.45	3.52
PBE0	3.11	6.55	1.44	2.34
PW6B95	3.01	6.01	0.92	2.81
B3LYP-D3	2.96	6.82	0.28	2.70
LC- $\omega$ PBEB95	2.89	4.29	0.78	3.62
LC-BOP-LRD[10,0]	2.56	3.63	0.43	3.51
LC-BOP-LRD[10,6]	2.56	3.50	0.49	3.54
BLYP-dDsC	2.45 (2.65)	3.71 (4.26)	0.62 (0.63)	3.05 (3.21)
LC- $\omega$ PBEB95-dDsC	2.39 (2.39)	4.15 (4.11)	0.66 (0.67)	2.65 (2.66)
LC- $\omega$ PBE-dDsC	2.37 (2.37)	4.82 (4.87)	0.43 (0.41)	2.38 (2.37)
PBE-dDsC	2.19 (2.22)	1.94 (1.94)	0.52 (0.57)	3.56 (3.58)
LC- $\omega$ PBELYP-dDsC	2.14 (2.04)	2.35 (2.05)	0.71 (0.59)	3.06 (3.08)
revPBE-dDsC	2.12 (1.92)	1.83 (1.89)	0.70 (0.59)	3.32 (2.90)
BP86-dDsC	2.03 (2.01)	2.44 (2.47)	0.81 (0.72)	2.68 (2.69)
TPSSM-dDsC	1.96 (1.96)	2.54 (2.61)	0.65 (0.63)	2.59 (2.56)
B3LYP-dDsC	1.67 (1.86)	2.43 (2.85)	0.48 (0.58)	2.11 (2.23)
BHHLYP-dDsC	1.66 (1.73)	1.76 (1.81)	0.48 (0.53)	2.47 (2.55)
PBE0-dDsC	1.59 (1.66)	1.98 (2.04)	0.42 (0.52)	2.22 (2.28)
M06-2X	1.41	2.94	0.40	1.26
PW6B95-dDsC	1.39 (1.39)	1.70 (1.67)	0.62 (0.66)	1.78 (1.76)
B2PLYP-D3	1.37	3.41	0.16	1.08
B97-dDsC	1.30 (1.32)	1.78 (1.82)	0.48 (0.47)	1.62 (1.65)

<sup>a</sup> Values in parentheses refer to the correction including coefficients up to  $C_{10}$  (dDsC10). All values are in  $\text{kcal mol}^{-1}$ . Results for B2PLYP-D3 and M06-2X are taken from refs 107 and 152.

Reaction energies and MADs for all methods tested are given in the Supporting Information, which also includes the corrected data with higher-order dispersion coefficients. Note that the effects of the higher-order terms strongly depend on the type of damping function. The TT-damping function applied herein “simulates” the missing higher-order dispersion terms by increasing the damping factor  $b$ .<sup>128</sup>

## COMPUTATIONAL METHODS

BLYP,<sup>71,74</sup> BP86,<sup>71,72</sup> PBE,<sup>77</sup> revPBE,<sup>129</sup> B3LYP,<sup>71,74–76</sup> and PBE0<sup>77,130</sup> computations were performed with a developmental



**Figure 1.** Set of illustrative examples of reactions poorly described by standard density functionals (e.g., B3LYP and B97) and corrected by dDsC. The reference values<sup>108,113,126,149</sup> are computed at the CCSD(T)/CBS level, except for 5, where SCS-MP3/CBS serves as the benchmark, and for 7, experimental values are used.<sup>153</sup> The DFT energies for 4–7 are computed with the def2-TZVP basis set.

version of ADF.<sup>131,132</sup> HF, BHHLYP,<sup>133</sup> Becke's hybrid B97<sup>78</sup> functional (that is to be distinguished from Grimme's GGA functional B97-D<sup>39</sup>), PW6B95,<sup>134</sup> LC- $\omega$ PBE<sup>79–81</sup> ( $\omega = 0.45$ ), LC- $\omega$ PBELYP ( $\omega = 0.45$ ), LC- $\omega$ PBEB95<sup>135</sup> ( $\omega = 0.45$ ), VV10 (rPW86<sup>136</sup>PBE<sup>77</sup>+nonlocal correction),<sup>82</sup> and vdW-DF-10 (rPW86<sup>136</sup>PW92<sup>137</sup>+nonlocal correction)<sup>83</sup> were performed in a developmental version of Q-Chem,<sup>138</sup> while LC-BOP,<sup>71,139–141</sup> LC-BOP-LRD,<sup>53,54</sup> and TPSSm<sup>142</sup> and all geometry optimizations were run with a modified version of GAMESS.<sup>143</sup> A patch for GAMESS 2010 (version 1 Oct 2010) will be available on our Web site. Due to SCF convergence problems, computations in GAMESS use the cc-pVTZ basis set<sup>144–146</sup> (augmented with diffuse functions, leading to aug-cc-pVTZ in order to minimize the BSSE for the WATER27 complexes and all but the benzene–indole complexes of the S22 test set), except for potassium and the heavier elements for which the def2-QZVP(-g) basis set was used. All Q-Chem computations were done with the def2-QZVP(-g)<sup>147</sup> basis set except for the clusters involving OH<sup>-</sup> from the WATER27 test set, for which the aug-cc-pVQZ basis set was used. In GAMESS and Q-Chem, the numerical integrations were performed on a fine 99/590 and 75/302 Euler-Maclaurin–Lebedev grid, respectively, with an integration threshold of 10<sup>-12</sup>. In ADF, the QZ4P basis set was used for all systems except for the OH<sup>-</sup>-containing WATER27 clusters, which were described by the ET-QZ3P-DIFFUSE basis set. All-electron computations in ADF for the HEAVY27 test set include the ZORA<sup>148</sup> relativistic corrections. The “dependency” and

“addDiffuseFit” keys were applied throughout and the integration accuracy set to 8. For the sake of clarity and brevity, only a selection of the tested functionals is included in the figures, but all of the statistics are collected in Table 1 and details given in the Supporting Information.

Geometries and reference values for the peptide conformational energies (4) and the cyclization reaction (5) are taken from ref 108 and refs 149 and 150, respectively. The Grubbs catalysts' (6 and 7) geometries and zero-point energy corrections are taken from ref 151.

The dDsC corrections are applied post-SCF, using atomic fragments computed on the fly with the same method and basis set as the molecular computation. All DFT-D3<sup>40</sup> and M06-2X<sup>27</sup> values are taken from the GMTKN30 Web page.<sup>107</sup>

## RESULTS AND DISCUSSION

The performance of dDsC is at first illustrated by Figure 1, which collects seven typical reactions for which a dispersion correction is essential. The first two reactions are taken from the S22 test set<sup>112</sup> and represent general  $\pi,\pi$ -stacking interactions (adenine–thymine base pair (1), which is unbound at the B3LYP level) and the phenol dimer (2) that features a combination of hydrogen-bond and other interactions often present in organic molecules. The isomerization reaction of  $\delta$ -valerolactone (3a) into 2,4-penandione (3b)<sup>126</sup> is characteristic of typical organic isomerization reactions and is also in the training set.

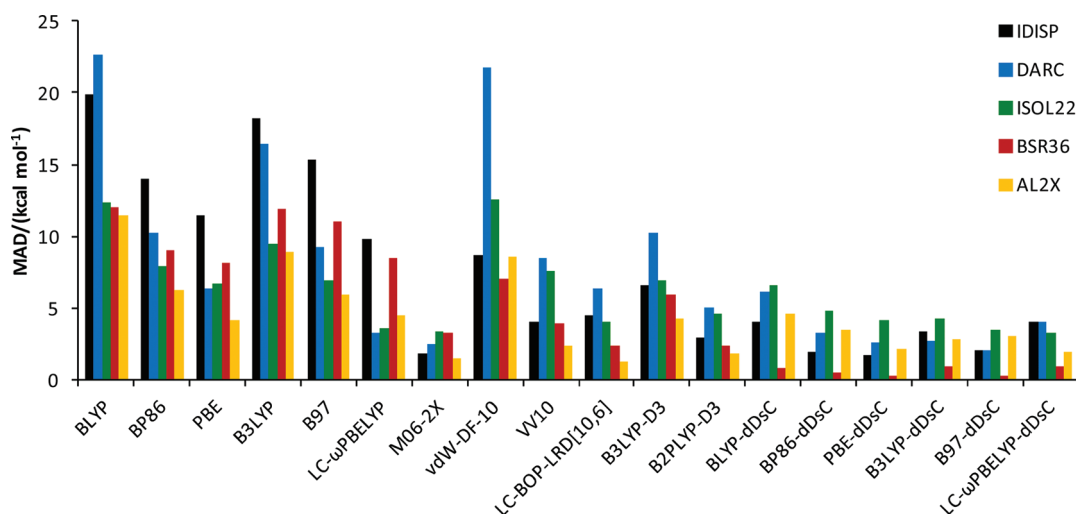


Figure 2. Mean absolute deviations for test sets dominated by intramolecular weak interactions.

The relative conformation energies of the two FGG tripeptides (4) is another example in which modeling of weak interactions is crucial to identifying the lower-lying conformer.<sup>149</sup> The cascade reaction leading to the formation of the steroid framework **5a** from the squalene precursor **5b** is a striking case with an error of almost 50 kcal mol<sup>-1</sup> at the B3LYP level.<sup>108</sup> Finally, the experimental<sup>153</sup> energy difference between the bond dissociation energies of PCy<sub>3</sub> from Grubbs' first (**6a**) and second generation (**7a**) catalysts<sup>154</sup> are qualitatively incorrect at standard density functional levels<sup>155</sup> but well reproduced when improving the treatment of medium-range correlation<sup>156</sup> or when using a dispersion correction.<sup>151</sup>

Reaction energies associated with a considerable change in molecular size and shape are challenging cases for density functional approximations. As discussed previously,<sup>61</sup> the problem may be associated with over-repulsiveness in the short range,<sup>109,157</sup> but missing weak interactions in the medium and long-ranges are the largest contributors to the errors.<sup>43,61,108,110,158</sup> By including reactions accounting for weak intramolecular interactions into the training set, our aim is to (i) obtain additional information regarding the proper form of the damping that is empirical in nature and (ii) devise a robust correction that improves both reaction energies and weak intermolecular interactions that are generally the only focus of empirical dispersion corrections.<sup>31,39,46,49,53,54</sup>

dDsC reduces the MAD of the parent functional for intramolecular interactions (see Figure 2) by a factor of 3–6, depending on the functional. The dramatically low (<1.0 kcal mol<sup>-1</sup>) MAD(BSR36) results from the highly systematic error in bond separation energies<sup>15,43,61</sup> along with the relatively large number (i.e., five) of such reactions included in the training set. The improvements for the intramolecular dispersion in hydrocarbons (IDISP) and the dimerizations of aluminum species (AL2X) as well as for the isomerizations of large organic molecules (ISOL22) highlight the high transferability of the density dependent scheme using the present parametrization. Long-range corrected functionals, such as LC- $\omega$ PBE, are among the best-uncorrected approximations (see Table 1 and the Supporting Information for more details). However, the remaining error is less systematic than that of standard functionals, and their combination with dDsC often leads to overcorrection.

LC- $\omega$ PBE-LYP-dDsC is the most accurate combination, but the variant does not present significant advantages over standard DFT-dDsC methods. The latter also clearly outperform the more sophisticated nonlocal van der Waals density functionals. The poorer performance of vdW-DF-10 as compared to VV10 is most likely related to the replacement of the local PW92 by the PBE correlation in VV10: The PBE correlation functional is known to capture intramolecular interactions involving weakly interacting densities that overlap reasonably well.<sup>61</sup> The changes in bond types of the AL2X, DARC, and ISOL22 test sets might be more accurate with the PBE than the PW92 correlation functional as well. LC-BOP-LRD[10,6] further lowers the MAD to 3.5 kcal mol<sup>-1</sup> in this category. With a MAD of 2.9 and 3.4 kcal mol<sup>-1</sup> over the five “intramolecular” test sets, M06-2X and B2PLYP-D3, respectively, improve considerably over the standard density functionals (e.g., MAD(B3LYP) = 12.2 kcal mol<sup>-1</sup>) but do not achieve the high accuracy of DFT-dDsC, where most functionals are corrected to a MAD of only about 2 kcal mol<sup>-1</sup>, with a minimum of 1.7 kcal mol<sup>-1</sup> for PW6B95-dDsC.

The improved energies for systems characterized by typical weak intermolecular interactions are collected in Figure 3. Most atom pairwise dispersion corrections and fully nonlocal van der Waals functionals are designed to improve the treatment of those interactions. Accordingly, the performance of methods such as B2PLYP-D3 is excellent, and VV10, vdW-DF-10, and LC-BOP-LRD[10,6] give relatively low errors as well. The remarkable performance of M06-2X is, on the other hand, illustrative of the success of extensive fitting. With an average MAD of 0.6 kcal mol<sup>-1</sup> (over 13 density functionals, excluding HF-dDsC), DFT-dDsC also performs well for diverse types of weak intermolecular interactions and relative conformational energies (see Table 1 and Supporting Information). The small errors obtained for the S22 test set (assessing pure dispersion to H-bonding) along with those on the heavy atom hydrides confirm the general accuracy of the density dependent dispersion scheme. Alkane dimers (ADIM6) are, however, overcorrected by dDsC. Our careful analysis suggests that ADIM6 is an exception rather than the result of an overfitting toward intramolecular interactions dominating the training set. Subtle changes in nonbonded interactions such as those dictating the relative conformational energies

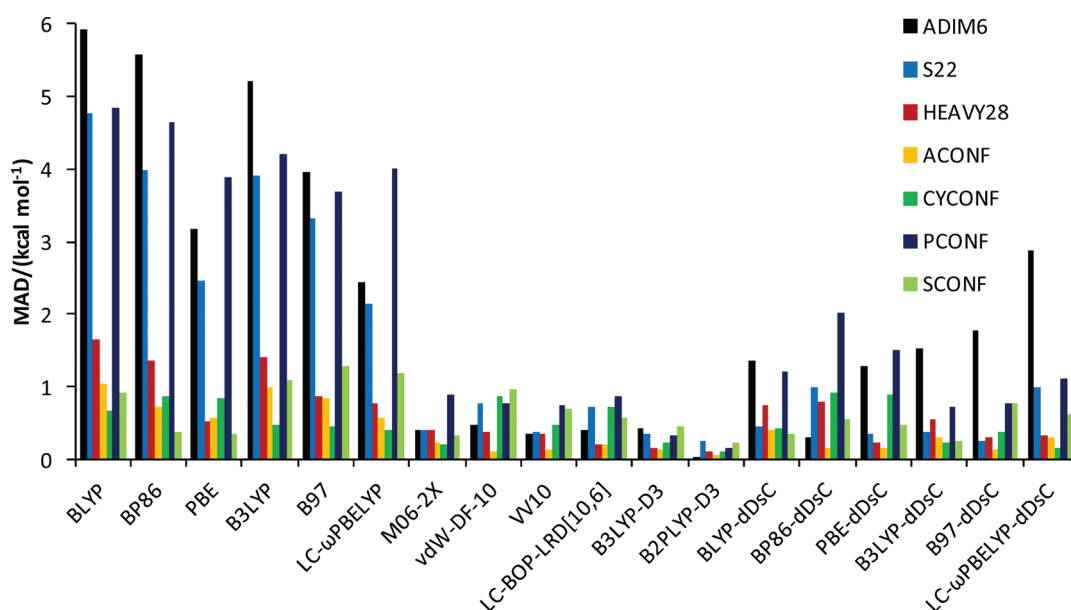


Figure 3. Mean absolute deviations for test sets featuring intermolecular weak interactions or relative conformational energies.

of alkanes (ACONF) are, for instance, well captured by dDsC, which shows that the strong correction needed for improving bond separation equations does not generally deteriorate longer-range interactions. To a much lesser extent, the D3 level also overcorrects alkane dimers, even though D3 is parametrized to perform well for these systems (see the detailed performance of D3 on the GMTKN Web site<sup>107</sup>). The peculiarity of the ADIM6 test set is further illustrated by the contrasting trend in the performance of MP2/CBS (MAD = 0.27 kcal mol<sup>-1</sup>) and SCS-MP2/CBS (MAD = 1.05 kcal mol<sup>-1</sup>), which is opposite that of the S22 test set.<sup>152</sup> The modest performance of dDsC for the Phe–Gly–Gly–peptide conformations (PCONF) is, to a large extent, influenced by the choice of “reference” conformer used in the relative energy computations. Standard functionals indeed identify the second lowest energy conformer instead of the correct conformation (at the CCSD(T) level) as the lowest energy one. The MADs are thus lowered by up to 50%, when considering the second lowest lying (0.14 kcal/mol higher according to the CCSD(T) reference values<sup>118</sup>) as the “reference compound”!

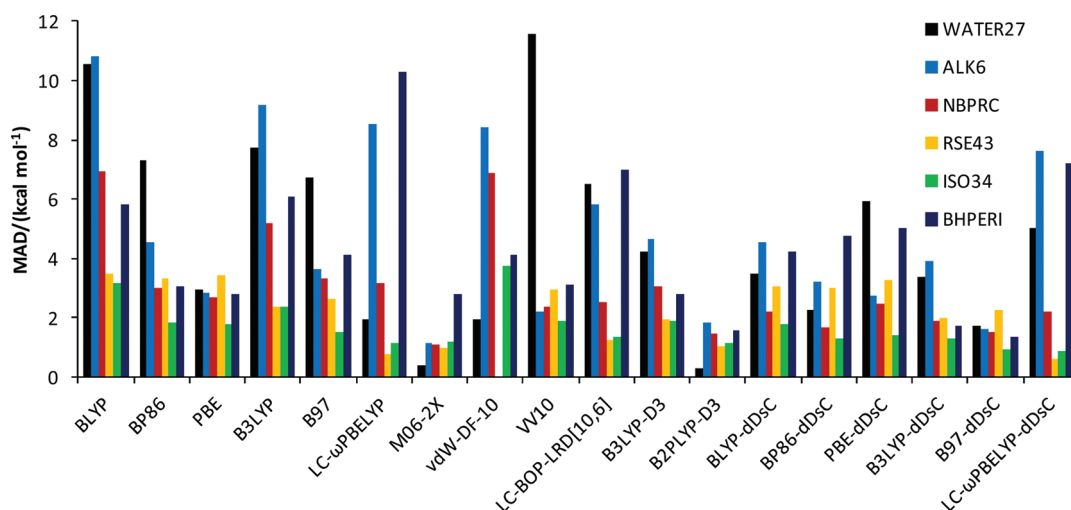
Several additional interesting features of Figure 3 can be better understood by considering the characteristics of the parent functional. For instance, the accurate treatment of the relative conformational energies of cysteine (CYCONF) relies on a balanced description between strong (e.g., OH···N) intramolecular hydrogen bonds (that dominate some of the conformers) and weaker interactions (e.g., NH···S present in other conformers). The good description of OH···N and NH···O hydrogen bonds by PBE and BP86 versus their underestimation of weak interactions bias the relative conformational energies and result in the poorer performance of PBE(-dDsC) and BP86(-dDsC) for CYCONF than for SCONF. The relative energies of sugar conformers, which are all dominated by strong hydrogen bonds, are indeed better described by these levels,<sup>117</sup> which do not benefit from the inclusion of a dispersion correction.

Figure 4 collects errors for the “mixed” category, regrouping six test sets, which are not all dominated by weak interactions but are nevertheless important for typical computational chemistry

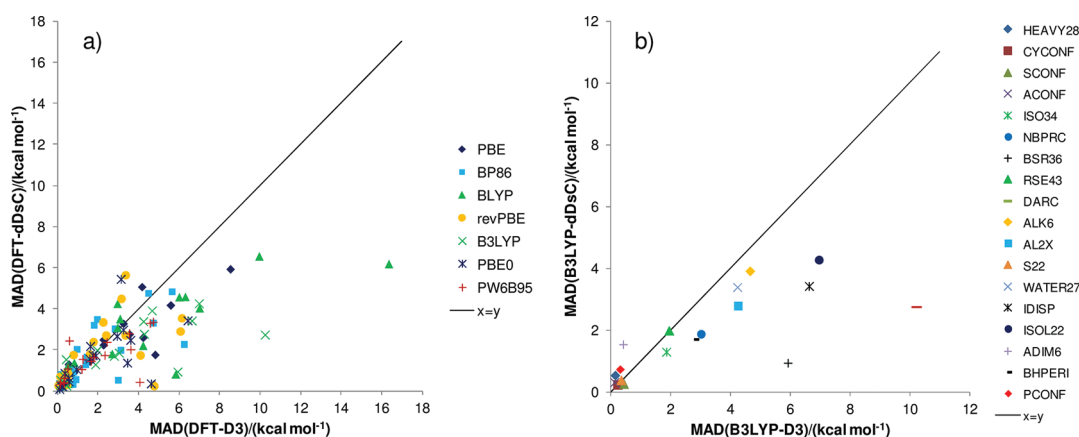
applications. The errors in radical stabilization energies (RSE43), isomerization energies of small molecules (ISO34), and the NBPRC test set, for instance, originate from subtle inaccuracies in, e.g., bond energies. The inaccurate treatment of barrier heights of pericyclic reactions (BHPERI) is generally attributed to the self-interaction error,<sup>159–163</sup> and to the delocalization error<sup>109</sup> (or the error in the repulsive wall<sup>61</sup>) that is also at the origin of the poor assessment of the related Diels–Alder reaction energies (see DARC in Figure 3). For “repulsive” functionals such as BLYP or B3LYP, the dispersion correction stabilizes the transition state and leads to a clear improvement. The barrier heights are, however, overcorrected with more attractive approximation such as PBE. The unexpected poor performance of LC- $\omega$ PBELYP (LC- $\omega$ PBE and LC- $\omega$ PBEB95 perform better in this case, with a MAD of about 6.7 kcal mol<sup>-1</sup> vs 10.3 kcal mol<sup>-1</sup>, but even BLYP (MAD = 5.8 kcal mol<sup>-1</sup>) outperforms the long-range corrected functionals) results from a strong overestimation of the barrier heights in line with that of HF (23.2 kcal mol<sup>-1</sup> and 10.6 kcal mol<sup>-1</sup> with HF-dDsC). The high error for BHPERI along with the general difficulty of systematically improving the LC- $\omega$ PBE functional group by a dispersion correction (*vide supra*) reflects the need for a better-devised long-range correction parameter  $\omega$ . A system dependence<sup>164–166</sup> could be a strategy that would, however, cause size-extensivity problems important for reaction energies. At higher computational costs, the more balanced description of range-separated local hybrids<sup>167</sup> represents another alternative. Note that M06-2X, with a MAD of 2.8 kcal mol<sup>-1</sup>, is also affected by the large amount of “exact” exchange (54%), while B97-dDsC (~19% “exact exchange”) performs best for these barrier heights (MAD = 1.3 kcal mol<sup>-1</sup>).

ALK6 played an important role in cross-validating the proposed density dependent dispersion correction: the three benzene–alkaline cation (Li<sup>+</sup>, Na<sup>+</sup>, K<sup>+</sup>) complexes are dominated by electrostatic and inductive interactions<sup>168</sup> and are thus well described by standard DFT levels. Such interactions are, however, problematic for classical dispersion correction schemes, which use dispersion coefficients and vdW radii corresponding (approximately) to the free (neutral) atoms, and not to the





**Figure 4.** Mean absolute deviations over test sets assessing various reaction energies and barrier heights for pericyclic reactions. For vdW-DF-10, the RSE43 set could not be computed since it is not defined for open-shell systems.



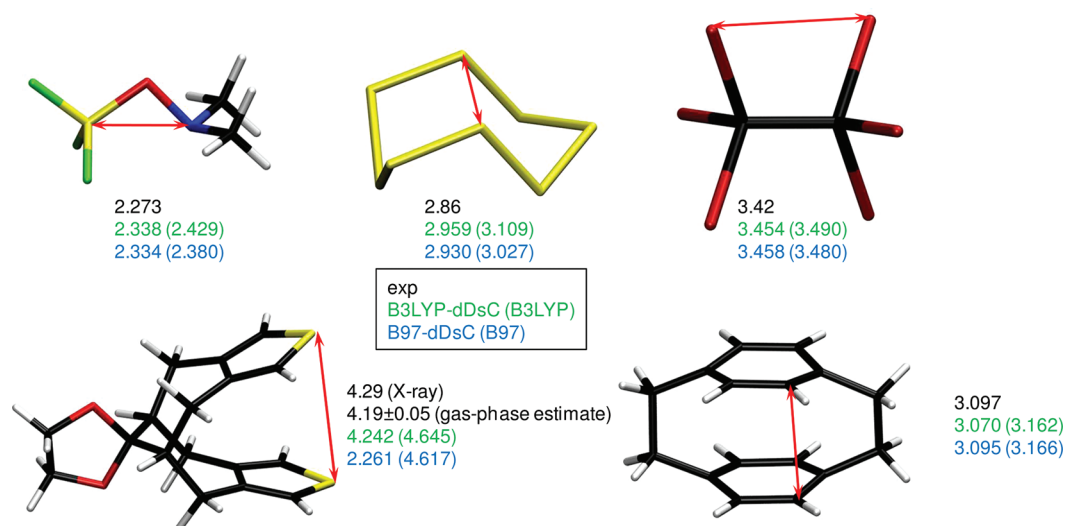
**Figure 5.** (a) Performance of DFT-dDsC versus DFT-D3 for seven functionals and the 18 selected test sets from the GMTKN30 database and (b) B3LYP-dDsC versus B3LYP-D3 with each test set as one point.

cations.<sup>40</sup> The other three systems in the test set are the decomposition of  $\text{Li}_8$ ,  $\text{Na}_8$ , and  $\text{K}_8$ , into their respective dimers. In our scheme, these clusters are characterized by relatively large dispersion coefficients and are almost as polarizable as free alkaline atoms. While most functionals underbind these clusters, our genuine damping factor,  $b_{ij,asym}$ , successfully avoids overcorrection due to its dependence on polarizability.

The overall description of test sets collected in the “mixed” category depends generally more strongly on the functional itself, than on the accuracy of the dispersion correction. For instance, the better performance of B2PLYP-D3 as compared to the dDsC corrected variants is due to B2PLYP, rather than to D3, as clearly illustrated by the comparison of B3LYP-D3 and B3LYP-dDsC (MADs of 2.7 and 2.1 kcal mol<sup>-1</sup>, respectively). Similarly, even though the LRD scheme (independently from the use of multi-center contributions, i.e., LRD[10,0] or LRD-[10,6]) improves the overall performance on the 18 test sets (3.32 vs 2.56 kcal mol<sup>-1</sup>), LC-BOP and LC-BOP-LRD[10,6], have almost the same MAD for these “mixed” test sets (3.52 and 3.54 kcal mol<sup>-1</sup>, respectively). The relatively large error of PBE-dDsC originates from the overcorrected PBE energies for WATER27 and BHPERI. A similar overcorrection is at the origin of the relatively

poor performance of VV10 (total MAD of 4.2 kcal mol<sup>-1</sup>). PBE-dDsC gives lower MAD than PBE-D3 for two reasons: (i) the ionic term in the damping function (eq 9) attenuates the correction for the strong and highly polarized hydrogen bonds of WATER27, and (ii) the polarizability-dependent damping factor prevents the energy overcorrection for the alkaline metal clusters (ALK6). Overall, B97-dDsC and PW6B95-dDsC (see Supporting Information) achieve MADs below 2.0 kcal mol<sup>-1</sup>, which illustrate that dDsC leads to improvements for this most challenging mixed category, albeit less impressive than for inter- and especially intramolecular (weak) interactions.

Figure 5 provides a detailed comparison of the MADs obtained with dDsC and the geometry-dependent D3 correction for seven functionals (Figure 5a) and the individual test sets (Figure 5b). The D3 correction performs better than dDsC in cases for which the latter has a tendency to overcorrect (e.g., ADIM6 or BHPERI with PBE) or for which the former scheme uses quasi-exact dispersion coefficients (HEAVY28). As expected, D3 also performs well for its targeted interactions (weak interactions between neutral molecules and relative conformational energies are in the training set<sup>40</sup>). On the other hand,



**Figure 6.** Geometrical structures of  $(\text{CH}_3)_2\text{NOSiF}_3$ ,<sup>172</sup>  $\text{S}_8$ ,<sup>2+,171</sup>  $\text{C}_2\text{Br}_6$  (first row),<sup>170</sup> RESVAN,<sup>174,175</sup> and [2.2]paracyclophane<sup>173</sup> (second row) with key nonbonded distances (in Ångstrom) indicated.

dDsC adjusts better to a given functional and provides a more robust performance, when considering both inter- and intramolecular interactions including challenging reaction energies (e.g., ISOL22, DARC, BSR36, IDISP, and AL2X).

The effect of dispersion corrections on thermochemistry has been thoroughly investigated. Geometries are usually less sensitive to the level of theory applied, but intramolecular nonbonded interactions are critical in certain cases. We thus compare the performance of two (un)corrected functionals, B3LYP and B97, for reproducing the geometry of five challenging molecules<sup>60,169</sup> for which experimental structures are available:  $\text{C}_2\text{Br}_6$ ,<sup>170</sup>  $\text{S}_8$ ,<sup>2+,171</sup>  $(\text{CH}_3)_2\text{NOSiF}_3$ ,<sup>172</sup> [2.2]paracyclophane,<sup>173</sup> and a bisthieno-fused molecule known under its CSD entry name RESVAN (see Figure 6).<sup>169,174–176</sup> B3LYP and B97 are overly repulsive for these intramolecular nonbonded contacts. The use of dDsC improves the geometries significantly, especially for the bisthieno-fused compound (RESVAN), mimicking stacked thiophene oligomers.

## CONCLUSIONS

The final parametrization and refinement of the density dependent dispersion correction, dDsC, introducing a simple atomic partitioning, computationally efficient dispersion coefficients, and advanced damping functions, considerably improves the performance of standard density functionals for various reaction energies and weakly interacting systems. With a MAD of  $1.3 \text{ kcal mol}^{-1}$  over the 18 test investigated sets, B97-dDsC performs slightly better than M06-2X and B2PLYP-D3 (MAD =  $1.4 \text{ kcal mol}^{-1}$  for both) but at a lower computational cost. The performance of B97-dDsC is especially impressive for the five intramolecular test sets (MAD =  $1.8 \text{ kcal mol}^{-1}$ ) for which M06-2X and B2PLYP-D3 are less satisfactory (MAD of 2.9 and  $3.4 \text{ kcal mol}^{-1}$ , respectively).

The correction is available for all elements of the periodic table. Due to its robust performance and general accuracy for various interactions, ranging from hydrocarbon reaction energies to heavy-atom hydride weak interaction energies, as well as for geometry optimization, we anticipate broad application of the dDsC scheme in diverse fields of computational

chemistry (e.g., organocatalysis, QM/MM hybrid schemes, prediction of crystal structures). The density dependence of both the dispersion coefficients and the damping function has been shown to be especially valuable for modeling oxygen reduction reactions by organic reducing agents,<sup>177</sup> the splitting of water by metallocenes,<sup>178</sup> as well as for the molecular receptors,<sup>179</sup> which all involve charged species.

## ASSOCIATED CONTENT

**S Supporting Information.** The optimal parameters  $a_0$  and  $b_0$  for all functionals tested herein and all reaction energies and MADs are listed. This material is provided free of charge via the Internet at <http://pubs.acs.org>

## AUTHOR INFORMATION

### Corresponding Author

\*E-mail: [clemence.corminboeuf@epfl.ch](mailto:clemence.corminboeuf@epfl.ch).

## ACKNOWLEDGMENT

C.C. acknowledges the Sandoz family foundation, the Swiss NSF Grant 200021\_121577/1, and EPFL for financial support. We are grateful to Q-Chem Inc. and SCM for providing the source code of Q-Chem and ADF, respectively.

## REFERENCES

- (1) Dabkowska, I.; Gonzalez, H. V.; Jurecka, P.; Hobza, P. *J. Phys. Chem. A* **2005**, *109*, 1131–1136.
- (2) Bashford, D.; Chothia, C.; Lesk, A. M. *J. Mol. Biol.* **1987**, *196*, 199–216.
- (3) Hollingsworth, M. D. *Science* **2002**, *295*, 2410–2413.
- (4) Coombes, D. S.; Price, S. L.; Willock, D. J.; Leslie, M. *J. Phys. Chem.* **1996**, *100*, 7352–7360.
- (5) Braga, D. *J. Chem. Soc., Dalton Trans.* **2000**, 3705–3713.
- (6) Knowles, R. R.; Jacobsen, E. N. *Proc. Natl. Acad. Sci. U.S.A.* **2010**, *107*, 20678–20685.
- (7) Kohn, W.; Sham, L. J. *Phys. Rev.* **1965**, *140*, A1133–A1138.
- (8) Kristyan, S.; Pulay, P. *Chem. Phys. Lett.* **1994**, *229*, 175–180.

- (9) Perez-Jorda, J. M.; Becke, A. D. *Chem. Phys. Lett.* **1995**, *233*, 134–137.
- (10) Hobza, P.; Sponer, J.; Reschel, T. J. *Comput. Chem.* **1995**, *16*, 1315–1325.
- (11) Zhang, Y.; Pan, W.; Yang, W. J. *Chem. Phys.* **1997**, *107*, 7921–7925.
- (12) Hehre, W. J.; Ditchfield, R.; Radom, L.; Pople, J. A. *J. Am. Chem. Soc.* **1970**, *92*, 4796–4801.
- (13) Pople, J. A.; Radom, L.; Hehre, W. J. *J. Am. Chem. Soc.* **1971**, *93*, 289–300.
- (14) Grimme, S. *Angew. Chem., Int. Ed.* **2006**, *45*, 4460–4464.
- (15) Wodrich, M. D.; Corminboeuf, C.; Schleyer, P. v. R. *Org. Lett.* **2006**, *8*, 3631–3634.
- (16) Wodrich, M. D.; Corminboeuf, C.; Schreiner, P. R.; Fokin, A. A.; Schleyer, P. v. R. *Org. Lett.* **2007**, *9*, 1851–1854.
- (17) Schreiner, P. R. *Angew. Chem., Int. Ed.* **2007**, *46*, 4217–4219.
- (18) von Lilienfeld, O. A.; Tavernelli, I.; Rothlisberger, U.; Sebastiani, D. *Phys. Rev. Lett.* **2004**, *93*, 153004.
- (19) von Lilienfeld, O. A.; Tavernelli, I.; Rothlisberger, U.; Sebastiani, D. *Phys. Rev. B* **2005**, *71*, 195119.
- (20) Lin, I.-C.; Coutinho-Neto, M. D.; Felsenheimer, C.; Lilienfeld, O. A. v.; Tavernelli, I.; Rothlisberger, U. *Phys. Rev. B* **2007**, *75*, 205131.
- (21) Mackie, I. D.; DiLabio, G. A. *J. Phys. Chem. A* **2008**, *112*, 10968–10976.
- (22) Nilsson Lill, S. O. *J. Phys. Chem. A* **2009**, *113*, 10321–10326.
- (23) Zhao, Y.; Lynch, B. J.; Truhlar, D. G. *J. Phys. Chem. A* **2004**, *108*, 2715–2719.
- (24) Xu, X.; Goddard, W. A. *Proc. Natl. Acad. Sci. U.S.A.* **2004**, *101*, 2673–2677.
- (25) Zhao, Y.; Truhlar, D. G. *J. Phys. Chem. A* **2004**, *108*, 6908–6918.
- (26) Zhao, Y.; Truhlar, D. G. *J. Phys. Chem. A* **2005**, *109*, 4209–4212.
- (27) Zhao, Y.; Truhlar, D. *Theor. Chem. Acc.* **2008**, *120*, 215–241.
- (28) Zhao, Y.; Truhlar, D. G. *Acc. Chem. Res.* **2008**, *41*, 157–167.
- (29) Meijer, E. J.; Sprik, M. J. *Chem. Phys.* **1996**, *105*, 8684–8689.
- (30) Wu, X.; Vargas, M. C.; Nayak, S.; Lotrich, V.; Scoles, G. *J. Chem. Phys.* **2001**, *115*, 8748–8757.
- (31) Wu, Q.; Yang, W. J. *Chem. Phys.* **2002**, *116*, 515–524.
- (32) Zimmerli, U.; Parrinello, M.; Koumoutsakos, P. *J. Chem. Phys.* **2004**, *120*, 2693–2699.
- (33) Grimme, S. *J. Comput. Chem.* **2004**, *25*, 1463–1473.
- (34) Conway, A.; Murrell, J. N. *Mol. Phys.* **1974**, *27*, 873–878.
- (35) Wagner, A. F.; Das, G.; Wahl, A. C. *J. Chem. Phys.* **1974**, *60*, 1885–1891.
- (36) Hepburn, J.; Scoles, G.; Penco, R. *Chem. Phys. Lett.* **1975**, *36*, 451–456.
- (37) Ahlrichs, R.; Penco, R.; Scoles, G. *Chem. Phys.* **1977**, *19*, 119–130.
- (38) Clementi, E.; Corongiu, G. *J. Phys. Chem. A* **2001**, *105*, 10379–10383.
- (39) Grimme, S. *J. Comput. Chem.* **2006**, *27*, 1787–1799.
- (40) Grimme, S.; Antony, J.; Ehrlich, S.; Krieg, H. *J. Chem. Phys.* **2010**, *132*, 154104.
- (41) Ducere, J.-M.; Cavallo, L. *J. Phys. Chem. B* **2007**, *111*, 13124–13134.
- (42) Olasz, A.; Vanommeslaeghe, K.; Krishtal, A.; Veszpremi, T.; Alsenoy, C. V.; Geerlings, P. *J. Chem. Phys.* **2007**, *127*, 224105.
- (43) Wodrich, M. D.; Jana, D. F.; Schleyer, P. v. R.; Corminboeuf, C. *J. Phys. Chem. A* **2008**, *112*, 11495–11500.
- (44) Murdachaew, G.; de Gironcoli, S.; Scoles, G. *J. Phys. Chem. A* **2008**, *112*, 9993–10005.
- (45) Jurecka, P.; Cerný, J.; Hobza, P.; Salahub, D. R. *J. Comput. Chem.* **2007**, *28*, 555–569.
- (46) Chai, J.-D.; Head-Gordon, M. *Phys. Chem. Chem. Phys.* **2008**, *10*, 6615–6620.
- (47) Krishtal, A.; Vanommeslaeghe, K.; Olasz, A.; Veszpremi, T.; Alsenoy, C. V.; Geerlings, P. *J. Chem. Phys.* **2009**, *130*, 174101.
- (48) Liu, Y.; Goddard, W. A. *Mater. Trans.* **2009**, *50*, 1664–1670.
- (49) Tkatchenko, A.; Scheffler, M. *Phys. Rev. Lett.* **2009**, *102*, 073005.
- (50) Steinmann, S. N.; Csonka, G.; Corminboeuf, C. *J. Chem. Theory Comput.* **2009**, *5*, 2950–2958.
- (51) Pernal, K.; Podeszwa, R.; Patkowski, K.; Szalewicz, K. *Phys. Rev. Lett.* **2009**, *103*, 4.
- (52) Podeszwa, R.; Pernal, K.; Patkowski, K.; Szalewicz, K. *J. Phys. Chem. Lett.* **2009**, *1*, 550–555.
- (53) Sato, T.; Nakai, H. *J. Chem. Phys.* **2009**, *131*, 224104.
- (54) Sato, T.; Nakai, H. *J. Chem. Phys.* **2010**, *133*, 194101.
- (55) Kannemann, F. O.; Becke, A. D. *J. Chem. Theory Comput.* **2009**, *5*, 719–727.
- (56) Kannemann, F. O.; Becke, A. D. *J. Chem. Theory Comput.* **2010**, *6*, 1081–1088.
- (57) Steinmann, S. N.; Corminboeuf, C. *J. Chem. Theory Comput.* **2010**, *6*, 1990–2001.
- (58) Grimme, S.; Diedrich, C.; Korth, M. *Angew. Chem., Int. Ed.* **2006**, *45*, 625–629.
- (59) Grimme, S.; Steinmetz, M.; Korth, M. *J. Chem. Theory Comput.* **2007**, *3*, 42–45.
- (60) Grimme, S.; Ehrlich, S.; Goerigk, L. *J. Comput. Chem.* **2011**, *32*, 1456–1465.
- (61) Steinmann, S. N.; Wodrich, M.; Corminboeuf, C. *Theor. Chem. Acc.* **2010**, *127*, 429–442.
- (62) Steinmann, S. N.; Corminboeuf, C. *Chimia* **2011**, *65*, 240–244.
- (63) Becke, A. D.; Johnson, E. R. *J. Chem. Phys.* **2005**, *122*, 154104.
- (64) Johnson, E. R.; Becke, A. D. *J. Chem. Phys.* **2005**, *123*, 024101.
- (65) Becke, A. D.; Johnson, E. R. *J. Chem. Phys.* **2005**, *123*, 154101.
- (66) Becke, A. D.; Johnson, E. R. *J. Chem. Phys.* **2006**, *124*, 014104.
- (67) Johnson, E. R.; Becke, A. D. *J. Chem. Phys.* **2006**, *124*, 174104.
- (68) Becke, A. D.; Johnson, E. R. *J. Chem. Phys.* **2007**, *127*, 154108.
- (69) Becke, A. D.; Johnson, E. R. *J. Chem. Phys.* **2007**, *127*, 124108.
- (70) Steinmann, S. N.; Corminboeuf, C. *J. Chem. Phys.* **2011**, *134*, 044117.
- (71) Becke, A. D. *Phys. Rev. A* **1988**, *38*, 3098.
- (72) Perdew, J. P.; Wang, Y. *Phys. Rev. B* **1986**, *33*, 8800–8802.
- (73) Perdew, J. P.; Yue, W. *Phys. Rev. B* **1989**, *40*, 3399.
- (74) Lee, C.; Yang, W.; Parr, R. G. *Phys. Rev. B* **1988**, *37*, 785.
- (75) Becke, A. D. *J. Chem. Phys.* **1993**, *98*, 5648–5652.
- (76) Stephens, P. J.; Devlin, F. J.; Chabalowski, C. F.; Frisch, M. J. *J. Phys. Chem.* **1994**, *98*, 11623–11627.
- (77) Perdew, J. P.; Burke, K.; Ernzerhof, M. *Phys. Rev. Lett.* **1996**, *77*, 3865–3868.
- (78) Becke, A. D. *J. Chem. Phys.* **1997**, *107*, 8554–8560.
- (79) Henderson, T. M.; Janesko, B. G.; Scuseria, G. E. *J. Chem. Phys.* **2008**, *128*, 194105.
- (80) Rohrdanz, M. A.; Martins, K. M.; Herbert, J. M. *J. Chem. Phys.* **2009**, *130*, 054112.
- (81) Weintraub, E.; Henderson, T. M.; Scuseria, G. E. *J. Chem. Theory Comput.* **2009**, *5*, 754–762.
- (82) Vydrov, O. A.; Voorhis, T. V. *J. Chem. Phys.* **2010**, *133*, 244103.
- (83) Lee, K.; Murray, E. D.; Kong, L.; Lundqvist, B. I.; Langreth, D. C. *Phys. Rev. B* **2010**, *82*, 081101.
- (84) Grimme, S. *J. Chem. Phys.* **2006**, *124*, 034108.
- (85) Tang, K. T.; Toennies, J. P. *J. Chem. Phys.* **1984**, *80*, 3726–3741.
- (86) Hirshfeld, F. L. *Theor. Chem. Acc.* **1977**, *44*, 129–138.
- (87) Angyan, J. G. *J. Chem. Phys.* **2007**, *127*, 024108.
- (88) Sheng, X. W.; Li, P.; Tang, K. T. *J. Chem. Phys.* **2009**, *130*, 174310.
- (89) Bohm, H.-J.; Ahlrichs, R. *J. Chem. Phys.* **1982**, *77*, 2028–2034.
- (90) Douketis, C.; Scoles, G.; Marchetti, S.; Zen, M.; Thakkar, A. J. *J. Chem. Phys.* **1982**, *76*, 3057–3063.
- (91) Tang, K. T.; Toennies, J. P.; Yiu, C. L. *Phys. Rev. Lett.* **1995**, *74*, 1546.
- (92) Misquitta, A. J.; Stone, A. J. *J. Chem. Theory Comput.* **2007**, *4*, 7–18.
- (93) Misquitta, A. J.; Stone, A. J.; Price, S. L. *J. Chem. Theory Comput.* **2007**, *4*, 19–32.

- (94) Misquitta, A. J.; Stone, A. J. *Mol. Phys.* **2008**, *106*, 1631–1643.
- (95) Tkatchenko, A.; Robert A. DiStasio, J.; Head-Gordon, M.; Scheffler, M. *J. Chem. Phys.* **2009**, *131*, 094106.
- (96) Thole, B. T. *Chem. Phys.* **1981**, *59*, 341–350.
- (97) Brinck, T.; Murray, J. S.; Politzer, P. *J. Chem. Phys.* **1993**, *98*, 4305–4306.
- (98) Miller, T. M. In *CRC Handbook of Chemistry and Physics*, 91th ed.; Taylor & Francis Group: London; pp 10–186.
- (99) Bader, R. F. W.; Carroll, M. T.; Cheeseman, J. R.; Chang, C. *J. Am. Chem. Soc.* **1987**, *109*, 7968–7979.
- (100) The use of a simple exponential decay was also investigated following one reviewer's suggestion. As shown in the Supporting Information, this function gives results very similar to those obtained with eq 8.
- (101) Mayer, I.; Salvador, P. *Chem. Phys. Lett.* **2004**, *383*, 368–375.
- (102) Mulliken, R. S. *J. Chem. Phys.* **1955**, *23*, 1841–1846.
- (103) Davidson, E. R.; Chakravorty, S. *Theor. Chem. Acc.* **1992**, *83*, 319–330.
- (104) Parr, R. G.; Ayers, P. W.; Nalewajski, R. F. *J. Phys. Chem. A* **2005**, *109*, 3957–3959.
- (105) Kong, J.; Gan, Z.; Proynov, E.; Freindorf, M.; Furlani, T. R. *Phys. Rev. A* **2009**, *79*, 042510.
- (106) Goerigk, L.; Grimme, S. *J. Chem. Theory Comput.* **2011**, *7*, 291–309.
- (107) GMTKN30. <http://toc.uni-muenster.de/GMTKN/GMTKN30/GMTKN30main.html> (accessed March 23, 2011).
- (108) Huenerbein, R.; Schirmer, B.; Moellmann, J.; Grimme, S. *Phys. Chem. Chem. Phys.* **2010**, *12*, 6940–6948.
- (109) Johnson, E. R.; Mori-Sanchez, P.; Cohen, A. J.; Yang, W. *J. Chem. Phys.* **2008**, *129*, 204112.
- (110) Krieg, H.; Grimme, S. *Mol. Phys.* **2010**, *108*, 2655–2666.
- (111) Schwabe, T.; Grimme, S. *Phys. Chem. Chem. Phys.* **2007**, *9*, 3397–3406.
- (112) Jurecka, P.; Sponer, J.; Cerny, J.; Hobza, P. *Phys. Chem. Chem. Phys.* **2006**, *8*, 1985–1993.
- (113) Takatani, T.; Hohenstein, E. G.; Malagoli, M.; Marshall, M. S.; Sherrill, C. D. *J. Chem. Phys.* **2010**, *132*, 144104–5.
- (114) Podeszwa, R.; Patkowski, K.; Szalewicz, K. *Phys. Chem. Chem. Phys.* **2010**, *12*, 5974–5979.
- (115) Gruzman, D.; Karton, A.; Martin, J. M. L. *J. Phys. Chem. A* **2009**, *113*, 11974–11983.
- (116) Goerigk, L.; Grimme, S. *J. Chem. Theory Comput.* **2010**, *6*, 107–126.
- (117) Csonka, G. I.; French, A. D.; Johnson, G. P.; Stortz, C. A. *J. Chem. Theory Comput.* **2009**, *5*, 679–692.
- (118) Řeha, D.; Valdés, H.; Vondrášek, J.; Hobza, P.; Abu-Riziq, A.; Crews, B.; de Vries, M. S. *Chem.—Eur. J.* **2005**, *11*, 6803–6817.
- (119) Wilke, J. J.; Lind, M. C.; Schaefer, H. F.; Csaszar, A. G.; Allen, W. D. *J. Chem. Theory Comput.* **2009**, *5*, 1511–1523.
- (120) Karton, A.; Tarnopolsky, A.; Lamere, J.-F.; Schatz, G. C.; Martin, J. M. L. *J. Phys. Chem. A* **2008**, *112*, 12868–12886.
- (121) Dinadayalane, T. C.; Vijaya, R.; Smitha, A.; Sastry, G. N. *J. Phys. Chem. A* **2002**, *106*, 1627–1633.
- (122) Guner, V.; Khuong, K. S.; Leach, A. G.; Lee, P. S.; Bartberger, M. D.; Houk, K. N. *J. Phys. Chem. A* **2003**, *107*, 11445–11459.
- (123) Ess, D. H.; Houk, K. N. *J. Phys. Chem. A* **2005**, *109*, 9542–9553.
- (124) Neese, F.; Schwabe, T.; Kossmann, S.; Schirmer, B.; Grimme, S. *J. Chem. Theory Comput.* **2009**, *5*, 3060–3073.
- (125) Grimme, S.; Kruse, H.; Goerigk, L.; Erker, G. *Angew. Chem., Int. Ed.* **2010**, *49*, 1402–1405.
- (126) Grimme, S.; Steinmetz, M.; Korth, M. *J. Org. Chem.* **2007**, *72*, 2118–2126.
- (127) Bryantsev, V. S.; Diallo, M. S.; van Duin, A. C. T.; Goddard, W. A. *J. Chem. Theory Comput.* **2009**, *5*, 1016–1026.
- (128) Tang, K. T.; Toennies, J. P. *J. Chem. Phys.* **2003**, *118*, 4976–4983.
- (129) Zhang, Y.; Yang, W. *Phys. Rev. Lett.* **1998**, *80*, 890.
- (130) Adamo, C.; Barone, V. *J. Chem. Phys.* **1999**, *110*, 6158–6170.
- (131) *ADF2010.02*; SCM, Theoretical Chemistry, Vrije Universiteit: Amsterdam, The Netherlands, 2010. <http://www.scm.com> (accessed October 2011).
- (132) Velde, G. t.; Bickelhaupt, F. M.; Baerends, E. J.; Guerra, C. F.; Gisbergen, S. J. A. v.; Snijders, J. G.; Ziegler, T. *J. Comput. Chem.* **2001**, *22*, 931–967.
- (133) Becke, A. D. *J. Chem. Phys.* **1993**, *98*, 1372–1377.
- (134) Zhao, Y.; Truhlar, D. G. *J. Phys. Chem. A* **2005**, *109*, 5656–5667.
- (135) Becke, A. D. *J. Chem. Phys.* **1996**, *104*, 1040–1046.
- (136) Murray, E. a. D.; Lee, K.; Langreth, D. C. *J. Chem. Theory Comput.* **2009**, *5*, 2754–2762.
- (137) Perdew, J. P.; Wang, Y. *Phys. Rev. B* **1992**, *45*, 13244–13249.
- (138) Shao, Y.; Molnar, L. F.; Jung, Y.; Kussmann, J.; Ochsenfeld, C.; Brown, S. T.; Gilbert, A. T. B.; Slipchenko, L. V.; Levchenko, S. V.; O'Neill, D. P.; DiStasio, R. A., Jr; Lochan, R. C.; Wang, T.; Beran, G. J. O.; Besley, N. A.; Herbert, J. M.; Lin, C. Y.; Voorhis, T. V.; Chien, S. H.; Sodt, A.; Steele, R. P.; Rassolov, V. A.; Maslen, P. E.; Korambath, P. P.; Adamson, R. D.; Austin, B.; Baker, J.; Byrd, E. F. C.; Dachsel, H.; Doerksen, R. J.; Dreuw, A.; Dunietz, B. D.; Dutoi, A. D.; Furlani, T. R.; Gwaltney, S. R.; Heyden, A.; Hirata, S.; Hsu, C.-P.; Kedziora, G.; Khalliulin, R. Z.; Klunzinger, P.; Lee, A. M.; Lee, M. S.; Liang, W.; Lotan, I.; Nair, N.; Peters, B.; Proynov, E. I.; Pieniazek, P. A.; Rhee, Y. M.; Ritchie, J.; Rosta, E.; Sherrill, C. D.; Simmonett, A. C.; Subotnik, J. E.; Woodcock, H. L., III; Zhang, W.; Bell, A. T.; Chakraborty, A. K.; Chipman, D. M.; Keil, F. J.; Warshel, A.; Hehre, W. J.; Schaefer, H. F., III; Kong, J.; Krylov, A. I.; Gill, P. M. W.; Head-Gordon, M. *Phys. Chem. Chem. Phys.* **2006**, *8*, 3172–3191.
- (139) Tsuneda, T.; Suzumura, T.; Hirao, K. *J. Chem. Phys.* **1999**, *110*, 10664–10678.
- (140) Iikura, H.; Tsuneda, T.; Yanai, T.; Hirao, K. *J. Chem. Phys.* **2001**, *115*, 3540–3544.
- (141) Song, J.-W.; Hirose, T.; Tsuneda, T.; Hirao, K. *J. Chem. Phys.* **2007**, *126*, 154105–7.
- (142) Perdew, J. P.; Ruzsinszky, A.; Tao, J.; Csonka, G. I.; Scuseria, G. E. *Phys. Rev. A* **2007**, *76*, 042506.
- (143) 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.; Su, S.; Windus, T. L.; Dupuis, M.; Montgomery, J. A., Jr. *J. Comput. Chem.* **1993**, *14*, 1347–1363.
- (144) Dunning, T. H., Jr. *J. Chem. Phys.* **1989**, *90*, 1007–1023.
- (145) Woon, D. E.; Dunning, T. H., Jr. *J. Chem. Phys.* **1993**, *98*, 1358–1371.
- (146) Wilson, A. K.; Woon, D. E.; Peterson, K. A.; Dunning, T. H., Jr. *J. Chem. Phys.* **1999**, *110*, 7667–7676.
- (147) Weigend, F.; Ahlrichs, R. *Phys. Chem. Chem. Phys.* **2005**, *7*, 3297–3305.
- (148) Lenthe, E. v.; Baerends, E. J.; Snijders, J. G. *J. Chem. Phys.* **1993**, *99*, 4597–4610.
- (149) Valdes, H.; Pluhackova, K.; Pitonak, M.; Rezac, J.; Hobza, P. *Phys. Chem. Phys.* **2008**, *10*, 2747–2757.
- (150) Rezac, J.; Jurecka, P.; Riley, K. E.; Cerny, J.; Valdes, H.; Pluhackova, K.; Berka, K.; Rezac, T.; Pitonak, M.; Vondrasek, J.; Hobza, P. *Collect. Czech. Chem. Commun.* **2008**, *73*, 1261–1270.
- (151) Minenkov, Y.; Occhipinti, G.; Jensen, V. R. *J. Phys. Chem. A* **2009**, *113*, 11833–11844.
- (152) Goerigk, L.; Grimme, S. *Phys. Chem. Chem. Phys.* **2011**, *13*, 6670–6688.
- (153) Torker, S.; Merki, D.; Chen, P. *J. Am. Chem. Soc.* **2008**, *130*, 4808–4814.
- (154) Sanford, M. S.; Love, J. A.; Grubbs, R. H. *J. Am. Chem. Soc.* **2001**, *123*, 6543–6554.
- (155) Tsipis, A. C.; Orpen, A. G.; Harvey, J. N. *Dalton Trans.* **2005**, 2849–2858.
- (156) Zhao, Y.; Truhlar, D. G. *Org. Lett.* **2007**, *9*, 1967–1970.
- (157) Song, J.-W.; Tsuneda, T.; Sato, T.; Hirao, K. *Org. Lett.* **2010**, *12*, 1440–1443.

- (158) Gonthier, J. F.; Wodrich, M. D.; Steinmann, S. N.; Corminboeuf, C. *Org. Lett.* **2010**, *12*, 3070–3073.
- (159) Stanton, R. V.; Kenneth M. Merz, J. *J. Chem. Phys.* **1994**, *100*, 434–443.
- (160) Zhang, Q.; Bell, R.; Truong, T. N. *J. Phys. Chem.* **1995**, *99*, 592–599.
- (161) Durant, J. L. *Chem. Phys. Lett.* **1996**, *256*, 595–602.
- (162) Baker, J.; Muir, M.; Andzelm, J. *J. Chem. Phys.* **1995**, *102*, 2063–2079.
- (163) Patchkovskii, S.; Ziegler, T. *J. Chem. Phys.* **2002**, *116*, 7806–7813.
- (164) Livshits, E.; Baer, R. *Phys. Chem. Chem. Phys.* **2007**, *9*, 2932–2941.
- (165) Stein, T.; Kronik, L.; Baer, R. *J. Am. Chem. Soc.* **2009**, *131*, 2818–2820.
- (166) Baer, R.; Livshits, E.; Salzner, U. *Annu. Rev. Phys. Chem.* **2010**, *61*, 85–109.
- (167) Haunschuld, R.; Scuseria, G. E. *J. Chem. Phys.* **2010**, *132*, 224106.
- (168) Marshall, M. S.; Steele, R. P.; Thanthiriwatte, K. S.; Sherrill, C. D. *J. Phys. Chem. A* **2009**, *113*, 13628–13632.
- (169) Pluhackova, K.; Grimme, S.; Hobza, P. *J. Phys. Chem. A* **2008**, *112*, 12469–12474.
- (170) Mandel, G.; Donohue, J. *Acta Crystallogr., Sect. B* **1972**, *28*, 1313–1316.
- (171) Cameron, T. S.; Deeth, R. J.; Dionne, I.; Du, H.; Jenkins, H. D. B.; Krossing, I.; Passmore, J.; Roobottom, H. K. *Inorg. Chem.* **2000**, *39*, 5614–5631.
- (172) Mitzel, N. W.; Losehand, U.; Wu, A.; Cremer, D.; Rankin, D. W. H. *J. Am. Chem. Soc.* **2000**, *122*, 4471–4482.
- (173) Grimme, S. *Chem.—Eur. J.* **2004**, *10*, 3423–3429.
- (174) Knoblock, K. M.; Silvestri, C. J.; Collard, D. M. *J. Am. Chem. Soc.* **2006**, *128*, 13680–13681.
- (175) Moellmann, J.; Grimme, S. *Phys. Chem. Chem. Phys.* **2010**, *12*, 8500–8504.
- (176) Sameera, W. M. C.; Maseras, F. *Phys. Chem. Chem. Phys.* **2011**, *13*, 10520–10526.
- (177) Olaya, A. J.; Ge, P.; Gonthier, J. F.; Pechy, P.; Corminboeuf, C.; Girault, H. H. *J. Am. Chem. Soc.* **2011**, *133*, 12115–12123.
- (178) Ge, P.; Olaya, A. J.; Todorova, T. K.; Corminboeuf, C.; Girault, H. H. *submitted*.
- (179) Rochat, S.; Steinmann, S. N.; Corminboeuf, C.; Severin, K. *Chem. Commun.* **2011**, *47*, 10584–10586.

# Influence of Triplet Instabilities in TDDFT

Michael J. G. Peach,\* Matthew J. Williamson,<sup>†</sup> and David J. Tozer\*

Department of Chemistry, Durham University, South Road, Durham, DH1 3LE United Kingdom

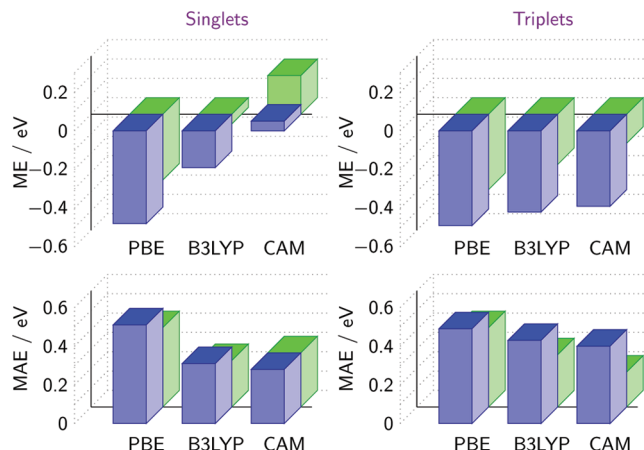
**S** Supporting Information

**ABSTRACT:** Singlet and triplet vertical excitation energies from time-dependent density functional theory (TDDFT) can be affected in different ways by the inclusion of exact exchange in hybrid or Coulomb-attenuated/range-separated exchange–correlation functionals; in particular, triplet excitation energies can become significantly too low. To investigate these issues, the explicit dependence of excitation energies on exact exchange is quantified for four representative molecules, paying attention to the effect of constant, short-range, and long-range contributions. A stability analysis is used to verify that the problematic TDDFT triplet excitations can be understood in terms of the ground state triplet instability problem, and it is proposed that a Hartree–Fock stability analysis should be used to identify triplet excitations for which the presence of exact exchange in the TDDFT functional is undesirable. The use of the Tamm–Dancoff approximation (TDA) significantly improves the problematic triplet excitation energies, recovering the correct state ordering in benzoquinone; it also affects the corresponding singlet states, recovering the correct state ordering in naphthalene. The impressive performance of the TDA is maintained for a wide range of molecules across representative functionals.

## 1. INTRODUCTION

Time-dependent<sup>1–4</sup> density functional theory<sup>5–8</sup> (TDDFT) in the adiabatic approximation is a widely used method for studying molecular electronic excited states. The accuracy of a TDDFT calculation is largely governed by the choice of exchange–correlation functional. Generalized gradient approximations (GGAs) have been largely superseded by hybrid functionals that incorporate a fixed amount of exact orbital exchange (hereafter denoted exact exchange), independent of the interelectron distance  $r_{12}$ . [Exact exchange in the DFT context is defined as the standard Hartree–Fock (HF) exchange energy expression, evaluated using the Kohn–Sham orbitals.] More recently, there has been enormous growth in the use of so-called Coulomb-attenuated or range-separated functionals<sup>9–18</sup> where the amount of exact exchange depends on  $r_{12}$ . The primary reason for this growth is that functionals where the amount of exact exchange increases with  $r_{12}$  have been shown to yield notably improved long-range, Rydberg and charge-transfer excitation energies, while maintaining good quality local excitations.<sup>12,14,16,17,19–26</sup> The majority of these studies have considered excitations to singlet excited states. The quality of excitations to triplet states with Coulomb-attenuated/range-separated functionals, of technological importance in phosphorescence in OLEDs, bioimaging, etc., is less well documented.<sup>27–29</sup>

Recent work by Thiel and co-workers<sup>30–32</sup> has provided a set of correlated wave function [complete active space self-consistent field with second order perturbation theory (CASPT2) and linear response coupled cluster with approximate perturbative triple excitations<sup>33</sup> (CC3)] benchmark results on small, closed-shell organic molecules, allowing comparison of low-lying local singlet and triplet vertical excitation energies within the same molecule. Among others, Silva-Junior et al.,<sup>34</sup> Jacquemin et al.,<sup>23,27</sup> Della Sala and Fabiano,<sup>35</sup> and Huix-Rotllant et al.<sup>36</sup> have assessed the performance of various DFT-based methods for this set. We have repeated the conventional TDDFT calculations of

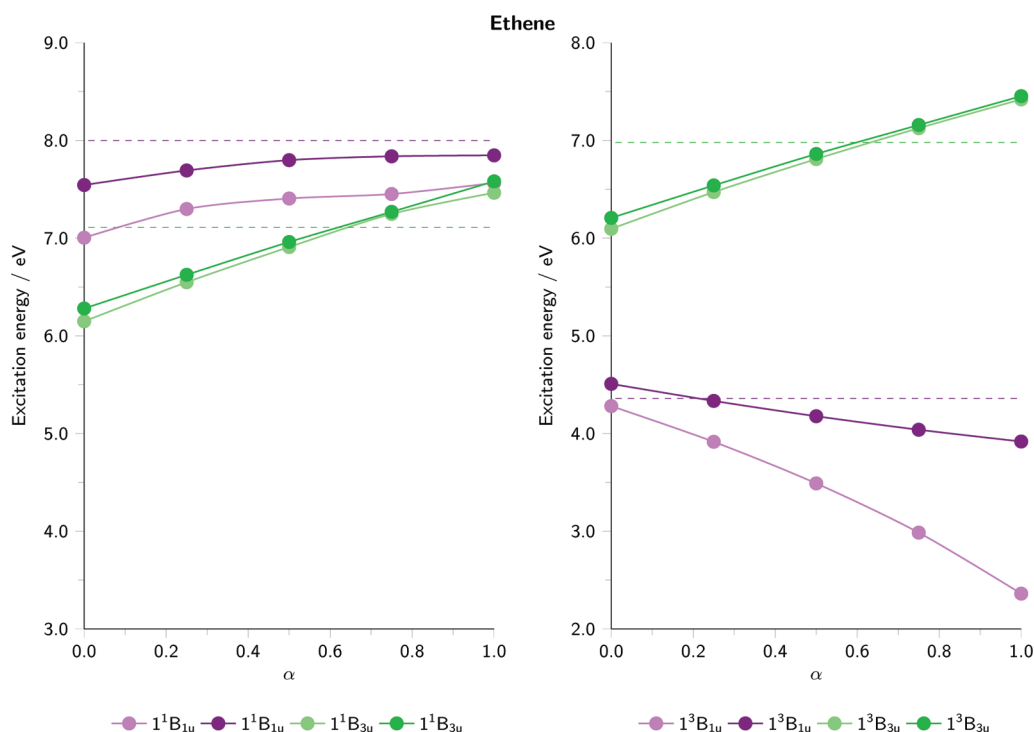


**Figure 1.** Mean errors (ME) and mean absolute errors (MAE), relative to the reference values of ref 32, for 57 singlet and 63 triplet vertical excitation energies. Blue bars represent conventional TDDFT errors; green bars represent TDA errors. CAM denotes CAM-B3LYP.

refs 23, 27, and 34, using the aug-cc-pVTZ basis set (which, unless otherwise stated, is used throughout this study), with the PBE<sup>37</sup> (GGA, no exact exchange), B3LYP<sup>38–43</sup> (hybrid, fixed 20% exact exchange), and CAM-B3LYP<sup>12</sup> (Coulomb-attenuated, with 19% exact exchange increasing with  $r_{12}$  to 65%) functionals at the same MP2/6-31G\* geometries. We consider 63 vertical triplet excitations and the 57 equivalent vertical singlet excitations for which reference CASPT2/CC3 reference values are available, using the Dalton<sup>44</sup> and Gaussian 09<sup>45</sup> programs. Mean and mean absolute errors, relative to the reference values, are presented as blue bars in Figure 1.

**Received:** September 16, 2011

**Published:** October 05, 2011



**Figure 2.** The variation of singlet (left panel) and triplet (right panel) excitation energies in ethene, as a function of the amount of exact exchange  $\alpha$ . The lighter version of the color represents the TDDFT results, the darker version, the TDA results. Dashed lines represent reference values.

The results for the singlet states illustrate the well-known trend: The PBE GGA functional underestimates the excitation energies, while increasing the amount of exact exchange (PBE  $\rightarrow$  B3LYP  $\rightarrow$  CAM-B3LYP) beneficially increases the excitation energies, reducing mean and mean absolute errors. For triplet states, PBE again underestimates the excitation energies, but the improvement upon increasing the amount of exact exchange is much less pronounced than for the singlet states. The reason for this different behavior is evident from an analysis of individual excitations—while in many cases the triplet excitation energy does (beneficially) increase with increasing exact exchange, in many other cases it *drops* significantly, leading to a degradation in accuracy. The latter behavior is not a consequence of low-overlap charge-transfer<sup>20</sup> failure.

It has long been known<sup>46–53</sup> that time-dependent Hartree–Fock theory (TDHF, 100% exact exchange) significantly underestimates triplet excitation energies when there is a triplet instability problem in the ground state wave function and that this underestimation can be largely overcome using configuration interaction singles (CIS). Given the similarity between the TDDFT and TDHF formalisms, we should anticipate similar problems in TDDFT, particularly as the amount of exact exchange increases, which could explain the observed underestimation of certain states. Bauernschmitt and Ahlrichs<sup>50</sup> and Hirata and Head-Gordon<sup>54</sup> presented early examples where hybrid functionals underestimate triplet excitation energies in systems known to have triplet instability problems. The latter authors also demonstrated that these errors are largely eliminated upon application of the Tamm–Dancoff approximation,<sup>55,56</sup> which is the TDDFT analogue of CIS.

In the present study, we explicitly quantify the influence of exact exchange on representative TDDFT singlet and triplet excitation energies and verify that the problematic triplet states

can be understood in terms of the triplet instability problem. Despite being highlighted in several studies,<sup>50,54,57–59</sup> this consequence of triplet instabilities is not widely appreciated in the TDDFT user community; it is, however, of increasing relevance due to the growth in the use of functionals containing large amounts of exact exchange. We propose that a stability analysis of the Hartree–Fock wave function should be used to identify triplet excitations for which the presence of exact exchange in the TDDFT functional is undesirable. By analogy with the TDHF/CIS case, and following ref 54, we then quantify the extent to which the TDDFT triplet problems can be overcome using the Tamm–Dancoff approximation. We also consider the effect of this approximation on singlet states, including state ordering in naphthalene, which is a challenging problem for approximate TDDFT. Finally, the full error analysis in Figure 1 is repeated using the Tamm–Dancoff approximation.

We commence in section 2 by quantifying the influence of exact exchange on singlet and triplet excitation energies for a representative set of molecules. Section 3 relates the observations to the triplet instability, and section 4 considers the Tamm–Dancoff approximation. Conclusions are presented in section 5.

## 2. EXCHANGE DEPENDENCE OF EXCITATION ENERGIES

To illustrate and quantify the influence of exact exchange in a systematic manner, we first consider the evolution of vertical excitation energies as a function of the fraction of exact exchange in a conventional global hybrid functional. Following Becke,<sup>60</sup> we define

$$E_{xc} = \alpha E_x^{\text{HF}}[\varphi] + (1 - \alpha) E_x^{\text{B}}[\rho, \nabla\rho] + E_c^{\text{LYP}}[\rho, \nabla\rho] \quad (1)$$

where the notation  $[\varphi]$ ,  $[\rho]$ , and  $[\nabla\rho]$  indicates explicit orbital, density, and density gradient dependence, respectively; B represents

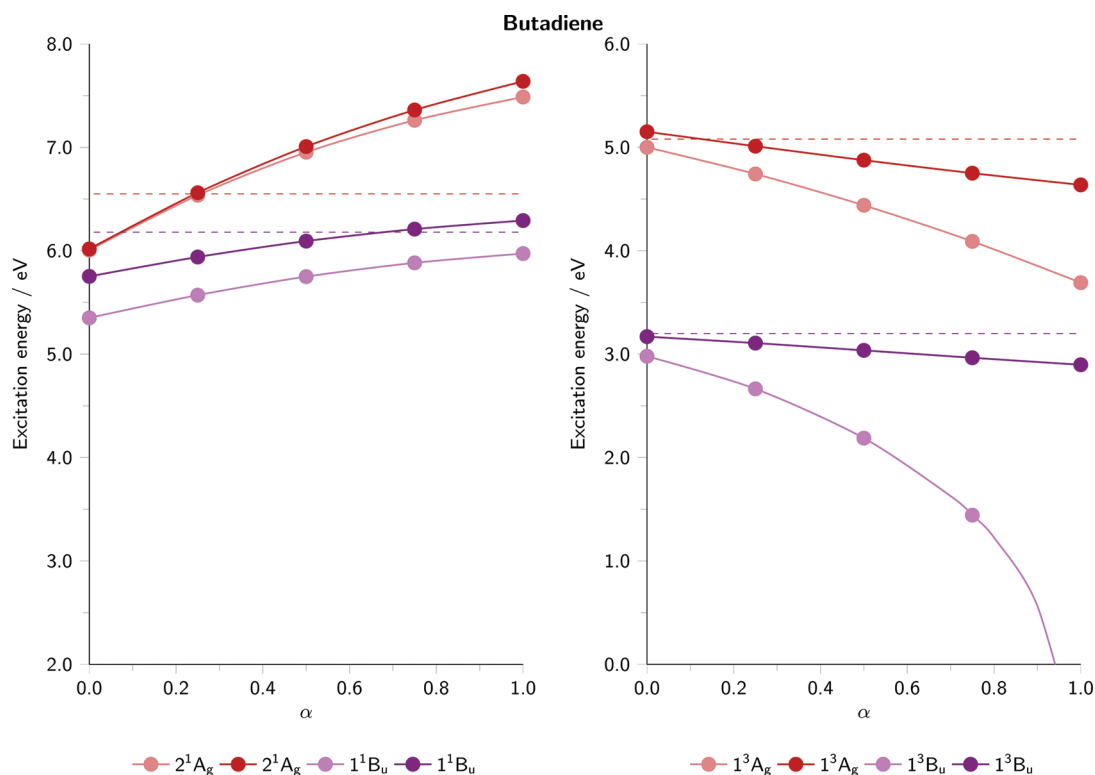


Figure 3. Excitation energies of butadiene (see caption to Figure 2). The TDDFT  ${}^3B_u$  excitation energy becomes imaginary at large  $\alpha$ .

Becke's 1988<sup>39</sup> gradient corrected exchange functional combined in equal proportions with Dirac/Slater<sup>61,62</sup> LDA exchange, and LYP represents the Lee–Yang–Parr<sup>40</sup> GGA correlation functional. From the benchmark set of Thiel and co-workers,<sup>32</sup> we consider four representative molecules: ethene, *E*-butadiene, *p*-benzoquinone, and naphthalene, using the same geometries as before. Additional results for formaldehyde and formamide are presented in the Supporting Information. To facilitate comparison with previous studies, we adopt the molecular orientation (and hence symmetry labels) of the earlier works.

In Figures 2–5, the left panel shows singlet excitation energies, while the right panel shows the equivalent triplet excitation energies (i.e., those that involve predominantly the same orbital transitions), as a function of the fraction of exact exchange  $\alpha$  in eq 1. The lighter solid line of each color represents conventionally evaluated TDDFT excitation energies, to be compared with the horizontal dashed lines that represent accurate reference values, taken from ref 32. For ethene, values were not available for all of the states we consider; comparison is instead made with the experimentally derived reference values used in ref 63, and the d-aug-cc-pVTZ basis set is used for the calculations. In all cases, the GGA ( $\alpha = 0$ ) singlet and triplet excitation energies underestimate the respective reference values.

First consider ethene in Figure 2. As  $\alpha$  increases, both of the singlet excitation energies increase, and each becomes more accurate (albeit at the expense of a less accurate relative energy). For the  ${}^3B_{3u}$  state, the variation with  $\alpha$  is nearly identical to the singlet transition; a significant amount of exact exchange is again optimal. By contrast, the  ${}^3B_{1u}$  state demonstrates markedly different behavior; the excitation energy becomes significantly less accurate with increasing  $\alpha$ , as it drops by nearly 2 eV between  $\alpha = 0$  and  $\alpha = 1$ .

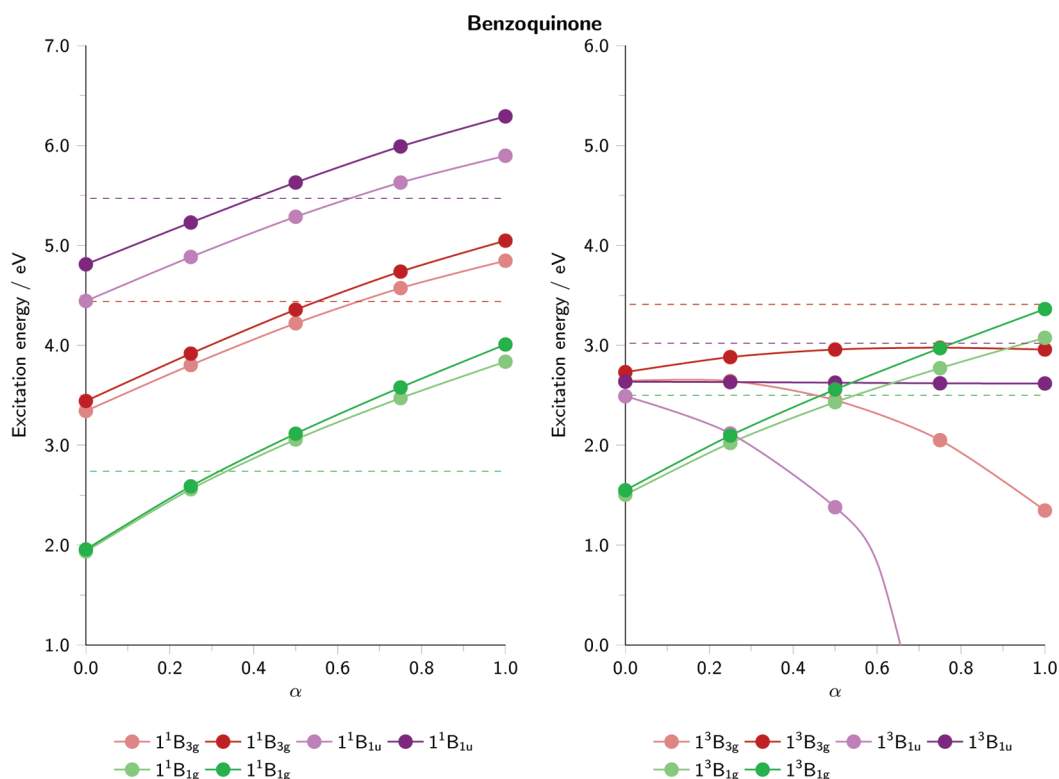
For butadiene in Figure 3, both of the singlet excitation energies again increase with  $\alpha$ , with notably different optimal values. However, both of the triplet excitation energies drop in energy, with each becoming significantly less accurate as  $\alpha$  increases. The  ${}^3A_g$  energy drops by over 1 eV while the  ${}^3B_u$  energy drops considerably more, eventually yielding an imaginary excitation energy (we only plot the real excitation energies).

Next, consider benzoquinone in Figure 4, where three states of each spin are considered. All three singlet excitation energies increase at a similar rate with  $\alpha$ , with modest amounts of exact exchange providing optimal results. By contrast, the three triplet states each exhibit a different dependence on  $\alpha$ . The  ${}^3B_{1g}$  energy behaves essentially identically to the singlet counterpart. The  ${}^3B_{3g}$  energy decreases, becoming less accurate, while the  ${}^3B_{1u}$  energy drops rapidly, becoming imaginary for large  $\alpha$ . This differential dependence means that the triplet state ordering is sensitive to the value of  $\alpha$ . The GGA calculation correctly places the  ${}^3B_{1g}$  state lowest in energy, but the ordering becomes incorrect as  $\alpha$  increases, with first the  ${}^3B_{1u}$  state and then the  ${}^3B_{3g}$  state dropping below the  ${}^3B_{1g}$ .

Finally, consider naphthalene in Figure 5, where the  $B_{2u}$  and  $B_{3u}$  states correspond to the  $L_a$  and  $L_b$  states in the usual Platt notation. Both of the singlet excitation energies increase in energy with  $\alpha$ , although no value of  $\alpha$  yields the correct ordering of the two states. This is a well-known problem in approximate TDDFT.<sup>20,64–68</sup> For the triplet states, the  ${}^3B_{3u}$  excitation energy increases gradually with  $\alpha$  and is accurately described, whereas the  ${}^3B_{2u}$  excitation energy decreases and becomes imaginary for large  $\alpha$ .

The results in Figures 2–5 were obtained using a global hybrid functional, where the amount of exact exchange is independent of  $r_{12}$ . In order to ascertain the relative importance of the long- and





**Figure 4.** Excitation energies of benzoquinone (see caption to Figure 2). The TDDFT  ${}^3B_{1u}$  excitation energy becomes imaginary at large  $\alpha$ .

short-range components of the exact exchange on the six identified “dropping” triplet states, we have performed additional calculations using Coulomb-attenuated/range-separated analogues of eq 1. For the long-range calculations, we considered a series of functionals with zero exact exchange at short  $r_{12}$ , increasing to  $\alpha$  exact exchange at large  $r_{12}$ . For the short-range calculations, we considered a series with  $\alpha$  exact exchange at short  $r_{12}$ , decreasing to zero exact exchange at large  $r_{12}$ . Both functionals used a standard error function partitioning with attenuation parameter  $\mu = 0.33 \text{ bohr}^{-1}$ . In all cases, the variation in excitation energy as a function of  $\alpha$  is smooth and monotonic. For the long-range functionals, the changes in excitation energies between  $\alpha = 0$  and  $\alpha = 1$  are  $-0.23$ ,  $-0.22$ ,  $-0.23$ ,  $-0.17$ ,  $+0.15$ , and  $-0.20$  eV, for ethene ( ${}^3B_{1u}$ ), butadiene ( ${}^3B_u$  and  ${}^3A_g$ ), benzoquinone ( ${}^3B_{1u}$  and  ${}^3B_{3g}$ ), and naphthalene ( ${}^3B_{2u}$ ), respectively. For the short-range functionals, the changes are more pronounced, at  $-1.48$ ,  $-1.59$ ,  $-0.93$ ,  $-2.49$  (imaginary beyond that point),  $-0.68$ , and  $-1.50$  eV. We conclude that long-range and short-range exact exchange each tend to cause these triplet excitation energies to decrease, with the effect of the latter (unsurprisingly) being more pronounced.

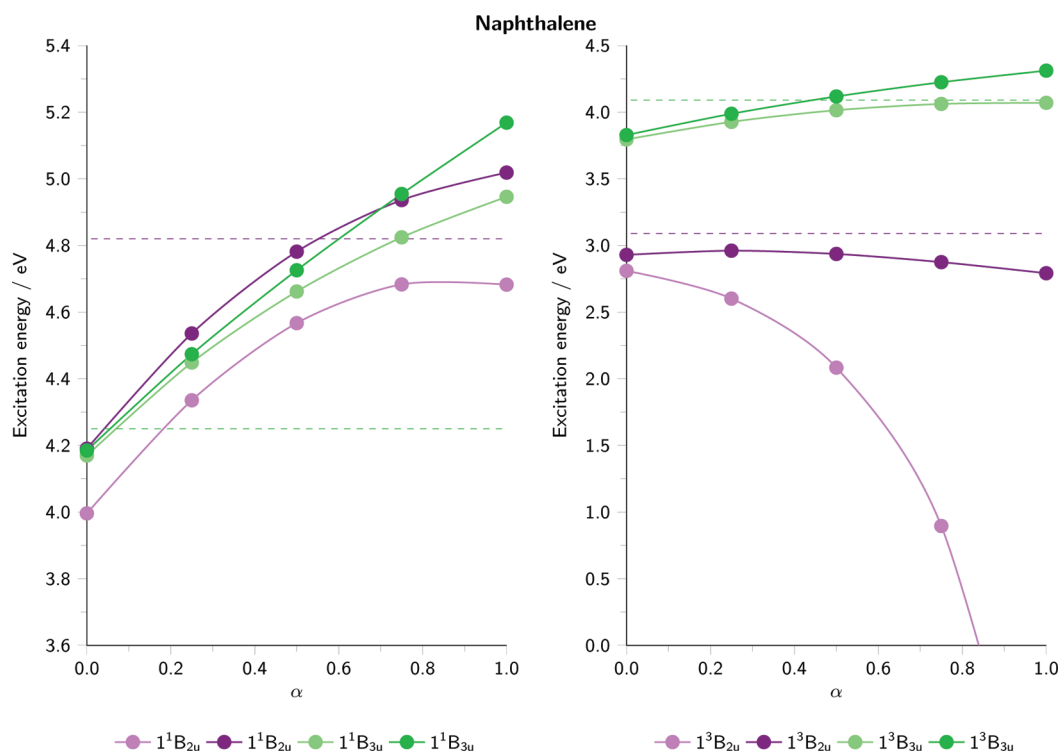
From this analysis, we would predict that both B3LYP (fixed exact exchange) and CAM-B3LYP (fixed- and long-range exact exchange) should underestimate the same six triplet excitation energies, and we have verified that this is indeed the case. Mean errors are  $-0.51$  and  $-0.64$  eV, respectively (compared to  $-0.38$  eV with PBE). The same would be true for any other molecule in the benchmark set where exact exchange causes the triplet excitation energy to drop. To understand why some, but not all, triplet excitation energies drop with exact exchange, we must consider the influence of the triplet instability on TDDFT results.

### 3. THE TRIPLET INSTABILITY PROBLEM

The triplet instability in Hartree–Fock theory is well-known.<sup>46,69</sup> Figure 6a presents potential energy curves for the prototypical molecule,  $H_2$ . The  ${}^1\Sigma_g^+$  spin-restricted Hartree–Fock (RHF) ground state energy becomes too high as the internuclear distance  $R$  increases, due to unphysical ionic components in the wave function. The repulsive  ${}^3\Sigma_u^+$  unrestricted Hartree–Fock (UHF) state does not contain any unphysical ionic components and so dissociates correctly. Consequently, instead of the  ${}^3\Sigma_u^+$  and  ${}^1\Sigma_g^+$  states becoming degenerate at large  $R$ , the energy of the  ${}^3\Sigma_u^+$  state drops below that of the  ${}^1\Sigma_g^+$  for  $R$  larger than  $\sim 3$  bohr. Also shown is the UHF ground state solution, which allows mixing of triplet state character into the singlet wave function. The UHF energy drops below the RHF energy beyond the Coulson–Fischer (CF)<sup>70</sup> point ( $\sim 2.3$  bohr in  $H_2$ ), and correct dissociation is obtained. The RHF solution is therefore unstable with respect to spin-symmetry breaking.

Computationally, this triplet instability manifests as a negative eigenvalue in the electronic Hessian, indicating that specific orbital rotations of an identified space–spin symmetry will lower the energy. Henceforth, we refer to the eigenvalues of this matrix as “stability measures”. Figure 6b presents the lowest  ${}^3\Sigma_u^+$  symmetry stability measure of the Hartree–Fock wave function for  $H_2$ , as a function of  $R$ . It reduces to zero at the CF point and becomes negative beyond. [The stability measures associated with a single determinant are simple to compute and can for instance be calculated in Gaussian<sup>45</sup> using the “stable” keyword, where IOp(9/41) controls the number computed.]

There are intrinsic similarities<sup>49,52</sup> between the equations used to determine the stability and the TDHF/TDDFT equations, and so triplet instabilities have significant implications for excited states determined using these methods. The eigenvectors of the



**Figure 5.** Excitation energies of naphthalene (see caption to Figure 2). The TDDFT  ${}^3B_{2u}$  excitation energy become imaginary at large  $\alpha$ . Note that the scale of this figure is different from that of Figures 2–4.

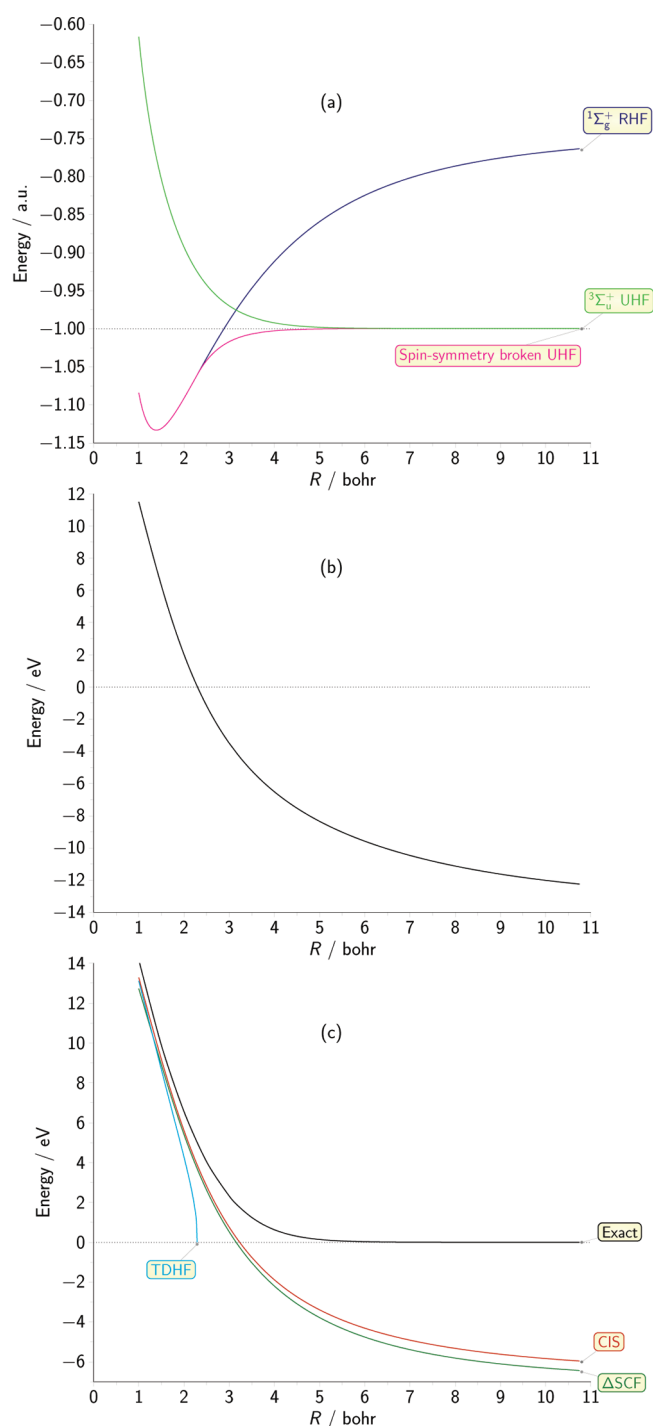
electronic Hessian have identifiable analogues among the orbital rotations associated with electronic excitations, and so it is generally possible to associate a stability measure with each excitation; the stability measure in Figure 6b corresponds to the lowest  ${}^1\Sigma_g^+ \rightarrow {}^3\Sigma_u^+$  excitation. Figure 6c presents this excitation energy as a function of  $R$ . The exact excitation energy approaches zero as  $R \rightarrow \infty$ . The “ $\Delta$ SCF” curve, obtained from the energy difference between the potential energy curves in Figure 6a, becomes increasingly negative at large  $R$ , reflecting the significant overestimation of the  ${}^1\Sigma_g^+$  energy. The influence of the triplet instability on the TDHF excitation energies is striking. The values are reasonable for small  $R$ , but as the CF point is approached, the values become increasingly underestimated, reaching zero at the CF point and becoming imaginary beyond. Analogous results (in the DFT context) have been presented by Casida et al.<sup>57</sup>

The unphysical TDHF excitation energies obtained for large  $R$  are exacerbated by the fact that while ground state HF theory is variational, TDHF is not; for an arbitrary state, the TDHF total electronic energy is no longer a rigorous upper bound on the exact energy. A simple way to restore the variational nature of the excited state energies is to use configuration–interaction singles (CIS) instead of TDHF theory. Figure 6c shows that the CIS excitation energies are close to those from  $\Delta$ SCF. The excitation energies become negative at large  $R$ , rather than imaginary, due to the Hermitian nature of the CIS matrix equations (see section 4).

The key result of this analysis is that as the triplet stability measure decreases toward zero, so the corresponding time-dependent triplet excitation energy also approaches zero, thereby increasingly underestimating the exact value; when the stability measure becomes negative (i.e., when there is a triplet instability), the excitation energy becomes imaginary. By contrast,

CIS is much less problematic. Returning to the TDDFT results in Figures 2–5, we have determined DFT stability measures for the  $\alpha = 0$  and  $\alpha = 1$  functionals for each of the triplet excitations; results are presented in Table 1. For the two states that (beneficially) increase significantly in energy with  $\alpha$ , the stability is large and increases between  $\alpha = 0$  and  $\alpha = 1$ . For the one state whose energy is approximately independent of the amount of exact exchange, the stability varies only slightly. For the three states that drop in energy, but do not become imaginary, the stability reduces significantly. For the three states whose energy drops and becomes imaginary, the stability again reduces significantly and becomes negative by  $\alpha = 1$ . The fact that the dropping triplets are associated with a significant reduction in the stability indicates that the drop—and the resultant underestimation from functionals such as B3LYP and CAM-B3LYP—can be understood in terms of the ground state triplet instability problem, consistent with refs 50 and 54. Analogous results for formaldehyde and formamide are presented in the Supporting Information. We note that an alternative explanation for the underestimated triplet state energies in ethene was recently proposed by Cui and Yang.<sup>28</sup>

Given that it is the inclusion of exact exchange that exacerbates the triplet instability problem, it is also pertinent to calculate the stabilities of these states for the Hartree–Fock wave function. Results are presented in Table 1, and the trend closely follows that of the  $\alpha = 1$  DFT results. [We have confirmed that in cases where the Hartree–Fock stability is large ( $>2$  eV), TDHF and CIS yield similar triplet excitation energies; when the stability is small but positive, TDHF excitation energies are notably smaller than CIS; when the stability is negative, TDHF excitation energies are imaginary, while the CIS values remain real.] This leads us to recommend that a Hartree–Fock stability analysis be



**Figure 6.** (a) HF electronic energy, (b) HF  $^3\Sigma_u^+$  stability measure, and (c)  $^1\Sigma_g^+ \rightarrow ^3\Sigma_u^+$  excitation energies, for  $H_2$  as a function of bond length  $R$ , using the d-aug-cc-pVTZ basis set.

undertaken when computing triplet excitations, to identify states for which the presence of exact exchange in the TDDFT functional is undesirable: The analysis in Table 1 and the Supporting Information (albeit on a limited number of molecules) suggests that in cases where the Hartree–Fock stability is less than  $\sim 2$  eV, the inclusion of exact exchange in the functional will lead to a decrease in excitation energy. If GGAs underestimate the triplet excitation energy (as they often do), then such a decrease will be

**Table 1.** Stability Measures for the DFT Functionals in eq 1 with  $\alpha = 0$ ,  $\alpha = 1$ , and for Hartree–Fock

molecule	state	$\alpha = 0$	$\alpha = 1$	HF
ethene	$^3B_{1u}$	3.22	0.81	0.05
	$^3B_{3u}$	6.04	7.30	6.61
butadiene	$^3B_u$	2.21	−0.16	−0.84
	$^3A_g$	4.06	1.88	1.16
benzoquinone	$^3B_{1u}$	1.89	−0.84	−1.41
	$^3B_{3g}$	2.21	0.33	−0.30
	$^3B_{1g}$	1.25	2.57	2.40
naphthalene	$^3B_{2u}$	2.24	−0.59	−1.25
	$^3B_{3u}$	3.47	2.99	2.66

detrimental. Of course, one could alternatively determine the stability of the DFT calculation directly, but test calculations suggest that the molecule-dependent amount of exact exchange introduced by Coulomb-attenuated/range-separated functionals yields a less-well-defined threshold.

#### 4. THE TAMM–DANCOFF APPROXIMATION IN TDDFT

The CIS approximation corresponds to setting  $B = 0$  in the TDHF generalized eigenvalue equations

$$\begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B} & \mathbf{A} \end{pmatrix} \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} = \omega \begin{pmatrix} \mathbf{1} & \mathbf{0} \\ \mathbf{0} & -\mathbf{1} \end{pmatrix} \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} \quad (2)$$

which is known as the Tamm–Dancoff<sup>55,56</sup> approximation (TDA) to TDHF. The TDDFT equations take exactly the same form (with different matrices  $\mathbf{A}$  and  $\mathbf{B}$ <sup>52</sup>), and so the TDA can equivalently be applied<sup>54</sup> to TDDFT by setting  $\mathbf{B} = \mathbf{0}$ ; see ref 71 for an earlier, related concept. Physically, the TDA corresponds to allowing only excitation between occupied–virtual orbital pairs (given by the eigenvector  $\mathbf{X}$ ) as opposed to conventional TDHF/TDDFT, where virtual–occupied de-excitation contributions ( $\mathbf{Y}$ ) are also allowed. The form of the TDA eigenvalue equation precludes the occurrence of imaginary excitation energies since  $\mathbf{A}$  is Hermitian. We note that there is sometimes concern<sup>52</sup> regarding the validity of transition intensities (oscillator strengths) computed from calculations involving the TDA as they do not satisfy the Thomas–Reiche–Kuhn sum rule.<sup>72–74</sup>

However, this is of no relevance to the calculation of non-spin–orbit coupled triplet transitions, and we therefore do not consider its implications in this study. The TDA is often used as an approximation to full TDDFT due to its relative computational simplicity. Results are often in excellent agreement with full TDDFT (the discrepancy is usually considerably smaller than between CIS and TDHF), but there are instances where TDA yields a better model of reality.<sup>53,54,58,66,75</sup>

Given that CIS is a significant improvement over TDHF when there is a failure associated with triplet instability problems, we now quantify the extent to which the TDA fixes the problematic TDDFT excitations of section 2. We return to Figures 2–5 and now consider the dark solid lines, which present results for TDA excitation energies, as a function of the amount of exact exchange

$\alpha$ . We use the NWChem 6.0 program<sup>76</sup> for the calculation of TDA excitation energies.

First consider ethene in Figure 2. The use of the TDA leads to a significant increase and improvement in the problematic  ${}^3B_{1u}$  excitation energies. By contrast, the  ${}^3B_{3u}$  excitation energies barely change. Notably, the TDA also leads to a shift by about +0.5 eV in the  ${}^1B_{1u}$  state energy—the singlet analogue of the problematic triplet state—resulting in a significant improvement. The  ${}^1B_{3u}$  state is barely affected.

For butadiene in Figure 3, the TDA leads to a significant increase in both of the problematic triplet states, greatly improving accuracy. The improvement is most pronounced for the  ${}^3B_u$  state, which was the most problematic. As with ethene, the corresponding singlet state energy is notably shifted and improved, while the  ${}^1A_g$  state is less affected.

For benzoquinone in Figure 4, the effect of the TDA again increases as the triplet instability problem becomes more severe from  ${}^3B_{1g}$  to  ${}^3B_{3g}$  to  ${}^3B_{1u}$ , leading to the correct state ordering for small values of  $\alpha$ ; the singlet states are again shifted to higher energy by proportionate amounts.

Finally, analogous observations are also made for naphthalene in Figure 5. The TDA leads to a dramatic improvement in the problematic triplet state. Significantly, the associated shift in the corresponding singlet state fixes the state ordering for most values of  $\alpha$  (although the energy difference remains poor), consistent with ref 66. This suggests that the origin of the incorrect state ordering is related to the triplet instability problem associated with the  ${}^3B_{2u}$  state (see ref 68 for an alternative discussion). We note that calculations with the CAM-B3LYP functional, which correctly predicts the state ordering with conventional TDDFT (by only 0.02 eV), is able to correctly increase the energy difference between the two states once the TDA is invoked (the difference becomes 0.18 eV).

Consistent with the findings of ref 54, the results of Figures 2–5 illustrate the benefit of using the TDA for calculating triplet excitation energies when there is a triplet instability problem. Perhaps less expected is the associated effect/improvement of the corresponding singlet states. We end this study by returning to the full assessment in Figure 1. We have repeated all calculations using the TDA, and the results are presented as green bars. The performance of the TDA is impressive, particularly for the triplet states. Indeed, the only error measure that discernibly degrades is the singlet CAM-B3LYP mean error.

## 5. CONCLUSIONS

In this study, we highlighted the fact that singlet and triplet vertical excitation energies in TDDFT can be affected in different ways by the inclusion of exact exchange in hybrid or Coulomb-attenuated/range-separated functionals. The improvement upon the addition of exact exchange is less pronounced for triplet states, which can be traced to the fact that some triplet excitation energies become significantly too low. We studied the explicit dependence of excitation energies on exact exchange for four representative molecules, illustrating the various behaviors and quantifying the effect of constant, short-, and long-range contributions.

We then used the  $H_2$  molecule to illustrate the effect of triplet instabilities on time-dependent excitation energies. As the triplet stability measure associated with an excitation decreases, so the corresponding triplet excitation energy increasingly underestimates the exact value, possibly becoming imaginary. By

determining DFT stability measures for the states of interest in the four representative molecules, we verified that the problematic TDDFT triplets can be understood in terms of the ground state triplet instability problem. We proposed that a Hartree–Fock stability analysis should be carried out to identify triplet excitations for which the presence of exact exchange in the TDDFT functional is undesirable.

We then considered the effect of the Tamm–Dancoff approximation in TDDFT. The use of the TDA significantly improves the problematic triplet states, recovering the correct state ordering in benzoquinone. It also affects the corresponding singlet states, recovering the correct state ordering in naphthalene, which is known to be a significant challenge for approximate TDDFT. The impressive performance of the TDA is maintained for the full assessment set, across representative functionals. We are presently expanding the current work to consider the effect of triplet instabilities and the TDA on singlet and triplet states for the more diverse set of molecules/excitations of ref 20 and for a more diverse set of functionals.

## ■ ASSOCIATED CONTENT

**S Supporting Information.** Plots of the exact exchange dependence of selected excitation energies in the formaldehyde and formamide molecules together with a table, listing the same information as Table 1, for these molecules. It also includes a plot of the Hartree–Fock stability measure plotted against the difference between the TDHF and CIS excitation energies for the molecules in Table 1, together with formaldehyde and formamide. Tables containing individual excitation energies for the data presented in Figure 1 are also included. This information is available free of charge via the Internet at <http://pubs.acs.org/>.

## ■ AUTHOR INFORMATION

### Corresponding Author

\*Fax: +44 191 384 4737. E-mail: M.J.G.Peach@Durham.ac.uk, D.J.Tozer@Durham.ac.uk.

### Present Addresses

<sup>†</sup>School of Physics and Astronomy, The University of Edinburgh, James Clerk Maxwell Building, Mayfield Road, Edinburgh, EH9 3JZ U. K.

## ■ ACKNOWLEDGMENT

M.J.G.P. and M.J.W. thank the EPSRC for financial support and for an undergraduate bursary, respectively.

## ■ REFERENCES

- (1) Runge, E.; Gross, E. K. U. *Phys. Rev. Lett.* **1984**, *52*, 997–1000.
- (2) Gross, E. K. U.; Ullrich, C. A.; Gossmann, U. J. In *Density Functional Theory*; Gross, E. K. U., Dreizler, R. M., Eds.; Plenum: New York, 1994; NATO ASI Series, pp 149–171.
- (3) Casida, M. E. In *Recent Advances in Density Functional Methods, Part I*; Chong, D. P., Ed.; World Scientific: Singapore, 1995; pp 155–192.
- (4) Marques, M. A. L.; Gross, E. K. U. *Time-Dependent Density Funct. Theory* **2003**, 144–184.
- (5) Hohenberg, P.; Kohn, W. *Phys. Rev.* **1964**, *136*, B864–B871.
- (6) Kohn, W.; Sham, L. J. *Phys. Rev.* **1965**, *140*, A1133–A1138.
- (7) Levy, M. *Proc. Natl. Acad. Sci. U.S.A.* **1979**, *76*, 6062–6065.

- (8) Parr, R. G.; Yang, W. *Density-Functional Theory of Atoms and Molecules*; Oxford University Press: Oxford, U. K., 1989; International Series of Monographs on Chemistry, pp 142–200.
- (9) Gill, P. M. W.; Adamson, R. D.; Pople, J. A. *Mol. Phys.* **1996**, *88*, 1005–1010.
- (10) Leininger, T.; Stoll, H.; Werner, H.-J.; Savin, A. *Chem. Phys. Lett.* **1997**, *275*, 151–160.
- (11) Iikura, H.; Tsuneda, T.; Yanai, T.; Hirao, K. *J. Chem. Phys.* **2001**, *115*, 3540–3544.
- (12) Yanai, T.; Tew, D. P.; Handy, N. C. *Chem. Phys. Lett.* **2004**, *393*, 51–57.
- (13) Baer, R.; Neuhauser, D. *Phys. Rev. Lett.* **2005**, *94*, 043002.
- (14) Peach, M. J. G.; Cohen, A. J.; Tozer, D. J. *Phys. Chem. Chem. Phys.* **2006**, *8*, 4543–4549.
- (15) Vydrov, O. A.; Heyd, J.; Krukau, A. V.; Scuseria, G. E. *J. Chem. Phys.* **2006**, *125*, 074106.
- (16) Chai, J.-D.; Head-Gordon, M. *J. Chem. Phys.* **2008**, *128*, 084106.
- (17) Rohrdanz, M. A.; Herbert, J. M. *J. Chem. Phys.* **2008**, *129*, 034107.
- (18) Besley, N. A.; Peach, M. J. G.; Tozer, D. J. *Phys. Chem. Chem. Phys.* **2009**, *11*, 10350–10358.
- (19) Tawada, Y.; Tsuneda, T.; Yanagisawa, S.; Yanai, T.; Hirao, K. *J. Chem. Phys.* **2004**, *120*, 8425–8433.
- (20) Peach, M. J. G.; Benfield, P.; Helgaker, T.; Tozer, D. J. *J. Chem. Phys.* **2008**, *128*, 044118.
- (21) Peach, M. J. G.; Le Sueur, C. R.; Ruud, K.; Guillaume, M.; Tozer, D. J. *Phys. Chem. Chem. Phys.* **2009**, *11*, 4465–4470.
- (22) Jacquemin, D.; Perpète, E. A.; Scuseria, G. E.; Ciofini, I.; Adamo, C. *J. Chem. Theory Comput.* **2008**, *4*, 123–135.
- (23) Jacquemin, D.; Wathelet, V.; Perpète, E. A.; Adamo, C. *J. Chem. Theory Comput.* **2009**, *5*, 2420–2435.
- (24) Rohrdanz, M. A.; Martins, K. M.; Herbert, J. M. *J. Chem. Phys.* **2009**, *130*, 054112.
- (25) Stein, T.; Kronik, L.; Baer, R. *J. Am. Chem. Soc.* **2009**, *131*, 2818–2820.
- (26) Peach, M. J. G.; Helgaker, T.; Salek, P.; Keal, T. W.; Lutnæs, O. B.; Tozer, D. J.; Handy, N. C. *Phys. Chem. Chem. Phys.* **2006**, *8*, 558–562.
- (27) Jacquemin, D.; Perpète, E. A.; Ciofini, I.; Adamo, C. *J. Chem. Theory Comput.* **2010**, *6*, 1532–1537.
- (28) Cui, G.; Yang, W. *Mol. Phys.* **2010**, *108*, 2745–2750.
- (29) Jacquemin, D.; Perpète, E.; Ciofini, I.; Adamo, C. *Theor. Chem. Acc.* **2011**, *128*, 127–136.
- (30) Schreiber, M.; Silva-Junior, M. R.; Sauer, S. P. A.; Thiel, W. *J. Chem. Phys.* **2008**, *128*, 134110.
- (31) Silva-Junior, M. R.; Sauer, S. P. A.; Schreiber, M.; Thiel, W. *Mol. Phys.* **2010**, *108*, 453–465.
- (32) Silva-Junior, M. R.; Schreiber, M.; Sauer, S. P. A.; Thiel, W. *J. Chem. Phys.* **2010**, *133*, 174318.
- (33) Christiansen, O.; Koch, H.; Jørgensen, P. *J. Chem. Phys.* **1995**, *103*, 7429–7441.
- (34) Silva-Junior, M. R.; Schreiber, M.; Sauer, S. P. A.; Thiel, W. *J. Chem. Phys.* **2008**, *129*, 104103.
- (35) Della Sala, F. D.; Fabiano, E. *Chem. Phys.* **2011**, DOI: 10.1016/j.chemphys.2011.05.020.
- (36) Huix-Rotllant, M.; Ipatov, A.; Rubio, A.; Casida, M. E. *Chem. Phys.* **2011**, DOI: 10.1016/j.chemphys.2011.03.019.
- (37) Perdew, J. P.; Burke, K.; Ernzerhof, M. *Phys. Rev. Lett.* **1996**, *77*, 3865–3868.
- (38) Becke, A. D. *J. Chem. Phys.* **1993**, *98*, 5648–5652.
- (39) Becke, A. D. *Phys. Rev. A* **1988**, *38*, 3098–3100.
- (40) Lee, C.; Yang, W.; Parr, R. G. *Phys. Rev. B* **1988**, *37*, 785–789.
- (41) Vosko, S. H.; Wilk, L.; Nusair, M. *Can. J. Phys.* **1980**, *58*, 1200–1211.
- (42) Stephens, P. J.; Devlin, F. J.; Chabalowski, C. F.; Frisch, M. J. *J. Phys. Chem.* **1994**, *98*, 11623–11627.
- (43) We use B3LYP with the VWN5 parametrization of the local correlation energy.
- (44) Dalton, a molecular electronic structure program, Release Dalton2011. See <http://daltonprogram.org/> (accessed October 2011).
- (45) Frisch, M. J.; Trucks, G. W.; Schlegel, H. B.; Scuseria, G. E.; Robb, M. A.; Cheeseman, J. R.; Scalmani, G.; Barone, V.; Mennucci, B.; Petersson, G. A.; Nakatsuji, H.; Caricato, M.; Li, X.; Hratchian, H. P.; Izmaylov, A. F.; Bloino, J.; Zheng, G.; Sonnenberg, J. L.; Hada, M.; Ehara, M.; Toyota, K.; Fukuda, R.; Hasegawa, J.; Ishida, M.; Nakajima, T.; Honda, Y.; Kitao, O.; Nakai, H.; Vreven, T.; Montgomery, J. A., Jr.; Peralta, J. E.; Ogliaro, F.; Bearpark, M.; Heyd, J. J.; Brothers, E.; Kudin, K. N.; Staroverov, V. N.; Kobayashi, R.; Normand, J.; Raghavachari, K.; Rendell, A.; Burant, J. C.; Iyengar, S. S.; Tomasi, J.; Cossi, M.; Rega, N.; Millam, J. M.; Klene, M.; Knox, J. E.; Cross, J. B.; Bakken, V.; Adamo, C.; Jaramillo, J.; Gomperts, R.; Stratmann, R. E.; Yazyev, O.; Austin, A. J.; Cammi, R.; Pomelli, C.; Ochterski, J. W.; Martin, R. L.; Morokuma, K.; Zakrzewski, V. G.; Voth, G. A.; Salvador, P.; Dannenberg, J. J.; Dapprich, S.; Daniels, A. D.; Farkas, Ö.; Foresman, J. B.; Ortiz, J. V.; Cioslowski, J.; Fox, D. J. *Gaussian 09*, Revision A.02; Gaussian Inc.: Wallingford, CT, 2009.
- (46) Seeger, R.; Pople, J. A. *J. Chem. Phys.* **1977**, *66*, 3045–3050.
- (47) Foresman, J. B.; Head-Gordon, M.; Pople, J. A.; Frisch, M. J. *J. Phys. Chem.* **1992**, *96*, 135–149.
- (48) Jamorski, C.; Casida, M. E.; Salahub, D. R. *J. Chem. Phys.* **1996**, *104*, 5134–5147.
- (49) Bauernschmitt, R.; Ahlrichs, R. *J. Chem. Phys.* **1996**, *104*, 9047–9052.
- (50) Bauernschmitt, R.; Ahlrichs, R. *Chem. Phys. Lett.* **1996**, *256*, 454–464.
- (51) Furche, F.; Ahlrichs, R. *J. Chem. Phys.* **2002**, *117*, 7433–7447.
- (52) Dreuw, A.; Head-Gordon, M. *Chem. Rev.* **2005**, *105*, 4009–4037.
- (53) Grimme, S.; Neese, F. *J. Chem. Phys.* **2007**, *127*, 154116.
- (54) Hirata, S.; Head-Gordon, M. *Chem. Phys. Lett.* **1999**, *314*, 291–299.
- (55) Tamm, I. *J. Phys. (USSR)* **1945**, *9*, 449–460.
- (56) Dancoff, S. M. *Phys. Rev.* **1950**, *78*, 382–385.
- (57) Casida, M. E.; Gutierrez, F.; Guan, J.; Gadea, F.-X.; Salahub, D.; Daudey, J.-P. *J. Chem. Phys.* **2000**, *113*, 7062–7071.
- (58) Cordova, F.; Doriol, L. J.; Ipatov, A.; Casida, M. E.; Filippi, C.; Vela, A. *J. Chem. Phys.* **2007**, *127*, 164111.
- (59) Lutnæs, O. B.; Helgaker, T.; Jaszunski, M. *Mol. Phys.* **2010**, *108*, 2579–2590.
- (60) Becke, A. D. *J. Chem. Phys.* **1993**, *98*, 1372–1377.
- (61) Dirac, P. A. M. *Proc. Cam. Phil. Soc.* **1930**, *26*, 376–385.
- (62) Slater, J. C. *Phys. Rev.* **1951**, *81*, 385–390.
- (63) Tozer, D. J.; Handy, N. C. *J. Chem. Phys.* **1998**, *109*, 10180–10189.
- (64) Grimme, S.; Parac, M. *ChemPhysChem* **2003**, *4*, 292–295.
- (65) Richard, R. M.; Herbert, J. M. *J. Chem. Theory Comput.* **2011**, *7*, 1296–1306.
- (66) Wang, Y.-L.; Wu, G.-S. *Int. J. Quantum Chem.* **2008**, *108*, 430–439.
- (67) Wong, B. M.; Hsieh, T. H. *J. Chem. Theory Comput.* **2010**, *6*, 3704–3712.
- (68) Kuritz, N.; Stein, T.; Baer, R.; Kronik, L. *J. Chem. Theory Comput.* **2011**, *7*, 2408–2415.
- (69) Lee, A. M.; Handy, N. C. *J. Chem. Soc., Faraday Trans.* **1993**, *89*, 3999–4003.
- (70) Coulson, C. A.; Fischer, I. *Philos. Mag.* **1949**, *40*, 386–393.
- (71) Grimme, S. *Chem. Phys. Lett.* **1996**, *259*, 128–137.
- (72) Thomas, W. *Naturwissenschaften* **1925**, *13*, 627–628.
- (73) Reiche, F.; Thomas, W. *Z. Phys.* **1925**, *34*, 510–525.
- (74) Kuhn, W. *Z. Phys.* **1925**, *33*, 408–412.
- (75) Hsu, C.-P.; Hirata, S.; Head-Gordon, M. *J. Phys. Chem. A* **2001**, *105*, 451–458.
- (76) Valiev, M.; Bylaska, E.; Govind, N.; Kowalski, K.; Straatsma, T.; Dam, H. V.; Wang, D.; Nieplocha, J.; Apra, E.; Windus, T.; de Jong, W. *Comput. Phys. Commun.* **2010**, *181*, 1477–1489.

# Assessment of Popular DFT and Semiempirical Molecular Orbital Techniques for Calculating Relative Transition State Energies and Kinetic Product Distributions in Enantioselective Organocatalytic Reactions

Sebastian Schenker,<sup>†,‡</sup> Christopher Schneider,<sup>‡</sup> Svetlana B. Tsogoeva,<sup>‡</sup> and Timothy Clark<sup>\*,†,§</sup>

<sup>†</sup>Computer-Chemie-Centrum der Friedrich-Alexander-Universität Erlangen-Nürnberg, Nögelsbachstrasse 25, 91052 Erlangen, Germany

<sup>‡</sup>Lehrstuhl I für Organische Chemie der Friedrich-Alexander-Universität Erlangen-Nürnberg, Henkestrasse 25, 91054 Erlangen, Germany

<sup>§</sup>Centre for Molecular Design, University of Portsmouth, Mercantile House, Portsmouth PO1 2EG, United Kingdom

 Supporting Information

**ABSTRACT:** The performance of computationally accessible levels of calculation for the transition states of organocatalytic reaction has been assessed. Reference post-Hartree–Fock single point energy calculations were used as standards for the gas-phase Born–Oppenheimer relative energies of pairs of alternative transition states that lead to the two product enantiomers. We show that semiempirical methods cannot even be relied on to yield qualitatively correct results. The geometries (optimized, for instance, with DFT) have a large impact on the results of high-level post-HF calculations, so that it is essential to use an adequate DFT technique and basis set. DFT can yield quantitatively correct results that are consistent with post-HF calculations if functionals that consider dispersion are used. Geometries for large systems show larger errors than those for smaller ones but are treated better by functionals such as M06-2X and w97Bxd that include dispersion implicitly or explicitly. Local correlation techniques introduce errors of comparable magnitude to those given by different levels of geometry optimization. We recommend RICCC2/TZVP//M06-2X/TZVP, RI-MP2/TZVP//M06-2X/TZVP, and M06-2X/TZVP//M06-2X/TZVP calculations in that order, depending on the size of the system.

## INTRODUCTION

**Aims.** Quantum mechanical calculations have long been important tools in mechanistic organic chemistry. This importance has increased with the availability of accurate modern DFT techniques. Valuable information about the atomistic details of reaction mechanisms that is not available from any other source can be obtained from calculations. Often, the structures and energies of transition states are of most interest. These details of transition states not only tell us qualitatively which reaction pathway is most likely but can also provide a theoretical estimate of the product distribution in the case of kinetic reaction control. Calculations have now advanced to the stage that they can be used to evaluate proposed reaction pathways or to predict the properties of unknown systems, as has been shown successfully for several enantioselective organocatalytic reactions.<sup>1,2</sup> It is therefore of paramount importance to be able to assess the reliability and probable error limits of different calculational methods for typical organic reactions systems. This is especially important as quantum mechanical calculations are now routinely applied to systems large enough that dispersion interactions between nonpolar residues become important.<sup>3–5</sup> Surprisingly, there have been relatively few systematic studies of the quality of commonly used methods. While several methodological benchmarks have been carried out that focus either on achiral activation barriers<sup>6–8</sup> or on relative conformational energies,<sup>9</sup> to our knowledge only two studies on relative enantiomeric or diastereomeric transition-state energies

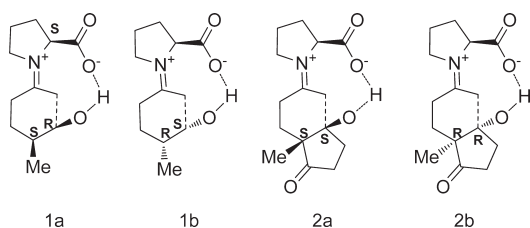
have been carried out to date, by Breslow and co-workers<sup>10</sup> and Simon and Goodman<sup>11</sup> (a concise review of DFT benchmark studies has been published by Ramos et al.<sup>12</sup>). The reaction yield and its enantio- or diastereoselectivity are the quantities of most interest to experimentalists. Whereas the yield is only marginally suitable for theoretical prediction (because it can be affected by many extraneous factors such as competing reactions), selectivities have long been the goal of predictive calculations at many levels.<sup>2,11</sup> We therefore concentrate on this aspect of the calculations, i.e., the difference between activation energies, rather than accurate prediction of their absolute values, because it provides an estimate of the quantity most needed in synthetic research.

Breslow and co-workers concentrated on single-electron transfer and radical reactions, and their primary aim was to optimize predictions of enantioselectivities. Their work considered only two parametrized DFT methods, B3LYP<sup>13</sup> and M06-2X,<sup>14</sup> and concentrated on the suitability of different solvent treatments for improving the quantitative prediction of reaction products. Predictions were good for all systems for which gas-phase UB3LYP/6-31G(d) energies alone gave good agreement with experimental findings, but the authors also noted the shortcomings of their approach for systems in which dispersion interactions are important. In these cases, even the newer of the two functionals, M06-2X,

**Received:** March 23, 2011

**Published:** October 03, 2011

Scheme 1. Transition States for the Aldol Reaction Systems 1 and 2



which treats dispersion effects implicitly, did not improve the results. However, the results of this study are unlikely to be applicable to the closed-shell electrophile/nucleophile reactivity usually involved in organocatalysis. As we were carrying out this study, Simon and Goodman reported a benchmark study on a large set of DFT methods, which they used to model known transition states that have been characterized in the literature.<sup>11</sup> They point out that most hybrid and meta-GGA functionals yield similar results for optimized structures. Because they did not carry out high-level reference calculations, they could not identify a recommended functional for optimizations but did focus on the vibrational free-energy corrections for the calculated structures and found that the resolution of the integration grid in the Hessian calculation is vital for accurate free-energy corrections. They also noted that the deviations in energies and geometries can be high throughout the spectrum of functionals.

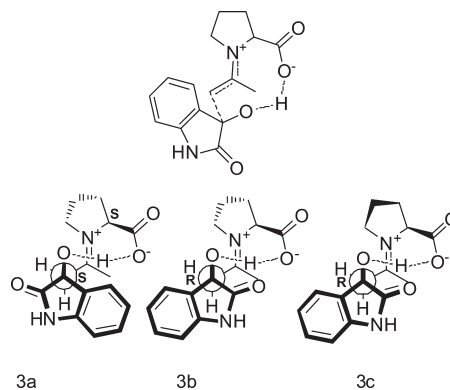
These two studies,<sup>10,11</sup> while instructive, are limited in that they focus on improving results that already agree well with experimental results. Our aim was to benchmark practical methods with the ultimate aim of predicting stereoselectivity, rather than to find the best corrections to apparently accurate gas-phase Born–Oppenheimer energies. We have therefore focused on a larger set of basis sets, DFT functionals, post-HF methods, and semiempirical Hamiltonians in order to identify practical techniques that agree well with the highest-level calculations. Here, we therefore first consider the fundamental problem of how to calculate good gas-phase data, which can then be combined with appropriate solvent and free-energy corrections to allow us to design versatile computational models for quantitative predictions for a large number of important reactions.

As our aim is to identify reliable computational predictions for routine studies, we focus on four known enantioselective and/or diastereoselective organocatalytic reactions from the work of Houk and Bahmanyar,<sup>15</sup> Tomassini et al.,<sup>16</sup> Papái et al.,<sup>17</sup> and Tsogoeva et al.<sup>18</sup> These are two proline-catalyzed aldol addition reactions with 40–47 atoms in the reaction system, a nitro-Michael reaction *via* an enamine intermediate with 67 atoms, and a thiourea-catalyzed nitro-Michael reaction *via* an amine intermediate with 81 atoms. These varied systems allow us to study the influence of different intermolecular interaction patterns on the quality of the results. While the small proline-catalyzed reactions are dominated by covalent and H-bond interactions, dispersion interactions play a more important role for the larger nitro-Michael reactions.

## METHODS

In order to test a wide spectrum of computationally economic techniques, we have investigated three common NDDO-based semiempirical molecular-orbital (MO) techniques, AM1,<sup>19</sup>

Scheme 2. Transition States for the Aldol Reaction System 3



PM3,<sup>20,21</sup> and PM6;<sup>22</sup> Hartree–Fock (HF) *ab initio* theory; and six popular DFT functionals. The six DFT methods are PBE<sup>23</sup> as a pure generalized gradient approximation (GGA) functional, B3LYP<sup>13</sup> as a hybrid-GGA functional, TPSS<sup>24</sup> and TPSSH<sup>24</sup> as meta-GGA functionals (pure and hybrid), and, as newer representatives, Head-Gordon and Chai’s more recent wB97xd functional, which includes an explicit empirical dispersion correction,<sup>25</sup> and M06-2X by Zhao and Truhlar,<sup>14</sup> which implicitly corrects for dispersion and for which very high accuracy for small organic transition states was reported.<sup>26</sup> All DFT calculations were carried out with the Pople 6-31G(d) double- $\zeta$ <sup>27,28</sup> and the Ahlrichs TZVP triple- $\zeta$ <sup>29</sup> basis sets.

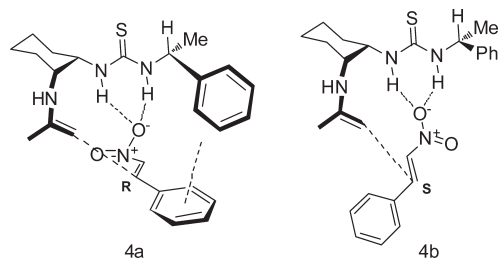
In order to provide benchmark results for comparison, we also performed *ab initio* post-HF single-point calculations. Since direct MP2<sup>30,31</sup> and CCSD<sup>32</sup> scale poorly,<sup>33</sup> we used RI-MP2<sup>34</sup> and RI-CC2<sup>35</sup> (note that RIC2 employs the CC2 approximation<sup>36</sup> in addition to RI). The resolution of identity (RI) approximation allows us to calculate single-point energies even for the largest systems with a polarized triple- $\zeta$  basis set. Geometry optimizations beyond RI-MP2 or single-point calculations with larger basis sets on systems of this size remain prohibitive. Unfortunately, CCSD(T), which is often referred to as the “gold standard” in quantum chemistry, is computationally still too costly to be used for the systems considered here.<sup>37</sup> The localized, and therefore theoretically linear scaling, LCCSD(T) ansatz of Werner and Schütz<sup>38–43</sup> was used to test the influence of triple excitations for the electronic energies. LCCSD(T) calculations were performed using Dunning’s cc-pVTZ basis set.<sup>44</sup> Different polarized triple- $\zeta$  basis sets were used for RI and LCCSD(T) calculations because of the authors’ recommended fitting basis sets. We could perform LCCSD(T) single-point calculations for all but the largest system.

## CALCULATED TRANSITION STATES IN DETAIL

**Aldol Reactions.** The transition states used in this study can be divided in two groups, the smaller aldol addition systems and the larger nitro-Michael additions. The first group consists of systems 1 and 2 (Scheme 1) from the work of Houk and Bahmanyar<sup>15</sup> and of system 3 (Scheme 2) from the work of Tomassini et al.<sup>16</sup> Systems 1 and 2 consist, in turn, of two sets of two diastereomeric transition states 1a,b and 2a,b (see Scheme 1).

System 3 consists of two diastereomeric transition states 3a,b and a third 3c, which only differs from 3b by a proline ring flip and is known to be very similar in energy to 3b<sup>16</sup> (see Scheme 2).

Scheme 3. Transition States for the Nitro-Michael System 4



In all cases, the major interactions between the catalyst and the reacting system are either covalent or hydrogen bonds. Both the C–C bond formation and a simultaneous proton transfer must be described correctly for this reaction. The size of the systems ranges from 41 (20 non-hydrogen) to 67 (33 non-hydrogen) atoms with molecular weights ranging from of  $225 \text{ g mol}^{-1}$  to  $301 \text{ g mol}^{-1}$  (i.e., 150–160 electrons; 384–408 basis functions with 6-31G(d) and 526–562 basis functions with TZVP).

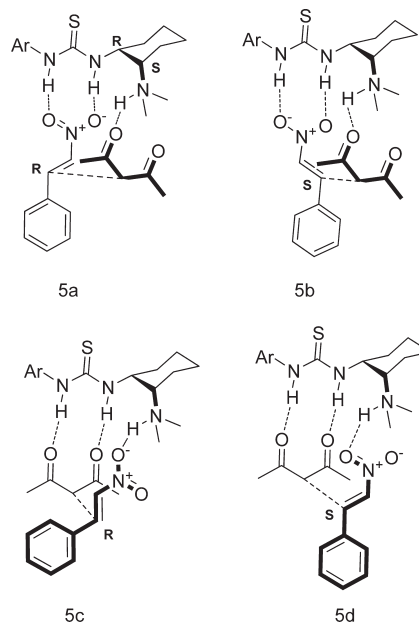
**Nitro-Michael Reactions.** For the larger systems, we used the transition states of two thiourea-catalyzed nitro-Michael reactions, published by Tsogoeva et al. (4)<sup>18</sup> and Papai et al. (5).<sup>17</sup> The nitro-Michael reaction 4 (Scheme 3) consists of the two diastereomeric transition states 4a and 4b. Stereocontrol is achieved mainly by covalent and H-bond interactions. For transition state 4a, we can also expect the formation of a  $\pi$ – $\pi$  interaction pattern between the phenyl groups of the nitrostyrene and the catalyst, in contrast to 4b, for which this interaction is not possible. The size of the system leads us to expect dispersion interactions between the reaction partners to be more important than for the aldol reactions. System 4 consists of 67 atoms (33 non-hydrogen) and has a molecular weight of  $464 \text{ g mol}^{-1}$  (250 electrons; 643 basis functions with 6-31G(d) and 871 basis functions with TZVP).

For 5 (Scheme 4), two pairs of diastereomeric transition states 5a and 5b and 5c and 5d were investigated. This system is mainly controlled by H-bond interactions, and we can also expect dispersion interactions to be important. In comparison to 4, system 5 is expected to show increased flexibility of the alignment of the two reactants at the catalytically active site. This was reported by Pápai et al., who observed low-energy normal modes for bending of the reactants in reactant complexes and the transition state.<sup>17</sup> TS 5 is the largest system considered, with 81 atoms (41 non-hydrogen) and a molecular weight of  $662 \text{ g mol}^{-1}$  (346 electrons; 831 basis functions with 6-31G(d) and 1123 basis functions with TZVP).

## RESULTS AND DISCUSSION

**Reference Energies.** As definitive a reference calculation as possible is necessary in order to be able to judge the quality of the results. The experimentally observed product distributions can only provide guidance, since the conditions of the gas-phase calculations do not necessarily correspond to the experimental ones. This is particularly important, as small effects (entropy, solvation) can control kinetic enantioselectivity, which reacts extremely sensitively to small changes in activation energies. In our case, gas-phase Born–Oppenheimer geometries and energies are the relevant target properties. We decided to use (RI-)MP2 and CCSD (as RI-CC2 or LCCSD) and local LCCSD(T)

Scheme 4. Transition States for the Nitro-Michael System 5



single-point calculations as our references. The systematic nature of these post-HF calculations allows us to approach the converged relative energies, so that they provide a good control for the DFT methods. As optimizations with the higher post-HF methods are too expensive, only single-point calculations were used systematically for all systems. However, we were able to optimize systems 1–4 using MP2. We expected a large influence of the optimization level on the post-HF energies, so that single-point energies were calculated for the optimized geometries obtained with each functional. The cheaper MP2 single points were carried out on all DFT double- $\zeta$  and triple- $\zeta$  geometries, while the more extensive RICCD and LCCSD(T) energy calculations were limited to DFT triple- $\zeta$  optimized geometries. All single-point calculations were carried out using augmented triple- $\zeta$  basis sets, which are sufficient to give accurate interaction energies and low basis set superposition errors (BSSE).<sup>45</sup> In the ideal case, the CCSD single-point energy differences should not vary strongly for geometries optimized at different levels. Although CCSD/TZVP energies are not definitive, they should not deviate strongly from the “correct” values, so that they provide at least a strong indication of the reliability of other techniques. The same is true for the LCCSD(T)/cc-pVTZ calculations, which are formally more accurate than the CCSD values (because they include a perturbational correction triple excitations) as long as the local-orbital approximation is applicable and accurate enough for our systems.

In order to test whether the local approximation is suitable for our purposes, The MP2 and CCSD single points can be compared with their local LMP2 and LCCSD equivalents. We can safely use the local approximations if the energies agree by significantly less than the change caused by the perturbational triples correction.

All energies discussed in this section are differences between the electronic energies of pairs of transition states optimized at the different DFT levels and are denoted  $\Delta\Delta E_{\text{TSA-TSB}}$ :

$$\Delta\Delta E_{\text{TSA-TSB}} = \Delta E_{\text{A}}^* - \Delta E_{\text{B}}^* \quad (1)$$



Table 1. Mean Electronic Energy Differences  $\Delta\Delta E_{\text{TSA-TSB}}$  with Standard Deviations (all values in kcal mol<sup>-1</sup>)<sup>a</sup>

	DFT	RI-MP2	RI-CC2	L-MP2	LCCSD	LCCSD(T)
1a–1b	0.81 ± 0.54	1.49 ± 0.11	1.50 ± 0.15	1.14 ± 0.33	0.54 ± 0.22	0.55 ± 0.25
2a–2b	3.56 ± 0.86	5.08 ± 0.25	5.36 ± 0.32	4.93 ± 0.12	4.84 ± 0.35	4.83 ± 0.33
3a–3b	0.11 ± 0.40	0.78 ± 0.53	0.68 ± 0.39	0.67 ± 0.83	0.28 ± 1.00	0.72 ± 0.89
3a–3c	1.02 ± 0.15	0.97 ± 0.04	1.35 ± 0.17	1.37 ± 0.24	0.79 ± 0.18	0.98 ± 0.11
3b–3c	1.14 ± 0.37	1.75 ± 0.54	2.03 ± 0.33	2.04 ± 0.73	1.07 ± 1.11	1.71 ± 0.93
4a–4b	3.14 ± 1.54	5.60 ± 1.07	4.87 ± 0.97	3.10 ± 0.77	4.15 ± 0.21	3.58 ± 0.24
5a–5b	3.79 ± 1.31	5.46 ± 1.01	6.05 ± 1.24			
5a–5c	3.27 ± 0.85	5.39 ± 0.89	5.77 ± 1.34			
5a–5d	0.17 ± 1.25	0.97 ± 1.24	1.71 ± 1.36			
5b–5c	7.06 ± 1.79	10.85 ± 0.55	11.82 ± 1.14			
5b–5d	3.96 ± 1.59	6.42 ± 1.05	7.76 ± 1.15			
5c–5d	3.09 ± 0.76	4.43 ± 0.61	4.06 ± 0.66			
		average standard deviation				
all data	0.95	0.66	0.77			
1–4 only	0.64	0.42	0.39	0.50	0.51	0.46

<sup>a</sup>The DFT column shows the statistics for fully optimized geometries using the different DFT methods (PBE, B3LYP, TPSS, TPSSH, w97XD, and M06-2X with the TZVP basis set). The other columns illustrate the effect of using the different DFT-optimized geometries on  $\Delta\Delta E_{\text{TSA-TSB}}$  at the given level of theory.

where  $\Delta E_A^*$  and  $\Delta E_B^*$  are the calculated Born–Oppenheimer energies for transition states A and B, respectively.

The reason for using this difference as the target value is explained above. We first compare the DFT energy differences with reference values and also the effect of using geometries optimized with different DFT functionals on the corresponding MP2 and CC2 single points. Table 1 shows the mean and the standard deviations of the  $\Delta\Delta E$  values from the mean for all pairs of comparable transition states for the DFT/MP2, DFT/CC2, and MP2/CC2 pairs. The mean values indicate how well each calculational level reproduces the values calculated at higher levels, and the standard deviations indicate the spread of the results caused by using different functionals for the DFT methods and the different DFT-optimized geometries for the post-HF calculations.

Table 1 shows that the differences in electronic energy for the pairs of transition states ( $\Delta\Delta E$ ) vary considerably between the methods considered. The variation between the DFT methods is the largest, as expected, because both the geometries and the functionals used vary. For the sets of comparable transition states, the  $\Delta\Delta E$  values vary on average by  $\pm 0.95$  kcal mol<sup>-1</sup>. For the *ab initio* single points, the deviation is still in the range of  $\pm 0.66$  (for RI-MP2) to  $\pm 0.77$  kcal mol<sup>-1</sup> for (RI-CC2). If the largest system is omitted, the deviation is smaller at  $\pm 0.39$  (for RI-CC2) to  $\pm 0.51$  kcal mol<sup>-1</sup> (for LCCSD). We can thus conclude that the DFT-optimized geometries cause non-negligible differences in the calculated single-point energies for reactivity studies on enantioselective systems. Remarkably, the deviations caused solely by the different geometries in the single-point calculations are almost as high as those found among the DFT methods themselves. Including the largest system, **5**, for which we expect the largest geometric deviations, leads to an increase in the deviation of the single-point energies. The choice of an appropriate economical level of calculation for geometry optimizations is therefore important and will be considered below.

It is known that the errors caused by the RI approximation with accurate fitting basis sets can be kept well below the basis set errors but increase linearly with the number of basis functions.<sup>46</sup> For the systems studied, MP2 optimizations without density

fitting were possible for systems **1a,b** and **2a,b**. An average difference from the RI-MP2 single points of 0.18 kcal mol<sup>-1</sup> is found, well below the standard deviation of 0.66 kcal mol<sup>-1</sup> between the RI-MP2 single-point energies (Table 1). We can therefore conclude that the error induced by the RI approximation is well below that introduced by optimizing the geometry at a more economical level of theory.

We also examined whether localized-orbital approximations introduce errors that would affect the accuracy of our predictions for the systems considered. The results are shown in Table 2 with the statistics for the comparison of RI-MP2 and MP2 calculations.

The local-RI and standard-RI results show significant deviations in the energies between LMP2 and MP2 and between LCCSD and RI-CC2. As we can see from Table 2, the deviation for the local approach is always comparable to or significantly larger than the additional perturbational triples corrections. We thus cannot expect to improve the results in our case using the local approach. The most reliable method for single-point energies in this case is thus RI-CC2.

In 11 of the 13 transition states, full optimization at the RI-MP2 level was possible. The root-mean-square difference of interatomic distances (geometric RMS) between all pairs of atoms was used to compare the geometries optimized at the different DFT levels to those optimized with MP2. As shown in Table 3, the best agreement with MP2 geometries is given by M06-2X, with average and maximum geometric RMSDs of only 0.11 Å and 0.27 Å, respectively. The next best functional in terms of the average RMSD is wB97xd (0.31 Å), but it gives the largest RMSD from the MP2 geometry of all (1.47 Å) for transition state **4b**. The same trend can be observed for energies. RI-MP2 single-point energies on DFT-optimized geometries show very good agreement with RI-MP2 optimized energies for M06-2X and wB97xd. RI-MP2 single points on optimizations with M06-2X lie within 2.05 kcal mol<sup>-1</sup> (maximum error) of the energies obtained for the RI-MP2-optimized geometries (mean error: 0.94 kcal mol<sup>-1</sup>). For pairs of transition states, the average error cancels out to only 0.16 kcal mol<sup>-1</sup>. wB97xd performs relatively well. wB97xd/TZVP//RI-MP2/TZVP single-point calculations

**Table 2.** Mean Absolute Deviation of  $\Delta\Delta E$  (kcal mol<sup>-1</sup>) for Local and Nonlocal Approaches for Pairs of Transition States Compared with the Effect of Higher Order Excitations<sup>a</sup>

	1a–1b	2a–2b	3a–3b	3a–3c	3b–3c	4a–4b	mean $\pm$ std dev.
LMP2-MP2	0.40	0.26	0.69	0.36	0.72	1.68	0.69 $\pm$ 0.52
LMP2-LCCSD(T)	0.58	0.30	0.46	0.48	0.69	0.73	0.54 $\pm$ 0.16
LCCSD-CC2	0.95	0.52	0.75	0.56	1.20	0.74	0.79 $\pm$ 0.25
LCCSD-LCCSD(T)	0.13	0.04	0.45	0.19	0.63	0.57	0.34 $\pm$ 0.25

<sup>a</sup> As for Table 1, the values given are the mean of the values calculated using all DFT levels for geometry optimization.

**Table 3.** RMS Difference of Interatomic Distances between Geometries Optimized at the DFT/TZVP and RI-MP2/TZVP Levels (all values in Å) and MP2 Single-Point Energy Differences on Both Geometries, Calculated for Individual Transition States and Pairs of Transition States

	RMS difference of interatomic distances [Å]					
	functional					
	PBE	TPSS	B3LYP	TPSSH	M06-2X	wB97xd
average	0.42	0.41	0.47	0.42	0.11	0.31
maximum	1.18	0.93	1.27	1.16	0.27	1.47
	energy deviation from the RI-MP2/TZVP: energy differences between pairs of transition states [kcal mol <sup>-1</sup> ]					
	PBE	TPSS	B3LYP	TPSSH	M06-2X	wB97xd
	average	0.29	0.48	0.77	0.86	0.16
maximum	0.90	1.75	2.72	1.72	0.92	0.36
individual TS	5.91	5.40	4.49	3.71	0.94	1.44
average						
maximum	11.26	10.54	9.17	7.52	2.05	2.82

**Table 4.** Mean Absolute Deviation of  $\Delta\Delta E$  Values (kcal mol<sup>-1</sup>) Calculated with MP2 and CC2 at the DFT-Optimized Geometries from Those Optimized with RI-MP2

	$\Delta\Delta E_{MP2} - \Delta\Delta E_{CC2}$					
	geometry					
	PBE	TPSS	B3LYP	TPSSH	M06-2X	wB97xd
mean absolute deviation	1.07	1.03	1.18	1.08	0.19	0.19
maximum deviation	3.60	2.97	3.15	2.71	0.54	0.51

show an average deviation in total energy from fully optimized RI-MP2/TZVP calculations of 1.44 kcal mol<sup>-1</sup>, which cancels out to 0.20 kcal mol<sup>-1</sup> for the  $\Delta\Delta E_{TSA-TSB}$  energy differences between pairs of transition states.

Table 4 shows that the differences  $\Delta\Delta E$  between pairs of transition states for MP2 and CC2 are also smallest for the geometries optimized with M06-2X (mean value: 0.19 kcal mol<sup>-1</sup>). The same value is found for the geometries optimized with the wB97xd functional. For all other DFT methods, both values increase rapidly. In combination with the excellent agreement of M06-2X geometries with MP2 geometries, these results suggest that geometries optimized with M06-2X are most suitable for higher-level single-point calculations. We have therefore used CC2/TZVP//M06-2X/TZVP energies as the reference for assessing the performance of further methods.

**Semiempirical MO Techniques.** Semiempirical MO techniques can handle very large systems easily<sup>47</sup> and would therefore

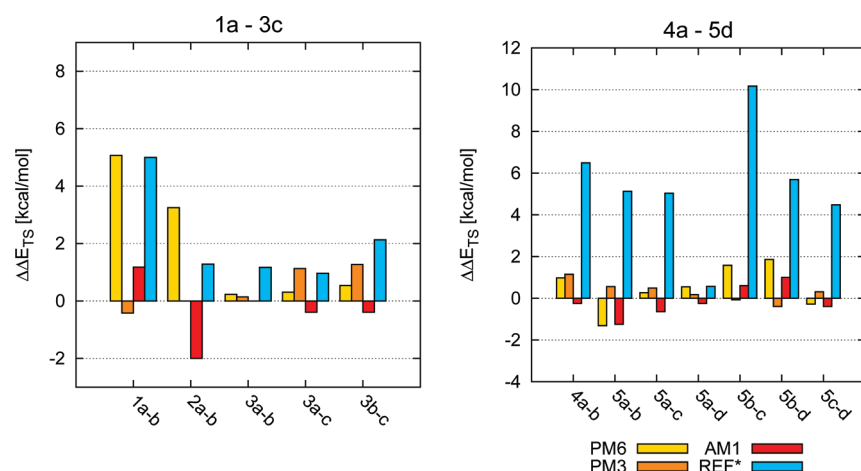
be extremely useful for fast scans to provide guidance for experimental studies if the results were reliable enough. Chart 1 shows the results obtained for the semiempirical MO techniques tested.

The results for the different systems vary for the semiempirical methods. While at least qualitatively correct descriptions are possible for the two nitro-Michael reactions, the case for the two aldol reactions is less encouraging. In contrast to a synchronous formation of the new C–C and O–H bonds found in the DFT optimizations, AM1 and PM3 find the reaction proceeds via a multistep path with separated C–C bond formation and proton-transfer steps, as can be shown by relaxed PES scans along the newly formed C–C and O–H bonds. The energies shown in Chart 1 and listed in Table S1 of the Supporting Information suggest that the qualitative prediction for the selectivity is often, but not reliably and consistently, correct.

Of the semiempirical methods tested, only PM6 gives results that are qualitatively consistent with the DFT calculations for the small systems 1–3. For the larger systems, none of the semiempirical methods give reliable energies. Remarkably, the geometries optimized with PM6 are quite accurate in some cases. The RMSD of all interatomic distances relative to the MP2-optimized geometries is 0.42 Å for 4a and 0.61 Å for 4b, comparable to or better than most DFT optimizations for these systems (see Table 6 for DFT results).

## ■ DFT-METHODS

**Energies.** The results shown above suggest that some DFT methods can provide results that are consistent with post-HF ab

Chart 1.  $\Delta\Delta E$  (kcal mol<sup>-1</sup>) Calculated with Semiempirical Molecular Orbital Theory for All Pairs of Transition States<sup>a</sup>

<sup>a</sup>Left: aldol reactions. Right: nitro-Michael reactions. AM1 and PM3 results are estimated from relaxed PES scans because these methods predict multistep reactions. Reference energies calculated at the CC2/TZVP//M06-2X/TZVP level.

Table 5. Mean Absolute Deviation of  $\Delta\Delta E$  (kcal mol<sup>-1</sup>), Calculated at DFT and MP2 Levels Based on DFT Geometries, from the Reference Level (RICC2/TZVP//M06-2X/TZVP) for All Pairs of Transition States

geometry optimization	energy calculation	PBE	TPSS	B3LYP	TPSSH	M06-2X	wB97XD
all systems							
DFT/6-31G(d)	RI-MP2/TZVP	0.83	1.12	1.01	1.47	0.72	0.44
	DFT/6-31G(d)	1.45	1.38	1.48	1.41	1.28	0.86
DFT/TZVP	RI-MP2/TZVP	0.68	0.67	1.02	0.99	0.20	0.74
	DFT/TZVP	1.88	2.20	2.32	2.15	1.06	0.75
aldol reactions (1–3)							
DFT/6-31G(d)	RI-MP2/TZVP	0.22	0.67	0.33	1.69	0.06	0.31
	DFT/6-31G(d)	1.02	1.10	1.21	1.14	0.37	0.47
DFT/TZVP	RI-MP2/TZVP	0.19	0.3	0.30	0.99	0.07	0.20
	DFT/TZVP	1.18	1.37	1.63	1.36	0.36	0.53
nitro-Michael reactions (4–5)							
DFT/6-31G(d)	RI-MP2/TZVP	1.18	1.38	1.39	1.35	1.10	0.50
	DFT/6-31G(d)	1.70	1.54	1.63	1.57	1.80	1.08
DFT/TZVP	RI-MP2/TZVP	0.95	0.88	1.43	0.99	0.26	1.06
	DFT/TZVP	2.28	2.68	2.72	2.60	1.46	0.88

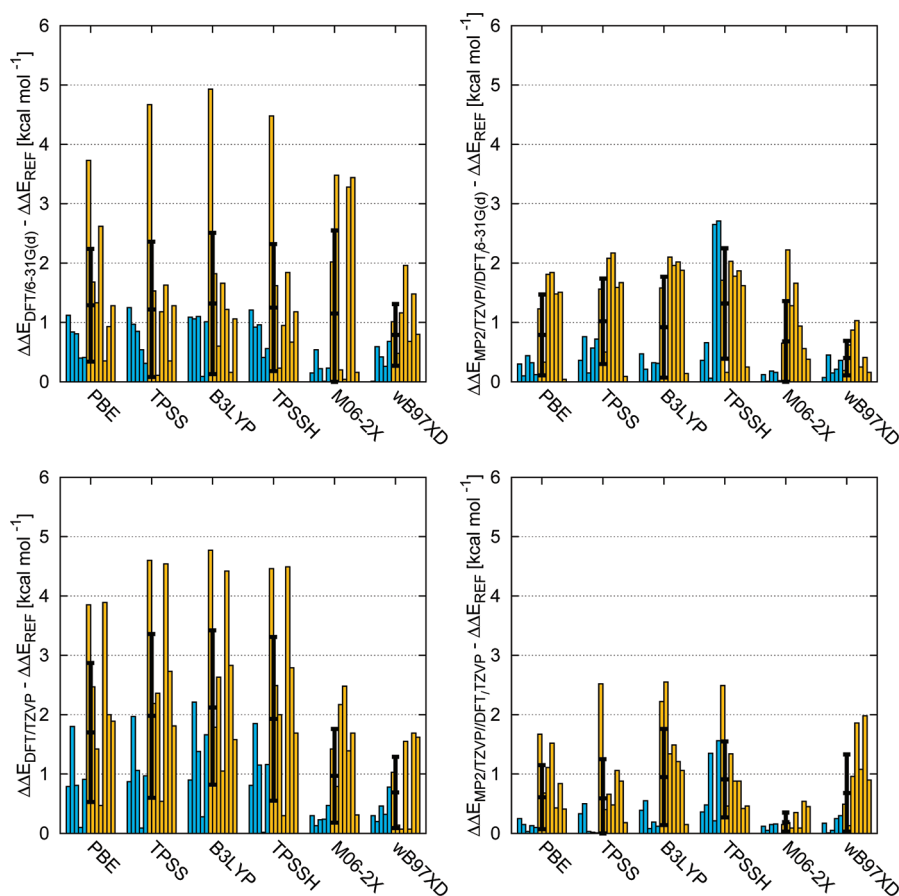
initio data. A question of practical relevance is how to achieve quantitatively correct results at the least computational expense. In this section, we investigate the performance of less computationally expensive methods relative to reference (CC2/TZVP//M06-2X/TZVP) data. We present the results of DFT optimizations with double- and triple- $\zeta$  basis sets and MP2/TZVP single points on the DFT-optimized geometries. The calculational protocol can be simplified by calculating MP2 in place of CC2 single points because MP2 is approximately 10 times faster than CC2 for the cases tested. A second possibility is to use the DFT results directly, which avoids the post-HF single point calculations, which are particularly expensive for large systems. If a further increase in efficiency is necessary for the largest systems, either smaller basis sets or less computationally expensive levels of DFT can be used. In particular, pure GGA and meta-GGA calculations can benefit from density fitting in RI approaches, which can lead to an acceleration of up to a factor of ten<sup>48</sup> (>30 in combination with multipole approximations<sup>49</sup>).

Reducing the level of the single-point calculations from CC2 to MP2 hardly affects the results, as shown in Table 5.

Table 5 and Chart 2 (bottom right graph) show that MP2/TZVP//M06-2X/TZVP calculations have an average deviation from the reference energies of only 0.20 kcal mol<sup>-1</sup> (with only one case with a deviation larger than 0.50 kcal mol<sup>-1</sup>, data shown in the Supporting Information). Using geometries optimized at other DFT levels leads to a significant degradation in performance, giving both larger average deviations and strong outliers. Only for the small systems (1–3) is PBE an economical alternative to M06-2X, with an average deviation of 0.19 kcal mol<sup>-1</sup>.

Surprisingly, using a different basis set hardly changes the situation. We expected less accurate energies for single-point calculations on DFT/6-31G(d) optimized geometries. The errors calculated for MP2 single points on DFT-optimized geometries with the TZVP (Chart 2, bottom-right) and 6-31G(d) (Chart 2, top-right) basis sets show that this is true in most

**Chart 2.** Absolute Deviation and Mean Absolute Deviation ( $\pm$  one standard deviation) of  $\Delta\Delta E$  (kcal mol $^{-1}$ ) Calculated at DFT and MP2 Levels on the Basis of DFT Geometries, from the Reference System (RICC2/TZVP//M06-2X/TZVP) for All Pairs of Transition States (blue, small systems; yellow, large systems)



cases. However, the smaller basis set gives lower errors (0.44 kcal mol $^{-1}$ ) for wb97xd. These are even better than those found with M06-2X geometries optimized with the smaller basis set (0.72 kcal mol $^{-1}$ ).

The computational protocol can be simplified further for very large systems by using DFT energies. Chart 2 shows that using DFT energies, rather than those from MP2 single points, leads to a significant increase in the deviations from the reference energies. The average deviation is larger than 1 kcal mol $^{-1}$  for all functionals except wb97xd, for which an error of 0.86 kcal mol $^{-1}$  is found with the double- $\zeta$  and 0.75 kcal mol $^{-1}$  with the triple- $\zeta$  basis set. The M06-2X functional performs marginally less well for energies (mean error = 1.06 kcal mol $^{-1}$  with the TZVP basis set) but might be preferred because it gives better geometries for MP2 single points.

**Geometries.** For the conventional DFT methods (PBE, B3LYP, TPSS, and TPSSH), energy differences calculated with the double and triple- $\zeta$  basis sets do not show any improvement for the triple- $\zeta$  basis, but rather the opposite. In order to explain these findings, the similarity of the geometries optimized at the different DFT level was studied. The results are shown in Table 6.

For the small systems 1 and 2, PBE, TPSS, B3LYP, and TPSSH give consistently better results with the smaller basis set. The opposite is true for M06-2X, and wb97xd performs similarly with the two basis sets. The results are more mixed for the largest

aldol reaction 3, although on balance, optimizations with the larger basis set are slightly better. The larger nitro-Michael reactions 4 and 5 give consistently better results with the larger basis set for all functionals except wb97xd, which performs similarly with the DZ and TZ basis sets. Statistically, M06-2X performs best, followed by wb97xd, as outlined above. However, both, but especially the latter, show a significant degradation in performance as the size of the system increases. Nevertheless, the two functionals that include dispersion implicitly or explicitly perform better for large systems than the others. Dispersion is the probable cause of this effect, although it is possible that MP2/TZVP is overestimating dispersion and therefore contributing to the error.

## COMPUTATIONAL METHODS

Semiempirical calculations were carried out using VAMP 10.0.<sup>50</sup> Transition states were characterized by calculating the normal vibrations within the harmonic approximation. Restricted potential energy surface scans for the systems 1–3 were carried out using MOPAC09<sup>51</sup> using the distance between the pairs of atoms defining the two newly formed covalent bonds as fixed reaction coordinates in steps of 0.1 Å.

DFT and HF calculations were carried out using Gaussian 09.<sup>52</sup> All transition states were fully optimized using the PBE,<sup>23</sup> TPSS,<sup>24</sup> B3LYP,<sup>13</sup> TPSSH,<sup>24</sup> M06-2X,<sup>14</sup> and wb97xd<sup>25</sup> functionals

**Table 6. Comparison of the DFT-Optimized Geometries with MP2/TZVP Optimized Geometries<sup>a</sup>**

functional	RMSD [Å] of all interatomic distances					
	PBE	TPSS	B3LYP	TPSSH	M06-2X	wB97xd
basis set	DZ/TZ	DZ/TZ	DZ/TZ	DZ/TZ	DZ/TZ	DZ/TZ
aldol reactions (small systems)						
1a	0.09/0.12	0.09/0.12	0.10/0.14	0.09/0.11	0.10/0.07	0.09/0.10
1b	0.15/0.17	0.17/0.19	0.18/0.2	0.16/0.18	0.06/0.05	0.14/0.14
2a	0.11/0.17	0.14/0.23	0.16/0.24	0.12/0.21	0.07/0.05	0.06/0.05
2b	0.10/0.11	0.09/0.10	0.10/0.11	0.08/0.09	0.04/0.03	0.05/0.04
3a	0.20/0.20	0.21/0.19	0.22/0.19	0.21/0.19	0.15/0.04	0.10/0.09
3b	0.15/0.13	0.17/0.15	0.15/0.18	0.16/0.13	0.05/0.07	0.07/0.07
3c	0.16/0.15	0.18/0.18	0.17/0.15	0.15/0.12	0.08/0.06	0.07/0.06
nitro-Michael reactions (large systems)						
4a	0.76/0.61	0.78/0.62	0.71/0.63	0.77/0.62	0.09/0.14	0.25/0.27
4b	0.64/0.66	0.66/0.68	0.72/0.75	0.65/0.67	0.33/0.22	0.41/0.67
5a	1.25/1.18	1.26/0.93	1.27/0.68	1.27/1.16	0.56/0.27	0.49/0.45
5b	0.99/0.88	1.03/0.89	1.09/0.90	1.09/0.97	1.03/0.23	1.69/1.47
mean	0.42/0.40	0.43/0.39	0.44/0.38	0.43/0.40	0.23/0.11	0.31/0.31
std. dev.	0.42/0.37	0.42/0.32	0.43/0.30	0.44/0.39	0.31/0.09	0.48/0.43

<sup>a</sup>RMSD is calculated as the root-mean square deviation between DFT and MP2 optimized structures for all interatomic distances.

in combination with the 6-31G(d)<sup>27,28</sup> and TZVP<sup>29</sup> basis sets. The transition states were characterized by calculating the normal vibrations within the harmonic approximation. All pure GGA functionals with TZVP basis set were calculated using the density-fitting approach as implemented in Gaussian 09.

RI-MP2 optimizations and RI post-HF single-point calculations were carried out with Turbomole 6.2.<sup>53</sup> All-electron MP2 and CCSD calculations used the RI (resolution of identity) approximated approaches RI-MP2 and RI-CC2 with the (def2-)TZVP<sup>54</sup> basis sets.

All localized correlation methods were carried out with MOLPRO 2010.1<sup>55</sup> using the Dunning cc-pTZV<sup>56</sup> basis set.

The geometrical RMSD between two structures was calculated as the unweighted RMS deviation of all interatomic distances.

## CONCLUSIONS

Semiempirical MO theory proved not to be suitable for calculating  $\Delta\Delta E$  values, as neither energies nor geometries are accurate enough to describe the reactions correctly.

In comparison, DFT methods perform better. The DFT results are qualitatively correctly rank-ordered for all cases examined, except for those with energy differences lower than 1 kcal mol<sup>-1</sup>. However, large quantitative differences exist. Only the wB97xd and M062X functionals, which consider dispersion either implicitly or explicitly, gave acceptable energies quantitatively consistent with MP2 and CCSD single points, whereby M06-2X performs best for our test set. CC2/TZVP//M06-2X/TZVP calculations were used as the most reliable reference calculations because M06-2X gives geometries that resemble those obtained with the post-HF ab initio techniques closely.

The gas-phase energy differences for single-point calculations converge to the correct result at the MP2 level of theory; additional explicit correlation hardly improves the results.

BSSE effects can be observed when the smaller (6-31G(d)) basis set is used. Remarkably, using the larger TZVP basis set only improved the results for the modern intrinsic or explicitly dispersion-corrected functionals. For the conventional DFT methods, more accurate energies are found for the smaller basis set. For these functionals, the BSSE contribution apparently cancels part of the missing dispersion energy. Conventional DFT methods give good geometries for small systems, but their performance degrades significantly with increasing system size. This trend is also found with M06-2X and wB97xd but is less pronounced, especially for M06-2X.

We were able to show that geometries optimized with the M06-2X density functional with a triple- $\zeta$  basis set can be used in subsequent post-HF single-point calculations to give very accurate gas-phase energies. Our work complements that of Friesner et al.<sup>10</sup> on the best techniques for calculating solvent effects and of Simon and Goodman<sup>11</sup> on accurate free-energy corrections to provide an accurate and economical calculational protocol for predicting the results of kinetically controlled organocatalytic reactions.

We therefore recommend RICC2/TZVP//M06-2X/TZVP, RI-MP2/TZVP//M06-2X/TZVP, and wB97xd/TZVP or M06-2X/TZVP, depending on the size of the system, for calculating energy differences between alternative transition states in stereoselective organocatalytic reactions.

## ASSOCIATED CONTENT

**S Supporting Information.** Detailed plot of optimized gas phase energies at the DFT MP2/DFT and CC2/DFT level (Chart S1), Table S1 containing semiempirical transition state energies plotted in Chart 1, gas-phase energies of the individual structures (Table S2–S10) and energy differences between pairs of transition states (Table S11), and geometric RMS and C–C bond length formation data (Table S12). This information is available free of charge via the Internet at <http://pubs.acs.org>.

## AUTHOR INFORMATION

### Corresponding Author

\*E-mail: [Tim.clark@chemie.uni-erlangen.de](mailto:Tim.clark@chemie.uni-erlangen.de).

## ACKNOWLEDGMENT

We thank the *Deutsche Forschungsgemeinschaft* (SPP 1179) for financial support and the *Leibnitz Rechenzentrum Munich* for computational time and technical support.

## REFERENCES

- (1) Shinisha, C.; Sunoj, R. Bicyclic proline analogues as organocatalysts for stereoselective aldol reactions: an in silico DFT study. *Org. Biomol. Chem.* **2007**, *5*, 1287–1294.
- (2) Houk, K.; Cheong, P. Computational prediction of small-molecule catalysts. *Nature* **2008**, *455*, 309–313.
- (3) Grimme, S. Accurate description of van der Waals complexes by density functional theory including empirical corrections. *J. Comput. Chem.* **2004**, *25*, 1463–1473.
- (4) Muller-Dethlefs, K.; Hobza, P. Noncovalent interactions: a challenge for experiment and theory. *Chem. Rev.* **2000**, *100*, 143–168.
- (5) Morgado, C.; Jurek, P.; Svozil, D.; Hobza, P.; Šponer, J. Reference MP2/CBS and CCSD(T) quantum-chemical calculations on stacked adenine dimers. Comparison with DFT-D, MP2. *5*, SCS

- (MI)-MP2, M06-2X, CBS (SCS-D) and force field descriptions. *Phys. Chem. Chem. Phys.* **2010**, *12*, 3522–3534.
- (6) Zhao, Y.; Truhlar, D. G. Density Functional Calculations of E2 and SN2 Reactions: Effects of the Choice of Density Functional, Basis Set, and Self-Consistent Iterations. *J. Chem. Theory Comput.* **2010**, *6*, 1104–1108.
- (7) Wheeler, S.; Moran, A.; Pieniazek, S.; Houk, K. Accurate Reaction Enthalpies and Sources of Error in DFT Thermochemistry for Aldol, Mannich, and -Aminoxylation Reactions. *J. Chem. Phys. A* **2009**, *113*, 10376–10384.
- (8) Zhao, Y.; González-García, N.; Truhlar, D. G. Benchmark Database of Barrier Heights for Heavy Atom Transfer, Nucleophilic Substitution, Association, and Unimolecular Reactions and Its Use to Test Theoretical Methods. *J. Chem. Phys. A* **2005**, *109*, 2012–2018.
- (9) Jiang, J.; Wu, Y.; Wang, Z.; Wu, C. Assessing the Performance of Popular Quantum Mechanics and Molecular Mechanics Methods and Revealing the Sequence-Dependent Energetic Features Using 100 Tetrapeptide Models. *J. Chem. Theory Comput.* **2010**, *6*, 1199–1209.
- (10) Schneebeli, S. T.; Hall, M. L.; Breslow, R.; Friesner, R. Quantitative DFT Modeling of the Enantiomeric Excess for Dioxirane-Catalyzed Epoxidations. *J. Am. Chem. Soc.* **2009**, *131*, 3965–3973.
- (11) Simon, L.; Goodman, J. M. How reliable are DFT transition structures? Comparison of GGA, hybrid-meta-GGA and meta-GGA functionals. *Org. Biomol. Chem.* **2011**, *9*, 689–700.
- (12) Sousa, S.; Fernandes, P.; Ramos, M. General performance of density functionals. *J. Chem. Phys. A* **2007**, *111*, 10439–10452.
- (13) Becke, A. Density-functional thermochemistry. III. The role of exact exchange. *Chem. Phys.* **1993**, *98*, 5648–5652.
- (14) Zhao, Y.; Truhlar, D. 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 functionals. *Theor. Chem. Acc.* **2008**, *120*, 215–241.
- (15) Bahmanyar, S.; Houk, K. The origin of stereoselectivity in proline-catalyzed intramolecular aldol reactions. *J. Am. Chem. Soc.* **2001**, *123*, 12911–12912.
- (16) Corrêa, R.; Garden, S.; Angelici, G.; Tomasini, C. A DFT and AIM Study of the Proline-Catalyzed Asymmetric Cross-Aldol Addition of Acetone to Isatins: A Rationalization for the Reversal of Chirality. *Eur. J. Chem. Chem.* **2008**, *2008*, 736–744.
- (17) Hamza, A.; Schubert, G.; Soós, T.; Pápai, I. Theoretical Studies on the Bifunctionality of Chiral Thiourea-Based Organocatalysts: Competing Routes to C–C Bond Formation. *J. Am. Chem. Soc.* **2006**, *128*, 13151–13160.
- (18) Yalalov, D.; Tsogoeva, S.; Schmatz, S. Chiral thiourea-based bifunctional organocatalysts in the asymmetric Nitro-Michael addition: a joint experimental-theoretical study. *Adv. Synth. Catal.* **2006**, *348*, 826–832.
- (19) Dewar, M.; Zebisch, E.; Healy, E.; Stewart, J. Development and use of quantum mechanical molecular models. 76. AM1: a new general purpose quantum mechanical molecular model. *J. Am. Chem. Soc.* **1985**, *107*, 3902–3909.
- (20) Stewart, J. Optimization of parameters for semiempirical methods I. Method. *J. Comput. Chem.* **1989**, *10*, 209–220.
- (21) Stewart, J. Optimization of parameters for semiempirical methods II. Applications. *J. Comput. Chem.* **1989**, *10*, 221–264.
- (22) Stewart, J. Optimization of parameters for semiempirical methods V: modification of NDDO approximations and application to 70 elements. *J. Mol. Model.* **2007**, *13*, 1173–1213.
- (23) Perdew, J.; Burke, K.; Ernzerhof, M. Generalized gradient approximation made simple. *Phys. Rev. Lett.* **1996**, *77*, 3865–3868.
- (24) Tao, J.; Perdew, J.; Staroverov, V.; Scuseria, G. Climbing the density functional ladder: Nonempirical meta-generalized gradient approximation designed for molecules and solids. *Phys. Rev. Lett.* **2003**, *91*, 146401.
- (25) Chai, J.; Head-Gordon, M. Long-range corrected hybrid density functionals with damped atom-atom dispersion corrections. *Phys. Chem. Chem. Phys.* **2008**, *10*, 6615–6620.
- (26) Xu, X.; Alecu, I. M.; Truhlar, D. G. How Well Can Modern Density Functionals Predict Internuclear Distances at Transition States? *J. Chem. Theory Comput.* **2011**, *7*, 1667–1676.
- (27) Rassolov, V.; Ratner, M.; Pople, J.; Redfern, P.; Curtiss, L. 6-31G\* basis set for third row atoms. *J. Comput. Chem.* **2001**, *22*, 976–984.
- (28) Rassolov, V.; Pople, J.; Ratner, M.; Windus, T. 6-31G basis set for atoms K through Zn. *J. Chem. Phys.* **1998**, *109*, 1223–1229.
- (29) Schäfer, A.; Huber, C.; Ahlrichs, R. Fully optimized contracted Gaussian basis sets of triple zeta valence quality for atoms Li to Kr. *J. Chem. Phys.* **1994**, *100*, 5829–5835.
- (30) Møller, C.; Plesset, M. Note on an approximation treatment for many-electron systems. *Phys. Rev.* **1934**, *46*, 618–622.
- (31) Head-Gordon, M.; Pople, J.; Frisch, M. MP2 energy evaluation by direct methods. *Chem. Phys. Lett.* **1988**, *153*, 503–506.
- (32) Cizek, J.; Paldus, J. Coupled Cluster Approach. *Phys. Scr.* **1980**, *21*, 251–254.
- (33) Cramer, C. *Essentials of computational chemistry: theories and models*; John Wiley & Sons Inc: New York, 2004; pp 203–248.
- (34) Weigend, F.; Häser, M. RI-MP2: first derivatives and global consistency. *Theor. Chem. Acc.* **1997**, *97*, 331–340.
- (35) Hättig, C.; Weigend, F. CC2 excitation energy calculations on large molecules using the resolution of the identity approximation. *J. Chem. Phys.* **2000**, *113*, 5154–5161.
- (36) Christiansen, O.; Koch, H.; Jørgensen, P. The second-order approximate coupled cluster singles and doubles model CC2. *Chem. Phys. Lett.* **1995**, *243*, 409–418.
- (37) Head-Gordon, M. In 10th American Conference on Theoretical Chemistry, Boulder, CO, 1999.
- (38) Schütz, M. Low-order scaling local electron correlation methods. V. Connected triples beyond (T): Linear scaling local CCSDT-1b. *J. Chem. Phys.* **2002**, *116*, 8772–8785.
- (39) Schütz, M.; Werner, H. Local perturbative triples correction (T) with linear cost scaling. *Chem. Phys. Lett.* **2000**, *318*, 370–378.
- (40) Schütz, M.; Werner, H. Low-order scaling local electron correlation methods. IV. Linear scaling local coupled-cluster (LCCSD). *J. Chem. Phys.* **2001**, *114*, 661–681.
- (41) Schütz, M. A new, fast, semi-direct implementation of linear scaling local coupled cluster theory. *Phys. Chem. Chem. Phys.* **2002**, *4*, 3941–3947.
- (42) Schütz, M. Low-order scaling local electron correlation methods. III. Linear scaling local perturbative triples correction (T). *J. Chem. Phys.* **2000**, *113*, 9986–10001.
- (43) Hampel, C.; Werner, H. Local treatment of electron correlation in coupled cluster theory. *J. Chem. Phys.* **1996**, *104*, 6286–6297.
- (44) Kendall, R.; Dunning, T., Jr.; Harrison, R. Electron affinities of the first row atoms revisited. Systematic basis sets and wave functions. *J. Chem. Phys.* **1992**, *96*, 6796–6806.
- (45) Jurečka, P.; Černý, J.; Hobza, P.; Salahub, D. R. Density functional theory augmented with an empirical dispersion term. Interaction energies and geometries of 80 noncovalent complexes compared with ab initio quantum mechanics calculations. *J. Comput. Chem.* **2007**, *28*, 555–569.
- (46) Skylaris, C. K.; Gagliardi, L.; Handy, N. C.; Ioannou, A. G.; Spencer, S.; Willetts, A. On the resolution of identity Coulomb energy approximation in density functional theory. *THEOCHEM* **2000**, *501–502*, 229–239.
- (47) Clark, T.; Stewart, J. J. P., MNDO-like Semiempirical Molecular Orbital Theory and its Application to Large Systems. In *Computational Methods for Large Systems*; Reimers, J. J., Ed.; Wiley: Chichester, U. K., 2011; pp 259–286.
- (48) Sundholm, D. Density functional theory calculations of the visible spectrum of chlorophyll a. *Chem. Phys. Lett.* **1999**, *302*, 480–484.
- (49) Sierka, M.; Hogekamp, A.; Ahlrichs, R. Fast evaluation of the Coulomb potential for electron densities using multipole accelerated resolution of identity approximation. *J. Chem. Phys.* **2003**, *118*, 9136–9148.
- (50) Clark, T.; Alex, A.; Beck, B.; Burckhardt, F.; Chandrasekhar, J.; Gedeck, P.; Horn, A.; Hutter, M.; Martin, B.; Rauhut, G.; Sauer, W.;

Schindler, T.; Steinke, T. *VAMP*, 10.0; Accelrys Inc.: San Diego, CA, 2007.

(51) Stewart, J. J. P. *MOPAC2009*; Stewart Computational Chemistry: Colorado Springs, CO, 2008.

(52) Frisch, M. J.; Schlegel, H. B.; Scuseria, G. E.; Robb, M. A.; Cheeseman, J. R.; Scalmani, G.; Barone, V.; Mennucci, B.; Petersson, G. A.; Nakatsuji, H.; Caricato, M.; Li, X.; Hratchian, H. P.; Izmaylov, A. F.; Bloino, J.; Zheng, G.; Sonnenberg, J. L.; Hada, M.; Ehara, M.; Toyota, K.; Fukuda, R.; Hasegawa, J.; Ishida, M.; Nakajima, T.; Honda, Y.; Kitao, O.; Nakai, H.; Vreven, T.; Montgomery, J. A., Jr.; Peralta, J. E.; Ogliaro, F.; Bearpark, M.; Heyd, J. J.; Brothers, E.; Kudin, K. N.; Staroverov, V. N.; Kobayashi, R.; Normand, J.; Raghavachari, K.; Rendell, A.; Burant, J. C.; Iyengar, S. S.; Tomasi, J.; Cossi, M.; Rega, N.; Millam, J. M.; Klene, M.; Knox, J. E.; Cross, J. B.; Bakken, V.; Adamo, C.; Jaramillo, J.; Gomperts, R.; Stratmann, R. E.; Yazyev, O.; Austin, A. J.; Cammi, R.; Pomelli, C.; Ochterski, J. W.; Martin, R. L.; Morokuma, K.; Zakrzewski, V. G.; Voth, G. A.; Salvador, P.; Dannenberg, J. J.; Dapprich, S.; Daniels, A. D.; Ö. Farkas, Foresman, J. B.; Ortiz, J. V.; Cioslowski, J.; Fox, D. J. *Gaussian 09*, Revision A.02; Gaussian, Inc.: Wallingford, CT, 2009.

(53) *TURBOMOLE*, V6.2; TURBOMOLE GmbH: Karlsruhe, Germany, 2010.

(54) Weigend, F.; Ahlrichs, R. Balanced basis sets of split valence, triple zeta valence and quadruple zeta valence quality for H to Rn: Design and assessment of accuracy. *Phys. Chem. Chem. Phys.* **2005**, *7*, 3297–3305.

(55) Werner, H.-J.; Knowles, P. J.; Manby, F. R.; Schütz, M.; Celani, P.; Knizia, G.; Korona, T.; Lindh, R.; Mitrushenkov, A.; Rauhut, G.; Adler, T. B.; Amos, R. D.; Bernhardsson, A.; Berning, A.; Cooper, D. L.; Deegan, M. J. O.; Dobbyn, A. J.; Eckert, F.; Goll, E.; Hampel, C.; Hesselmann, A.; Hetzer, G.; Hrenar, T.; Jansen, G.; Köppl, C.; Liu, Y.; Lloyd, A. W.; Mata, R. A.; May, A. J.; McNicholas, S. J.; Meyer, W.; Mura, M. E.; Nicklass, A.; Palmieri, P.; Pflüger, K.; Pitzer, R.; Reiher, M.; Shiozaki, T.; Stoll, H.; Stone, A. J.; Tarroni, R.; Thorsteinsson, T.; Wang, M.; Wolf, A. *MOLPRO*, 2010.1; Cardiff University: Cardiff, Wales, 2010.

(56) Dunning, T. H. Correlation Consistent Basis Sets. *J. Chem. Phys.* **1989**, *90*, 1007–1023.

# Time-Reversible Velocity Predictors for Verlet Integration with Velocity-Dependent Right-Hand Side

Jiří Kolafa<sup>\*,†</sup> and Martin Lísal<sup>†,§</sup>

<sup>†</sup>Department of Physical Chemistry, Institute of Chemical Technology, Prague, Technická 5, 166 28 Praha 6, Czech Republic

<sup>‡</sup>E. Hála Laboratory of Thermodynamics, Institute of Chemical Process Fundamentals of the ASCR, v. v. i., 165 02 Prague 6, Czech Republic

<sup>§</sup>Department of Physics, J. E. Purkinje University, 400 96 Ústí n. Lab., Czech Republic

 Supporting Information

**ABSTRACT:** Time-reversible velocity predictors (TRVPs) with increasing orders of the time-reversibility error are developed to be used with the Verlet integrator for equations of motion with the right-hand side depending on velocities. The method performs outside a possible SHAKE algorithm to constrain bond lengths and does not require repeated SHAKE iterations nor RATTLE. We have tested the TRVPs with the Nosé–Hoover thermostat on four model systems (coupled harmonic and anharmonic oscillators, liquid argon, SPC/E water, and a small peptide), comparing them to the Gear integrator with the Lagrangian formulation of constraint dynamics, the Martyna, Tuckerman, Tobias, and Klein (MTTK) method, and the velocity iteration method. The TRVP method performs similarly to the iteration method. In addition, we discuss three methodology improvements: (i) We tested several formulas for the kinetic energy compatible with the Verlet/SHAKE algorithm and found that the leapfrog velocities are usually the best; (ii) we proposed two modifications of the MTTK method; and (iii) we suggest that thermostats directly controlling the translational kinetic temperature may give more accurate values of some thermodynamic quantities.

## 1. INTRODUCTION

The Newton equations of motion for atomistic systems relate accelerations to forces which are functions of positions only, not velocities. Although many smart methods have been proposed, a great deal of existing molecular dynamics (MD) code relies on the simplest Verlet integrator.<sup>1,2</sup>

Extended Lagrangian methods, as the Nosé–Hoover thermostat and Andersen barostat, add a velocity-dependent term to the equations of motion; hence, using the Verlet integrator and its clones (leapfrog, velocity Verlet, Beeman) directly is no longer possible because the velocities at time  $t$  are known after the forces at time  $t$  have been evaluated. Alternative integration methods include the predictor–corrector methods,<sup>3</sup> of these the Gear integrators<sup>4</sup> are most popular. The integrators, based on the Trotter decomposition of the Liouville operator,<sup>5</sup> may be viewed as an extension of the Verlet integrator and may be combined with SHAKE and RATTLE. If one wants to adhere to the Verlet scheme, either iterations can be used to obtain the velocities or some approximation of these.<sup>6</sup>

The time reversibility error leading to a drift in the total energy (Hamiltonian) which should be conserved is worst in the Gear methods. The MTTK method<sup>5</sup> (see Section 2.7.3) is time reversible, although not symplectic: There is no drift in the total energy, but the mean quadratic error grows as the square root of time. In the iteration methods, there is a small drift decreasing with an increasing number of iterations.

In this paper we propose a velocity predictor so that iterations can be avoided. The predictor is of the second order (as the Verlet method), and the main requirement for its construction is time reversibility. During extensive testing of the methods, we found several improvements of the simulation methodology.

## 2. THEORY

**2.1. Notation and Kinetic Temperature.** Let us consider a system of  $N$  atoms with masses  $m_i$ , positions described by vectors  $\vec{r}_i$ , and velocities by  $\dot{\vec{r}}_i$ ,  $i = 1, \dots, N$ . The force acting on particle  $i$  is denoted as  $\vec{f}_i$ . The kinetic temperature is then defined by formula:

$$T_{\text{kin}} = \frac{1}{fk} \sum_{i=1}^N m_i \dot{\vec{r}}_i^2 \quad (1)$$

where  $k$  is the Boltzmann constant, and  $f = Nd + f_{\xi} - f_c$  is the number of degrees of freedom. In this formula,  $d$  denotes the space dimensionality,  $f_{\xi}$  the number of additional degrees of freedom with quadratic term for the kinetic energy, and  $f_c$  denotes the number of constraints, typically bond lengths and conserved quantities (momentum, angular momentum, and also total energy).

**2.2. Nosé–Hoover Thermostat.** We use the Nosé–Hoover thermostat as a model of equations of motion with the right-hand side containing velocities. The equations of motion of this system at temperature  $T$  are<sup>7,8</sup>

$$\ddot{\vec{r}}_i = \frac{\vec{f}_i}{m_i} - \dot{\vec{r}}_i \dot{\xi} \quad (2)$$

$$\ddot{\xi} = \frac{1}{\tau^2} \left( \frac{T_{\text{kin}}}{T} - 1 \right) \quad (3)$$

**Received:** February 15, 2011

**Published:** August 31, 2011



Here,  $\xi$  is the additional dynamic variable and  $\tau$  the typical correlation time of the thermostat. The particle accelerations depend on all velocities (including  $\xi$  which we consider a velocity, although it is not a velocity in the original Lagrangian formalism<sup>7</sup>), whereas the acceleration of the additional variable depends on the real degrees of freedom only. The equations of the Andersen barostat<sup>1,2</sup> (not considered here) have the same structure.

During integration the following total energy (derived from the Hamiltonian<sup>7</sup>) is conserved

$$E_{\text{NH}} = E_{\text{kin}} + E_{\text{pot}} + fkT \left( \xi + \frac{\tau^2 \dot{\xi}^2}{2} \right) \quad (4)$$

where  $E_{\text{kin}}$  is the kinetic energy and  $E_{\text{pot}}$  the potential (configurational) energy.

**2.3. Verlet Method.** The Verlet method for eqs 2 and 3 is respectively

$$\vec{r}_i(t+h) = 2\vec{r}_i(t) - \vec{r}_i(t-h) + \left[ \frac{\vec{f}_i(t)}{m_i} - \vec{r}_i(t)\dot{\xi}(t) \right] h^2 \quad (5)$$

$$\xi(t+h) = 2\xi(t) - \xi(t-h) + \frac{1}{\tau^2} \left[ \frac{T_{\text{kin}}(t)}{T} - 1 \right] h^2 \quad (6)$$

where  $h$  is the time step. Various approximations for the unknown velocities  $\vec{r}_i(t)$  and  $\xi(t)$  as well as the kinetic temperature  $T_{\text{kin}}(t)$  (depending on the velocities) will be discussed below.

The Verlet method can be rewritten in the leapfrog form using definition:

$$\dot{q}(t+h/2) = \frac{q(t+h) - q(t)}{h} \quad (7)$$

where  $q$  stands for any coordinate component or  $\xi$ . It may be even advantageous in a computer code to use difference  $q(t+h) - q(t)$  instead of  $\dot{q}(t+h/2)$ . All these variants yield identical trajectories and need not be distinguished, provided that eq 7 is treated as a formal definition of symbol  $\dot{q}(t+h/2)$  (which approximates the velocity at time  $t+h/2$  up to the order of  $\mathcal{O}(h^2)$ ).

**2.4. Time-Reversible Velocity Predictor.** We propose to calculate the unknown velocities  $\dot{q}^p(t)$  from the knowledge of the history (previous positions) by the following predictor:

$$\dot{q}^p(t) = \frac{1}{h} \sum_{i=0}^{k+1} A_i q(t-ih) \quad (8)$$

where  $k$  stands for the (additional) predictor length. For  $k=0$  the predictor uses only information known within the Verlet algorithm at time  $t$ , namely  $q(t)$  and  $q(t-h)$ .

To determine constants  $A_j$ ,  $i=0, \dots, k+1$ , we will use the method described elsewhere.<sup>8</sup> Let us Taylor expand the right-hand side of eq 8:

$$\sum_{i=0}^{k+1} A_i q(t-ih) = \frac{1}{j!} \sum_{j=0}^{\infty} X_j q^{(j)} h^j$$

where

$$X_j = \sum_{i=0}^{k+1} (-i)^j A_i$$

Equation 8 should give the velocity  $\dot{q}^p(t)$  correct up to the second order (the order of the Verlet method). The following three

equations must be thus satisfied

$$X_0 = \sum_{i=0}^{k+1} A_i = 0 \quad (9)$$

$$X_1 = - \sum_{i=0}^{k+1} i A_i = 1 \quad (10)$$

$$X_2 = \sum_{i=0}^{k+1} i^2 A_i = 0 \quad (11)$$

There is no solution for  $k=0$ , therefore the minimum predictor length is  $k=1$ . The first error term in eq 8 is then  $X_3 h^2$ ; it is even and therefore it does not cause time irreversibility. The next term,  $X_4 h^3$ , causes time irreversibility of the third order  $\mathcal{O}(h^3)$ ; it means that running a simulation with doubled time step multiplies the energy drift eight times.

If we consider  $k > 1$ , we can achieve a better time reversibility by nullifying the terms at odd powers of  $h$ :

$$X_{2j} = \sum_{i=0}^{k+1} i^{2j} A_i = 0, \quad j = 1, 2, \dots, k \quad (12)$$

The solution of eqs 9–12 is

$$\begin{aligned} A_0 &= \frac{2k+1}{k+1} \\ A_1 &= -2(2k+1) \frac{1}{k+2} \\ A_2 &= +2(2k+1) \frac{k}{(k+2)(k+3)} \\ A_3 &= -2(2k+1) \frac{k(k-1)}{(k+2)(k+3)(k+4)} \\ &\vdots \end{aligned}$$

or in a compact form

$$A_i = (-1)^i (1 - \delta_{0i}/2) \binom{2k+2}{k+1-i} / \binom{2k}{k}$$

where  $\delta$  stands for the Kronecker delta.

**2.4.1. Proof.** To prove the above statement, let us first consider expressions for  $X_{2j}$ ,  $j=0, 1, \dots, k$ . They are composed of terms  $A_i i^{2j}$ ,  $i > 0$ , which we write as

$$A_i i^{2j} = \frac{1}{2} [A_i (-i)^{2j} + A_i (+i)^{2j}]$$

The equation for  $X_{2j}$ ,  $j > 0$ , then becomes

$$\begin{aligned} 2 \binom{2k}{k} (-1)^{k+1} X_{2j} &= + \binom{2k+2}{0} (-k-1)^{2j} \\ &- \binom{2k+2}{1} (-k)^{2j} + \binom{2k+2}{2} (-k+1)^{2j} \\ &- \dots \binom{2k+2}{2k+2} (k+1)^{2j} \end{aligned} \quad (13)$$

This is the operator of the  $(2k+2)$ -th difference applied to function  $f(j) = (j-n-1)^{2j}$ , and therefore the result is 0 for  $2k+2 > 2j$ , i.e.,  $j \leq k$ .<sup>3</sup> (The difference operator is defined by  $\Delta f(j) = f(j) - f(j-1)$ ). A degree of a polynomial is decreased by one by

applying this operator. The sum in eq 13 then equals  $\Delta^{2k+2}f(j)$  because the powers of the difference operator contain binomial coefficients with alternating signs, as can be shown by using the Pascal triangle.)

It remains to calculate  $X_1$ . After inserting for  $A_i$ , eq 10 becomes

$$\begin{aligned} \binom{2k}{k} X_1 &= \binom{2k}{k} \\ &= \binom{2k+2}{k} - 2 \binom{2k+2}{k-1} \\ &\quad + 3 \binom{2k+2}{k-3} - + \dots \end{aligned} \quad (14)$$

To prove it, we recursively expand the binomial coefficient:

$$\begin{aligned} \binom{2k}{k} &= \binom{2k+1}{k} - \binom{2k}{k-1} \\ &= \binom{2k+1}{k} - \binom{2k+1}{k-1} + \binom{2k}{k-2} \\ &\quad \vdots \\ &= \binom{2k+1}{k} - \binom{2k+1}{k-1} + \binom{2k+1}{k-2} \\ &\quad - + \dots \end{aligned}$$

and then apply the same recursive expansion to every term in the last expression.

**2.4.2. Final formula for TRVPs.** It may be useful to express  $\dot{q}^p(t)$  using the first differences. The advantages include smaller rounding errors and easier conversion to the Gear-type methods. The algorithm is then closer to the leapfrog form:

$$h\dot{q}^{\text{PR}}(t) = \sum_{i=0}^k B_i [q(t-ih) - q(t-[i+1]h)] \quad (15)$$

It holds  $B_0 = A_0$  and

$$B_j = (-1)^j (2k+1) \frac{k(k-1)\cdots(k+1-j)}{(k+1)(k+2)\cdots(k+1+j)}$$

or recursively

$$B_0 = \frac{2k+1}{k+1} \quad (16)$$

$$B_j = -B_{j-1} \times \frac{k+1-j}{k+1+j}, \quad j > 0 \quad (17)$$

which can be easily coded.

The drift in the total energy is of the order of  $h^{2k+1}$ .

**2.5. Velocity Estimators and Kinetic Temperature.** Several approximations of velocities can be used to calculate the kinetic temperature, eq 1. The simplest possibility is the difference formula which is equivalent to the so-called velocity Verlet (VV) algorithm:

$$\begin{aligned} \dot{\vec{r}}_i^{\text{VV}}(t) &= \frac{\dot{\vec{r}}_i(t-h/2) + \dot{\vec{r}}_i(t+h/2)}{2} \\ &= \frac{\vec{r}_i(t+h) - \vec{r}_i(t-h)}{2h} \end{aligned} \quad (18)$$

The “harmonic approximation” (HA):

$$\begin{aligned} \dot{\vec{r}}_i^{\text{HA}}(t)^2 &= \dot{\vec{r}}_i(t-h/2) \cdot \dot{\vec{r}}_i(t+h/2) \\ &= \frac{[\vec{r}_i(t+h) - \vec{r}_i(t)] \cdot [\vec{r}_i(t) - \vec{r}_i(t-h)]}{h^2} \end{aligned} \quad (19)$$

gives exactly constant total energy when the Verlet method is applied to a harmonic oscillator.

The “leap-frog approximation” (LF) makes an average of the approximated kinetic energies at midpoints:

$$\begin{aligned} \dot{\vec{r}}_i^{\text{LF}}(t)^2 &= \frac{\dot{\vec{r}}_i(t-h/2)^2 + \dot{\vec{r}}_i(t+h/2)^2}{2} \\ &= \left[ \frac{\vec{r}_i(t+h) - \vec{r}_i(t)}{h} \right]^2 + \left[ \frac{\vec{r}_i(t) - \vec{r}_i(t-h)}{h} \right]^2 \end{aligned} \quad (20)$$

Finally, one can use the predicted velocities,  $\dot{\vec{r}}_i^{\text{PR}}(t)$ , see eq 15, to calculate the kinetic temperature. This choice, denoted hereafter as PR, is independent of the particular form of the right-hand side. In contrast, options VV, HA, and LF require a particular right-hand side, eqs 2 and 3 or similar, as algorithmized below.

**2.6. TRVP Algorithm.** The proposed method can be easily combined with the SHAKE algorithm to maintain constraints (typically bond lengths). It can be written in several equivalent forms. Because of reduced numerical errors, we store the history of differences  $q(t) - q(t-h)$ ,  $q(t-h) - q(t-2h)$ , etc., rather than the positions. One step from  $t$  to  $t+h$  of the combined predicted velocity Nosé–Hoover Verlet integration with optional SHAKE method is then summarized below:

- (1) Calculate the potential energy and forces  $\vec{f}_i(t)$  from known positions  $\vec{r}_i(t)$ .
- (2) Predict velocities  $\dot{q}$  ( $q = \{\vec{r}_i, \xi\}$ ) from known history  $q(t) - q(t-h)$ ,  $q(t-h) - q(t-2h)$ , ...,  $q(t-kh) - q(t-[k+1]h)$ , using eqs 15, 16, and 17.
- (3) Perform one step of the Verlet method, eq 5, to get  $\vec{r}_i(t+h)$ .
- (4) Run the SHAKE algorithm;  $\vec{r}_i(t+h)$  are modified.
- (5) Calculate new differences  $\vec{r}_i(t+h) - \vec{r}_i(t)$ .
- (6) Calculate the kinetic temperature by one of four available eqs 18–20 and 15; below we will recommend eq 20.
- (7) Perform one step of the Verlet method, eq 6, for  $\xi$  to get  $\xi(t+h)$ .
- (8) Calculate the total energy at time  $t$ .
- (9) Advance time,  $t := t+h$ .

The algorithm does not need iterations (except those inside SHAKE). At the integration start the history needed for Step 2 is not known, and shorter predictors ( $k=0,1,2,\dots$ ) must be used. Since a typical MD run includes equilibration, this will rarely be a problem.

**2.7. Other Methods.** We compare the proposed TRVP method with several known methods.

**2.7.1. Gear Integration and Lagrangian Constraint Dynamics.** The Gear integration method<sup>1,4</sup> is based on storing the history in the form of a vector of higher derivatives at time  $t$ ,  $(q(t), h\dot{q}(t), (h^2/2)\ddot{q}(t), \dots)$ , from which the positions and velocities at time  $t+h$  are predicted by the Taylor expansion. The integration step thus has an easy access to velocities, and the method is straightforward unless constraints are to be satisfied.

For systems with constrained bonds we use the Lagrangian formulation of the constrained dynamics.<sup>10</sup> It is based on evaluation of the constraint forces by Lagrange multipliers; the set of linear equations is solved by the conjugate gradient method.<sup>11</sup> Since the constraints are satisfied only up to the order of the method and rounding errors, they are corrected by the same method at every step.

**2.7.2. Verlet with Iterated Velocities.** As already mentioned, the Verlet-family integrators suffer from a chicken-or-egg problem: The velocities needed in an integration step are known after the integration step is finished, i.e., the velocities are given implicitly by a set of equations. These equations can be solved by iterations.<sup>6</sup> The algorithm reads as<sup>12</sup>

- (1) Calculate the potential energy and forces  $\vec{f}_i(t)$  from known positions  $\vec{r}_i(t)$ .
- (2) Initial approximation:  
Version 1:  $\vec{r}_i(t) := [\vec{r}_i(t) - \vec{r}_i(t-h)]/h$ ,  $\dot{\xi}(t) := [\dot{\xi}(t) - \dot{\xi}(t-h)]/h$ .  
Version 2:  $\vec{r}_i(t) := 0$ ,  $\dot{\xi}(t) := 0$ .
- (3) Repeat the following loop:
  - (a) Perform one step of the Verlet method, eq 5, to get  $\vec{r}_i(t+h)$  (using velocities  $\vec{v}_i(t)$ ).
  - (b) Run the SHAKE algorithm;  $\vec{r}_i(t+h)$  are modified. (Omitted in the first iteration of Version 2.)
  - (c) Calculate new velocities  $\vec{v}_i(t) := [\vec{r}_i(t+h) - \vec{r}_i(t-h)]/(2h)$ .
  - (d) Calculate the kinetic temperature by one of eqs 18–20.
  - (e) Perform one step of the Verlet method, eq 6, to get  $\dot{\xi}(t+h)$ .
  - (f) Calculate the new velocity  $\dot{\xi}(t) := [\dot{\xi}(t+h) - \dot{\xi}(t-h)]/(2h)$ .
- (4) Calculate the total energy at time  $t$ .
- (5) Advance time,  $t := t+h$ .

The number of iterations (Step 3) may be either fixed or controlled by a predefined accuracy of the velocities. Note that the forces are calculated once per step, however, SHAKE iterations must be repeated (although less iterations are then needed). One pass (iteration) of Version 1 is equivalent to TRVP ( $k=0$ ), which is not accurate enough.

**2.7.3. MTTK Method.** A smart method keeping the Verlet scheme for positions but replacing the integrator of  $\xi$  has been proposed.<sup>5</sup> The algorithm is<sup>12</sup>

- (1)  $\dot{\xi}(t+h/4) := \dot{\xi}(t) + (h/4)a(t)$ , where  $a(t) = [T_{\text{kin}}(t)/T - 1]/\tau^2$ .
- (2)  $\vec{r}_i(t) := \vec{r}_i(t)$ ,  $\text{sym}[-(h/2)\dot{\xi}(t+h/4)]$ .
- (3)  $\dot{\xi}(t+h/2) := \dot{\xi}(t+h/4) + (h/4)a^*(t)$ , where  $a^*(t)$  is calculated from the velocities calculated in the previous step.
- (4)  $\vec{r}_i(t+h/2) := \vec{r}_i(t) + (h/2)\vec{f}_i/m_i$ .
- (5)  $\vec{r}_i(t+h) := \vec{r}_i(t) + h\vec{v}_i(t+h/2)$ .
- (6) Run (the first part of) the RATTLE algorithm.
- (7) Calculate forces  $\vec{f}_i(t+h)$  from  $\vec{r}_i(t+h)$ .
- (8)  $\vec{r}_i(t+h) := \vec{r}_i(t+h/2) + (h/2)\vec{f}_i(t+h)/m_i$ .
- (9) Run (the second part of) the RATTLE algorithm.
- (10)  $\dot{\xi}(t+3h/4) := \dot{\xi}(t+h/2) + (h/4)a(t+h)$ , where  $a(t+h)$  is calculated from the RATTLE velocities.
- (11)  $\vec{r}_i(t+h) := \vec{r}_i(t+h)/\text{sym}[(h/2)\dot{\xi}(t+3h/4)]$ .
- (12)  $\dot{\xi}(t+h) := \dot{\xi}(t+3h/4) + (h/4)a^*(t+h)$ , where  $a^*(t+h)$  is calculated from the velocities calculated in the previous step.

In the original work the function

$$\text{sym}(x) = \exp(x) \quad (21)$$

is derived by the Trotter decomposition of the Liouville operator. We will denote this version by letter  $e$ .

**2.7.4. Modification of MTTK Method.** It follows from reading the MTTK algorithm in the bottom-up direction that the method is time reversible for any function  $\text{sym}(x)$ . However,  $\text{sym}(x) \approx 1+x$  is required to maintain the order of the method. We thus propose two computationally less expensive functions:

$$\text{sym}_+(x) = 1+x \quad (22)$$

$$\text{sym}_-(x) = 1/(1-x) \quad (23)$$

The speed gain is marginal for large atomic systems but may be significant for simple enough systems. We will denote these versions by symbols  $+$  and  $-$ , respectively.

**2.7.5. Berendsen Thermostat.** The Berendsen (friction) thermostat<sup>1,2</sup> is used for comparison. The equations of motion are modified by a friction term:<sup>11</sup>

$$\ddot{\vec{r}}_i = \frac{\vec{f}_i}{m_i} - \frac{\ln(T_{\text{kin}}/T)}{2\tau} \dot{\vec{r}}_i$$

In a simulation of a dilute system (ideal gas), the temperature relaxes to the thermostat value  $T$  exponentially with the correlation time of  $\tau$ . In the simplest implementation to the Verlet scheme, velocities  $\vec{v}_i(t+h/2)$  are multiplied by factor  $\exp[-\ln(T_{\text{kin}}(t)/T)h/(2\tau)]$  after every step.

**2.8. Simulation Details.** **2.8.1. Potential Cutoff.** In molecular models of argon and water we use a smoothly truncated Lennard-Jones potential given by the formula:<sup>11</sup>

$$u_{\text{LJ}}(r) \approx \begin{cases} 4\epsilon[(\sigma/r)^{12} - (\sigma/r)^6] & \text{for } r < C_1 \\ A(r^2 - C_2)^2 & \text{for } C_1 < r < C_2 \\ 0 & \text{for } C_2 < r \end{cases}$$

where  $C_1$  and  $A$  are calculated from the cutoff  $C_2$  so that both the potential and the forces are continuous. Standard cutoff corrections for the potential energy and pressure are calculated using the usual assumption that the radial distribution function is unity beyond  $C_1$ .

The SPC/E water model<sup>13</sup> contains partial charges. In order to avoid additional errors inherent to more sophisticated methods (Ewald summation, reaction field) and also to gain speed, we use a simple truncated formula<sup>14</sup> to approximate the electrostatic forces. The  $1/r$  term in the Coulomb energy is replaced by

$$\frac{1}{r} \approx \begin{cases} 1/r - s & \text{for } r < C_1 \\ (r - C_2)^3(A + Br) & \text{for } C_1 < r < C_2 \\ 0 & \text{for } C_2 < r \end{cases} \quad (24)$$

where  $C_1 = 0.7C_2$ . The shift  $s$  and parameters  $A$  and  $B$  are determined so that the potential, forces, and the derivative of forces are continuous. The electrostatic force is thus neglected beyond the cutoff, shifted at short separations, and smoothly interpolated in between.

**2.8.2. Mechanical Quantities.** The averaged potential energy,  $E_{\text{pot}}$  also called residual internal energy, can be regarded as the most important and simplest mechanical quantity.

The second mechanical quantity of interest is pressure, which is calculated from the virial of force. This is straightforward for liquid argon. For constrained models there are two possibilities. To evaluate the instantaneous pressure in the TRVP, Gear, and Berendsen methods, we use the atom-based formula:

$$P = \frac{1}{3V} (2E_{\text{kin}} + \sum_{i < j} r_{ij} f_{ij}) \quad (25)$$

where we use notation  $|\vec{r}| = r$ . The sum runs over all interacting pairs of particles (with the distance calculated by the nearest image convention). Forces  $f_{ij}$  include the constraint forces calculated within the SHAKE procedure or given by Lagrange multipliers. This pressure thus depends on velocities.

To evaluate the pressure in the iteration and MTTK methods, we use the molecule-based formula:

$$P = \frac{1}{3V} \left[ 2E_{\text{tr}} + \sum_{n < m} \sum_i \sum_j \frac{(\vec{R}_{nm} \cdot \vec{r}_{ni, mj})}{R_{nm}} f_{ni, mj} \right] \quad (26)$$

where  $\vec{R}_n$  is the position of a reference point (for water we use oxygen) in molecule  $n$ . Position of atom  $i$  at molecule  $n$  is denoted as  $\vec{r}_{ni}$  and  $\vec{r}_{ni, mj} = \vec{r}_{mj} - \vec{r}_{ni}$  (with the nearest-image rule). Instead of the total kinetic energy, only the translational part is used

$$E_{\text{tr}} = \frac{1}{2} \sum_{n=1}^N M_n V_n^2$$

where  $M_n$  ( $V_n$ ) is the mass (velocity) of the center-of-mass of molecule  $n$ . Both methods give the same results up to the order of the integrator.

The virial of electrostatic forces is the same (but sign) as the electrostatic energy. However, if the electrostatic forces are approximated by eq 24, this relation is only approximate. We use a direct evaluation of the virial because it has been shown<sup>14</sup> that it gives a better approximation of the true pressure.

**2.8.3. Control Quantities.** The drift in the total energy  $E_{\text{NH}}$ , eq 4, is an auxiliary control quantity which is sensitive to time irreversibility. We calculate it by linear regression from the  $E_{\text{NH}}(t)$  dependence.

The quality of the canonical distribution of particle velocities can be generally monitored by moments. All odd moments are zero because of symmetry. The second moment, variance of velocity (averaged over particles), is proportional to the kinetic temperature. The next nonzero moment is the kurtosis:

$$\text{kurtosis} = \frac{\langle v^4 \rangle}{\langle v^2 \rangle^2} - 3 \quad (27)$$

Here  $v$  represents any component of the velocity and  $\langle \cdot \rangle$  includes averaging over all equivalent degrees of freedom. The kurtosis is zero for the Gaussian distribution.

Since we use different approximations of velocity, eqs 15 and 18–20, it is “natural” to use the same definition also for kurtosis. This is easy because only a squared velocity appears in eq 27. For some purposes (e.g., the velocity autocorrelation function), the velocity versions HA and LF are not directly applicable, and therefore either the simple leapfrog  $\vec{r}(t - h/2)$  or VV must be used instead of the “natural” definition. It then makes sense to determine kurtosis also for these alternate definitions.

The variance of the kinetic temperature is an important quantity related to heat capacity. Unlike kurtosis, which is a single-particle

property, the variance of temperature depends on the whole configuration. From the Maxwell–Boltzmann distribution one can easily calculate that

$$\text{var} T_{\text{kin}} = \frac{2T^2}{f} \quad (28)$$

It is important in simulations of complex systems that both slow and fast degrees of freedom are well equilibrated. For systems of small molecules the translational degrees of freedom can be easily separated from rotations. We thus define the translational temperature:

$$T_{\text{tr}} = \frac{2}{k(3N - 3)} E_{\text{tr}} \quad (29)$$

The number of degrees of freedom is  $3N - 3$ , taking into account momentum conservation. The rotational temperature is then

$$T_{\text{rot}} = T_{\text{kin}} - T_{\text{tr}}$$

Both temperatures should be the same.

**2.8.4. Diffusivity.** The diffusivity, as the simplest example of kinetic quantity, was calculated by the Einstein formula from squared displacements averaged over all particles, coordinates, and overlapping blocks 100 ps long with the coverage (overlap) factor three times; correlations of the consecutive data were taken into account in error estimation.<sup>15</sup> First, the dependence of the mean squared displacement on time calculated over shorter blocks and averaged over all simulations was plotted, and the initial nonlinear part determined. For argon we omitted the first 2 ps and for water 20 ps, and the rest was fitted to a linear function. The obtained diffusivity is not corrected for finite size errors.

**2.8.5. Gyration Tensor and Shape Anisotropy.** Conformation changes of a single molecule can be studied by the gyration tensor and related shape descriptors.<sup>16</sup> We use the mass-weighted tensor:

$$S_{ab} = \frac{1}{M} \sum_i m_i r_{i,a} r_{i,b} \quad (30)$$

where indices  $a, b$  run over coordinates  $x, y, z$ , index  $i$  labels atoms, the atom coordinates are with respect to the center of mass, and  $M = \sum_i m_i$  is the molecule mass. The radius of gyration is given by the tensor trace,  $R_{\text{gyr}}^2 = S_{xx} + S_{yy} + S_{zz}$ . To obtain shape descriptors, the tensor must be diagonalized so that in the new coordinates (main axes of inertia) it writes as  $S = \text{diag}(\lambda_{xx}, \lambda_{yy}, \lambda_{zz})$ . The relative shape anisotropy  $\kappa$  is then defined by

$$\kappa^2 = \frac{\lambda_{xx}^2 + \lambda_{yy}^2 + \lambda_{zz}^2 - \lambda_{xx}\lambda_{yy} - \lambda_{yy}\lambda_{zz} - \lambda_{zz}\lambda_{xx}}{(\lambda_{xx} + \lambda_{yy} + \lambda_{zz})^2} \quad (31)$$

Note that  $0 \leq \kappa \leq 1$  and that  $\lambda_{xx} + \lambda_{yy} + \lambda_{zz} = R_{\text{gyr}}^2$  because of trace invariance.

### 3. NUMERICAL TESTS

**3.1. Ring of Oscillators.** The first testing system is small—a ring of  $N = 6$  particles of unit mass in one direction. The neighboring particles interact via either harmonic or anharmonic terms:

$$\begin{aligned} U_h(r) &= Kx^2/2 \\ U_a(r) &= K(x^2/2 + x^4/4) \end{aligned}$$

where  $K$  is the force constant. We used values  $K_{i,i+1} = (N + i)/N$ , where  $i = 0, \dots, 5$  and  $i = 6$  is identical to 0 so that a ring is created. The momentum is conserved; since the numerical noise becomes detectable in long runs, we reset the total momentum to zero after every integration step. The number of degrees of freedom is  $f = N - 1$ .

For each method we ran three numerical tests with time steps of  $h = 1/16, 1/32$ , and  $1/64$ , and the correlation time  $\tau = 0.3$ ; in addition, results for  $\tau = 0.2$  and  $0.5$  are presented in the Supporting Information. Each run lasted  $2^{25}$  time units.

The inspection of trajectories shows that the harmonic oscillator simulations exhibit rather poor ergodicity. It is not surprising because this system is similar to the path-integral molecular dynamics, which is known to require a Nosé–Hoover chain for ergodic thermostating. Both harmonic and anharmonic simulations with the longest time step using a predictor also for energy (eq 15) are unstable and crash with infinite values. The MTTK methods are more (but not 100%) prone to bad ergodicity. Better time reversibility and longer  $\tau$  improve ergodicity.

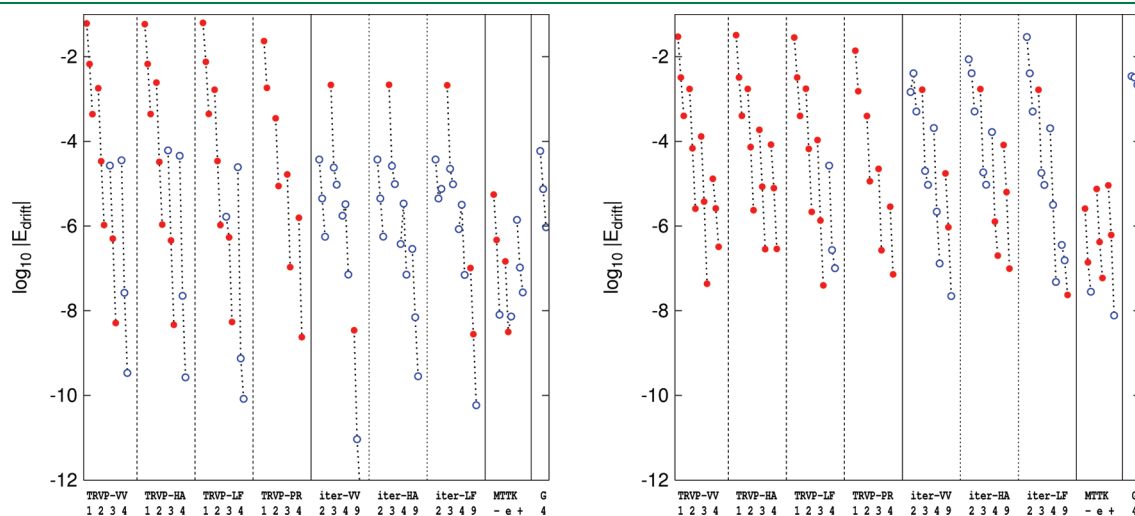
The drift in the total energy is shown in logarithmic scale (as  $\log_{10}|E_{\text{drift}}|$ ) in Figure 1. The drift diminishes if: (i) the time step decreases, (ii) the predictor length or the number of iterations increases, and also (iii)  $\tau$  increases (see the Supporting Information). The TRVP and iteration methods are similar if we use a sufficient predictor length or equivalently increase the number of iterations (increasing the predictor length as well as the number of iterations by 1 decreases the irreversibility error by  $\mathcal{O}(h^2)$ ). The MTTK method is uniformly good even for the longest time steps where the above methods may be unstable. The original Trotter-based version with  $\text{sym} = \text{exp}$  is the best for the harmonic oscillators but worst for the anharmonic ones. The Gear method exhibits the biggest drift, especially for the anharmonic oscillators. We should, however, bear in mind that the drift is only a control quantity and never a final result of interest.

The averaged kinetic temperature (see the Supporting Information) equals the nominal thermostat value with a better precision than the estimated statistical uncertainty for all methods using the same single kinetic temperature in the right-hand side, namely TRVP, iteration, and Gear methods. The original MTTK method, eq 21, gives the average temperature distinguishable from the nominal value only for the longest time step. The modified MTTK methods, eqs 22 and 23, lead to second-order differences; the biggest difference of 0.0025 for  $h = 1/16$  and  $\tau = 0.2$  drops to  $2.5 \times 10^{-5}$  for  $h = 1/64$  and  $\tau = 0.5$ .

The variance of temperature is  $\text{var}T_{\text{kin}} = 0.4$  for both oscillator rings. The results are shown in Figure 2 along with standard errors (68% confidence level) estimated from blocks. The iteration method with only two iterations fails to yield the canonical distribution at all, the longest-time step version fails even with four iterations (whereas three iterations give satisfactory results); however, the fully iterated version is satisfactory. The Gear method does not give the canonical distribution either. The predictor method is slightly better, and the longer predictor results are satisfactory. The MTTK method is uniformly good. The results are much better for the more ergodic anharmonic ring, unless the time step is too long and at the same time the reversibility poor. The best results are for the MTTK method and iterations with the HA kinetic energy, then for predictions with the VV kinetic energy. The Gear method converges poorly and needs a short time step.

As seen from Figure 3, the kurtosis is sensitive to ergodicity problems. For the more ergodic anharmonic oscillator, the MTTK method works well for all time steps, the TRVP method needs  $k \geq 2$ , and the iteration method  $i = 3$  or more iterations to converge well. Different velocity definitions give results differing in the second-order term. Particularly, the LF kinetic temperature but VV velocity is the best with the TRVP method.

The potential energy (see the Supporting Information) gives similar results. The harmonic results are scattered witnessing about poor ergodicity. The MTTK method works well, and the



**Figure 1.** The drift in the total energy (as logarithm of absolute value; negative drifts are open blue circles, positive solid red) for a ring of harmonic (left) and anharmonic (right) oscillators integrated by various methods. The triplets connected by dotted lines are from left for time steps  $1/16, 1/32$ , and  $1/64$  time units (if the simulation with  $1/16$  failed, only a doublet is shown). Label ‘TRVP’ denotes the proposed Nosé–Hoover integration with velocity predictor, the numbers below denote the value of  $k$ . Label ‘iter’ denotes the iteration method (Version 1), the numbers below denote the number of iterations. Symbols ‘VV’, ‘HA’, ‘LF’, and ‘PR’ refer to the kinetic temperature version, eqs 15 and 18–20. MTTK is the Martyna et al.<sup>5</sup> method, and the symbol below defines function  $\text{sym}()$ , eqs 21–23. Label ‘G’ denotes the Gear method ( $m = 4$ ). Error bars are not shown for clarity, they become apparent for  $|E_{\text{drift}}| < 10^{-8}$ .

TRVP and iteration methods need several iterations or a longer predictor and are not reliable with long time steps. The anharmonic case, more similar to a typical many-particle atomistic simulation, shows a clearer picture—the systems are ergodic. The best results are obtained with the LF formula for the kinetic energy. The iteration method is the best, followed by the predicted velocities (with  $k \geq 2$ ) and then the MTTK methods.

**3.2. Liquid Argon.** The second testing system is liquid argon modeled by the Lennard-Jones potential with parameters  $\varepsilon/k = 119.8$  K and  $\sigma = 3.405$  Å. The simulated temperature was 143.76 K (reduced temperature 1.2) and density  $1.3443$  g cm $^{-3}$  (reduced number density 0.8). We used  $N = 200$  atoms in the standard periodic cubic box, smooth potential cutoff  $C_2 = 10$  Å, time steps  $h = 20, 10,$  and  $5$  fs, and  $\tau = 0.1$  ps; more results with  $\tau = 0.3$  and  $0.5$  ps are presented in the Supporting Information. Trajectory length of each point was 200 ns.

We compare the propose predictor method with the iteration method controlled by the error limit in velocity of  $10^{-6}$  reduced units per particle, the Gear integrator ( $m = 4$ ), and also the Berendsen thermostat (with both the Verlet integrator as well as Gear).

Figure 4 collects the basic thermodynamic results—internal energy and pressure. Error bars were omitted because they are comparable to symbol sizes. It is seen that all methods converge well, however, using the LF kinetic temperature gives the best results, even better than the MTTK method. This observation applies to all methods (TRVP, iterations as well as Gear and also the Berendsen thermostat). The Gear values are shifted because of the finite-size ensemble error; this difference will diminish for larger systems. The TRVP method with the shortest predictor  $k = 1$  is sufficient; only if the kinetic energy was calculated from the predicted velocities (version PR), a longer predictor would be needed, but this version is not recommended anyway.

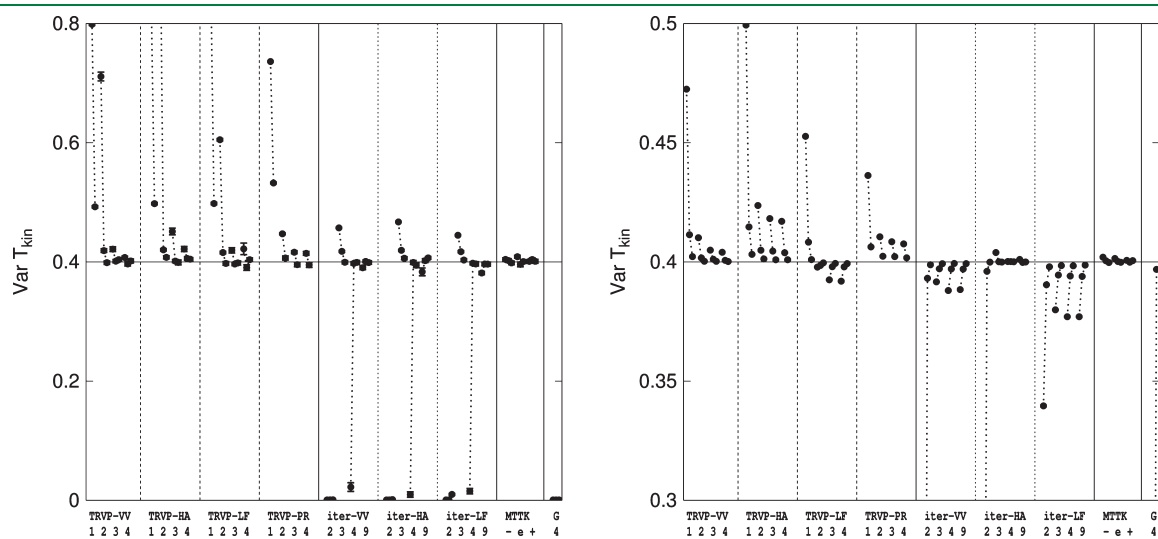


Figure 2. The variance of temperature for a ring of oscillators. See Figure 1 for symbol explanation. Note the different scales of the vertical axes.

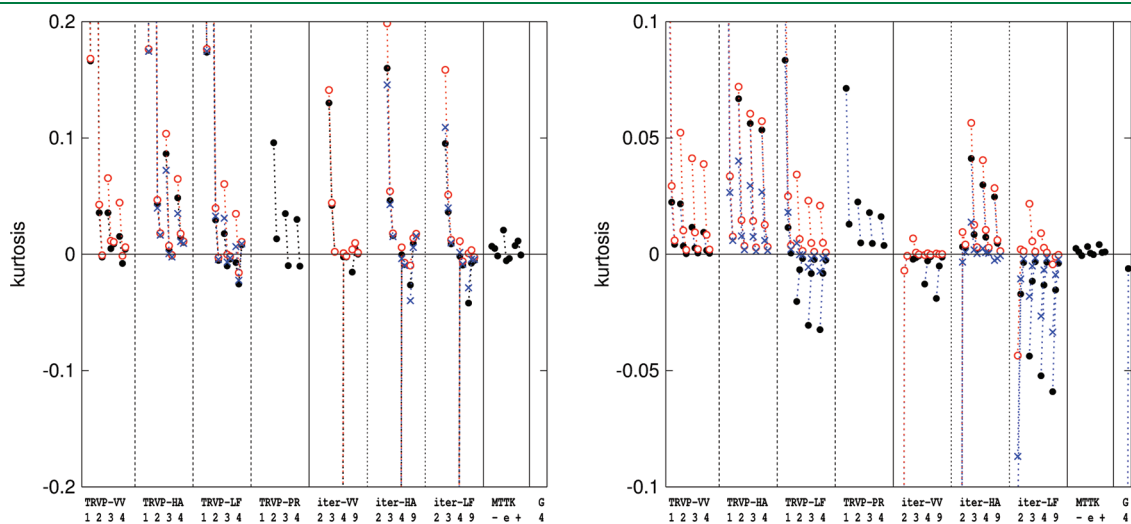
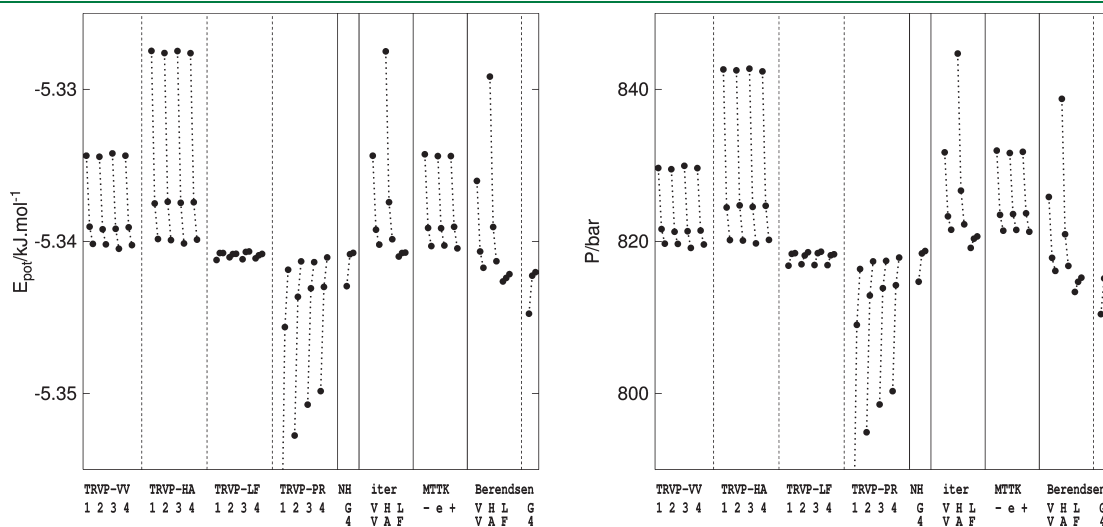


Figure 3. Kurtosis of the velocity distribution averaged over all 6 oscillators. Left: harmonic, right: anharmonic. Black solid circles: “natural” velocity, red open circles: single leapfrog velocity  $[r(t) - r(t - h)]/h$ , blue cross: VV  $([r(t + h) - r(t - h)]/(2h))$ , if differs from the “natural” velocity. See Figure 1 for symbol explanation. Note the different scales of the vertical axes.

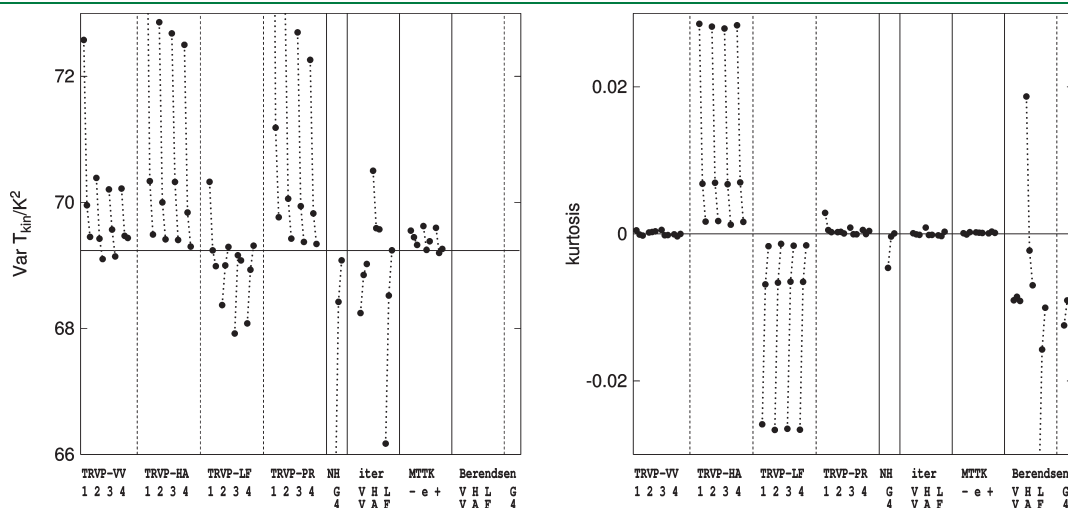
Figure 5 shows, along with the temperature variance, also kurtosis of the distribution of velocities; note that the kurtosis is calculated from the “natural” velocity definition for the TRVP method (see Section 2.8.3), whereas from the VV formula for iterations. It is seen that all the Nosé–Hoover methods converge to the correct values (at least within error bars, which are for the variance several symbol sizes). The Berendsen values are out of the graph because this method is not canonical. The TRVP method gives better temperature variance with  $k = 2$  than  $k = 1$ , however, increasing the predictor length to  $k > 2$  gives only a marginal improvement. The MTK results are the best even for large time steps. The VV velocity (eq 18) works best with TRVPs while the HA velocity (eq 19) with iterations.

The results for diffusivity, Figure 6, are subject of larger errors than for mechanical quantities. All methods and both thermostats converge well, only the TRVPs with the harmonic and predicted temperatures perform worse with long time steps. Both the TRVPs and iterations with the VV and LF temperatures are acceptable. Surprisingly, the Gear integration (both with the Nosé–Hoover and Berendsen thermostats) wins the comparison.

**3.3. Liquid Water.** The third testing system consists of  $N = 200$  SPC/E<sup>13</sup> water molecules in a cubic box simulated at ambient conditions (density  $997 \text{ kg m}^{-3}$ , temperature  $300 \text{ K}$ ). Both the Lennard-Jones and electrostatics cutoffs were set to  $C_2 = 9 \text{ \AA}$ . We further utilized  $h = 2, 1, \text{ and } 0.5 \text{ fs}$  and  $\tau = 0.1 \text{ ps}$ , and the runs took  $50 \text{ ns}$ . Selected results



**Figure 4.** Averaged potential energy (left) and pressure (right) for liquid argon. The triplets connected by dotted lines correspond (from left) to time steps 20, 10, and 5 fs. Label ‘TRVP’ denotes the proposed Nosé–Hoover integration with velocity predictor, the numbers below denote the value of  $k$ . Label ‘iter’ denotes the iteration method (Version 2) with the number of iterations controlled by precision. Symbols ‘VV’, ‘HA’, ‘LF’, and ‘PR’ refer to the kinetic temperature version. MTK is the Martyna et al.<sup>5</sup> method and the symbol below defines function  $\text{sym}()$ , eqs 21–23. Label ‘G’ denotes the Gear method ( $m = 4$ ), either with Nosé–Hoover (NH) or Berendsen thermostat. Error bars are comparable to symbol sizes.



**Figure 5.** Variance of temperature and kurtosis of the velocity distribution for liquid argon. The horizontal lines represent the respective theoretical values. See Figure 4 for symbol explanation.

with  $\tau = 0.3$  and  $0.5$  ps are presented in the Supporting Information.

The results for the potential energy, Figure 7, are similar as for argon: the TRVP and iteration methods are equivalent. It is more important to choose the kinetic energy formula. The LF version is the best, the VV version is worse and of the same accuracy as the MTTK method (using internally the same velocity approximation). The Gear integrator performs surprisingly well.

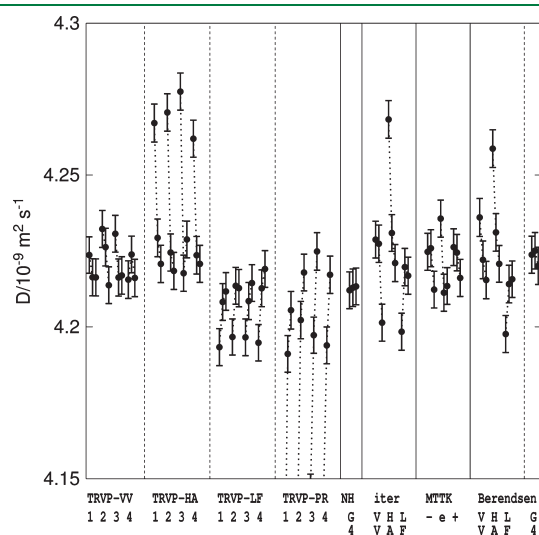
The results for pressure are different, for TRVPs (and also Berendsen thermostat) the VV kinetic temperature overperforms the LF one. However, this result can be explained by a “dynamic” atom-based algorithm, eq 25, used here for pressure. Of course, both formulas give results differing by a term proportional to  $h^2$ .

The variance of temperature, Figure 8, is in spite of 50 ns runs and a relatively small system, a subject of large errors. Essential all methods, perhaps with the exception of the

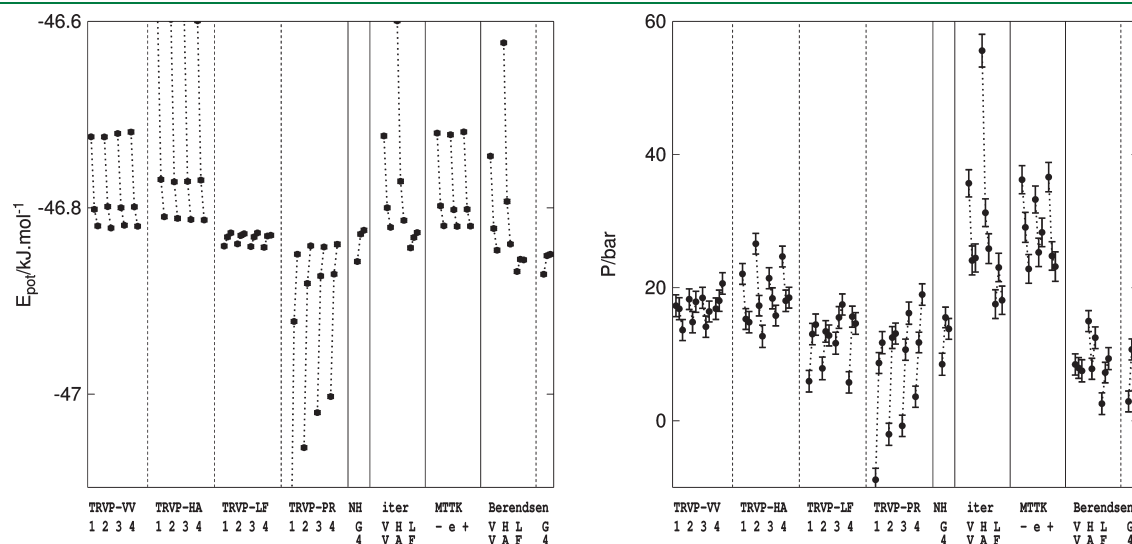
harmonic approximation for the kinetic temperature (and of course the Berendsen thermostat), give satisfactory results. The kurtosis is worse for the HA and LF temperatures and long time steps.

The equipartition of translational and rotational kinetic energies, Figure 9, depends on the time step and the kinetic temperature formula, but not on the integration method: the LF version is the best, followed by VV, whereas HA and PR are not so good. The diffusivity, Figure 10, behaves in the same way, only the data are more noisy. To a great extent also the potential energy and (with the “static” eq 26) also pressure follow the same pattern. This observation is not surprising because most thermostats keep the averaged (translational and rotational) temperature constant, whereas the Verlet integration errors increase the translational temperature. The center-of-mass velocities are then bigger, and the diffusivity increases. Similarly, the energy and pressure depend mainly on the interparticle energy. One may ponder that Nosé–Hoover and Berendsen thermostats with only the translational kinetic energy could give (some) thermodynamic quantities more accurately.

**3.4. Peptide.** As the last testing system we chose a small peptide in vacuum, hexalaanine modified by acetyl at the N-terminus and *N*-methylamide at the C-terminus ( $\text{CH}_3\text{—CO—}[\text{NH—CHCH}_3\text{—CO}]_6\text{—NH—CH}_3$ ) because: (i) a complex set of bonds provides a comprehensive test of constraint dynamics; (ii) it has a sufficiently rich conformational space to check for possible ergodicity problems; and (iii) a relatively small number of atoms facilitates sampling of all important conformations. The molecule was modeled by the GROMOS96<sup>17</sup> force field version 43A1 (with united-atom approximation for the  $\text{CH}_3$  and  $\text{CH}$  groups). However, for simplicity of coding we did not apply nonbonded fixes (exceptions) for 1–4 nonbonded interactions and included them in full. This change has a minor impact on the peptide properties. We also replaced the torsion term around the peptide bond by a locally equivalent harmonic potential and prohibited thus the *cis* conformations at all, because in trial runs with the original dihedral potential, we detected the *trans*–*cis* conversion at a hundred

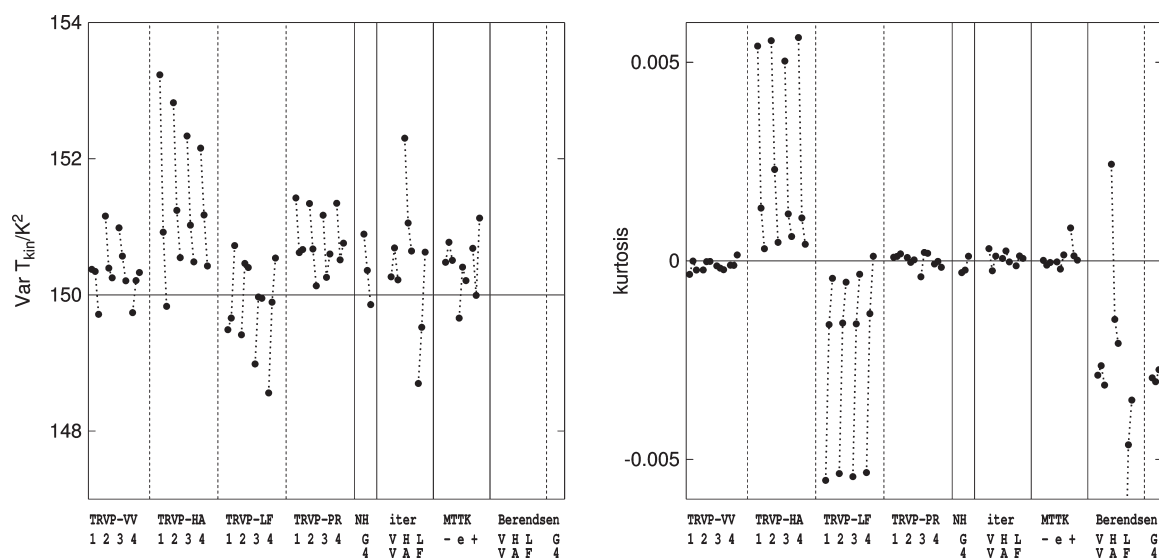


**Figure 6.** Diffusivity of liquid argon. See Figure 4 for symbol explanation. The error bars denote estimated standard errors.

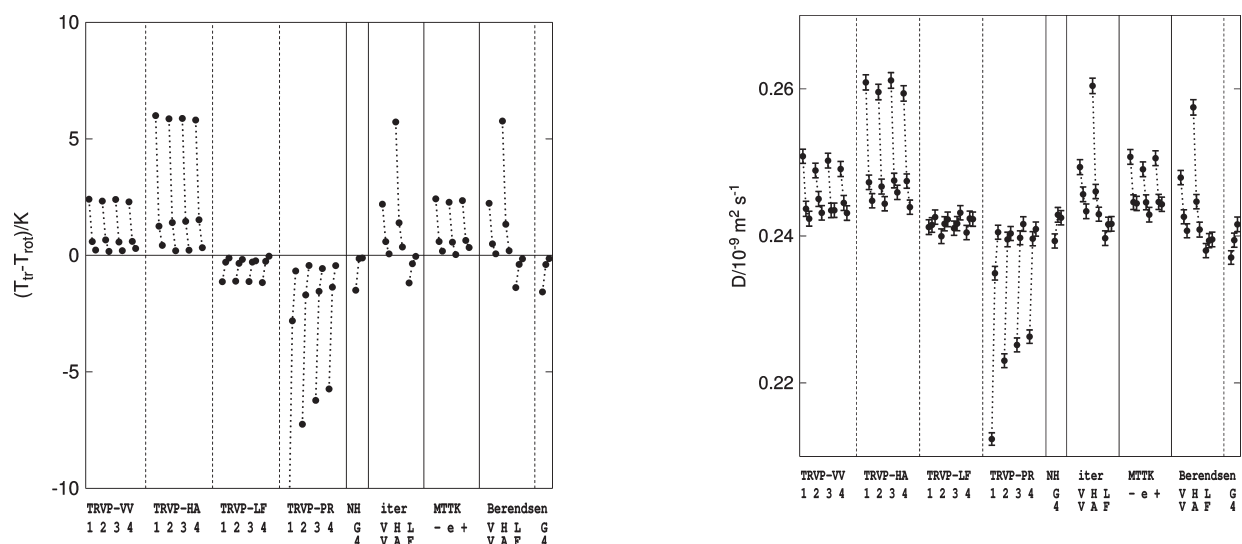


**Figure 7.** Averaged potential energy (left) and pressure (right) for liquid SPC/E water. The triplets correspond to time steps of 2, 1, and 0.5 fs, for other symbols see Figure 4.





**Figure 8.** Variance of temperature and kurtosis of the velocity distribution for liquid water. The horizontal lines represent the theoretical values. See Figure 7 for symbol explanation.



**Figure 9.** Equipartition of the translational and rotational temperatures for liquid water. See Figure 7 for symbol explanation.

nanosecond scale, which spoiled the statistics. Note that our aim is not to obtain an accurate realistic representation of this peptide but a model system reasonably close to realistic biomolecules, yet simple enough to allow for accurate comparison of methods.

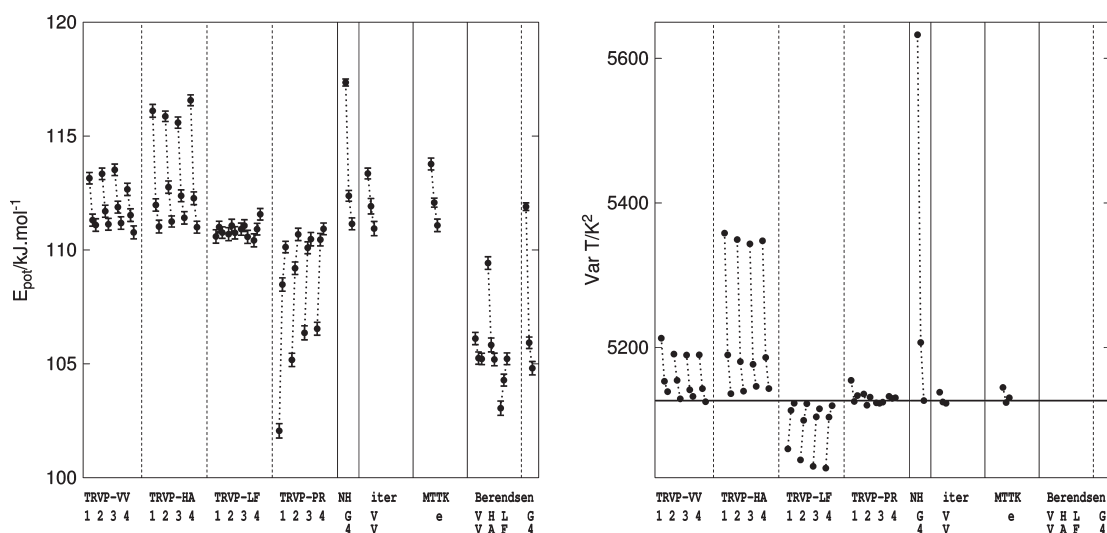
The peptide with constrained bond lengths was simulated at 450 K for 1  $\mu$ s using three time steps and the Nosé–Hoover or Berendsen correlation times of  $\tau = 0.1$  ps. The TRVP method and Berendsen thermostat were implemented using MACSIMUS,<sup>11</sup> for the iterated velocity, and for the MTTK method we used DL\_POLY;<sup>12</sup> only the VV version of the kinetic temperature, eq 18, and the original MTTK method are available. The momentum and angular momentum were periodically reset to exact zero and do not contribute to the number of degrees of freedom.

**Figure 10.** Diffusivity of liquid water. See Figure 7 for symbol explanation. The error bars denote estimated standard errors.

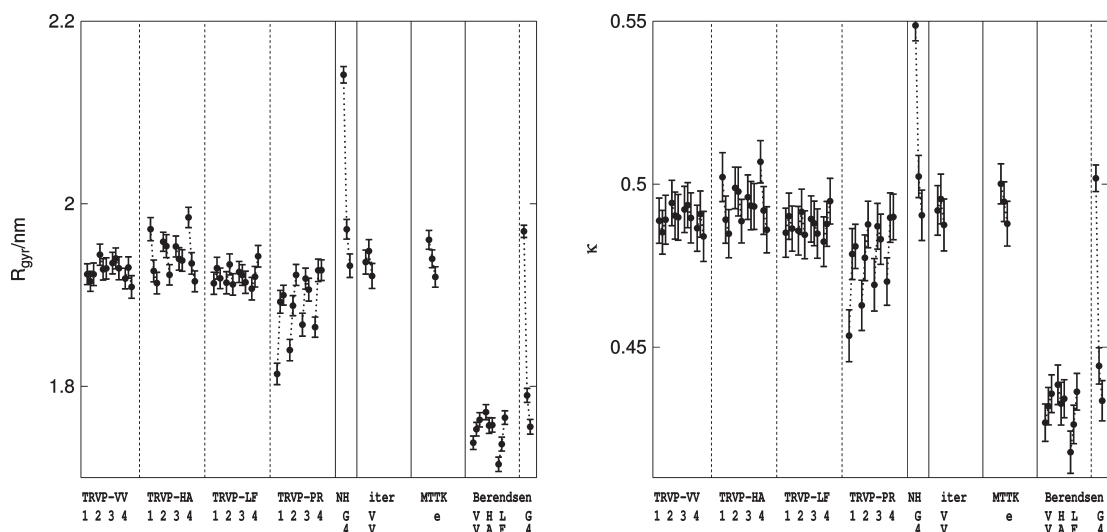
The internal energy is the least noisy variable of interest, see Figure 11 left. All canonical methods converge as the time step decreases, whereas the Berendsen value is off as expected for this small system of only 79 degrees of freedom. The leapfrog version of the kinetic energy definition is again the best.

The variance of temperature, see Figure 11 right, converges to the correct value similarly as for other investigated systems; because of a small number of degrees of freedom, the results are less noisy than for liquids. The values for the noncanonical friction thermostat are off the graph scale.

The averaged radius of gyration and averaged shape anisotropy show a similar pattern of convergence as the potential energy; only the data are more noisy (Figure 12). The TRVP method with both major kinetic energy definitions, VV and LF, performs well even for the longest time step and so do the benchmark methods, iterations (with VV), and MTTK. The Gear integrator is satisfactory with the shortest time step only.



**Figure 11.** Averaged potential energy (left) and variance of temperature (right) for modified hexaalanine in vacuum at 450 K. The triplets correspond to time steps of 2, 1, and 0.5 fs, for other symbols see Figure 4. The horizontal line is the theoretical value of temperature variance, eq 28.



**Figure 12.** Radius of gyration (left) and relative shape anisotropy  $\kappa$  (right) for modified hexaalanine in vacuum at 450 K. The triplets correspond to time steps of 2, 1, and 0.5 fs, for other symbols see Figure 4.

## 4. CONCLUSIONS

**4.1. TRVP Method.** The proposed TRVP method, see Section 2.6, is the main result of this work. It has been successfully applied to the Nosé–Hoover thermostat and Verlet integrator for many-particle atomistic systems, including those with constrained bond lengths. The method quality and speed are similar to the iteration method, and the decision which method to use may thus rather depend on algorithm subtleties. The TRVP method can be coded outside the existing Verlet + SHAKE algorithm, which may simplify the code design. Only one set of SHAKE iterations is needed, and thus measuring some quantities (e.g., the pressure tensor) does not interfere with repeated calculations. On the other hand, the TRVP method requires more memory than the iteration method. Often the shortest predictor ( $k = 1$ ) is sufficient, although  $k = 2$  is never worse.

The MTTK method is for large atomistic systems comparable to both the above simple Verlet-based methods. As it is more

complex, the kinetic energy is calculated four times per step, and for constrained systems it requires the RATTLE algorithm that is more CPU consuming than the SHAKE algorithm. For small systems, it performs better than both the TRVPs and iterations, avoiding many ergodicity problems; on the other hand, it is slower.

**4.2. Velocity Estimates in Verlet Schemes.** The TRVP and iteration methods are of the second order in the time step. At the same time there are several formulas, eqs 15 and 18–20 available to approximate velocities and in turn the kinetic temperature. They yield different coefficients at the second-order error terms for different quantities. We found that the LF eq 20 gives the overall best results for quantities of interest (energy, pressure, diffusivity), although the quality of the canonical distribution is worse. The method accuracy may depend on the details of the algorithm, e.g., an atom-based formula for pressure for constrained systems is better with the VV kinetic energy, eq 18,

probably because this formula depends on the constrained forces, which depend on velocities.

**4.3. Modification of MTTK Algorithm.** We proposed two modifications of the MTTK algorithm, see eqs 22 and 23, which avoid evaluating the exponential functions and therefore run faster. The speedup is insignificant for large atomistic systems but may be worth considering in some special cases, like Nosé–Hoover chains for simple systems. The kinetic temperature derived from the velocities available within the algorithm is (inconsistently) off by a second-order term in the time step, however, other observable quantities are not affected. One might consider alternating  $\text{sym}_+$  and  $\text{sym}_-$  in the Nosé–Hoover chain in order to decrease the kinetic temperature error. Compromise functions  $\text{sym}()$  better approximating  $\exp()$  than eqs 22 and 23 are also possible.

**4.4. Translational Temperature.** In the simulations of liquid water we found that the values of many quantities (energy, pressure, diffusivity) depend mainly on the translational temperature, eq 29. Using this temperature in the thermostat may compensate for large time step errors of the Verlet integration. A similar correction for equipartition errors might be possible for the internal vibrations instead of rotations.

## ■ ASSOCIATED CONTENT

**S Supporting Information.** Extended sets of results are available: runs with different correlation times  $\tau$  (except for the peptide); more quantities (e.g., the kinetic temperatures, electrostatic energy for water); time profiles of the total energy for the ring of oscillators; and more shape descriptors, the end-to-end distance, and correlation times for the model peptide. This material is available free of charge via the Internet at <http://pubs.acs.org>.

## ■ AUTHOR INFORMATION

### Corresponding Author

\*E-mail: [jiri.kolafa@vscht.cz](mailto:jiri.kolafa@vscht.cz).

## ■ ACKNOWLEDGMENT

This work was supported by the Czech Science Foundation, projects P208/10/1724 (J.K.) and 104/08/0600 (M.L.) and the Internal Grant Agency of the J. E. Purkinje University, grant no. 53222 15 0006 01 (M.L.).

## ■ REFERENCES

- (1) Allen, M. P.; Tildesley, D. J. *Computer Simulation of Liquids*; Clarendon Press: Oxford, U.K., 1986.
- (2) Frenkel, D.; Smit, B. *Understanding Molecular Simulation*; Academic Press: San Diego, CA, 2002.
- (3) Ralston, A. *A First Course in Numerical Analysis*; McGraw-Hill: New York, 1965.
- (4) Gear, C. W. *Numerical Initial Value Problems in Ordinary Differential Equations*; Prentice Hall: Upper Saddle River, NJ, 1971.
- (5) Martyna, G.; Tuckerman, M. E.; Tobias, D. J.; Klein, M. L. *Mol. Phys.* **1996**, *87*, 1117–1157.
- (6) Toxvaerd, S. *Mol. Phys.* **1991**, *72*, 159–168.
- (7) Nosé, S. *Mol. Phys.* **1984**, *52*, 255–268.
- (8) Hoover, W. G. *Phys. Rev. A* **1985**, *31*, 1695–1697.
- (9) Kolafa, J. *J. Comput. Chem.* **2004**, *25*, 335–342.
- (10) de Leeuw, S. W.; Perram, J. W.; Petersen, H. G. *J. Stat. Phys.* **1990**, *61*, 1203–1222.

(11) *MACSIMUS*; Institute of Chemical Technology: Prague, Czech Republic, 2011; <http://www.vscht.cz/fch/software/macsimus>, (accessed Aug 19, 2011).

(12) Todorov, I.; Smith, W. *The DL POLY 3 User Manual*; STFC Daresbury Laboratory: Cheshire, U.K., 2009.

(13) Berendsen, H. J. C.; Grigera, J. R.; Straatsma, T. P. *J. Phys. Chem.* **1987**, *91*, 6269–6271.

(14) Kolafa, J.; Moučka, F.; Nezbeda, I. *Collect. Czech. Chem. Commun.* **2008**, *73*, 481–506.

(15) Picálek, J.; Kolafa, J. *J. Mol. Liq.* **2007**, *134*, 29–33.

(16) Mattice, W. L.; Suter, U. W. *Conformational Theory of Large Molecules*; Wiley Interscience: Hoboken, NJ, 1994.

(17) van Gunsteren, W. F.; Billeter, S. R.; Eising, A. A.; Hünenberger, P. H.; Krueger, P.; Mark, A. E.; Scott, W. R. P.; Tironi, I. G. *The GROMOS96 Manual and User Guide*; BIOMOS: Zurich, Switzerland, 1996.

# On-the-Fly Numerical Surface Integration for Finite-Difference Poisson–Boltzmann Methods

Qin Cai,<sup>†,‡</sup> Xiang Ye,<sup>‡,§</sup> Jun Wang,<sup>‡</sup> and Ray Luo<sup>\*,†,‡</sup><sup>†</sup>Department of Biomedical Engineering and <sup>‡</sup>Department of Molecular Biology and Biochemistry, University of California, Irvine, California<sup>§</sup>Department of Physics, Shanghai Normal University, Shanghai, China

**ABSTRACT:** Most implicit solvation models require the definition of a molecular surface as the interface that separates the solute in atomic detail from the solvent approximated as a continuous medium. Commonly used surface definitions include the solvent accessible surface (SAS), the solvent excluded surface (SES), and the van der Waals surface. In this study, we present an efficient numerical algorithm to compute the SES and SAS areas to facilitate the applications of finite-difference Poisson–Boltzmann methods in biomolecular simulations. Different from previous numerical approaches, our algorithm is physics-inspired and intimately coupled to the finite-difference Poisson–Boltzmann methods to fully take advantage of its existing data structures. Our analysis shows that the algorithm can achieve very good agreement with the analytical method in the calculation of the SES and SAS areas. Specifically, in our comprehensive test of 1555 molecules, the average unsigned relative error is 0.27% in the SES area calculations and 1.05% in the SAS area calculations at a grid spacing of 1/2 Å. In addition, a linear correlation analysis was found to improve the accuracy of the coarse-grid SES areas, with the average unsigned relative error reduced to 0.13%. These validation studies indicate that the proposed algorithm can be applied to biomolecules over a broad range of sizes and structures. Finally, the numerical algorithm can also be adapted to evaluate the surface integral of either a vector field or a scalar field defined on the molecular surface for additional solvation energetics and force calculations.

## INTRODUCTION

Water is crucial to the functions of biological molecules such as nucleic acids and proteins. The solute–solvent interactions can be accurately modeled by explicit solvent models in biomolecular simulations. Nevertheless, extra computational cost has to be paid to handle thousands to millions of extra degrees of freedom in the explicit solvent. This is because a system of higher dimensions requires more sampling to achieve equilibrium and to cover a sufficient number of biologically interesting conformations. Alternatively, the solvent molecules can be treated implicitly as in the implicit solvent models to capture average solvent behaviors. Most modern implicit solvent models require the definition of an interface that separates the solute in atomic detail from the solvent approximated as a continuum medium. The solvation free energy and force are both sensitively dependent on the interface location and presentation.

The interface is often based on a molecular surface definition. Well-known molecular surface definitions are the solvent accessible surface (SAS), the solvent excluded surface (SES), and the van der Waals surface (VDWS). The SAS was defined by Lee and Richards as a union of atomic spheres with radii augmented by the probe radius for a given molecule.<sup>1</sup> The SES was introduced by Richards in 1977.<sup>2</sup> It is more complex, with two types of surfaces, the contact surface that consists of solvent exposed portions of van der Waals spheres and the re-entrant surface that consists of the inward-facing surface of the solvent probe sphere as it rolls over the molecule. Finally, the VDWS represents the molecular interior as the union of the atomic spheres with van der Waals radii. In contrast to the hard sphere definition of atomic volumes in the above surface definition, a smoothly varying

dielectric boundary using the Gaussian(-like) density approach has also been reported.<sup>3,4</sup>

Recent analyses have shown that the numerical Poisson–Boltzmann methods with the SES definition are reasonable in the calculation of electrostatic solvation energetics and forces with respect to explicit solvent simulations.<sup>3–7</sup> Unfortunately, the SES is not differentiable with respect to atomic positions, making it difficult to adopt in molecular dynamics simulations.<sup>5</sup> The SAS definition has difficulty reproducing the electrostatic energetics in the explicit solvent models due to its much enlarged atomic cavities, but the SAS definition has been used to estimate the nonpolar hydration energy with great success.<sup>6–14</sup> The SAS is also more efficient than the SES and differentiable with respect to atomic positions.<sup>7,15–20</sup> The VDWS definition is both efficient and smooth over time, making it advantageous to molecular simulations. However, it has been pointed out that the VDWS definition has the tendency to assign a much higher value to the apparent protein interior dielectric constant, as in the pK<sub>a</sub> calculations which may or may not be wanted.<sup>21–26</sup> On the other hand, Zhou and co-workers have shown that the electrostatic free energies using the VDWS definition as the dielectric boundary are in better accord with the experimental data and the explicit solvent simulation results for a series of tests, including a single mutation to the folded protein,<sup>27</sup> protein–protein complexes,<sup>28,29</sup> and protein–RNA complexes.<sup>30</sup>

Given a surface definition, an important issue is how to evaluate the surface area or surface integration for biomolecular

Received: June 8, 2011

Published: September 27, 2011

simulations. Due to the efficiency of the SAS definition, significant prior efforts have been invested to develop analytical methods<sup>15,16,18,20,31–45</sup> and numerical algorithms.<sup>46–51</sup> The analytical SAS methods can be further divided into two types, exact<sup>15,16,18,20,32–38</sup> and approximate.<sup>31,39–45</sup> A comprehensive review of the SAS methods can be found in ref 45. In contrast, the SES methods are less common. The first analytical algorithm on the SES area was proposed by Connolly.<sup>32</sup> It was implemented into AMS and later other programs such as PQMS<sup>52</sup> and Amber/MOLSURF.<sup>53</sup> Due to the complexity of the SES definition and the resulting self-intersecting regions of the surface, there had been no improvements on the analytical algorithm until Sanner et al.'s work in 1996.<sup>54</sup> Sanner et al. proposed a new method to detect the self-intersecting singularities with the assistance of the reduced surface elements.<sup>54</sup> Although the new method may successfully compute the areas for most molecules, it does encounter errors in certain molecular structures and has to be restarted with modified atomic radii.<sup>54</sup> A different analytical approach was derived from the  $\alpha$  shape theory in computational geometry.<sup>55,56</sup> It uses Delaunay complexes and their filtrations to describe the topological structure of a molecule, which can facilitate the computation of the surface area by removing redundant terms in the direct inclusion–exclusion method.<sup>38</sup> In contrast, more methods exist to compute the SES area numerically.<sup>47,57–66</sup> Available programs include GEOPOL,<sup>47</sup> MASKER,<sup>66</sup> MOLSURF,<sup>65</sup> and USURF.<sup>64</sup> Numerical methods use geometrical objects, such as dots, triangles, cubes, and polyhedrons, to tessellate the surface or fill the interior volume. They more easily circumvent the pathological singularity cases but inevitably sacrifice accuracy.

In this work, we propose a numerical algorithm to compute the molecular surface area or surface integration for both the SES and SAS definitions. Different from previous numerical approaches of a geometrical nature, our algorithm is physics-inspired. The algorithm is implemented as an integral part of our PBSA program<sup>67–70</sup> and is based upon the finite-difference Poisson–Boltzmann (FDPB) methods where a grid labeling step, i.e., mapping the molecular surface to the grid points, is necessary before computing the surface area. The pioneer work on the grid labeling was proposed by You and Bashford<sup>71</sup> and later improved upon by Rocchia et al. in efficiency.<sup>72</sup> Our mapping strategy followed Rocchia et al.'s basic idea but tried to record additional information, such as the locations of the intersection points of the grid and the analytical surface, for the subsequent determination of dielectric constants.<sup>73</sup> In the following, we first go over the development of our algorithm. This is followed by detailed numerical tests to validate the accuracy, convergence, and timing of the numerical algorithm.

## METHODS

**Overview of the Algorithm.** Our algorithm is based on the observation that a molecular surface, either in the SES definition or in the SAS definition, can be treated as a union of partial spheres of different centers and radii. The SAS is strictly a union of solvent exposed portions of extended van der Waals spheres, i.e., spheres of van der Waals radii augmented by the solvent probe radius. The SES is more complex with two types of surfaces—the contact surface that consists of solvent exposed portions of van der Waals spheres and the re-entrant surface that consists of inward-facing portions of the solvent probe sphere as it rolls over the molecule. The re-entrant surface can be further classified into

two types according to how many atoms the solvent probe is in contact with simultaneously: (1) saddle surfaces if two atoms are in contact with the probe and (2) spherical triangles if three atoms are in contact with the probe. A re-entrant surface formed by the probe's concurrent contact with more than three atoms can always be divided into multiple spherical triangles. Contact surfaces and re-entrant spherical triangles are partial spheres, but re-entrant saddle surfaces are clearly not. However, the latter can be considered as consisting of small partial spheres of the solvent probe at different probe sites located on discretized solvent accessible arcs from the numerical point of view, for example in numerical surface definitions proposed for FDPB.<sup>71,72</sup> Apparently, in the numerical surface definitions, the number of the solvent probe sites has to be finite. Therefore, at each site, the probe is responsible for an area on the approximate saddle surface in the shape of a partial sphere.<sup>71,72</sup> Given the numerical representation of the molecular surface, if it is possible to compute the surface area/surface integration for each partial sphere, no matter how small it may be, the total surface area/surface integration is simply the sum of the contributions from all of the partial spheres.

In the following, we first describe the terminologies used in the finite-difference discretization. Then, we validate our aforementioned assumption that the saddle surfaces in the SES definition can be numerically treated as unions of partial spheres. This is followed by how to compute the surface area with our physics-inspired strategy and its numerical implementation for the finite-difference discretization. Finally, its extension to the surface integrations of both scalar and vector fields is discussed.

**Finite-Difference Discretization.** Without a loss of generality, we focus on Poisson's equation in this study since the Boltzmann term is nonzero only outside the Stern layer, which is typically set 2 Å away from the molecular surface. The partial differential equation

$$\nabla \cdot \epsilon \nabla \phi = -\rho \quad (1)$$

establishes a relation between the charge density ( $\rho$ ) and the electrostatic potential ( $\phi$ ) given a predefined dielectric distribution function ( $\epsilon$ ) for a solvated molecule.

A commonly used numerical method to solve Poisson's equation is the finite-difference method where a uniform Cartesian grid is used to discretize a rectangular box containing the molecule. The grid points are numbered as  $(i, j, k)$ ,  $i = 1, \dots, xm$ ;  $j = 1, \dots, ym$ ; and  $k = 1, \dots, zm$ , where  $xm$ ,  $ym$ , and  $zm$  are the numbers of grid points along the  $x$ ,  $y$ , and  $z$  axes, respectively. The spacing between two neighboring grid points is uniformly set to be  $h$ . With the finite-difference discretization, eq 1 can be written as

$$\begin{aligned} & \epsilon_{i-1/2,j,k}(\phi_{i-1,j,k} - \phi_{i,j,k}) + \epsilon_{i+1/2,j,k}(\phi_{i+1,j,k} - \phi_{i,j,k}) \\ & + \epsilon_{i,j-1/2,k}(\phi_{i,j-1,k} - \phi_{i,j,k}) + \epsilon_{i,j+1/2,k}(\phi_{i,j+1,k} - \phi_{i,j,k}) \\ & + \epsilon_{i,j,k-1/2}(\phi_{i,j,k-1} - \phi_{i,j,k}) + \epsilon_{i,j,k+1/2}(\phi_{i,j,k+1} - \phi_{i,j,k}) \\ & = -q_{i,j,k}/h \end{aligned} \quad (2)$$

Use of eq 2 requires the dielectric constant  $\epsilon$  to be defined at the grid edge centers between two neighboring grid points (denoted as grid index  $\pm 1/2$  above). In this study, the dielectric constant is defined in the following way: if the grid edge center is in the solute, the dielectric constant is equal to the solute permittivity otherwise,

the dielectric constant is equal to the solvent permittivity. It also requires mapping the point charges onto the grid points. The solution is the potentials on the grid points. More detailed implementation information specific to this study can be found in our recent publications.<sup>5,67-70</sup>

**Spherical Representation of Saddle Surfaces.** Suppose two atomic spheres of radii  $r_1$  and  $r_2$ , respectively, are in contact with the solvent probe sphere of radius  $r_p$ . The distance between the two atoms is  $d$ . In a typical finite-difference scheme, only a finite number of solvent probe sites are used to represent the solvent accessible circle formed by the two atoms. Here, the resolution of the discretized solvent accessible circle in radian is denoted by  $\psi$ . Each probe contributes a small partial sphere to the spherical representation of the saddle surface between the two atoms, and the area of each partial sphere can be computed with the following integral, provided that the centers of the two atoms are on the  $z$  axis:

$$S_\psi = \int_{-b}^a \left( 2 \arcsin \frac{l \sin(\psi/2)}{\sqrt{r_p^2 - z^2}} - \psi \right) r_p dz \quad (3)$$

where  $l$  is the radius of the discretized circle of probe sites and  $l = \{(r_p + r_1)^2 - [(d/2) + (\beta/d)]^2\}^{1/2}$ ,  $a = r_p[(d/2) + (\beta/d)]/(r_p + r_1)$ ,  $b = r_p[(d/2) + (\beta/d)]/(r_p + r_2)$ , and  $\beta = (2r_p + r_1 + r_2)(r_1 - r_2)/2$ . To evaluate the error of the above approximation of the corresponding analytical saddle surface area, we expand eq 3 with the Taylor series at  $\psi = 0$  and obtain

$$\begin{aligned} S_\psi &= \psi r_p l \left( \arcsin \frac{a}{r_p} + \arcsin \frac{b}{r_p} \right) - \psi r_p (a + b) \\ &+ \frac{\psi^3 r_p l}{24} \left( \frac{l^2}{r_p^2} \left( \frac{a}{\sqrt{r_p^2 - a^2}} + \frac{b}{\sqrt{r_p^2 - b^2}} \right) - \arcsin \frac{a}{r_p} \right. \\ &\left. - \arcsin \frac{b}{r_p} \right) + O(\psi^5) \end{aligned} \quad (4)$$

where the first two terms are the exact area of the saddle surface. Therefore, the third term is the leading error of the spherical representation of the saddle surface. Omitting the higher-order terms in eq 4, the relative error can be computed as

$$\begin{aligned} \Delta_\psi &= \frac{\psi^2}{24} \left( \frac{l^2}{r_p^2} d - l \arcsin \frac{a}{r_p} - l \arcsin \frac{b}{r_p} \right) \\ &/ \left( l \arcsin \frac{a}{r_p} + l \arcsin \frac{b}{r_p} - a - b \right) \end{aligned} \quad (5)$$

Equation 5 demonstrates that the spherical representation of the saddle surface is second-order accurate. It can be shown that the relative error of the approximate saddle surface area reaches the maximum when one atom is buried in the other, and the relative error becomes

$$\begin{aligned} \Delta_{\psi, \max} &= \frac{\psi^2 (r_1 r_2 + r_1 r_p + r_2 r_p) (r_1 + r_p) (r_2 + r_p)}{12 r_p^2 (2 r_1 r_2 + r_1 r_p + r_2 r_p)} \\ &= \frac{\psi^2 (mn + m + n) (1 + m) (1 + n)}{12 (2mn + m + n)} \end{aligned} \quad (6)$$

where  $m = r_1/r_p$  and  $n = r_2/r_p$ . In typical molecular simulations,  $r_p = 1.4 \text{ \AA}$ , and  $r_1$  and  $r_2$  are between 1 and 2  $\text{\AA}$ . Since  $\Delta_{\psi, \max}$  is monotonically increasing as either  $m$  or  $n$  increases, the maximum relative error can be estimated as

$$\Delta_{\psi, \max} = 0.35 \psi^2, \text{ if } m = 1.43 \text{ and } n = 1.43 \quad (7)$$

As the distance  $d$  increases, the saddle surface reduces to two self-intersecting parts, to which eqs 3–5 can no longer be applied. Here, the singularity in the SES definition is not considered because it is much less common than the saddle surface. However, it can be shown that the maximum relative error of using discrete probe arcs in the singular parts is actually smaller than the error given in eq 7.

Next, we further assess the error quantitatively in the framework of the PBSA program. The density of the solvent probe sites in the PBSA program is determined by the *arces* option, defined as the arc length between two neighboring probe sites. The relation of *arces* and the probe site resolution  $\psi$  introduced here is

$$\frac{2\pi l}{\text{arces}} = \frac{2\pi}{\psi} \text{ or } \psi = \frac{\text{arces}}{l} \quad (8)$$

Provided that the self-intersecting region is not considered,  $l = r_p = 1.4$  at the upper bound of  $d$ . On the other side, as the two atoms approach each other,  $l$  decreases and reaches zero when one atom is buried in the other. To estimate the error in the spherical representation of the saddle surface in realistic situations, we further assume that the lower bound of  $d$  also gives  $l = r_p = 1.4$ . This corresponds to an interatomic distance larger than the difference of the radii of the two atoms by at most 0.15  $\text{\AA}$ . This assumption is reasonable according to the atomic radii and bond lengths in the AMBER force fields.<sup>74</sup> The default *arces* in the PBSA program is 1/16  $\text{\AA}$ , corresponding to

$$\Delta_{\psi, \max} = 0.35 \psi^2 \leq 0.35 \left( \frac{1/16}{1.4} \right)^2 = 6.98 \times 10^{-4} \quad (9)$$

Therefore, the spherical representation of saddle surfaces is an appropriate approximation in the calculation of the molecular surface area. Note that the maximum error in eq 9 is underestimated when a smaller probe is used because the curvature will need more grid points to resolve, which is a general problem in the finite-difference scheme.

Now that both the SAS and SES can be treated as the union of partial spheres with good accuracy, we are ready to introduce our new algorithm to compute the surface areas of partial spheres.

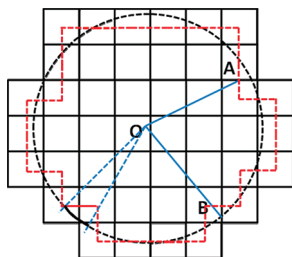
**Field-View Method.** Suppose there is a point charge  $Q$  at the center of a sphere of radius  $R$  in the vacuum. The flux  $\Phi$  of the electric displacement through any closed surface containing the charge is given by the integral form of Gauss's law, which is equal to the total free charge inside, i.e.,

$$\Phi = Q \quad (10)$$

By definition we have

$$\Phi = \iint_S \epsilon_0 \vec{E} \cdot \hat{n} dS \quad (11)$$

where  $\epsilon_0$  is the vacuum permittivity,  $\vec{E}$  is the vacuum electric field, and  $\hat{n}$  is the normal direction of the infinitesimal surface element  $dS$ . Given the central symmetry of the radial field due to a point charge at the spherical center, the flux on a spherical surface at



**Figure 1.** A 2-D diagram of the finite-difference discretization of a sphere (black dashed circle). Black solid lines are grid edges, whose intersection points denote grid points. O is the center of the sphere. A is the center of a square surface element or a grid edge. B is an intersection point of the sphere and a grid edge. Red dashed lines denote square surface elements at grid edge centers. The red dashed lines can also be viewed as the finite-difference approximation of the spherical surface. The black solid arc represents a spherical surface element, and the red solid line represents a square surface element that subtends the same solid angle, so they have the same flux passing through them.

radius  $R$  can be computed with

$$\Phi = \epsilon_0 E S \quad (12)$$

where  $S$  is the total surface area of the sphere. This is also true for a spherical surface element  $\delta S$ :

$$\delta\Phi = \epsilon_0 E \delta S = \frac{Q \delta S}{4\pi R^2} \quad (13)$$

Instead of directly computing the flux through the spherical surface, we compute the flux via the finite-difference data structure by exploiting the conservation of electric flux. Figure 1 is a 2-D illustration of the finite-difference discretization of a sphere. Consider the closed surface represented by the red dashed lines in Figure 1. Note that each red dashed line represents a square surface element (area =  $h^2$ , with  $h$  being the finite-difference spacing) located at the center of a grid edge that is intersected by the sphere. Since the square surface elements at the grid edge centers form a closed surface, the flux through the closed surface is equal to the flux through the spherical surface. For each square surface element, there is always a spherical surface element that subtends the same solid angle, so that they have the same flux passing through them (see Figure 1 for the two types of surfaces with the same flux: spherical surface element (black solid arc) and square surface element (red solid line)). This is the key relation that is to be exploited below to infer the surface area of the sphere or the spherical surface element.

The flux through each square surface element, as well as the closed surface formed by the square surface elements, can be computed with the double integrals as follows:

$$\begin{aligned} \delta\Phi_i &= \iint_{S_i} \epsilon_0 \vec{E}_i \cdot \hat{n}_i \, dS = \iint_{S_i} \epsilon_0 E_i \cos \theta_i \, dS \\ \delta\Phi_j &= \iint_{S_j} \epsilon_0 \vec{E}_j \cdot \hat{n}_j \, dS = \iint_{S_j} \epsilon_0 E_j \cos \theta_j \, dS \\ \delta\Phi_k &= \iint_{S_k} \epsilon_0 \vec{E}_k \cdot \hat{n}_k \, dS = \iint_{S_k} \epsilon_0 E_k \cos \theta_k \, dS \\ \Phi &= \sum_{N_i} \delta\Phi_i + \sum_{N_j} \delta\Phi_j + \sum_{N_k} \delta\Phi_k \end{aligned} \quad (14)$$

where the subscripts  $i$ ,  $j$ , and  $k$  are the indices in the  $x$ ,  $y$ , and  $z$  directions, respectively;  $\delta\Phi_i$ ,  $\delta\Phi_j$ , and  $\delta\Phi_k$  are the fluxes

through the square surface elements at the grid edge centers in the  $x$ ,  $y$ , and  $z$  directions, respectively;  $N_i$ ,  $N_j$ , and  $N_k$  are the numbers of the square surface elements in the  $x$ ,  $y$ , and  $z$  directions, respectively;  $\vec{E}_i$ ,  $\vec{E}_j$ , and  $\vec{E}_k$  are the electrical fields at the grid edge centers (or the centers of the square surface elements) in the  $x$ ,  $y$ , and  $z$  directions, respectively; and  $\hat{n}_i$ ,  $\hat{n}_j$ , and  $\hat{n}_k$  are the unit normal vectors of the square surface elements in the  $x$ ,  $y$ , and  $z$  directions, respectively. Finally,  $\theta_i$  is the angle between  $\vec{E}_i$  and  $\hat{n}_i$ ,  $\theta_j$  is the angle between  $\vec{E}_j$  and  $\hat{n}_j$ , and  $\theta_k$  is the angle between  $\vec{E}_k$  and  $\hat{n}_k$ .

On the  $i$ th square surface element in the  $x$  direction, angle  $\theta_i$  between the direction of the radial electric field and the normal direction of the square surface element can be computed with

$$\cos \theta_i = \frac{|x_i|}{r_i} \quad (15)$$

where  $x_i$  is the  $x$  coordinate of the center of the  $i$ th square surface element in the  $x$  direction, and  $r_i$  is the distance between the center of the square surface element and the center of the sphere (see the blue solid line OA in Figure 1, given the center of the sphere at the origin). Similarly in the  $y$  and  $z$  directions, we have

$$\cos \theta_j = \frac{|y_j|}{r_j}, \cos \theta_k = \frac{|z_k|}{r_k} \quad (16)$$

Thus, the electric flux through the  $i$ th square surface element in the  $x$  direction at  $(x_i, y_i, z_i)$  is

$$\begin{aligned} \delta\Phi_i &= \int_{z_i - h/2}^{z_i + h/2} \int_{y_i - h/2}^{y_i + h/2} \epsilon_0 E_i(x_i, y', z') \frac{|x_i|}{r_i(x_i, y', z')} \, dy' \, dz' \\ &= \int_{-h/2}^{h/2} \int_{-h/2}^{h/2} \frac{Q|x_i|}{4\pi[x_i^2 + (y_i + u)^2 + (z_i + v)^2]^{3/2}} \, du \, dv \end{aligned} \quad (17)$$

Applying the Taylor series expansion to eq 17 at the center of the square surface element, i.e.,  $(x_i, y_i, z_i)$ , it can be simplified as

$$\begin{aligned} \delta\Phi_i &= \frac{Q|x_i|}{4\pi} \int_{-h/2}^{h/2} \int_{-h/2}^{h/2} \left( \frac{1}{r_i^3} + \frac{3}{2} \left( \frac{5y_i^2 u^2 + 5z_i^2 v^2 - u^2 + v^2}{r_i^7} - \frac{u^2 + v^2}{r_i^5} \right) \right. \\ &\quad \left. + O\left(\frac{h^4}{r_i^4}\right) \right) \, du \, dv = \frac{Q|x_i|h^2}{4\pi r_i^3} \left( 1 + \left( \frac{3}{8} - \frac{5x_i^2}{8r_i^2} \right) \frac{h^2}{r_i^2} \right. \\ &\quad \left. + O\left(\frac{h^4}{r_i^4}\right) \right) \end{aligned} \quad (18)$$

where  $r_i = (x_i^2 + y_i^2 + z_i^2)^{1/2}$ . All of the odd functions in the expansion in eq 18 disappear due to the symmetry of the integral interval. The double integrals in the  $y$  and  $z$  directions in eq 14 can be expanded in the same way. Thus

$$\delta\Phi_j = \frac{Q|y_j|h^2}{4\pi r_j^3} \left( 1 + \left( \frac{3}{8} - \frac{5y_j^2}{8r_j^2} \right) \frac{h^2}{r_j^2} + O\left(\frac{h^4}{r_j^4}\right) \right) \quad (19)$$

$$\delta\Phi_k = \frac{Q|z_k|h^2}{4\pi r_k^3} \left( 1 + \left( \frac{3}{8} - \frac{5z_k^2}{8r_k^2} \right) \frac{h^2}{r_k^2} + O\left(\frac{h^4}{r_k^4}\right) \right) \quad (20)$$

Substitution of eqs 18, 19, and 20 into eq 14 and omission of fourth or higher order terms give

$$\begin{aligned}\delta\Phi_i &= \frac{Q|x_i|h^2}{4\pi r_i^3}(1 + \beta_i) \\ \delta\Phi_j &= \frac{Q|y_j|h^2}{4\pi r_j^3}(1 + \beta_j) \\ \delta\Phi_k &= \frac{Q|z_k|h^2}{4\pi r_k^3}(1 + \beta_k) \\ \Phi &= \frac{Qh^2}{4\pi} \left( \sum_{N_i} \frac{|x_i|}{r_i^3}(1 + \beta_i) + \sum_{N_j} \frac{|y_j|}{r_j^3}(1 + \beta_j) \right. \\ &\quad \left. + \sum_{N_k} \frac{|z_k|}{r_k^3}(1 + \beta_k) \right) \quad (21)\end{aligned}$$

where  $\beta_i = ((3/8) - (5x_i^2)/(8r_i^2))(h^2/r_i^2)$ ,  $\beta_j = ((3/8) - (5y_j^2)/(8r_j^2))(h^2/r_j^2)$ , and  $\beta_k = ((3/8) - (5z_k^2)/(8r_k^2))(h^2/r_k^2)$  are the coefficients of the second-order terms. Comparison between eqs 12 or 13 and 21 shows that the spherical surface element areas and the total surface area of the sphere can be written as

$$\begin{aligned}\delta S_i &= \frac{|x_i|h^2R^2}{r_i^3}(1 + \beta_i) \\ \delta S_j &= \frac{|y_j|h^2R^2}{r_j^3}(1 + \beta_j) \\ \delta S_k &= \frac{|z_k|h^2R^2}{r_k^3}(1 + \beta_k) \\ S &= h^2R^2 \left( \sum_{N_i} \frac{|x_i|}{r_i^3}(1 + \beta_i) \right. \\ &\quad \left. + \sum_{N_j} \frac{|y_j|}{r_j^3}(1 + \beta_j) + \sum_{N_k} \frac{|z_k|}{r_k^3}(1 + \beta_k) \right) \quad (22)\end{aligned}$$

Since the SES and the SAS are both composed of partial spheres (the SAS is the union of atomic spheres and the SES consists of atomic spheres and probe spheres), the above algorithm can be directly used to evaluate the SES and SAS areas numerically. In molecular applications, eq 22 becomes

$$\begin{aligned}\delta S_i &= \frac{|x_i|h^2R_i^2}{r_i^3}(1 + \beta_i) \\ \delta S_j &= \frac{|y_j|h^2R_j^2}{r_j^3}(1 + \beta_j) \\ \delta S_k &= \frac{|z_k|h^2R_k^2}{r_k^3}(1 + \beta_k) \\ S &= h^2 \left( \sum_{N_i} \frac{|x_i|R_i^2}{r_i^3}(1 + \beta_i) + \sum_{N_j} \frac{|y_j|R_j^2}{r_j^3}(1 + \beta_j) \right. \\ &\quad \left. + \sum_{N_k} \frac{|z_k|R_k^2}{r_k^3}(1 + \beta_k) \right) \quad (23)\end{aligned}$$

where  $x_i$ ,  $y_j$ , and  $z_k$  are the relative coordinates of the  $i$ th,  $j$ th, and  $k$ th square surface element centers (or the grid edge centers) in the  $x$ ,  $y$ , and  $z$  directions, respectively, and  $R_i$ ,  $R_j$ , and  $R_k$  are the radii of the atomic/probe spheres that intersect the  $i$ th,  $j$ th, and

$k$ th grid edges in the  $x$ ,  $y$ , and  $z$  directions, respectively. All relative coordinates are with respect to the atomic/probe sphere centers. We term the algorithm the field-view method in this study.

The second-order terms in eq 23 show that the zeroth-order truncation is second-order accurate and converges quadratically. In theory, inclusion of the second-order terms results in a convergence of fourth order, but actually it is not guaranteed to achieve more accurate result in the re-entrant region at coarse grid spacings. This is because at coarse grid spacings, the leading error comes from the spherical representation of saddle surfaces, which will be shown in the Results and Discussion.

**Surface Integration.** Given the spherical surface element areas in eq 23, it is straightforward to evaluate a surface integral on the SES or SAS numerically. With the field-view method, the surface integral of any vector field  $\vec{A}$  defined on the surface can be written as

$$\begin{aligned}\iint_S \vec{A} \cdot d\vec{S} &\cong \sum_{N_i} A_i \delta S_i \cos \gamma_i + \sum_{N_j} A_j \delta S_j \cos \gamma_j \\ &+ \sum_{N_k} A_k \delta S_k \cos \gamma_k \cong h^2 \left( \sum_{N_i} \frac{A_i |x_i| R_i^2 (1 + \beta_i) \cos \gamma_i}{r_i^3} \right. \\ &\quad \left. + \sum_{N_j} \frac{A_j |y_j| R_j^2 (1 + \beta_j) \cos \gamma_j}{r_j^3} + \sum_{N_k} \frac{A_k |z_k| R_k^2 (1 + \beta_k) \cos \gamma_k}{r_k^3} \right) \quad (24)\end{aligned}$$

where  $\delta S_i$ ,  $\delta S_j$ , and  $\delta S_k$  are the spherical surface element areas introduced in the field-view method;  $A_i$ ,  $A_j$ , and  $A_k$  are the field magnitudes on the spherical surface elements  $\delta S_i$ ,  $\delta S_j$ , and  $\delta S_k$ , respectively;  $\gamma_i$ ,  $\gamma_j$ , and  $\gamma_k$  are the angles between the vector field  $\vec{A}$  and the normal directions of the spherical surface elements  $\delta S_i$ ,  $\delta S_j$ , and  $\delta S_k$ , respectively. For any scalar field  $A$  defined on the surface, eq 24 becomes

$$\begin{aligned}\iint_S A \, dS &\cong \sum_{N_i} A_i \delta S_i + \sum_{N_j} A_j \delta S_j + \sum_{N_k} A_k \delta S_k \\ &\cong h^2 \left( \sum_{N_i} \frac{A_i |x_i| R_i^2 (1 + \beta_i)}{r_i^3} + \sum_{N_j} \frac{A_j |y_j| R_j^2 (1 + \beta_j)}{r_j^3} \right. \\ &\quad \left. + \sum_{N_k} \frac{A_k |z_k| R_k^2 (1 + \beta_k)}{r_k^3} \right) \quad (25)\end{aligned}$$

For example, if the finite-difference method is used to solve Poisson's equation or the Poisson–Boltzmann equation, eqs 24 and 25 can be used to compute the hydrophobic solvation free energy that is often modeled as being linearly proportional to the surface area.

**Computational Details.** The field-view method was implemented into the Amber/PBSA program.<sup>53,67,69</sup> Implicit solvent models under the finite-difference scheme require the molecular surface be mapped to the grid points using a grid labeling algorithm.<sup>73</sup> After mapping, the dielectric constant was set to be 1 in the solute region and 80 in the solvent region. On the dielectric boundary, the dielectric constant was set to be either 1 or 80 depending on whether the grid edge center was in the solute or solvent. The solvent probe was set to be 1.4 Å.

All of the molecules in our tests come from the PBSA test set, which includes 579 proteins,<sup>53</sup> 364 nucleic acids,<sup>70</sup> and 622 protein–protein complexes (see the Appendix). However, 10 of them have been left out due to the failures in the reference surface area program Amber/MOLSURF.<sup>53</sup> The MOLSURF program



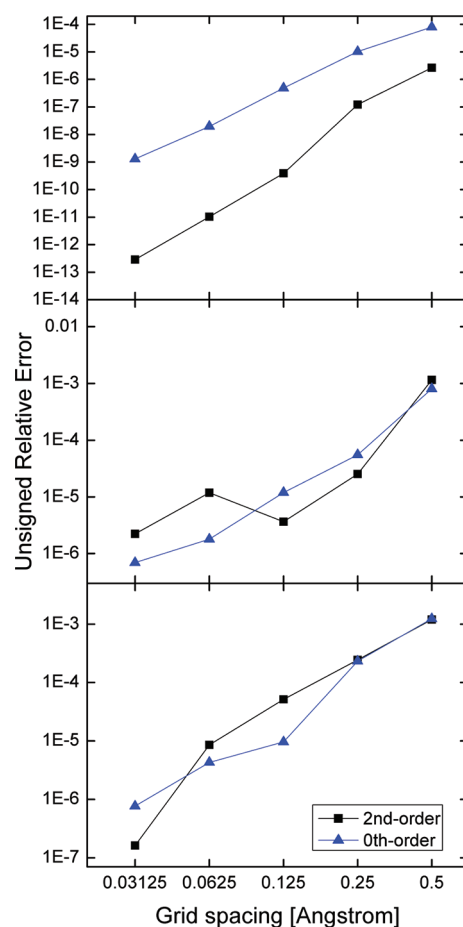
was implemented by Beroza according to Connolly's analytical algorithm.<sup>32</sup> When possible, the orientation and the origin of the finite-difference grid were randomized 100 times for each finite-difference run to reduce the numerical uncertainty of reported values.

## RESULTS AND DISCUSSION

**Simple Geometries: Agreement with Exact Analytical Solutions.** To assess the accuracy of the proposed algorithm, we first used a single sphere, double spheres, and triple spheres as test cases because their analytical surface areas can be readily calculated. The algorithm was tested at five different grid spacings, 1/2 Å, 1/4 Å, 1/8 Å, 1/16 Å, and 1/32 Å, to study its convergence behavior. The *arces* option (the arc length between two neighboring probe sites) was set to be as fine as 0.01 Å to reduce the error introduced by limited resolution of solvent accessible arcs.

Figure 2 shows the unsigned relative errors of the numerical SES areas. It can be seen that at the finest grid spacing, the relative errors of the numerical SES areas reach as low as  $10^{-9}$  for the zeroth-order truncation and  $10^{-13}$  for the second-order truncation in the test case of the single sphere, but the relative errors only reach  $10^{-6}$  to  $10^{-7}$  in the test cases of the double spheres and triple spheres. This is due to the finite resolution of the numerical representation of the re-entrant regions as a union of partial spheres. For example, eq 5 can be used to estimate the unsigned relative error of the re-entrant region in the double spheres, which is  $4.33 \times 10^{-6}$  and consistent with the numerical analysis. The standard deviations in the numerical SES areas (not shown in the plot) are larger than or comparable with the errors of the means at all grid spacings in these simple test cases.

It is also interesting to note that the second-order truncation exhibits better accuracy at the tested grid spacing (1/2 Å) in the single sphere case only. The reason is that there are not many square surface elements in the re-entrant region at coarse grid spacings (such as the tested 1/2 Å), so that only a fraction of solvent probe sites are employed even if a very high resolution is used in the discretized solvent accessible arcs. As a consequence, the error due to the spherical representation of saddle surfaces is dominant. For example, the number of square surface elements in the re-entrant region of the double-sphere test case is around 60, much smaller than the number of available solvent probe sites, about 1300, at the default computation setting in the Amber/PBSA program. In general, the more square surface elements that exist in a saddle surface, the more solvent probe sites that can be used. This is because a solvent probe contributes to the approximate saddle surface via a spherical surface element that shares the same flux as a square surface element, and the square surface element is exclusively assigned to the solvent probe at a certain site. According to eq 7, when the area of the re-entrant region reaches its maximum (the radii of the two atoms are both 2 Å; the distance is about 4.4 Å) and at the same time the number of square surface elements reaches its maximum (about 150), the lowest possible error is about 0.06%. In contrast, the truncation error of each surface element is about 1–2% for the zeroth-order truncation and about 0.04–0.06% for the second-order truncation. These error estimations were obtained by comparison between the truncated Taylor series and the exact analytical expression for a surface element area (eq 17). Although the truncation error of each surface element at the zeroth truncation level (1–2%) is larger than the error in the

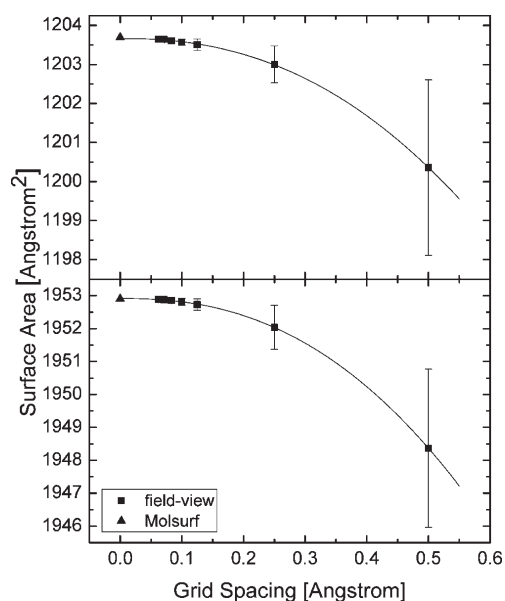


**Figure 2.** Unsigned relative errors of the numerical SES areas of the simple geometries at successively fine grid spacings. Top: A single sphere. The radius of the sphere is 1.5 Å. Middle: Double spheres. The radii of the two spheres are both 1.5 Å, and the distance between the centers of the two spheres is 4 Å. Bottom: Triple spheres. The radii of the three spheres are all 1.5 Å, and the distance between one pair of spherical centers is 4 Å. The distances between the other two pairs of spherical centers are both 3.606 Å. The results are obtained from 100 area calculations with randomized grid orientations.

spherical representation of saddle surfaces (can be as low as 0.06%), the errors of the surface elements tend to cancel each other, leading to a much smaller overall error in the total surface area. As shown in Figure 2 (top), the overall error in the surface area of a single sphere is  $\sim 0.01\%$  at a grid spacing of 1/2 Å (the error cancellation may not be so dramatic where there is no central symmetry). As a result, the error in the spherical representation of saddle surfaces becomes the leading error term.

Thus, at a grid spacing of 1/2 Å, numerical surface areas computed by the zeroth-order truncation are similarly accurate to those computed by higher-order truncations. For example, in the double-sphere model (Figure 2 (middle)), the error of the zeroth-order truncation is similar to the second-order truncation ( $\sim 0.1\%$ ). Therefore, the leading error shown here is actually from the spherical representation of saddle surfaces. Thus, in the following tests, only the zeroth-order truncation was used to compute surface areas.

At fine grid spacings when there are abundant square surface elements in the re-entrant region, more solvent probe sites can be

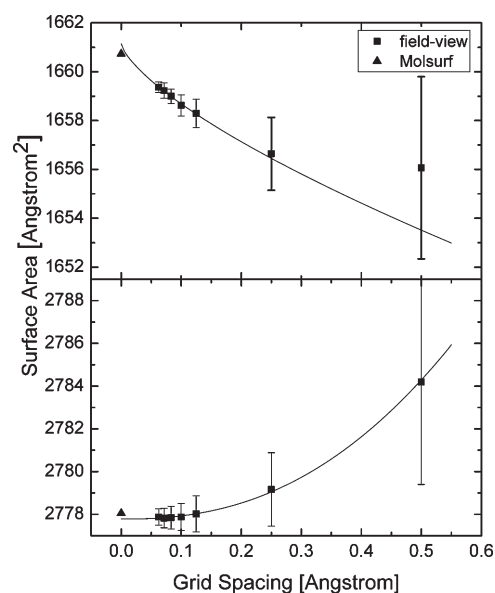


**Figure 3.** Convergence of the numerical SES areas versus grid spacings for 1BRV and 1FN2, respectively. Top: 1BRV (the field-view method converges to  $1203.66 \text{ \AA}^2$ , the power order is 2.30, and the MOLSURF result is  $1203.69 \text{ \AA}^2$ ). Bottom: 1FN2 (the field-view method converges to  $1952.91 \text{ \AA}^2$ , the power order is 2.37, and the MOLSURF result is  $1952.90 \text{ \AA}^2$ ). The uncertainty bars are estimated as the standard deviations from 100 FDPB calculations with randomized grid orientations. The uncertainty bars for finer grid spacings are too small to be seen.

utilized to resolve the re-entrant surface. Under this circumstance, the zeroth-order truncation may have larger error than the spherical representation of saddle surfaces. For example, at a grid spacing of  $1/8 \text{ \AA}$ , the number of square surface elements in the saddle surface of the double spheres increases to more than 1000, compared to around 60 at a grid spacing of  $1/2 \text{ \AA}$ . Consequently, at a grid spacing of  $1/8 \text{ \AA}$ , the second-order truncation performs better than the zeroth-order truncation in the double sphere test case (see Figure 2 (middle)). In this case, the higher-order terms generally help improve the accuracy of surface areas. Therefore, we suggest that higher-order truncations be used at fine grid spacings for better accuracy and the resolution of solvent probe sites be dependent on grid spacings for the best performance and efficiency.

**Convergence Tests on Realistic Biomolecules.** Next, we analyzed the convergence behavior of the numerical algorithm with two realistic but small biomolecules: 1BRV (268 atoms) and 1FN2 (348 atoms). Both the SES and SAS areas were computed at successively fine grid spacings, and the results are shown in Figures 3 and 4, respectively. The *arces* option was set to be  $1/16 \text{ \AA}$  in this test, and only the zeroth-order truncation was used. The surface areas computed by the Amber/MOLSURF program that implements Connolly's algorithm were used as a reference.<sup>32,53</sup>

As shown in Figure 3, the performance of the field-view method in the calculation of SES areas of realistic biomolecules is promising. The unsigned relative error is around 0.2% at the grid spacing of  $1/2 \text{ \AA}$  often used in FDPB calculations of biomolecules. We also extrapolated the "converged" value using the nonlinear curve fitting method (with the formula  $y = a + bh^c$ , where  $a$  is the predicted surface area when the grid spacing  $h$  goes to zero and  $c$  is the convergence rate). Obviously, the converged surface areas are highly consistent with the Amber/MOLSURF

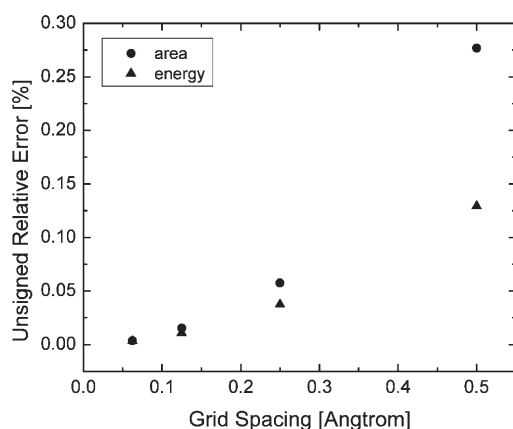


**Figure 4.** Convergence of the numerical SAS areas versus grid spacings for 1BRV and 1FN2, respectively. Top: 1BRV (the field-view method converges to  $1661.15 \text{ \AA}^2$ , the power order is 0.70, and the MOLSURF result is  $1660.73 \text{ \AA}^2$ ). Bottom: 1FN2 (the field-view method converges to  $2777.78 \text{ \AA}^2$ , the power order is 2.36, and the MOLSURF result is  $2778.05 \text{ \AA}^2$ ). Note that for 1BRV, only the data at fine grid spacings can be used in the extrapolation of the surface area.

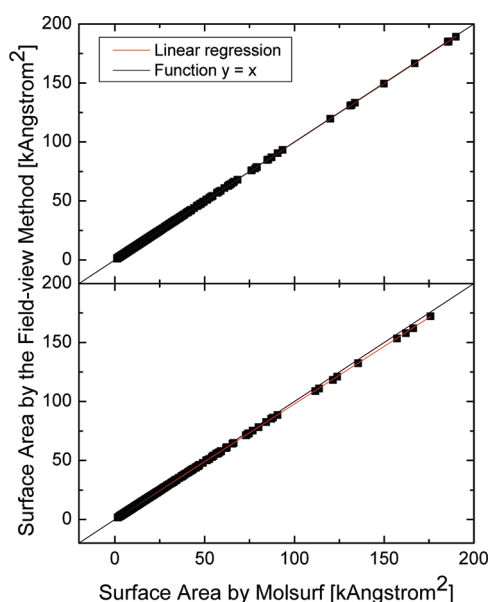
results. Interestingly, the convergence rate is somewhat higher than quadratic. This is because the error at a grid spacing of  $1/2 \text{ \AA}$  is mainly due to the spherical representation of saddle surfaces, larger than the truncation error of the zeroth-order truncation, while at fine grid spacings, the truncation error dominates, which is quadratic, as analyzed in the Methods section. Figure 3 also shows that the standard deviations for both molecules over 100 random grid orientations are smaller than the errors of the means at all grid spacings.

Figure 4 shows that the field-view method performs less well in the calculation of SAS areas. Although the unsigned relative error is still at the same low level, the convergence rate is no longer guaranteed to be quadratic. The reason is probably that the finite-difference-based numerical algorithm cannot resolve the cusps inside grid cells that are ubiquitous all over the molecular surface in the SAS definition. The situation can be mitigated at finer grid spacings, but it still remains significant at the finest tested grid spacing of  $1/16 \text{ \AA}$ . The standard deviation of numerical SAS areas is also larger than that of numerical SES areas. Although the SES definition also has cusps or self-intersecting parts, they are much rarer and likely not the major source of errors in the numerical calculation. Indeed, the unsigned relative errors in SES areas computed with the field-view method are comparable between the three analytical test cases discussed above and the two biomolecules. Of course for the biomolecules, the errors are larger due to the higher proportion of the re-entrant surface.

To better appreciate the convergence quality of the new surface area integration, we obtained the numerical surface area and the numerical reaction field energy in the same context of FDPB calculations, and the comparison of their convergence behaviors is shown in Figure 5. The reaction field energy  $\Phi_{RC}$  was calculated by the product of atomic charges and polarization



**Figure 5.** Convergence rates of the numerical SES surface area and the numerical reaction field energy of 1BRV at successively fine grid spacings. Each result is obtained from 100 calculations with randomized grid orientations.



**Figure 6.** Correlation between the numerical surface areas of 1555 molecules computed at a grid spacing of 1/2 Å and the analytical surface areas computed with MOLSURF. Top: SES areas—AURE, 0.27%; slope, 0.99696; R-square, 1.00000. Bottom: SAS areas—AURE, 1.05%; slope, 0.97753; R-square, 0.99999.

charges, i.e.,

$$\Phi_{RC} = \frac{1}{2} \sum_n \sum_m \frac{q_m^{\text{atom}} q_n^{\text{pol}}}{r_{mn}} \quad (26)$$

where  $r_{mn}$  is the distance between the atomic charge ( $q_m^{\text{atom}}$ ) and the finite-difference polarization charge ( $q_n^{\text{pol}}$ ). Here,  $q_n^{\text{pol}}$  is computed by a finite volume integral within a finite-difference grid cell:

$$q_n^{\text{pol}} = h(6\phi_n - \sum_{l=1}^6 \phi_{n,l}) \quad (27)$$

where  $\phi_n$  is the potential at the center of the finite-difference grid cell,  $\phi_{n,l}$   $l = 1, 2, \dots, 6$  are the potentials at the centers of the six

**Table 1. Timing Analysis for FDPB Calculations and the Proposed “On-the-Fly” Surface Area Calculation of 2MRB<sup>a</sup>**

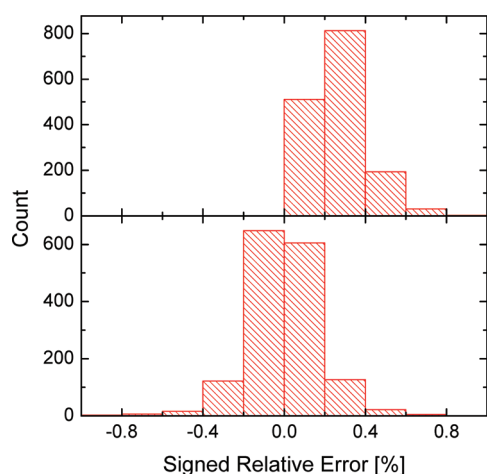
	SES	SAS
FDPB probe generation	52.80	N/A
FDPB grid labeling	38.17	9.63
FDPB solver	509.11	506.75
FDPB force calculation	5.76	5.01
surface area calculation	0.10	0.08
MOLSURF surface area calculation	3.38	3.87

<sup>a</sup>The timing for MOLSURF surface area calculation is also shown as a reference. Note that the reported times of FDPB calculations (in second) are for 100 calculations with randomized grid orientations. The MOLSURF calculation is also repeated 100 times.

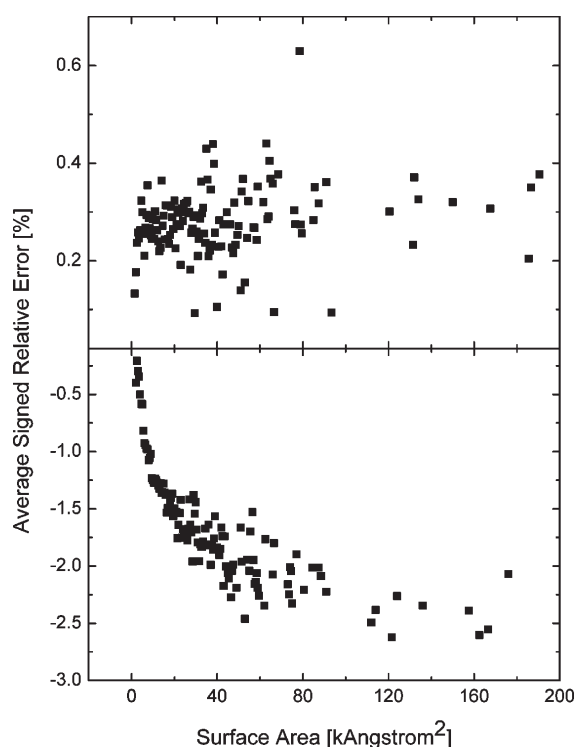
neighboring grid cells, and all of the potentials are obtained from the finite-difference solution of the Poisson–Boltzmann equation. It has been shown that the relative error in the surface area by the field-view method is  $O(h^2)$ . The error in the numerical reaction field energy primarily originates from the polarization charge. Note that the potential on the dielectric boundary from the FDPB method is first-order accurate, but the polarization charge is second-order accurate because the finite volume integration used to compute the polarization charges cancels out the leading error in the potential. To verify this analysis, we performed power curve fitting on the two data sets in Figure 5. The analysis shows that the power of the energy error curve is 1.79, and the power of the area error curve is 2.31 (the power curve here is a little different from that in Figure 3 (top) due to the use of different sets of data points in the curve fitting). In addition to the common quadratic convergence rate, the relative errors of the two calculations are also on the same order of magnitude.

**Consistency Tests.** Next, we tested the field-view method with a diversified set of 1555 molecules and molecular complexes to demonstrate its applicability and robustness in realistic biomolecular applications. All test cases were run at a grid spacing of 1/2 Å without random grid orientations. The correlation between the numerical molecular surface areas and the Amber/MOLSURF results is plotted in Figure 6. As in the previous tests, the numerical SES areas computed with the field-view method are overall more consistent with those with Amber/MOLSURF than the numerical SAS areas. It should be noted that the correlation coefficients are both very close to one, indicating that the numerical algorithm is consistent with the analytical algorithm regardless of the sizes or the structures of the tested molecules. The average unsigned relative error (AURE) of the numerical SES areas is 0.27%, and the AURE of the numerical SAS areas is 1.05%. The analysis demonstrates that the field-view method agrees well with the analytical method at the tested coarse grid spacing often used in biomolecular applications of FDPB.

**Timing Analysis.** It is interesting to analyze the timings for the proposed numerical algorithm. Table 1 lists the detailed timing analysis of a FDPB calculation with the field-view method turned on. The tested molecule was 2MRB (377 atoms), and the FDPB calculation was repeated 100 times with randomized grid orientations and origins. As expected, the time spent in computing the surface area is negligible compared to those used by other FDPB components in the calculation, supporting its “on-the-fly” application in FDPB calculations.

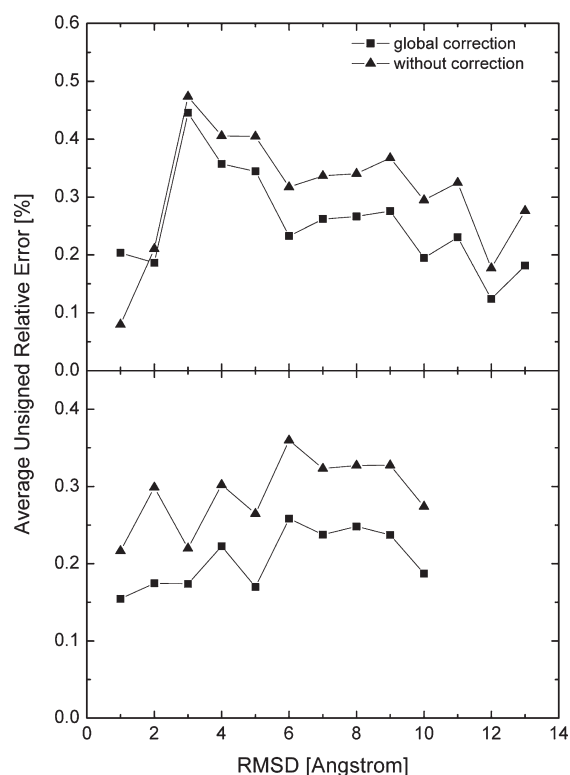


**Figure 7.** Distributions of the signed relative errors of the numerical SES areas computed at a grid spacing of  $1/2$  Å. Top: before correction, AURE, 0.27%. Bottom: after correction, AURE, 0.13%.



**Figure 8.** Correlation between the numerical surface areas and their signed relative errors. All of the numerical surface areas were computed at a grid spacing of  $1/2$  Å. Each data point represents the average signed relative error of the surface areas within a range of  $0.5$  kÅ<sup>2</sup>. Top: SES areas. Bottom: SAS areas.

**Limitations and Possible Remedies.** Despite the high-level consistency between the coarse-grid numerical calculations and analytical calculations, convergence errors do exist at coarse grid spacings. Interestingly, the convergence errors in the numerical SES areas are systematic and can be identified with a linear regression analysis. The distributions of the signed relative errors with respect to the analytical values before and after the linear correction are shown in Figure 7. It can be seen that the AURE is reduced after the correction by half from 0.27% to 0.13%.

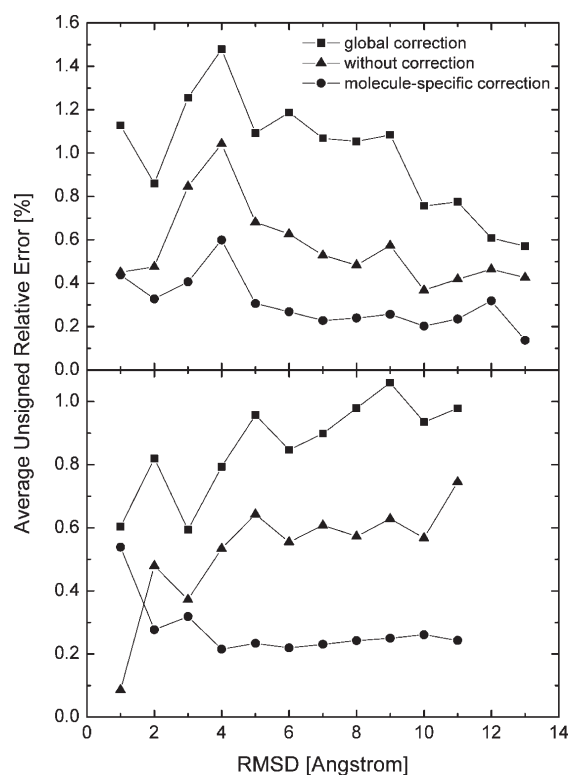


**Figure 9.** Average unsigned relative errors in the numerical SES areas with the global correction and without any correction versus the backbone root-mean-square deviation (RMSD) for the two tested peptides. Each data point represents the average unsigned relative error in the numerical surface areas for structures with the backbone RMSD within a range of  $1$  Å with respect to the crystal structure. Top: hairpin. Bottom: helix (grid spacing:  $1/2$  Å).

We did not apply the linear correction strategy to the numerical SAS areas on the basis of the following observations. Figure 8 shows the trend of the signed relative errors of the SES and SAS areas. If there is a systematic error that can be corrected with a linear correction, the signed relative errors should go to a constant as the surface area increases. If the linear correction is in the form of  $y_i = \beta_1 x_i + \beta_0 + e_i$ , where  $y_i$  is the analytical surface area,  $x_i$  is the numerical surface area,  $\beta_0$  and  $\beta_1$  are the regression estimators, and  $e_i$  is the residual error, the signed relative error  $\varepsilon_i$  of the numerical result  $x_i$  with respect to the analytical result  $y_i$  can be written as

$$\begin{aligned} \varepsilon_i &= \frac{x_i - y_i}{y_i} = \frac{(y_i - \beta_0 - e_i)/\beta_1 - y_i}{y_i} \\ &= \left( \frac{1}{\beta_1} - 1 \right) - \frac{\beta_0 + e_i}{\beta_1 y_i} \end{aligned} \quad (28)$$

The second term is negligible, and the signed relative error becomes constant if the residual error  $e_i$  does not have a positive correlation with  $y_i$ . Otherwise, the linear correction is unsuitable for the numerical method. It can be seen from Figure 8 that the signed relative errors of the numerical SES areas finally become stable at about 0.3%, whereas the signed relative errors of the numerical SAS areas are always negative and keep going down. This suggests that the numerical SAS areas probably need molecule-specific corrections.



**Figure 10.** Average unsigned relative errors in the numerical SAS areas with the global correction, without any correction, and with molecule-specific corrections versus the backbone RMSD for the two tested peptides. Each data point represents the average unsigned relative error in the numerical surface areas for structures with the backbone RMSD within a range of 1 Å with respect to the crystal structure. Top: hairpin. Bottom: helix (grid spacing: 1/2 Å).

We present two examples to substantiate the above claims, each consisting of folded and unfolded structures extracted from molecular simulations of two peptides. First, the global corrections were applied to the numerical SES areas and numerical SAS areas of those structures, as shown in Figures 9 and 10, respectively. It is encouraging to note reduced errors in the numerical SES areas in both folded and unfolded structures even if these tested peptides and structures are outside the training molecules that were used to obtain the linear correction equation. In contrast, the global correction, parametrized from the numerical SAS areas of the 1555 molecules, does not work for the numerical SAS areas of the two tested peptides. Next, molecule-specific corrections were applied to the numerical SAS areas. Specifically, for each molecule, one-fifth of the structures from the trajectory were picked as the training set for parametrization of the linear correction, and the remaining structures were used as the test set to validate the linear correction. The results are also shown in Figure 10. It can be seen that the molecule-specific correction not only reduces the error but achieves more uniform performance over different conformations. The overall better applicability of the molecular-specific correction to more extended structures is mainly because the peptides unfolded shortly after the simulations started and stayed longer in the denatured status, resulting in higher population of extended structures in the training set.

The cusps in the SAS definition introduce singularity and numerical difficulty to the finite-difference approaches, especially at coarse grid spacings. The SES definition also has pathological

self-intersecting parts that cannot be fully resolved with any finite difference method. These situations can be alleviated by using a finer grid, but the computational cost grows cubically with the inverse of the grid spacing. The above linear corrections aim at reducing errors with much less computational overhead, but the final solution is to develop a new and physically reasonable molecular surface definition that is both smooth and analytical everywhere. It should be noted that the linear correction is not required as in the application of any numerical method, i.e., the FDPB solution of the Poisson–Boltzmann equation, where nobody conducts any correction in biomolecular applications. Of course, no such global correction is possible for FDPB.

## CONCLUSION

We have developed a new numerical algorithm, the field-view method, to calculate the SES or SAS areas under the finite-difference scheme, which is based on the observation that a molecular surface, either in the SES definition or in the SAS definition, can be treated as a union of partial spheres of different centers and radii. To compute the surface area of a spherical surface element, the algorithm exploits the central symmetry of the radial field of a test point charge placed at the spherical center. Specifically, the flux through the spherical surface element is proportional to its surface area in the radial field of the point charge. The new algorithm computes the flux on the finite-difference grid surface element that subtends the same solid angle as the spherical surface element by exploiting the conservation of electric flux. Finally, the summation of the surface areas of the spherical surface elements gives a good approximation for the molecular surface area. Utilization of the finite-difference data structure leads to the new algorithm's particular suitability for the FDPB calculations. The algorithm can also be easily adapted to evaluate the surface integral of either a vector field or a scalar field defined on the SES or SAS.

Two major error sources can be identified in the field-view method: spherical representation of saddle surfaces and truncation in the Taylor expansion used in the flux calculation. The first error is influenced by how many solvent probe sites are used to generate the re-entrant surface as the probe rolls over the atoms. The second error, i.e., the truncation error, is influenced by how many terms are used in the Taylor expansion. The convergence rate of the zeroth-order is quadratic, and that of the second order is in the fourth power. We suggest using the zeroth-order truncation to compute molecular surface areas because the leading error comes from the first source at the typical coarse grid spacing of 1/2 Å often used for FDPB in biomolecular applications.

The quadratic convergence of the field-view method (the zeroth-order truncation) has been demonstrated both theoretically and numerically. The proposed algorithm can achieve very good agreement with the analytical method as far as surface area calculations are concerned. The additional consistency test using a large test set of 1555 molecules and complexes shows that the average unsigned relative error of the SES areas is 0.27% and that of the SAS areas is 1.05% at the typical coarse grid spacing of 1/2 Å, indicating that the proposed algorithm can be applied to biomolecules and complexes over a broad range of sizes and structures. The timing test further shows that the algorithm takes little additional time in the context of FDPB calculations. More interestingly, it was found that a systematic correction can improve the accuracy of the numerical SES areas calculated at coarse grid spacings. For example, the average unsigned relative

error of the numerical SES areas by the algorithm can be reduced from 0.27% to 0.13%. In contrast, the numerical SAS areas can only be improved by molecule-specific corrections. Considering the poor numerical behavior of the SAS, we only recommend the algorithm for SES surface integrations and area calculations.

At fine grid spacings, the zeroth-order truncation may have larger error than the spherical representation of saddle surfaces. In this case, the higher-order terms generally help improve the accuracy of surface areas. Therefore, we have added the option to use higher-order terms in the algorithm so that higher-order truncations can be used at fine grid spacings for better accuracy. Apparently the computation time increases after adopting higher-order terms, but the surface area calculation uses much less time than other components in FDPB calculations. Moreover, higher-order truncations may be advantageous to computing molecular surface integrations due to lack of symmetry in an arbitrary field and little chance of significant error cancellation.

## APPENDIX: TEST SET OF 622 PROTEIN–PROTEIN COMPLEXES

The following molecular structures, in both the Amber format and the pqr format, can be downloaded from <http://rayl.bio.uci.edu/rayl/#Database>: 1A9X, 1APY, 1AQW, 1B0N, 1B5F, 1BBZ, 1BLX, 1C1Y, 1CCW, 1CG5, 1CKA, 1CLV, 1CQ4, 1CSB, 1CZQ, 1CZY, 1D2Z, 1D3B, 1D4T, 1DDV, 1DGW, 1DKZ, 1DNU, 1DOW, 1DPJ, 1DTD, 1E1H, 1E6I, 1E6Y, 1EER, 1EEX, 1EF1, 1EG4, 1EGP, 1ELR, 1ELW, 1EMU, 1EPT, 1EUUV, 1EVH, 1F2T, 1F3U, 1F47, 1F60, 1FIP, 1FLT, 1FM0, 1FS1, 1FVU, 1G1S, 1G6G, 1G73, 1G8K, 1GCV, 1GJ7, 1GK9, 1GL2, 1GL4, 1GO3, 1GUX, 1GVN, 1GYB, 1H0H, 1H2S, 1H32, 1H6K, 1H6W, 1H9O, 1HBN, 1HFE, 1HLE, 1HTR, 1I7Q, 1IHJ, 1J2X, 1J34, 1JAT, 1JBO, 1JD5, 1JDH, 1JDP, 1JEK, 1JGK, 1JLT, 1JMX, 1JNR, 1JSD, 1JSM, 1JW6, 1JWI, 1JY2, 1JYO, 1K2X, 1K5N, 1K8K, 1KQF, 1KSH, 1KVE, 1KYF, 1L2W, 1L6X, 1LB6, 1LM8, 1LQV, 1LSH, 1LUC, 1LVM, 1M1N, 1M2T, 1M4S, 1M93, 1MA3, 1MFG, 1MHW, 1MIZ, 1MJU, 1MSO, 1MTP, 1MTY, 1MZW, 1NOW, 1N12, 1N13, 1N1J, 1N62, 1N7F, 1N7S, 1NH0, 1NKZ, 1NQ7, 1NRJ, 1NTV, 1NVM, 1NX1, 1O6L, 1OAL, 1OAO, 1OAO, 1OAX, 1OBY, 1OK7, 1OO0, 1OR0, 1OR7, 1OU8, 1OV3, 1P57, 1P5V, 1PB5, 1PDQ, 1PFB, 1PK1, 1PK6, 1PQ1, 1PXV, 1PYO, 1PYU, 1Q1A, 1Q3L, 1Q40, 1Q7L, 1QAV, 1QGE, 1QOP, 1QTN, 1R0R, 1R17, 1R1Q, 1R4P, 1R8O, 1R8S, 1RBD, 1RDQ, 1REQ, 1REW, 1RM6, 1RXZ, 1RYP, 1SSD, 1SSP, 1S6C, 1SB2, 1SC3, 1SCT, 1SE0, 1SEM, 1SHA, 1SR4, 1SSH, 1SVF, 1SVZ, 1T0F, 1T0H, 1T0P, 1T15, 1T3Q, 1T61, 1T6G, 1T6O, 1TA3, 1TAF, 1TQY, 1TZY, 1U00, 1U0S, 1U7B, 1U8T, 1UGH, 1UGP, 1UGX, 1UHE, 1UJ0, 1UMD, 1UPK, 1UPT, 1UTI, 1UVQ, 1UW4, 1V74, 1V7P, 1VC3, 1VLE, 1VRA, 1W2W, 1W6S, 1W70, 1W7J, 1W85, 1W9E, 1WAS, 1WDC, 1WDD, 1WHS, 1WMH, 1WQJ, 1WUI, 1WVE, 1WXC, 1XEW, 1XG0, 1XG2, 1XK4, 1XKP, 1XU1, 1Y43, 1Y5I, 1Y7L, 1YAR, 1YCS, 1YDI, 1YFN, 1YMT, 1YPH, 1YQW, 1YRO, 1YTV, 1YUC, 1YUK, 1YWO, 1Z0J, 1Z0K, 1Z3E, 1ZSY, 1Z6O, 1Z9O, 1ZAV, 1ZGX, 1ZHH, 1ZUD, 1ZUK, 1ZV8, 2A3I, 2A50, 2A5T, 2A9K, 2AD6, 2AIJ, 2AIR, 2AKA, 2APO, 2AQ2, 2AQ9, 2ARP, 2ASU, 2B1X, 2B3G, 2B9H, 2BBA, 2BBK, 2BCG, 2BCN, 2BEQ, 2BEZ, 2BFD, 2BGR, 2BKR, 2BKY, 2BL0, 2BLF, 2BMO, 2BO9, 2BPT, 2BR9, 2BS2, 2BUR, 2BW3, 2BZ6, 2BZ8, 2C1D, 2CCH, 2CIO, 2CJS, 2CKL, 2CNZ, 2CWG, 2CZV, 2D00, 2D1X, 2D7C, 2DE6, 2DF6, 2DG5, 2DJF, 2DKO, 2DRM, 2DS2, 2DS8, 2DYO, 2DYR, 2DZE, 2E2D, 2E4M, 2EJF, 2EKE,

2EQ7, 2EQ8, 2ES4, 2F4M, 2F69, 2F91, 2F9I, 2F9N, 2FCW, 2FF4, 2FFU, 2FGR, 2FHZ, 2FLU, 2FMM, 2FOJ, 2FOM, 2FP7, 2FTX, 2FU5, 2FYM, 2G2S, 2G2U, 2G30, 2G5L, 2GAG, 2GBW, 2GGV, 2GH0, 2GHT, 2GIA, 2GL9, 2GPH, 2GPO, 2GSM, 2GUZ, 2GW4, 2H1C, 2H4P, 2H6F, 2H7Z, 2H88, 2H9A, 2HEY, 2HMH, 2HO2, 2HPE, 2HPL, 2HQH, 2HQS, 2HT9, 2HUE, 2HYS, 2I3S, 2IG0, 2INC, 2IUH, 2IVF, 2IZX, 2J12, 2J32, 2J6F, 2J7P, 2J7Y, 2J8C, 2J9A, 2J9U, 2JDI, 2JE6, 2JGB, 2JJS, 2JK9, 2JKH, 2KIN, 2LTN, 2NL9, 2NNU, 2NPT, 2NS1, 2NW2, 2O02, 2O4J, 2O4X, 2O5G, 2O8M, 2O9V, 2OBH, 2ODE, 2OGX, 2OIZ, 2OKR, 2OMZ, 2OQ1, 2OVH, 2OX0, 2OXG, 2OZN, 2P0W, 2P1M, 2P1T, 2P4S, 2P54, 2P58, 2PA8, 2PBI, 2PBK, 2PI2, 2PQR, 2PTT, 2PU9, 2PUY, 2PV2, 2Q0O, 2Q5W, 2QA9, 2QAC, 2QDY, 2QFA, 2QIY, 2QKH, 2QM6, 2QME, 2QWO, 2R2S, 2R7G, 2RHI, 2RHK, 2RI7, 2RKY, 2RMC, 2UUF, 2UWJ, 2UYZ, 2V1T, 2V2F, 2V36, 2V3S, 2V3Z, 2V52, 2V6X, 2V89, 2V8C, 2V9T, 2VGO, 2VIF, 2VLQ, 2VN6, 2VNF, 2VOF, 2VOL, 2VPB, 2VR3, 2VSM, 2VT1, 2VWF, 2VZG, 2W0P, 2W3O, 2W9R, 2WJN, 2WWX, 2YVJ, 2Z30, 2Z3Q, 2Z5B, 2Z8P, 2Z9I, 2ZA4, 2ZD1, 2ZD7, 2ZFD, 2ZMI, 2ZON, 2ZS0, 2ZSI, 2ZVV, 2ZYZ, 2ZZD, 3A1G, 3B5N, 3BC1, 3BEJ, 3BFQ, 3BH7, 3BOM, 3BP6, 3BQO, 3BRL, 3BSS, 3BU3, 3BWU, 3BX4, 3BXM, 3BZY, 3C4M, 3C6W, 3C7B, 3C9A, 3CAL, 3CF4, 3CJS, 3CLS, 3CPT, 3CV0, 3CWW, 3D1K, 3D1M, 3D32, 3D3B, 3D44, 3D9N, 3D9T, 3DAC, 3DBO, 3DD7, 3DDC, 3DGP, 3DKS, 3DLQ, 3DRA, 3DS4, 3DSS, 3DWG, 3DXE, 3DY0, 3E1R, 3EBB, 3ECH, 3EGV, 3EHU, 3EJ9, 3EJB, 3EMH, 3EMW, 3EP6, 3EQS, 3ERY, 3ET3, 3EXE, 3F02, 3F1P, 3F4Y, 3F6Q, 3F75, 3F9X, 3FAP, 3FDT, 3FGR, 3FHV, 3FIV, 3FJU, 3FP2, 3FPN, 3G2S, 3G5O, 3G9A, 3G9K, 3GE3, 3GJ3, 3GL6, 3GLR, 3GV4, 3H11, 3H6P, 3H7H, 3H87, 3H8K, 3H91, 3HDS, 3HEI, 3HHS, 3HHT, 3HNA, 3HPW, 3HQR, 3HTU, 3HXI, 3JQL, 3JRV, 3JVK, 3KB3, 3KDF, 3KDJ, 3KNB, 3PCC, 4UBP, 6RLX, 6TMN.

## AUTHOR INFORMATION

### Corresponding Author

\*Fax: (949) 824-9551. E-mail: rluo@uci.edu.

## ACKNOWLEDGMENT

This work is supported in part by NIH/NIGMS [GM079383 & GM093040].

## REFERENCES

- (1) Lee, B.; Richards, F. M. *J. Mol. Biol.* **1971**, *55*, 379.
- (2) Richards, F. M. *Annu. Rev. Biophys. Biol.* **1977**, *6*, 151.
- (3) Im, W.; Beglov, D.; Roux, B. *Comput. Phys. Commun.* **1998**, *111*, 59.
- (4) Grant, J. A.; Pickup, B. T.; Nicholls, A. *J. Comput. Chem.* **2001**, *22*, 608.
- (5) Lu, Q.; Luo, R. *J. Chem. Phys.* **2003**, *119*, 11035.
- (6) Eisenberg, D.; McLachlan, A. D. *Nature* **1986**, *319*, 199.
- (7) Wesson, L.; Eisenberg, D. *Protein Sci.* **1992**, *1*, 227.
- (8) Cramer, C. J.; Truhlar, D. G. *Science* **1992**, *256*, 213.
- (9) Gallicchio, E.; Kubo, M. M.; Levy, R. M. *J. Phys. Chem. B* **2000**, *104*, 6271.
- (10) Gallicchio, E.; Zhang, L. Y.; Levy, R. M. *J. Comput. Chem.* **2002**, *23*, 517.
- (11) Levy, R. M.; Zhang, L. Y.; Gallicchio, E.; Felts, A. K. *J. Am. Chem. Soc.* **2003**, *125*, 9523.
- (12) Gallicchio, E.; Levy, R. M. *J. Comput. Chem.* **2004**, *25*, 479.
- (13) Su, Y.; Gallicchio, E. *Biophys. Chem.* **2004**, *109*, 251.
- (14) Tan, C.; Tan, Y. H.; Luo, R. *J. Phys. Chem. B* **2007**, *111*, 12263.

- (15) Richmond, T. J. *J. Mol. Biol.* **1984**, *178*, 63.
- (16) Perrot, G.; Cheng, B.; Gibson, K. D.; Vila, J.; Palmer, K. A.; Nayeem, A.; Maigret, B.; Scheraga, H. A. *J. Comput. Chem.* **1992**, *13*, 1.
- (17) Sridharan, S.; Nicholls, A.; Sharp, K. A. *J. Comput. Chem.* **1995**, *16*, 1038.
- (18) Fraczkiwicz, R.; Braun, W. *J. Comput. Chem.* **1998**, *19*, 319.
- (19) Bryant, R.; Edelsbrunner, H.; Koehl, P.; Levitt, M. *Discrete Comput. Geom.* **2004**, *32*, 293.
- (20) Hayryan, S.; Hu, C. K.; Skrivaneck, J.; Hayryan, E.; Pokorny, I. *J. Comput. Chem.* **2005**, *26*, 334.
- (21) Yang, A. S.; Gunner, M. R.; Sampogna, R.; Sharp, K.; Honig, B. *Proteins* **1993**, *15*, 252.
- (22) Yang, A. S.; Honig, B. *J. Mol. Biol.* **1993**, *231*, 459.
- (23) You, T. J.; Bashford, D. *Biophys. J.* **1995**, *69*, 1721.
- (24) Antosiewicz, J.; McCammon, J. A.; Gilson, M. K. *Biochemistry* **1996**, *35*, 7819.
- (25) Alexov, E. G.; Gunner, M. R. *Biophys. J.* **1997**, *72*, 2075.
- (26) Alexov, E. *Proteins* **2003**, *50*, 94.
- (27) Vijayakumar, M.; Zhou, H. X. *J. Phys. Chem. B* **2001**, *105*, 7334.
- (28) Dong, F.; Vijayakumar, M.; Zhou, H. X. *Biophys. J.* **2003**, *85*, 49.
- (29) Dong, F.; Zhou, H. X. *Proteins* **2006**, *65*, 87.
- (30) Qin, S. B.; Zhou, H. X. *Biopolymers* **2007**, *86*, 112.
- (31) Wodak, S. J.; Janin, J. *Proc. Natl. Acad. Sci. U. S. A.* **1980**, *77*, 1736.
- (32) Connolly, M. L. *J. Appl. Crystallogr.* **1983**, *16*, 548.
- (33) Eisenhaber, F.; Argos, P. *J. Comput. Chem.* **1993**, *14*, 1272.
- (34) Edelsbrunner, H. *Discrete Comput. Geom.* **1995**, *13*, 415.
- (35) Gogonea, V.; Osawa, E. *J. Comput. Chem.* **1995**, *16*, 817.
- (36) Augspurger, J. D.; Scheraga, H. A. *J. Comput. Chem.* **1996**, *17*, 1549.
- (37) Gibson, K. D.; Scheraga, H. A. *Mol. Phys.* **1988**, *64*, 641.
- (38) Liang, J.; Edelsbrunner, H.; Fu, P.; Sudhakar, P. V.; Subramaniam, S. *Proteins* **1998**, *33*, 1.
- (39) Street, A. G.; Mayo, S. L. *Fold. Des.* **1998**, *3*, 253.
- (40) Weiser, J.; Shenkin, P. S.; Still, W. C. *J. Comput. Chem.* **1999**, *20*, 217.
- (41) Weiser, J.; Shenkin, P. S.; Still, W. C. *Biopolymers* **1999**, *50*, 373.
- (42) Vasilyev, V.; Purisima, E. O. *J. Comput. Chem.* **2002**, *23*, 737.
- (43) Cavallo, L.; Kleinjung, J.; Fraternali, F. *Nucleic Acids Res.* **2003**, *31*, 3364.
- (44) Guvench, O.; Brooks, C. L. *J. Comput. Chem.* **2004**, *25*, 1005.
- (45) Rychkov, G.; Petukhov, M. *J. Comput. Chem.* **2007**, *28*, 1974.
- (46) Shrake, A.; Rupley, J. A. *J. Mol. Biol.* **1973**, *79*, 351.
- (47) Silla, E.; Villar, F.; Nilsson, O.; Pascualahuir, J. L.; Tapia, O. *J. Mol. Graphics* **1990**, *8*, 168.
- (48) Still, W. C.; Tempczyk, A.; Hawley, R. C.; Hendrickson, T. *J. Am. Chem. Soc.* **1990**, *112*, 6127.
- (49) Legrand, S. M.; Merz, K. M. *J. Comput. Chem.* **1993**, *14*, 349.
- (50) Eisenhaber, F.; Lijnzaad, P.; Argos, P.; Sander, C.; Scharf, M. *J. Comput. Chem.* **1995**, *16*, 273.
- (51) Masuya, M.; Doi, J. *J. Mol. Graphics* **1995**, *13*, 331.
- (52) Connolly, M. L. *J. Mol. Graphics* **1993**, *11*, 139.
- (53) Case, D. A.; Darden, T. A.; Cheatham, T. E., III; Simmerling, C. L.; Wang, J.; Duke, R. E.; Luo, R.; Walker, R. C.; Zhang, W.; Merz, K. M.; Roberts, B.; Wang, B.; Hayik, S.; Roitberg, A.; Seabra, G.; Kolosváry, L.; Wong, K. F.; Paesani, F.; Vanicek, J.; Liu, J.; Wu, X.; Brozell, S. R.; Steinbrecher, T.; Gohlke, H.; Cai, Q.; Ye, X.; Wang, J.; Hsieh, M. -J.; Cui, G.; Roe, D. R.; Mathews, D. H.; Seetin, M. G.; Sagui, C.; Babin, V.; Luchko, T.; Gusarov, S.; Kovalenko, A.; Kollman, P. A. *AMBER 11*; University of California: San Francisco, CA, 2010.
- (54) Sanner, M. F.; Olson, A. J.; Spehner, J. C. *Biopolymers* **1996**, *38*, 305.
- (55) Edelsbrunner, H.; Kirkpatrick, D. G.; Seidel, R. *IEEE Trans. Inf. Theory* **1983**, *29*, 551.
- (56) Edelsbrunner, H.; Mücke, E. P. *ACM Trans. Graph.* **1994**, *13*, 43.
- (57) Greer, J.; Bush, B. L. *Proc. Natl. Acad. Sci. U. S. A.* **1978**, *75*, 303.
- (58) Finney, J. L. *J. Mol. Biol.* **1978**, *119*, 415.
- (59) Pearl, L. H.; Honegger, A. *J. Mol. Graphics* **1983**, *1*, 9.
- (60) Muller, J. J. *J. Appl. Crystallogr.* **1983**, *16*, 74.
- (61) Pavlov, M. Y.; Fedorov, B. A. *Biopolymers* **1983**, *22*, 1507.
- (62) Connolly, M. L. *J. Appl. Crystallogr.* **1985**, *18*, 499.
- (63) Pascualahuir, J. L.; Silla, E.; Tomasi, J.; Bonaccorsi, R. *J. Comput. Chem.* **1987**, *8*, 778.
- (64) Moon, J. B.; Howe, W. J. *J. Mol. Graphics* **1989**, *7*, 109.
- (65) Zauhar, R. J.; Morgan, R. S. *J. Comput. Chem.* **1990**, *11*, 603.
- (66) Bystroff, C. *Protein Eng.* **2002**, *15*, 959.
- (67) Luo, R.; David, L.; Gilson, M. K. *J. Comput. Chem.* **2002**, *23*, 1244.
- (68) Cai, Q.; Wang, J.; Zhao, H. K.; Luo, R. *J. Chem. Phys.* **2009**, *130*, 145101.
- (69) Wang, J.; Luo, R. *J. Comput. Chem.* **2010**, *31*, 1689.
- (70) Cai, Q.; Hsieh, M. J.; Wang, J.; Luo, R. *J. Chem. Theory Comput.* **2010**, *6*, 203.
- (71) You, T.; Bashford, D. *J. Comput. Chem.* **1995**, *16*, 743.
- (72) Rocchia, W.; Sridharan, S.; Nicholls, A.; Alexov, E.; Chiabrera, A.; Honig, B. *J. Comput. Chem.* **2002**, *23*, 128.
- (73) Wang, J.; Luo, R. Manuscript in Preparation.
- (74) Case, D. A.; Cheatham, T. E.; Darden, T.; Gohlke, H.; Luo, R.; Merz, K. M.; Onufriev, A.; Simmerling, C.; Wang, B.; Woods, R. J. *J. Comput. Chem.* **2005**, *26*, 1668.

# Pairwise Long-Range Compensation for Strongly Ionic Systems

Seyit Kale<sup>†</sup> and Judith Herzfeld<sup>\*,‡</sup>

<sup>†</sup>Graduate Program in Biophysics and Structural Biology, Brandeis University, Waltham, Massachusetts 02454, United States

<sup>‡</sup>Department of Chemistry, Brandeis University, Waltham, Massachusetts 02454, United States

**ABSTRACT:** We propose a pairwise compensation method for long-range electrostatics, as an alternative to traditional infinite lattice sums. The approach represents the third generation in a series beginning with the shifted potential corresponding to counterions surrounding a cutoff sphere. That simple charge compensation scheme resulted in pairwise potentials that are continuous at the cutoff, but forces that are not. A second-generation approach modified both the potential and the force such that both are continuous at the cutoff. Here, we introduce another layer of softening such that the derivative of the force is also continuous at the cutoff. In strongly ionic liquids, this extension removes structural artifacts associated with the earlier pairwise compensation schemes and provides results that compare well with Ewald sums.

## INTRODUCTION

Accurate treatment of long-range electrostatics is crucial for the reliability of molecular simulations.<sup>1</sup> The slow convergence of the  $1/r$  term precludes termination at a practical cutoff distance,<sup>2–4</sup> as is typically done, e.g., for van der Waals interactions.<sup>5,6</sup> A widely accepted solution imposes an infinitely repeating lattice that allows the slowly converging Coulomb sum to be separated into a sum that converges rapidly in real space and another that converges rapidly in reciprocal space.<sup>7</sup> Collectively known as Ewald or lattice summations,<sup>8–11</sup> these methods rely on the suitability of the infinite lattice.

A more intuitive alternative has been proposed by Wolf et al. in a study of Madelung energies in perfect crystals.<sup>12</sup> In their approach, the electrostatic pair potentials are shifted by their value at the cutoff distance:

$$U_{\text{SP}}(r_{ij}) = \begin{cases} U(r_{ij}) - U(r_c) & r_{ij} \leq r_c \\ 0 & r_{ij} > r_c \end{cases} \quad (1)$$

$$F_{\text{SP}}(r_{ij}) = \begin{cases} -\frac{dU(r_{ij})}{dr} & r_{ij} \leq r_c \\ 0 & r_{ij} > r_c \end{cases} \quad (2)$$

where  $U$  is the original potential,  $dU/dr$  is its derivative with respect to distance,  $r_{ij}$  is the distance between particles  $i$  and  $j$ , and  $r_c$  is the distance cutoff, typically chosen in the range of 9–12 Å. Physically, this adjustment amounts to placing counterions on the cutoff sphere. Mathematically, this adjustment corresponds to the previously published shifted potential (SP), which achieves continuity at the cutoff for potentials of any form.<sup>13</sup>

In this scheme, the force (eq 2) does not “feel” the charge neutralization and remains discontinuous at the cutoff (Figure 1). Wolf et al.<sup>12</sup> addressed this issue by applying damping, and Zahn et al.<sup>14</sup> subsequently revised Wolf’s damping to achieve energy conservation in MD simulations. However, damping introduces an additional arbitrary parameter. A more straightforward approach to energy conservation is the shifted force (SF),<sup>13,15</sup>

which meets two continuity requirements at the cutoff boundary, i.e., for both the potential and its derivative (similar to the boundary conditions applied to the Poisson–Boltzmann calculation in reaction field methods, but without the need to assume a uniform continuum with known dielectric constant beyond the cutoff<sup>16–18</sup>):

$$U_{\text{SF}}(r_{ij}) = \begin{cases} U(r_{ij}) - U(r_c) - (r_{ij} - r_c) \frac{dU(r_c)}{dr} & r_{ij} \leq r_c \\ 0 & r_{ij} > r_c \end{cases} \quad (3)$$

$$F_{\text{SF}}(r_{ij}) = \begin{cases} -\frac{dU(r_{ij})}{dr} + \frac{dU(r_c)}{dr} & r_{ij} \leq r_c \\ 0 & r_{ij} > r_c \end{cases} \quad (4)$$

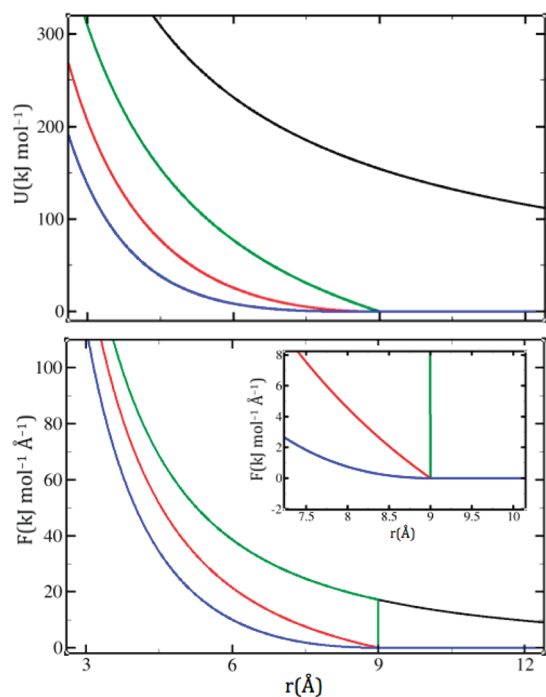
Rigorous tests on this shifted force method have been reported by Fennell and Gezelter.<sup>19</sup> On the basis of comparisons with Ewald energies and forces in SPC/E water<sup>20</sup> and in high temperature molten salts, they conclude that force shifting can be a viable alternative to lattice sums. More recently, Toxvaerd and Dyre reported that SF potentials permit smaller cutoffs in weakly bound systems.<sup>6</sup>

However, problems arise when we apply SF to an unusual ionic liquid of “molecules” that are inherently polarizable and reactive. In this model, molecules comprise charged and independently mobile atomic cores that are surrounded by fully charged and independently mobile valence electron pairs.<sup>21</sup> According to this “LEWIS” construct, a water molecule, e.g., is composed of seven independently mobile particles: a +6 charged oxygen core, two +1 charged protons, and four –2 charged electron pairs (Figure 2, inset). These charges are an order of magnitude larger than the partial charges of typical empirical force fields. In this unusual ionic liquid, SF does

Received: June 9, 2011

Published: September 15, 2011





**Figure 1.** Three different levels of shifting on a purely electrostatic potential  $U = 1/r$  (top) and its force  $F = 1/r^2$  (bottom), as they approach a cutoff of 9 Å. The unmodified potential and force are shown in black, the SP in green, the SF in red, and the SFG in blue. In b, the inset shows a magnified view of the cutoff region. Note that in both SF and SFG the energies go smoothly to zero, whereas the force goes smoothly to zero only in SFG.

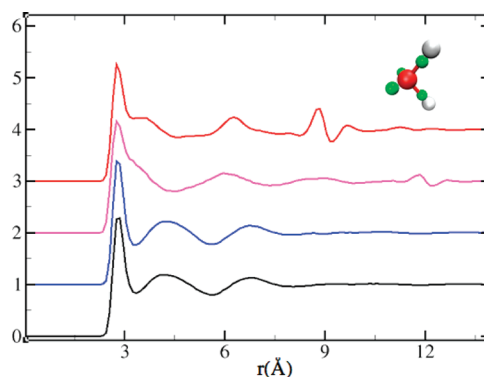
better than SP but still produces significant structural artifacts (Figure 2). To address this problem, we take the approach one step further, by shifting the gradient of the force to make it continuous at the cutoff and adjusting the force and potential accordingly:

$$U_{\text{SFG}}(r_{ij}) = \begin{cases} U(r_{ij}) - U(r_c) - (r_{ij} - r_c) \frac{dU(r_c)}{dr} - \frac{1}{2} (r_{ij} - r_c)^2 \frac{d^2U(r_c)}{dr^2} & r_{ij} \leq r_c \\ 0 & r_{ij} > r_c \end{cases} \quad (5)$$

$$F_{\text{SFG}}(r_{ij}) = \begin{cases} -\frac{dU(r_{ij})}{dr} + \frac{dU(r_c)}{dr} + (r_{ij} - r_c) \frac{d^2U(r_c)}{dr^2} & r_{ij} \leq r_c \\ 0 & r_{ij} > r_c \end{cases} \quad (6)$$

Here, SFG stands for *shifted force gradient*, and  $d^2U/dr^2$  denotes the second derivative of potential with respect to distance. Note that the second derivatives appear as constants, and they need to be evaluated only once per type of interaction. Since the force remains the exact derivative of the potential, energy is still conserved in molecular dynamics (MD) simulations.<sup>6</sup>

We should note that a general shifting scheme was discussed long ago by Levitt et al.<sup>22</sup> in the form of a truncated Taylor series



**Figure 2.** The smooth particle mesh Ewald (SPME) oxygen–oxygen radial distribution function of LEWIS water<sup>21</sup> compared to the (vertically translated) results obtained using SFG with  $r_c = 9$  Å in blue; SF with  $r_c = 12$  Å in magenta; and SF with  $r_c = 9$  Å in red. The inset shows one molecule of LEWIS water: the oxygen ion is rendered in red, protons in white, and electron pairs in green.

expansion:

$$U_{n\text{-shifted}}(r_{ij}) = \begin{cases} U(r_{ij}) - U(r_c) - \sum_{m=1}^n \frac{1}{m!} (r_{ij} - r_c)^m \frac{d^m U(r_c)}{dr^m} & r_{ij} \leq r_c \\ 0 & r_{ij} > r_c \end{cases} \quad (7)$$

Levitt et al.<sup>22</sup> explored both the  $n = 1$  case (corresponding to SF) and the  $n = 2$  case (corresponding to SFG) and concluded that the former provides a better electrostatics description in weakly to mildly ionic systems, such as aqueous solutions of biological macromolecules. They also argued that force shifting, i.e.,  $n = 1$ , has little influence on a hydrogen bond that is modeled by partial charges on atom centers. Here, we show that in cases of extreme ionicity, as encountered in novel or rare model systems, these conclusions may be challenged, and  $n > 1$  orders of shifting can become necessary for proper long-range electrostatics.

## COMPUTATIONAL DETAILS

All MD simulations were performed with Gromacs software<sup>23</sup> version 4.5.3 and analyzed using Gromacs and VMD<sup>24</sup> (version 1.8.7). Potentials were introduced as user tabulated functions with a distance increment of  $\Delta r = 0.005$  Å. Ewald sums were calculated using Gromacs' smooth particle mesh Ewald (SPME), implemented with fourth order spline interpolation and a tolerance of  $10^{-5}$ . In these runs, the non-Coulombic terms were handled as in eqs 5 and 6 with  $r_c = 9$  Å. Since these terms decay rapidly, their effect in the long-range was small compared to the Coulomb term. Following the approach of Fennell and Gezelter,<sup>19</sup> we take the Ewald results as our reference.

The sodium chloride simulations used the Charmm27 force field,<sup>26,27</sup> such that the interaction between ions  $i$  and  $j$  is

$$U_{ij} = \frac{q_i q_j}{4\pi\epsilon_0 r_{ij}} + \epsilon_{ij}^{\min} \left[ \left( \frac{R_{ij}^{\min}}{r_{ij}} \right)^{12} - 2 \left( \frac{R_{ij}^{\min}}{r_{ij}} \right)^6 \right] \quad (8)$$

where  $\epsilon_{ij}^{\min} = (\epsilon_i^{\min} \epsilon_j^{\min})^{1/2}$  and  $R_{ij}^{\min} = (R_i^{\min} + R_j^{\min})/2$  and the parameter values are listed in Table 1.

Simulations were run at 2000 K where the model potentials predict a molten liquid that still exhibits extensive structural order.

**Table 1. Charges and Lennard-Jones Parameters Used in NaCl Simulations<sup>25–27</sup>**

	$q$ (e)	$R^{\min}$ (Å)	$\epsilon^{\min}$ (kcal/mol)
Na <sup>+</sup>	1.00	2.7275	−0.0469
Cl <sup>−</sup>	−1.00	3.8164	−0.0300

**Table 2. Charges and Lennard-Jones Parameters Used in ThCl<sub>4</sub> Simulations (as found in the Gromacs 4.5.3 release)**

	$q$ (e)	$\rho$ (Å)	$\epsilon$ (kJ/mol)
Th <sup>4+</sup>	4.00	3.30000	0.209200
Cl <sup>−</sup>	−1.00	4.41724	0.492833

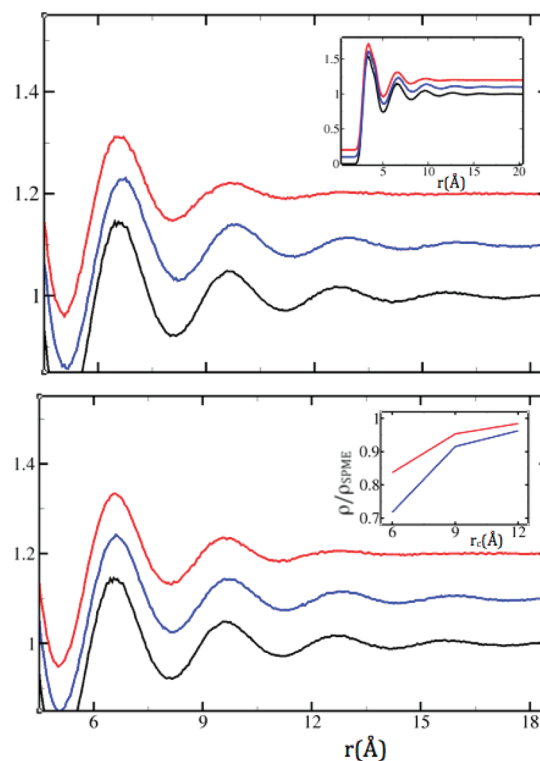
All simulations were in the NPT ensemble. The temperature was maintained using stochastic velocity rescaling<sup>28</sup> with a time constant of 0.1 ps. Pressures were maintained at 1 atm using an isotropic Berendsen barostat<sup>29</sup> with a time constant of 10 ps and compressibility of  $4.5 \times 10^{-5} \text{ bar}^{-1}$ . All NaCl runs began with conjugate gradient energy minimization of a perfect simple cubic  $16 \times 16 \times 16$  crystal with a lattice spacing of 2.0 Å. The integration time step was 2 fs, and neighbor lists were updated very frequently (every five steps) to avoid possible artifacts. Neighbor list radii were 2 Å longer than the cutoffs. Each simulation ran for 1 ns, and the last 800 ps were used for analysis. Coordinates were saved every other picosecond.

Thorium(IV) tetrachloride simulations used the OPLS-AA force field<sup>30,31</sup> such that the interaction between ions  $i$  and  $j$  is given by

$$U_{ij} = \frac{q_i q_j}{4\pi\epsilon_0 r_{ij}} + 4\epsilon_{ij} \left[ \left( \frac{\rho_{ij}}{r_{ij}} \right)^{12} - \left( \frac{\rho_{ij}}{r_{ij}} \right)^6 \right] \quad (9)$$

where  $\epsilon_{ij} = (\epsilon_i \epsilon_j)^{1/2}$  and  $\rho_{ij} = (\rho_i \rho_j)^{1/2}$  and the parameter values are listed in Table 2. Note that both the Lennard-Jones expression (eq 9) and the combination rules of OPLS are slightly different from the CHARMM potential (eq 8). The temperature was maintained at 1000 K. Both NVT and NPT simulations ( $P = 1000$  atm) were run for comparison. NPT runs began with energy minimization of 1000 ThCl<sub>4</sub> molecules arranged in a  $10 \times 10 \times 10$  box with a lattice spacing of 8.0 Å and an intramolecular Th–Cl distance of 2.8 Å. NVT runs began with the final positions, velocities, and volume of the constant pressure SPME simulation. Commensurate with larger interparticle separations than for NaCl, longer cutoffs (15 Å, 18 Å, 21 Å, and 24 Å) were used, neighbor list radii were 3 Å longer than the cutoffs, and the Ewald radius was 21 Å. Trajectories were propagated for 10 ns and coordinates recorded every 20 ps. The first 5 ns of NPT and the first 2 ns of NVT runs were excluded from analysis.

LEWIS water was simulated under similar MD conditions, except for a shorter time step of 0.2 fs and a lower temperature of 300 K. The 9 Å cutoff SPME run used a cubic box of 500 water molecules (edge length  $\sim 24.7$  Å). The 12 Å cutoff SF run used a larger box (edge length  $\sim 36$  Å) of 1500 molecules. Due to the relatively small time step and low temperature, the neighbor lists were updated only every 100 steps. The large water box was run for 250 ps, while all others were run for 1 ns. The first 200 ps of each run were excluded from analysis.



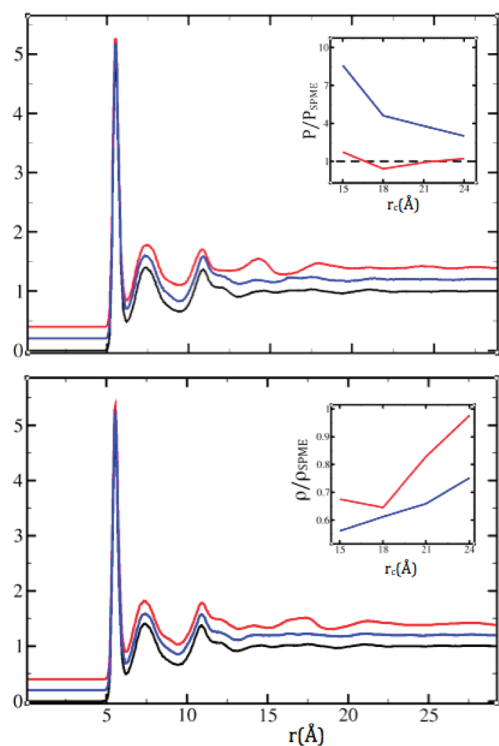
**Figure 3.** The SPME Na<sup>+</sup>–Na<sup>+</sup> radial distribution function (black) compared to the (vertically translated) results obtained using SF (red) and SFG (blue), in the NPT ensemble with  $r_c = 9$  Å (top) and  $r_c = 12$  Å (bottom). The inset in the top panel shows the full radial distribution functions. The inset in the lower panel shows the convergence to 1 of the ratios of the predicted densities from NPT simulations to the SPME value ( $2.486 \text{ g/cm}^3$ ), with increasing cutoff radius (6, 9, and 12 Å). Peaks beyond the cutoff are better preserved with SFG. This occurs at the expense of a mild outward shift that is reduced with the longer cutoff.

## RESULTS

Our LEWIS model for water exhibits dramatic artifacts in the O<sup>6+</sup>–O<sup>6+</sup> correlations with the SF method. The prominent signature of SF is the presence of an artificially dense shell just inside the cutoff with a depletion layer just beyond (Figure 2). The associated peak and valley in the radial distribution are distinctive and become more prominent with smaller cutoffs. In the case of a 9 Å cutoff, the artifact peak is higher than the physical second and third neighbor peaks. In fact, the latter are pushed inward and become ordered by the artificial layer so that the hydrogen bond network is restructured. For a larger cutoff of 12 Å, the problem is less dramatic, yet the artifact is still significant. On the other hand, SFG resolves this abnormality already for  $r_c = 9$  Å and reproduces the Ewald structure to a reasonable accuracy. For the present model, SPME takes about 8–9% longer on a parallel machine with 16 virtual cores.

In molten NaCl, neither SF nor SFG causes a distinct cutoff layer analogous to the LEWIS water artifact. However, SF harshly suppresses order beyond the cutoff, whereas SFG reproduces it (Figure 3). On the other hand, both the SF and SFG softening methods underestimate density with decreasing cutoff, SFG more so than SF (Figure 3, bottom panel inset).

As charge magnitudes increase, the artifact around the SF cutoff re-emerges. In ThCl<sub>4</sub>, the effect is very similar to that in the LEWIS



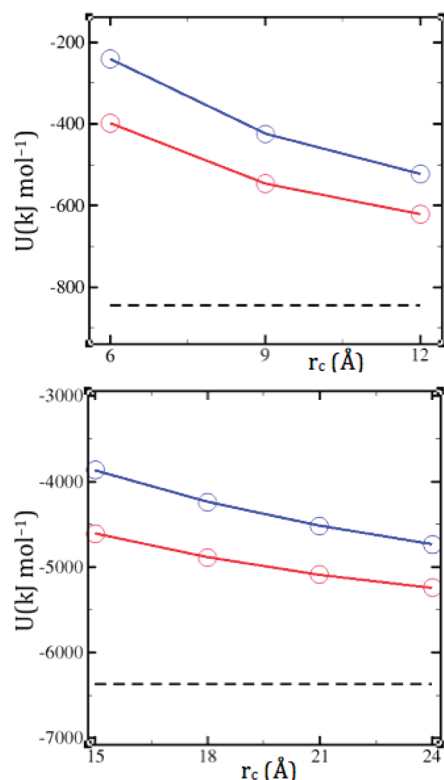
**Figure 4.** Simulations of  $\text{ThCl}_4$ . The SPME  $\text{Th}^{4+}-\text{Th}^{4+}$  radial distribution function (black) compared to the (vertically translated) results obtained using SF (red) and SFG (blue), in the NVT ensemble with  $r_c = 15 \text{ \AA}$  (top) and  $r_c = 18 \text{ \AA}$  (bottom). The artifact layer in SF is still present with the longer cutoff, but less distinct. Insets show the dependence on the cutoff distance (15, 18, 21, and 24  $\text{\AA}$ ) of the pressure in NVT simulations (top) and the density in NPT simulations (bottom) relative to the SPME values (590.5 bar and  $3.011 \text{ g/cm}^3$ , respectively).

water case, in the sense that it appears only in the  $\text{Th}^{4+}-\text{Th}^{4+}$  correlations and becomes more pronounced and localized at shorter cutoffs (Figure 4). SFG lifts this abnormality in the structure. However, it does so at the expense of an underestimation of the density in NPT simulations or an overestimation of the pressure in NVT simulations (Figure 4 insets). Nevertheless, at constant volume, SFG predicts a structure that is virtually identical to SPME already at a cutoff below the Ewald radius (Figure 4, bottom panel).

## DISCUSSION

This work provides evidence that, even in extreme systems, a pairwise compensation scheme can reproduce results similar to those obtained with conventional infinite lattice sums that are typically more CPU-intensive and more difficult to parallelize. When used in conjunction with neighbor lists and cell domain decomposition, pairwise methods can also offer linear scaling with the number of particles  $N$ .<sup>13,19</sup> Currently, most mainstream lattice-sum algorithms scale as  $N \log(N)$ ,<sup>9–11</sup> which makes pairwise sums advantageous as systems grow in size.

We characterize three artifacts, two associated with SF and one with SFG. While SF can be a viable solution for weakly ionic liquids,<sup>6</sup> it results in an artificial layer just inside the cutoff as charges increase in magnitude. In our highly ionic water model, this layer appears in the homoionic correlations between +6 charged particles, and in molten  $\text{ThCl}_4$  it appears between +4 charged



**Figure 5.** Total energy per salt molecule (kJ/mol) as a function of cutoff distance in NPT simulations of molten NaCl (top) and molten  $\text{ThCl}_4$  (bottom). Cohesion is reduced by both SF (red) and SFG (blue), as compared to SPME (dashed black line).

particles. It does not appear in the other correlations of these liquids, or in any of the correlations in molten NaCl where the ionic charges are smaller. However, SF results in a different structural artifact in NaCl; i.e., long-range order is lost for small cutoffs. In contrast, SFG provides a reliable liquid structure in NVT simulations of  $\text{ThCl}_4$ . On the other hand, in NPT simulations, SFG causes more outward shifted correlations and greater underestimation of the density than SF. However, it is notable that in our strongly ionic water model, SFG obtains the Ewald density already at a cutoff of 9  $\text{\AA}$ . While this radius is small for an ionic liquid, it is still  $\sim 20$ -fold larger than the smallest (i.e., intramolecular) ion separation of  $\sim 0.3\text{--}0.5 \text{ \AA}$  in this system. Ions of molten  $\text{ThCl}_4$  are significantly less densely distributed, with nearest neighbor distances varying between 2.7  $\text{\AA}$  and 5.5  $\text{\AA}$ . The cutoffs considered in this system, while large in magnitude, remain small multiples of typical interion separations.

The observed outward correlation shifts, and related density underestimations (Figure 3, bottom panel inset, and Figure 4, bottom panel inset), can be rationalized by considering the corrected potentials in Figure 1. To the extent that the potentials that hold the system together are attenuated, less cohesion is expected (Figure 5). Thus, the greater potential softening in the  $n = 2$  correction than in the  $n = 1$  correction is consistent with the greater correlation shifts and density underestimations. Constant volume ensembles can circumvent this issue at the expense of elevated pressures (Figure 4, top panel inset). Another alternative may be reoptimization of the force field for use with the specific long-range correction, as has been done, e.g., for Ewald compatibility of water models.<sup>32,33</sup>

We arrived at SFG compensation (eqs 5 and 6), to address our own needs for a novel, highly ionic model of water. However, the demonstrated advantages in more conventional ionic systems indicate that the approach may be of wider benefit for the computational community.

## AUTHOR INFORMATION

### Corresponding Author

\*E-mail: herzfeld@brandeis.edu.

## ACKNOWLEDGMENT

We thank Mason Kramer for technical advice and for his initial exploration of SP. This work was supported by NIH grant R01EB001035. Additional computational support was provided by the Brandeis HPC.

## REFERENCES

- (1) Sagui, C.; Darden, T. A. Molecular dynamics simulations of biomolecules: Long-range electrostatic effects. *Annu. Rev. Biophys. Biomol. Struct.* **1999**, *28*, 155–179.
- (2) Schreiber, H.; Steinhauser, O. Molecular Dynamics Studies of Solvated Polypeptides - Why the Cutoff Scheme Does Not Work. *Chem. Phys.* **1992**, *168*, 75–89.
- (3) Schreiber, H.; Steinhauser, O. Cutoff Size Does Strongly Influence Molecular Dynamics Results on Solvated Polypeptides. *Biochemistry* **1992**, *31*, 5856–5860.
- (4) Auffinger, P.; Beveridge, D. L. A Simple Test for Evaluating the Truncation Effects in Simulations of Systems Involving Charged Groups. *Chem. Phys. Lett.* **1995**, *234*, 413–415.
- (5) Brooks, B. R.; Bruccoleri, R. E.; Olafson, B. D.; States, D. J.; Swaminathan, S.; Karplus, M. CHARMM - a Program for Macromolecular Energy, Minimization, and Dynamics Calculations. *J. Comput. Chem.* **1983**, *4*, 187–217.
- (6) Toxvaerd, S.; Dyre, J. C. Communication: Shifted forces in molecular dynamics. *J. Chem. Phys.* **2011**, *134*, 081102.
- (7) Sagui, C.; Darden, T. Multigrid methods for classical molecular dynamics simulations of biomolecules. *J. Chem. Phys.* **2001**, *114*, 6578–6591.
- (8) Ewald, P. The Berechnung optischer und elektrostatischer Gitterpotentiale. *Ann. Phys.* **1921**, *369*, 253–287.
- (9) Toukmaji, A. Y.; Board, J. A. Ewald summation techniques in perspective: A survey. *Comput. Phys. Commun.* **1996**, *95*, 73–92.
- (10) Darden, T.; York, D.; Pedersen, L. Particle Mesh Ewald - an N Log(N) Method for Ewald Sums in Large Systems. *J. Chem. Phys.* **1993**, *98*, 10089–10092.
- (11) Essmann, U.; Perera, L.; Berkowitz, M. L.; Darden, T.; Lee, H.; Pedersen, L. G. A Smooth Particle Mesh Ewald Method. *J. Chem. Phys.* **1995**, *103*, 8577–8593.
- (12) Wolf, D.; Keblicinski, P.; Phillpot, S. R.; Eggebrecht, J. Exact method for the simulation of Coulomb systems by spherically truncated, pairwise 1/r summation. *J. Chem. Phys.* **1999**, *110*, 8254–8282.
- (13) Allen, M. P.; Tildesley, D. J. *Some tricks of the trade. In Computer Simulation of Liquids*; Oxford University Press: New York, 1987; pp 145–146.
- (14) Zahn, D.; Schilling, B.; Kast, S. M. Enhancement of the Wolf damped Coulomb potential: Static, dynamic, and dielectric properties of liquid water from molecular simulation. *J. Phys. Chem. B* **2002**, *106*, 10725–10732.
- (15) Stoddard, S. D.; Ford, J. Numerical Experiments on Stochastic Behavior of a Lennard-Jones Gas System. *Phys. Rev. A* **1973**, *8*, 1504–1512.
- (16) Barker, J. A. Reaction Field, Screening, and Long-Range Interactions in Simulations of Ionic and Dipolar Systems. *Mol. Phys.* **1994**, *83*, 1057–1064.
- (17) Alper, H.; Levy, R. M. Dielectric and Thermodynamic Response of a Generalized Reaction Field Model for Liquid-State Simulations. *J. Chem. Phys.* **1993**, *99*, 9847–9852.
- (18) Tironi, I. G.; Sperb, R.; Smith, P. E.; van Gunsteren, W. F. A Generalized Reaction Field Method for Molecular Dynamics Simulations. *J. Chem. Phys.* **1995**, *102*, 5451–5459.
- (19) Fennell, C. J.; Gezelter, J. D. Is the Ewald summation still necessary? Pairwise alternatives to the accepted standard for long-range electrostatics. *J. Chem. Phys.* **2006**, *124*, 234104.
- (20) Berendsen, H. J. C.; Grigera, J. R.; Straatsma, T. P. The Missing Term in Effective Pair Potentials. *J. Phys. Chem.* **1987**, *91*, 6269–6271.
- (21) Kale, S.; Herzfeld, J.; Dai, S.; Blank, M. Lewis-inspired representation of dissociable water in clusters and Grothuss chains. *J. Biol. Phys.* **2011**, DOI: 10.1007/s10867-011-9229-5.
- (22) Levitt, M.; Hirshberg, M.; Sharon, R.; Daggett, V. Potential-Energy Function and Parameters for Simulations of the Molecular-Dynamics of Proteins and Nucleic-Acids in Solution. *Comput. Phys. Commun.* **1995**, *91*, 215–231.
- (23) van der Spoel, D.; Lindahl, E.; Hess, B.; Groenhof, G.; Mark, A. E.; Berendsen, H. J. C. Gromacs: Fast, Flexible, and Free. *J. Comput. Chem.* **2005**, *26*, 1701–1718.
- (24) Humphrey, W.; Dalke, A.; Schulten, K. VMD: Visual molecular dynamics. *J. Mol. Graphics* **1996**, *14*, 33–38.
- (25) Beglov, D.; Roux, B. Finite Representation of an Infinite Bulk System - Solvent Boundary Potential for Computer Simulations. *J. Chem. Phys.* **1994**, *100*, 9050–9063.
- (26) MacKerell, A. D.; Bashford, D.; Bellott, M.; Dunbrack, R. L.; Evanseck, J. D.; Field, M. J.; Fischer, S.; Gao, J.; Guo, H.; Ha, S.; Joseph-McCarthy, D.; Kuchnir, L.; Kuczera, K.; Lau, F. T. K.; Mattos, C.; Michnick, S.; Ngo, T.; Nguyen, D. T.; Prodhom, B.; Reiher, W. E.; Roux, B.; Schlenkrich, M.; Smith, J. C.; Stote, R.; Straub, J.; Watanabe, M.; Wiorkiewicz-Kuczera, J.; Yin, D.; Karplus, M. All-atom empirical potential for molecular modeling and dynamics studies of proteins. *J. Phys. Chem. B* **1998**, *102*, 3586–3616.
- (27) MacKerell, A. D.; Banavali, N.; Foloppe, N. Development and current status of the CHARMM force field for nucleic acids. *Biopolymers* **2001**, *56*, 257–265.
- (28) Bussi, G.; Donadio, D.; Parrinello, M. Canonical sampling through velocity rescaling. *J. Chem. Phys.* **2007**, *126*, 014101.
- (29) Berendsen, H. J. C.; Postma, J. P. M.; van Gunsteren, W. F.; Dinola, A.; Haak, J. R. Molecular Dynamics with Coupling to an External Bath. *J. Chem. Phys.* **1984**, *81*, 3684–3690.
- (30) Jorgensen, W. L.; Maxwell, D. S.; TiradoRives, J. Development and testing of the OPLS all-atom force field on conformational energetics and properties of organic liquids. *J. Am. Chem. Soc.* **1996**, *118*, 11225–11236.
- (31) Jensen, K. P.; Jorgensen, W. L. Halide, ammonium, and alkali metal ion parameters for modeling aqueous solutions. *J. Chem. Theory Comput.* **2006**, *2*, 1499–1509.
- (32) Price, D. J.; Brooks, C. L. A modified TIP3P water potential for simulation with Ewald summation. *J. Chem. Phys.* **2004**, *121*, 10096–10103.
- (33) Rick, S. W. A reoptimization of the five-site water potential (TIP5P) for use with Ewald sums. *J. Chem. Phys.* **2004**, *120*, 6085–6093.

# Adaptive-Partitioning Redistributed Charge and Dipole Schemes for QM/MM Dynamics Simulations: On-the-fly Relocation of Boundaries that Pass through Covalent Bonds

Soroosh Pezeshki and Hai Lin\*

Chemistry Department, University of Colorado Denver, Denver, Colorado 80217-3364, United States

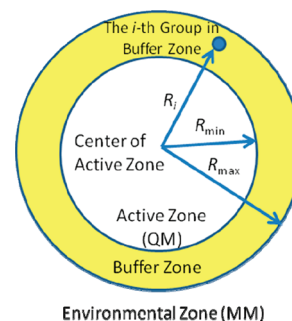
**S** Supporting Information

**ABSTRACT:** Recently, Heyden, Lin, and Truhlar (*J. Phys. Chem. B* **2007**, *111*, 2231–2241) formalized the adaptive-partitioning schemes for quantum mechanical and molecular mechanical (QM/MM) molecular dynamics simulations. The adaptive-partitioning schemes permit on-the-fly reclassification of atoms/groups as part of the QM or MM subsystems during dynamics simulations. Test simulations of argon atoms in a periodic box with dual-level MM potentials in the microcanonical ensemble demonstrated that the adaptive-partitioning schemes conserved energy and momentum, which is critical to ensure correct sampling of configuration spaces of desired ensembles. In this work, we reported the extension of the adaptive-partitioning schemes to deal with groups that are molecular fragments. The newly developed adaptive-partitioning redistributed charge scheme and adaptive-partitioning redistributed charge and dipole schemes allow on-the-fly relocation of the QM/MM boundaries that cut through covalent bonds during dynamics simulations. Test QM/MM simulations with a variety of QM levels of theory in the microcanonical ensembles demonstrated that the new schemes conserve energy and momentum.

## 1. INTRODUCTION

The past decade has witnessed rapid growth in the applications of combined quantum mechanical and molecular mechanical (QM/MM)<sup>1–15</sup> methods. But conventional QM/MM methods have noticeable limits.<sup>16</sup> One of those limits is the prohibition of on-the-fly reclassification of atoms/groups as part of the QM or MM subsystems during molecular dynamics (MD) simulations. For many systems with localized active sites, such a limit is not a concern, but for other systems with nonlocalized active sites, such as ion solvation and transport, defect propagation in materials, and diffusion on catalytic surfaces, such a limit would require the use of large-size QM subsystems, which is very expensive or impractical at all.

Recently, there have been increasing interests in the development of new QM/MM schemes that go beyond this limit.<sup>17–22</sup> Several algorithms have been published on the dynamically partitioning of atoms/groups into QM or MM subsystems in MD simulations.<sup>17–20</sup> When an atom changes its QM or MM identity, the potential energy and forces will show sudden changes. Note that not just the force on the given atom is changing but also the forces on all the other atoms are changing. The discontinuities in energy and forces could lead to numerical instability in MD simulations.<sup>19</sup> Moreover, it could prevent correct sampling of configuration space of the desired ensemble.<sup>19</sup> To cope with those discontinuities, a narrow buffer zone (also called switching shell) is identified between the QM subsystem (also called active zone) and the MM subsystems (also called environmental zone). As illustrated in Figure 1, the active zone is usually defined as a sphere of the inner radius  $R_{\min}$  centered at a given primary atom (or a given spatial location defined by the cartesian coordinates), and the buffer zone is defined as a shell within the inner and outer radii  $R_{\min}$  and  $R_{\max}$ . Various smoothing algorithms are applied to remove the



**Figure 1.** Illustration of the active, buffer, and environmental zones.

discontinuities in the potential energy and/or forces when atoms/groups enter or leave the buffer zone. The smoothing functions, which take forms ranging from simple polynomials to complicated functions, usually depend on the distances ( $R_i$ ) between the active-zone center and the atoms/groups in the buffer zone.

The “hot-spot” method<sup>17</sup> developed by Rode and co-workers in 1996 is one of those examples. Although this method does not conserve momentum and the energy is not evaluated or defined, for simulations in the canonical (*NVT*) ensemble, the kinetic energy is approximately constant, and the numerical instabilities are reduced. The hot-spot method has been applied in a series of studies on the ion solvation, e.g.,  $\text{Li}^+$  in ammonia,<sup>17</sup>  $\text{Ca}^{2+}$ ,<sup>23</sup>  $\text{Na}^+$  and  $\text{K}^+$ ,<sup>24</sup>  $\text{Cu}^{2+}$ ,<sup>25,26</sup>  $\text{Mn}^{2+}$ ,<sup>27</sup>  $\text{Ni}^{2+}$ ,<sup>28</sup>  $\text{Li}^+$ ,<sup>29</sup>  $\text{V}^{2+}$ ,<sup>30</sup>  $\text{Fe}^{2+}$  and  $\text{Fe}^{3+}$ ,<sup>31</sup>  $\text{F}^-$  and  $\text{Cl}^-$ ,<sup>32,33</sup>  $\text{NO}_3^-$ ,<sup>34</sup> and  $\text{HCOO}^-$  in water.<sup>35</sup> In 2002, in discussion of the hot-spot method, Kerdcharoen and Morokuma suggested the ONIOM-XS method,<sup>18</sup> which does

**Received:** July 27, 2011

**Published:** September 19, 2011

conserve momentum. The ONIOM-XS scheme was tested with  $\text{Li}^+$  ion<sup>18</sup> and  $\text{Ca}^{2+}$  ion<sup>36</sup> in liquid ammonia in the isothermal–isobaric (*NPT*) ensemble. However, the ONIOM-XS scheme does not conserve energy in microcanonical (*NVE*) simulations if two or more groups are present in the buffer zone. In 2007, Heyden, Lin, and Truhlar<sup>19</sup> proposed two adaptive-partitioning (AP) schemes, the permuted AP and the sorted AP schemes. Those two schemes were tested by MD simulations of argon atoms in a periodic box using dual-level MM (MM/MM) potentials—the interactions between argon atoms in the active zone were described by a Morse potential and the interactions in the environmental zone by a Lennard-Jones potential. Simulations in the *NVE* ensemble demonstrated that the AP schemes conserved energy and momentum much better than the other schemes in comparison. Buló et al.<sup>20</sup> introduced in 2009 the so-called difference-based adaptive solvation potential, which is related to the sorted AP scheme with a different way in constructing smoothing functions. They also put forward an interesting idea, the continuous-force scheme, which, although not energy conserving, retains a related conserved quantity (energy corrected by a book-keeping term obtained by integration of forces over the trajectory). The methods were tested by dual-level MM simulations in the *NVE* ensemble of a water-in-water model system and of an acetonitrile-in-water system.

In views of the above methods, it is clear that simulations in the *NVE* ensemble is a stricter test for the algorithms, because coupling with a heat bath in the *NVT* and *NPT* simulations will help to keep the kinetic energy approximately constant, reducing numerical instabilities. In the thermodynamics limit, calculations in different ensembles should converge; we have observed in ref 19 that the methods conserving energy and momentum produced very similar radial distribution functions in the *NVE* and *NVT* simulations of the Ar system, while methods with poor conservation of energy and momentum did not. Schemes that do not conserve energy or momentum might still be useful in *NVT* and *NPT* simulations, if proper care is taken to avoid/minimize artifacts in the results.

All methods discussed so far only deal with groups that are whole molecules, such as water and ammonia. A question remains unanswered: What to do if one wants to define a fragment of a molecule as a group? For example, can we define the backbone or side chain of a residue in a protein to be a group? Doing so requires that one must be able to dynamically relocate the QM/MM boundary that passes through covalent bonds during MD simulations. Solving the problem of “fragmental group” for QM/MM dynamics simulations is not only very interesting but also has practical uses. For example, an enzyme active site is usually modeled at the QM level in a QM/MM setup. During MD simulations, side chains of residues may flip away, cofactors may bind, products may be released, and solvent molecules may diffuse in. It will be beneficial to dynamically vary the contents of the QM subsystem in response to those conformational changes. Another example is the simulations of an ion or a molecule transport through membrane channel proteins, where it is desirable to include the ion or molecule and its first solvation shell into the QM subsystem. With the new development here, one could construct a moving QM subsystem centered at the given ion or molecule; when the ion or molecule passes by a residue or lipid, one can add the residue or part of the lipid into the QM subsystem or delete it from the QM subsystem as needed.

In this work, we report an extension of the adaptive-partitioning schemes in ref 19 to handle the fragmental groups. To our knowledge, this is the first implementation of the schemes that allow

on-the-fly relocation of the QM/MM boundaries that cut through covalent bonds in MD simulations. The development results in two new QM/MM schemes, namely the adaptive-partitioning redistributed charge (AP-RC) and the adaptive-partitioning redistributed charge and dipole (AP-RCD) schemes, which are described in Section 2. Test calculations are present in Section 3. The results are analyzed in Section 4, and discussion is given in Section 5.

## 2. METHODOLOGY

**2.1. Adaptive-Partitioning Treatments.** The algorithms of the adaptive-partitioning QM/MM have been given in detail in ref 19. Briefly, in the permuted AP scheme, the potential energy is expressed in a many-body expansion manner:

$$\begin{aligned}
 V = & V^A + \sum_{i=1}^N P_i(V_i^A - V^A) \\
 & + \sum_{i=1}^{N-1} \sum_{j=i+1}^N P_i P_j (V_{ij}^A - [V^A + \sum_{r=i,j}^N (V_r^A - V^A)]) \\
 & + \sum_{i=1}^{N-2} \sum_{j=i+1}^{N-1} \sum_{k=j+1}^N P_i P_j P_k (V_{ijk}^A - (V^A + \sum_{r=i,j,k}^N (V_r^A - V^A))) \\
 & + \sum_{(p,q)=(i,j),(i,k),(j,k)}^{N-1,N} (V_{p,q}^A - (V^A \\
 & + \sum_{r=i,j}^N (V_r^A - V^A))) + \dots
 \end{aligned} \quad (1)$$

where  $V^A$  is the energy determined with the groups in the active zone at the QM level,  $V_i^A$  with all active-zone groups and the  $i$ -th buffer-zone group at the QM level,  $V_{ij}^A$  with all active-zone groups, the  $i$ -th buffer-zone group, and the  $j$ -th buffer-zone group at the QM level, ...  $V_{1,2,\dots,N}^A$  with all active-zone groups and all  $N$  buffer-zone groups at the QM level, and  $P_i$  is the smoothing function of the  $i$ -th buffer-zone group in terms of the dimensionless reduced radial coordinate  $\alpha_i$ :

$$\begin{aligned}
 P_i(\alpha_i) &= -6\alpha_i^5 + 15\alpha_i^4 - 10\alpha_i^3 + 1 \\
 \alpha_i &= \frac{R_i - R_{\min}}{R_{\max} - R_{\min}} \quad \text{for } R_{\min} < R_i < R_{\max}
 \end{aligned} \quad (2)$$

eq 1 can be conveniently rewritten as

$$\begin{aligned}
 V = & V^A \left( 1 - \sum_{i=1}^N P_i + \sum_{i=1}^{N-1} \sum_{j=i+1}^N P_i P_j \right. \\
 & \left. - \sum_{i=1}^{N-2} \sum_{j=i+1}^{N-1} \sum_{k=j+1}^N P_i P_j P_k + \dots \right) \\
 & + \sum_{i=1}^N P_i V_i^A \left( 1 - \sum_{j \neq i}^N P_j + \sum_{j \neq i}^{N-1} \sum_{k=j+1 \neq i}^N P_j P_k - \dots \right) \\
 & + \sum_{i=1}^{N-1} \sum_{j=i+1}^N P_i P_j V_{ij}^A \left( 1 - \sum_{k \neq i,j}^N P_k + \dots \right) + \dots
 \end{aligned} \quad (3)$$

or

$$\begin{aligned}
 V = & V^A \prod_{i=1}^N (1 - P_i) + \sum_{i=1}^N P_i V_i^A \prod_{j \neq i}^N (1 - P_j) \\
 & + \sum_{i=1}^{N-1} \sum_{j=i+1}^N P_i P_j V_{i,j}^A \prod_{k \neq j \neq i}^N (1 - P_k) + \dots
 \end{aligned} \quad (4)$$

Note the constraint that the sum of the smoothing functions is always 1.

$$\begin{aligned}
 & \prod_{i=1}^N (1 - P_i) + \sum_{i=1}^N P_i \prod_{j \neq i}^N (1 - P_j) \\
 & + \sum_{i=1}^{N-1} \sum_{j=i+1}^N P_i P_j \prod_{k \neq j \neq i}^N (1 - P_k) + \dots \\
 & = \prod_{i=1}^N ((1 - P_i) + P_i) = 1
 \end{aligned} \quad (5)$$

In total,  $2^N$  QM calculations are to be performed. All derivatives of the potential energy with respect to the coordinates vary smoothly up to the same order for which the smoothing functions  $P_i$  vary continuously. Since  $0 < P_i < 1$ , the energy contributions of the terms in the series in eq 1 decrease rapidly, it may be advisable to truncate the series. The truncation significantly reduces the number of embedded QM calculations, but it also introduces small (but controllable) discontinuities in the energy and the derivatives. Our test calculations showed that the discontinuities are insignificant if the series is truncated at the fifth order. That is, only up to five groups in the buffer zone are included in the QM calculations, and the fifth or higher order terms of  $P_i$  are neglected.

In the sorted-AP scheme,<sup>19</sup> the groups in the buffer zone are sorted in a canonical order with respect to  $R_i$  from the smallest to largest. The QM calculations begin with the active zone only, and the buffer-zone groups are added one at a time according to the increasing distance, leading to in total  $N + 1$  calculations. The potential energy in sorted AP is given by

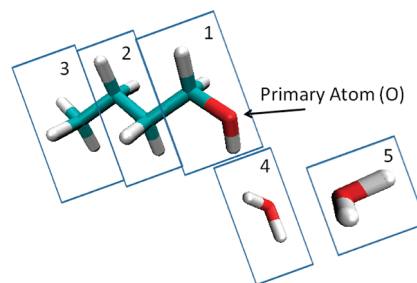
$$V = \sum_{i=0}^N (\Phi_i V_{1,2,\dots,N}^A \prod_{j=i+1}^N (1 - \Phi_j)) \quad (6)$$

or

$$\begin{aligned}
 V = & V^A \prod_{j=1}^N (1 - \Phi_j) + \Phi_1 V_1^A \prod_{j=2}^N (1 - \Phi_j) \\
 & + \Phi_2 V_{1,2}^A \prod_{j=3}^N (1 - \Phi_j) + \dots
 \end{aligned} \quad (7)$$

with the smoothing function

$$\begin{aligned}
 \Phi_i &= (1 - \chi_i)^{-3} \\
 \chi_i &= \sum_{j=1}^{i-1} \frac{1 - P_j}{P_j - P_i} + \frac{1 - P_i}{P_i} + \sum_{j=i+1}^N \left( P_j \frac{1 - P_i}{P_i - P_j} \right)
 \end{aligned} \quad (8)$$



**Figure 2.** Groups in the test model system of a butanol in complex with two water molecules: group 1 is the  $-\text{CH}_2\text{OH}$  fragment, group 2 is the  $-\text{CH}_2\text{CH}_2-$  fragment, group 3 is the  $-\text{CH}_3$  fragment, and groups 4 and 5 are two water molecules, respectively. The O atom in group 1 is the primary atom of the QM subsystem, i.e., the center of the active zone.

and

$$\Phi_0 \equiv 1 \quad (9)$$

The chosen smoothing function  $\Phi_i$  ensures that the energy and the gradient stay constant when two groups in the buffer change the rank.

Comparing eqs 4 and 7 reveals the relationship between the permuted-AP and sorted-AP. The  $n$ -th order contributions in eq 4 comprise  $W(n, N) = N!/n!(N - n)!$  number of terms, and the sorted-AP scheme keeps only one term—the term with the  $n$  buffer-zone groups that are closest to the active-zone center (presumably the biggest contribution)—with very complicated smoothing function  $\Phi_i$ . Note that the sum of smoothing functions is also 1 in the sorted-AP scheme. It is also interesting to note that, if one truncates eq 3 at the first order, keeps only the constant and the linear terms of  $P_i$ , and makes the assumption that  $V_1^A = V_2^A = \dots = V_N^A = (1/N)V_{1,2,\dots,N}^A = (1/N)V^{A+B}$ , then one will arrive at the energy expression for ONIOM-XS.<sup>18</sup>

**2.2. Fragmental groups.** Next, we will focus on the issues that are particularly related to the treatment of the fragmental groups. To deal with the fragmental groups in calculations, one must solve three problems. Here we illustrate the problems by using a small model system: A butanol molecule in complex with two water molecules. As shown in Figure 2, the butanol molecule is divided into three fragmental groups:  $-\text{CH}_2\text{OH}$ ,  $-\text{CH}_2\text{CH}_2-$ , and  $-\text{CH}_3$  as groups 1–3, while each water molecule forms a group (groups 4 and 5). The O atom in group 1 is set to be the primary atom in the QM subsystem. The distance  $R_i$  between the primary atom and the  $i$ -th group is calculated as the distance between the primary atom and the center of mass of the group (or a delegate atom of the group if needed).

The first problem is how to deal with the dangling bond at the QM/MM boundary. For example, if the boundary is passing through the C–C bond connecting groups 1 and 2 and if group 1 is the QM group, then group 1 will have a dangling bond at the boundary. Various schemes are available in the literature,<sup>2,3,16,37–52</sup> and here we adopt the redistributed-charge (RC) and redistributed-charge and dipole (RCD) schemes by Lin and Truhlar,<sup>49</sup> which are classic mechanical mimics to the quantum mechanical description of the charge distribution near the QM/MM boundary by the generalized hybrid orbital (GHO) scheme by Gao and co-workers.<sup>40,53</sup> In both the RC and RCD schemes, the QM subsystem is capped by a hydrogen link atom, and the MM point charge at the M1 atom (the MM atom that directly bounded to a QM atom) is evenly

redistributed to the midpoints at the M1–M2 bonds, where M2 is the MM atom that directly bond to the M1 atom. The capped group 1 becomes CH<sub>3</sub>OH. The redistributed charges and the MM charges at the M2 atom are further modified in the RCD scheme to preserve the M1–M2 bond dipoles. Despite their simplicities, the RC and RCD schemes have been found to yield reasonably good accuracy in energies and geometries in QM/MM calculations.<sup>49</sup> More details about the RC and RCD schemes can be found in ref 49 and are not repeated here.

The second problem is how to define the zeros of QM and MM energies for the fragmental groups, which are to be subtracted from the raw total energy of the QM/MM calculations. The absolute energies given by the QM calculations are several orders of magnitude larger than the MM counterparts; properly setting the zeros of energy to be subtracted is therefore critical to the calculations of energies and forces in the simulations. The zero of QM (or MM) energy for a group that is a whole molecule is straightforward to obtain, which we set to the QM (or MM) energy of the isolated molecule at its geometry optimized at the given QM level of theory (or MM force field). For a group that is part of a molecule, the situation is more complicated, and the zero of energy depends on how this group is linked to other groups of the molecule. For the butanol molecule in Figure 2, group 1 is always in the active zone, and its zero of energy is given by the capped group 1, i.e., CH<sub>3</sub>OH:

$$E_0(\text{group1}) = E(\text{CH}_3\text{OH}) \quad (10)$$

If group 2 is present in the active or buffer zone, it will be included in some of the calculations together with group 1. In those calculations, group 2 will merge with group 1 to form a “super group” –CH<sub>2</sub>CH<sub>2</sub>CH<sub>2</sub>OH. Consequently, the link atom will cap the merged supergroup and will be located at the boundary now between groups 2 and 3. Therefore, there will be one capped supergroup CH<sub>3</sub>CH<sub>2</sub>CH<sub>2</sub>OH instead of two separately capped groups CH<sub>3</sub>OH and CH<sub>3</sub>CH<sub>3</sub> in the calculations. Accordingly, the zero of energy for group 2 is obtained as the energy difference between the capped supergroup and the capped group 1, each at its optimized geometries:

$$E_0(\text{group2}) = E(\text{CH}_3\text{CH}_2\text{CH}_2\text{OH}) - E(\text{CH}_3\text{OH}) \quad (11)$$

Similarly, the zero of energy for group 3 is obtained as the energy difference between the supergroup by merging groups 1 to 3 and the supergroup by merging groups 1 and 2:

$$E_0(\text{group3}) = E(\text{CH}_3\text{CH}_2\text{CH}_2\text{CH}_2\text{OH}) - E(\text{CH}_3\text{CH}_2\text{CH}_2\text{OH}) \quad (12)$$

Apparently, the zero of energy of a fragmental group must be obtained in accord to how the group is merged with other groups in the active and buffer zones. For example, if the C atom in group 3 is set to be the primary atom of the QM subsystem, we will have

$$E_0(\text{group1}) = E(\text{CH}_3\text{CH}_2\text{CH}_2\text{CH}_2\text{OH}) - E(\text{CH}_3\text{CH}_2\text{CH}_3) \quad (13)$$

$$E_0(\text{group2}) = E(\text{CH}_3\text{CH}_2\text{CH}_3) - E(\text{CH}_4) \quad (14)$$

$$E_0(\text{group3}) = E(\text{CH}_4) \quad (15)$$

This then gives rise to the third problem, which is a technical issue, that the computer program must be able to automatically relocate the boundary and find the correct zero of energy based

on the topology of the system. This becomes quite cumbersome if a group is linked to many other groups via covalent bonds.

When determining the zero of energy for a fragmental group, it was not necessary to include groups that were present in the active and/or buffer zones but did not covalently connect to the fragmental group. Inclusion of those groups changed the zero of energy slightly but seemed to have negligible effects on the energy and momentum conservations in the MD simulations.

The total zero of energy for the whole system,  $E_0(\text{sys})$ , is the sum of the zeros of energy for all groups, according to eq 4 for the permuted-AP method and eq 8 for the sorted-AP method. A group in the buffer zone has dual (QM and MM) characteristics, so its contribution to the total zero of energy varies when its distance to the active-zone center changes. As a result, the total zero of energy for the whole system can change significantly (a few to a few hundreds of hartree, depending on the system) and rapidly (in a few tens of femtoseconds) during simulations, presenting a challenge in maintaining numerical precision and stabilities in long-time simulations, especially the NVE simulations. (Note that the gradients due to the smoothing functions depend on the difference between the QM and MM energies.) Such drastic variations in the zero of energy were not present in previous dual-level MM simulations,<sup>19,20</sup> since the zeros of MM energy of a group, such as a water molecule, are usually rather small. The numerical stability also relies critically on the availability of highly accurate gradients. Therefore, for QM calculation, a tight SCF convergence is desired. In particular, for density functional theory (DFT) calculations, fine grids for numerical integration are recommended.

The newly developed AP-RC and AP-RCD schemes have been implemented in a new version of the QMMM program.<sup>54</sup>

### 3. COMPUTATION

The energy and momentum conservations by the AP-RC and AP-RCD methods were tested by MD simulations on two model systems. The first model system was the butanol molecular in complex with two water molecules, as described in Figure 2, for which detailed analysis would be carried out. The second and larger model system was an extension of the first model system: a butanol molecule solvated in one 30 × 30 × 30 Å periodic water box of 813 water molecules, and the butanol molecule was divided into the same three fragmental groups as in the first model system. The MM force field was OPLS-AA<sup>55–60</sup> for butanol and TIP3P<sup>61</sup> for water. Table S1 in the Supporting Information lists the employed MM parameters. The QM levels of theory included the semiempirical method AM1,<sup>62</sup> the hartree-fock (HF) method,<sup>63</sup> the DFT model B3LYP,<sup>64–66</sup> and the post-HF method MP2.<sup>67</sup> The 6-31G(d) basis set<sup>68–72</sup> was employed for the HF, B3LYP, and MP2 calculations. As will be seen in the Results Section, basically all tested QM methods performed equally well in conserving energy and momentum of the systems. Due to high computational costs, simulations for the second model system were only carried out with AM1. In addition to the adaptive-partitioning QM/MM MD simulations, we also performed simulations at the pure-MM, the pure-QM (for the first model system only), and the fixed-partitioning QM/MM level. In the fixed-partitioning QM/MM simulations, group 1 (the –CH<sub>2</sub>OH group) was the QM subsystem while the other groups belonged to the MM subsystem. Table 1 summarizes the calculations that have been done in this work.



Table 1. List of Test Calculations<sup>a</sup>

entry	model system	ensemble	description	QM level	boundary treatment
1.0	first	NVE	pure-MM	n/a	n/a
1.1	first	NVE	pure-QM	AM1	n/a
1.2	first	NVE	fixed-partition QM/MM	AM1	RC
1.3	first	NVE	fixed-partition QM/MM	AM1	RCD
1.4	first	NVE	permuted-AP QM/MM	AM1	RC
1.5	first	NVE	permuted-AP QM/MM	AM1	RCD
1.6	first	NVE	permuted-AP QM/MM	HF/6-31G(d)	RC
1.7	first	NVE	permuted-AP QM/MM	HF/6-31G(d)	RCD
1.8	first	NVE	Permuted-AP QM/MM	B3LYP/6-31G(d)	RC
1.9	first	NVE	permuted-AP QM/MM	B3LYP/6-31G(d)	RCD
1.10	first	NVE	permuted-AP QM/MM	MP 2/6-31G(d)	RC
1.11	first	NVE	permuted-AP QM/MM	MP 2/6-31G(d)	RCD
2.0	second	NVE	pure-MM	n/a	n/a
2.1	second	NVE	fixed-partition QM/MM	AM1	RC
2.2	second	NVE	fixed-partition QM/MM	AM1	RCD
2.3	second	NVE	permuted-AP QM/MM	AM1	RC
2.4	second	NVE	permuted-AP QM/MM	AM1	RCD
2.5	second	NVT	pure-MM	n/a	n/a
2.6	second	NVT	fixed-partition QM/MM	AM1	RC
2.7	second	NVT	fixed-partition QM/MM	AM1	RCD
2.8	second	NVT	permuted-AP QM/MM	AM1	RC
2.9	second	NVT	permuted-AP QM/MM	AM1	RCD

<sup>a</sup>The first model system is a butanol molecule in complex with two water molecules (Figure 2), and the second model system is a butanol molecule in a  $30 \times 30 \times 30 \text{ \AA}$  periodic box of 813 water molecules. In the MD simulations, each trajectory was propagated using the velocity Verlet algorithm with 0.5 fs time steps for 10 000 steps. In the fixed-partitioning QM/MM simulations, group 1 (the  $-\text{CH}_2\text{OH}$  fragment of the butanol molecule) was set to be the QM subsystem, while all other groups were set to be the MM subsystem. In the adaptive-partitioning QM/MM simulations, the active zone centered at the O atom of the butanol with a radius of 3.05  $\text{\AA}$ , the buffer zone had a thickness of 0.5  $\text{\AA}$ , the butanol was divided into three groups as shown in Figure 2, and each water molecule forms a group alone. The MM force field was OPLS-AA for butanol and TIP3P for water.

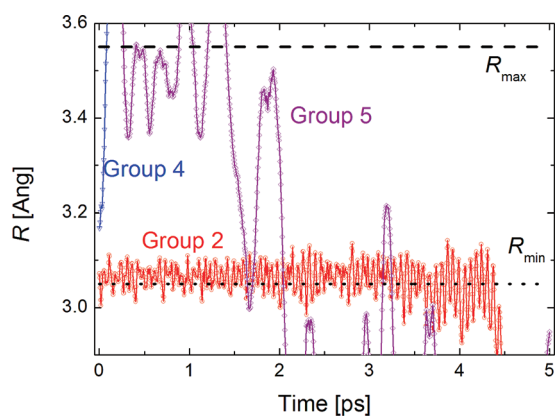
The QM/MM MD simulations were carried out by employing the new version of the QMMM program,<sup>54</sup> which for the QM calculations invoked the ORCA<sup>73</sup> program (for AM1) or the Gaussian03<sup>74</sup> program (for HF, B3LYP, and MP2) and for the MM calculations invoked the TINKER<sup>75</sup> program. Only the simulations with the permuted-AP method are presented here. The sorted-AP method was found to be less satisfactory in conserving energy in the NVE simulations due to numerical instabilities caused by two groups switching positions in the buffer zone,<sup>19,20</sup> and the results are omitted from this work. In the permuted-AP simulations, we did not observe noticeable difference in energy and momentum conservation between simulations employing the full energy expression and an expression truncated at the fifth order. Consequently, all the simulations presented here were done with the truncations at the fifth order. The trajectories were propagated by using the velocity Verlet algorithm,<sup>76</sup> and the time steps were set to 0.5 fs. Smaller time steps, such as 0.2 and 0.1 fs, had also been tested, and we found that they did equally well as the 0.5 fs time step in the energy and momentum conservations for the systems tested here; therefore, those results are not shown. The trajectories were recorded every 10 steps. Although we focused on the NVE ensemble, we also carried out the NVT simulations for the second model system, where the temperature was set to 300 K and was controlled by a Berendsen thermostat<sup>77</sup> with a coupling constant of 2 fs. The active zone was centered at the O atom of the butanol molecule with a radius of  $R_{\min} = 3.05 \text{ \AA}$ , and the thickness of the buffer zone was 0.5  $\text{\AA}$  ( $R_{\max} = 3.55 \text{ \AA}$ ).

The  $R_{\min} = 3.05 \text{ \AA}$  was chosen such that group 2 entered and left the buffer zone frequently during MD simulations. The zeros of energies for the groups in butanol and for the water group are listed in Table S2 in the Supporting Information.

We note that the purpose of the test calculations is to verify the energy and momentum conservations by the AP-RC and AP-RCD schemes rather than to achieve good agreement with experimental results or with MD simulations at the pure-MM or pure-QM levels of theory. Such agreements are unlikely to achieve without optimizing selected parameters describing the interactions between the QM and MM subsystems;<sup>14,78,79</sup> examples of those parameters include the parameters for the van der Waals interactions between QM and MM atoms<sup>80</sup> and the partial atomic charge parameters that enter the effective QM Hamiltonians (as one-electron operators for the electrostatic interaction between the nuclei and electrons of the QM subsystem and the partial atomic charges of the MM subsystem).<sup>16,51,52,81,82</sup>

## 4. RESULTS

**4.1. First Model System.** First, we look at the NVE simulations by permuted-AP RC with AM1 as the QM level of theory, i. e., entry 1.4 in Table 1. An overview of the on-the-fly boundary relocation is provided by Figure 3. The initial geometry was such that group 1 was in the active zone, groups 2 and 4 were in the buffer zone, and groups 3 and 5 in the environmental zone. After the simulation started, the distance between group 2 and

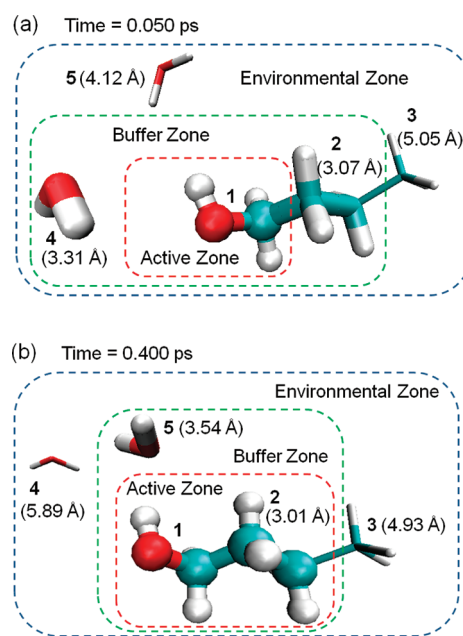


**Figure 3.** Respective distances between the active zone center (the O atom in butanol) and groups 2, 4, and 5 for the first model system in the NVE simulations by permuted-AP RC with AM1 as the QM level of theory, i.e., entry 1.4 in Table 1. The buffer zone region was between the dotted ( $R_{\min} = 3.05 \text{ \AA}$ ) and dashed ( $R_{\max} = 3.55 \text{ \AA}$ ) lines. Group 1 always stayed inside the active zone, while group 3 was outside the buffer zone during the simulation, and they are not plotted here.

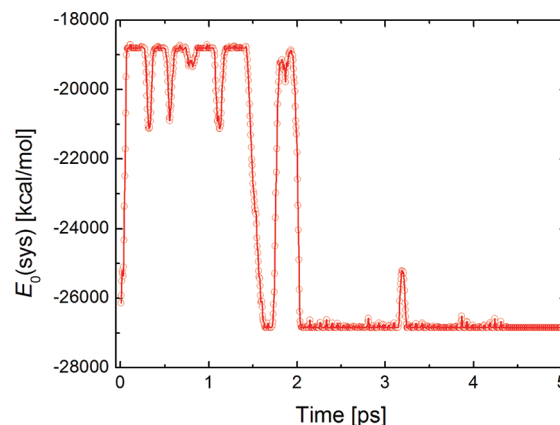
the active-zone center fluctuated around  $3.07 \text{ \AA}$  due to the vibrations of butanol, in particular, due to the bending of the C–C–C and O–C–C angles and the torsion of the O–C–C–C dihedral. The fast fluctuation let group 2 enter and leave the buffer zone quickly and constantly, making the model system a demanding test for the AP-RC (and AP-RCD) treatment. Soon after  $t = 0.075 \text{ ps}$ , group 4 moved out of the buffer zone and stayed in the environmental zone. Group 5 then stepped into the buffer zone and after entered and left the buffer zone a few times, it eventually moved into the active zone and formed a hydrogen bond with the hydroxyl group of butanol. During the last stage of the simulation ( $t > 4.380 \text{ ps}$ ), group 2 walked into the active zone and resided there until the end of the simulation; visualization of the trajectory revealed that the butanol molecule underwent a noticeable conformational change from the trans to the gauche form for the bond axes of  $\text{CH}_3\text{--CH}_2\text{--CH}_2\text{--CH}_2\text{OH}$ . The conformational change shortened the distances between the active-zone center and groups 2 and 3 (by  $0.4$  and  $1.2 \text{ \AA}$ , respectively).

Figure 4 shows two snapshots of the trajectory at simulation time (a)  $t = 0.050 \text{ ps}$  and (b)  $t = 0.400 \text{ ps}$ . In Figure 4a, group 1 ( $\text{--CH}_2\text{OH}$ ) was in the active zone, while groups 2 ( $\text{--CH}_2\text{CH}_2\text{--}$ ) and 4 ( $\text{H}_2\text{O}$ ) were in the buffer zone. In Figure 4b, group 2 had entered the active zone, group 4 moved into the environmental zone, and group 5 ( $\text{H}_2\text{O}$ ) entered the buffer zone from the environmental zone. Group 2 then merged with group 1, producing a supergroup in the active zone:  $\text{--CH}_2\text{CH}_2\text{CH}_2\text{OH}$ . Clearly, the QM/MM boundary had shifted outward from at  $t = 0.050$  to  $t = 0.400 \text{ ps}$ .

The exchange of groups between the active, buffer, and environmental zones had led to significant variations in the zero of energy for the system  $E_0(\text{sys})$ . Figure 5 shows such changes in the zero of energy as a function of simulation time. Those large variations (a few thousands of kcal/mol) were caused by the water molecules due to their significant changes in  $R$ . The contributions by group 2 were smaller, usually less than  $100 \text{ kcal/mol}$ . As pointed out earlier,  $E_0(\text{sys})$  needs to be subtracted from the raw total energy of a QM/MM calculation. Obtaining accurate  $E_0(\text{sys})$  is therefore crucial to the success of the adaptive-partitioning QM/MM simulations.

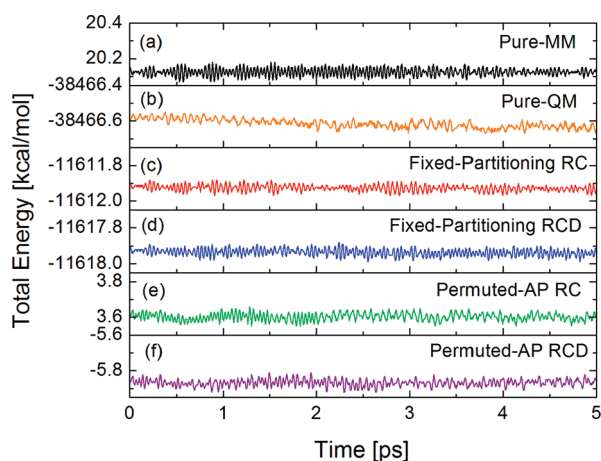


**Figure 4.** Snapshots at (a)  $t = 0.050 \text{ ps}$  and (b)  $t = 0.400 \text{ ps}$  from the NVE simulation by permuted-AP RC for the first model system. The QM level of theory was AM1. Groups in the active zone are shown as balls and sticks, in the buffer zone as licorice (thick), and in the environmental zone as licorice (thin). The distance between the primary atom (O in butanol) and a group (center of mass) is listed in parentheses next to the given group. The corresponding distances for group 1, which are not shown in the figure, are  $0.65 \text{ \AA}$  in both (a) and (b).



**Figure 5.** The zero of energy of the first model system vs simulation time for the NVE simulation by permuted-AP RC with AM1 as the QM level of theory, i.e., entry 1.4 in Table 1.

The conservation of energy is readily examined by plotting the total energies of the system as a function of simulation time, which was done in Figure 6 for the pure-MM (entry 1.0), pure-QM (entry 1.1), fixed-partitioning QM/MM (entries 1.2 and 1.3), and permuted-AP QM/MM (entries 1.4 and 1.5) simulations; the simulations were carried out with AM1 as the QM method. As can be seen, the performances by all methods are very similar. The fluctuations in the total energy were generally smaller than  $0.1 \text{ kcal/mol}$ . The adaptive-partitioning QM/MM schemes did not yield any notable long-term drift in the total energy during the  $5 \text{ ps}$  simulations, despite the large variations in the zeros of the energy as discussed above. As a quick and rough



**Figure 6.** Total energies of the first model system in the *NVE* simulations at (a) the pure-MM, (b) the pure-QM, (c) and (d) the fixed-partitioning QM/MM, and (e) and (f) the adaptive-partitioning QM/MM levels. The QM theory was AM1. The calculations are entries 1.0–1.5 in Table 1. For the adaptive-partitioning QM/MM calculations, the zero of energy of the system  $E_0(\text{sys})$  has been subtracted from the raw total energy.

estimation, we compared the total energies averaged over the first picosecond ( $\bar{E}_{1 \text{ ps}}$ ) and over the last picosecond ( $\bar{E}_{5 \text{ ps}}$ ). The energy difference  $\Delta\bar{E} = |\bar{E}_{1 \text{ ps}} - \bar{E}_{5 \text{ ps}}|$  was 0.001 kcal/mol in the permuted-AP RC and 0.006 kcal/mol in the permuted-AP RCD simulations. For comparison,  $\Delta\bar{E}$  was 0.001 kcal/mol for the pure-MM simulation, 0.050 kcal/mol for pure-QM, 0.009 kcal/mol for fixed-partitioning RC, and 0.010 kcal/mol for fixed-partitioning RCD, respectively. Although 5 ps is not a very long time, the above numbers do indicate that the permuted-AP RC and RCD schemes conserve energy and momentum reasonably well. Adaptive-partitioning QM/MM with higher levels of QM theory HF, B3LYP, and MP2 was tested in entries 1.6–1.11, respectively, and their total energies were displayed in Figure S1, Supporting Information. Again, very satisfactory performance has been observed. The values of  $\Delta\bar{E}$  were found to be small in all six simulations: for HF, 0.007 kcal/mol in permuted-AP RC and less than 0.001 kcal/mol in permuted-AP RCD; the corresponding values were 0.007 and 0.005 kcal/mol for B3LYP and 0.004 and 0.003 kcal/mol for MP2.

**4.2. Second Model System.** The total energies of the *NVE* simulations for the second model system (entries 2.0–2.4) were plotted in Figure S2, Supporting Information, and the temperatures were shown in Figure S3, Supporting Information. Both the energies and temperature were found very stable, although the fluctuations were larger than those in the smaller first model system. The performances by the permuted-AP RC and permuted-AP RCD schemes seemed comparable to those by the pure-MM and by the fixed-partitioning RC and RCD schemes. We found that  $\Delta\bar{E}$  was 0.1 kcal/mol for the pure-MM simulation, 0.1 kcal/mol for fixed-partitioning RC, 0.2 kcal/mol for fixed-partitioning RCD, 0.2 kcal/mol for permuted-AP RC, and 0.5 kcal/mol for permuted-AP RCD, respectively. For the *NVT* simulations (entries 2.5–2.9), the total energies and temperatures were illustrated in Figures S4 and S5, Supporting Information, respectively. The energy and momentum were reasonably well conserved, although the fluctuations in the total energy were larger than those in the *NVE* simulation. As expected, the fluctuations in the temperatures were smaller than those in the

*NVE* simulations. In a similar way to  $\Delta\bar{E}$ , we computed for the *NVT* simulations the temperatures averaged over the first picosecond  $\bar{T}_{1 \text{ ps}}$  and over the last picosecond  $\bar{T}_{5 \text{ ps}}$ , and the difference  $\Delta\bar{T} = |\bar{T}_{1 \text{ ps}} - \bar{T}_{5 \text{ ps}}|$  was found to be 0.008, 0.005, <0.001, 0.050, and 0.040 K for entries 2.5–2.9, respectively.

## 5. DISCUSSION

The computational costs of the adaptive-partitioning QM/MM methods depend on the number of groups in the buffer zone.<sup>19,20</sup> The permuted-AP treatment is the most rigorous algorithm and provided the most satisfactory conservation of energy and momentum in our tests. Unfortunately, its computational cost scales unfavorably as  $2^N$ , where  $N$  is number of groups in the buffer zone. One obvious way to lower the costs is to reduce the thickness of the buffer zone. For an ion solvated in water, an active zone with diameter of 9 Å and a buffer zone of 1 Å thickness will give rise to approximately 10 water molecules in the buffer zone, leading to  $2^{10} = 1024$  QM calculations for one time step. However, if the buffer zone can be narrowed down to 0.5 Å, there will only be 5 water molecules in the buffer zone, requiring  $2^5 = 32$  calculations, which is a lot less. Based on what we have found in our test calculations, the difference in the qualities of energy and momentum conservations was insignificant between a buffer of 0.5 Å thickness and one of 1.0 Å. The 0.5 Å option appeared to be a good choice. The problem of computational cost scaling seems less severe when modeling ion transport through channel proteins and when modeling conformational changes in enzyme active sites; in those examples, the numbers of buffer-zone groups are small (likely less than 5). Another way to reduce the computational cost of the permuted-AP treatments is to truncate the many-body expansion-like energy expression to a given order, as described in the Methodology Section. The results obtained in this work demonstrated that the energy and momentum were conserved reasonably well with the truncation at the fifth order. Finally, we emphasize that all the QM calculations for a given time step are parallel in nature, making large-scale parallel computation very feasible, especially on supercomputers with lots of computational nodes and CPUs. The wall time will in principle be determined by the QM calculations with the largest number of buffer groups.

The permuted-AP RC and RCD schemes belong to the electrostatic embedding QM/MM schemes,<sup>39</sup> where the QM subsystem is embedded in a background of atomic partial charges of the MM atoms. As a result, the QM wave function is polarized by the MM subsystem, providing a more realistic description than that by the mechanical embedding,<sup>39</sup> where the QM subsystem is computed in the gas phase. However, the gradient calculations for the embedded-QM subsystem are very expensive, especially when a large number of background point charges enter the effective QM Hamiltonian. Using cutoff to reduce the number of background point charges in the embedded-QM calculations should lower the computational costs. Another way to reduce the costs is to use the frozen density approximation that neglects the changes in the polarization of the density caused by the varying MM coordinates for a number of steps<sup>83</sup> and the density reduced to point charge approximation that uses the electrostatic potential (ESP) fitted charges<sup>84,85</sup> to replace the full density in the energy and gradient calculations of the electrostatic interactions between the QM and MM atoms.<sup>83,85–89</sup> It will be interesting to see how the gradients computed employing those approximations, although will not be as rigorous as the gradients

computed here with full density updated every step, affect the conservation of energy and momentum in MD simulations.

The current study is one of the steps toward our goal in developing the open-boundary QM/MM methods,<sup>16</sup> which will be a combination of the flexible-boundary treatments<sup>16,52</sup> and the adaptive-partitioning schemes. The flexible-boundary QM/MM methods aim to go beyond another limitation<sup>16</sup> in conventional QM/MM methods to permit partial charge transfer across the QM/MM boundaries. In the flexible-boundary treatments, both the QM and MM subsystems can have fractional numbers of charge, which in principle provides a more realistic picture for the charge distributions within the entire system. The marriage of the flexible-boundary treatments and the adaptive-partitioning schemes will make it possible for the QM and MM subsystem to dynamically exchange partial charges as well as atoms/groups during MD simulations. The open-boundary QM/MM would (at least in principle) facilitate the adoption of a relatively small QM subsystem in QM/MM molecular dynamics simulations, which will assist the use of high-level QM theory and/or long simulation time and could potentially lead to new insights.

## ■ ASSOCIATED CONTENT

**S Supporting Information.** Table S1 lists the MM parameters used for model systems in this work, and Table S2 tabulates the zeros of energy for the groups. We plotted the total energies of the system versus simulation time for entries 1.6–1.11 in Figure S1, for entries 2.0–2.4 in Figure S2, and for entries 2.5–2.9 in Figure S4. The temperatures versus simulation time were displayed for entries 2.0–2.4 in Figure S3 and for entries 2.5–2.9 in Figure S5. This material is available free of charge via the Internet at <http://pubs.acs.org>.

## ■ AUTHOR INFORMATION

### Corresponding Author

\*E-mail: [hai.lin@ucdenver.edu](mailto:hai.lin@ucdenver.edu).

## ■ ACKNOWLEDGMENT

This work is supported by Research Corporation (CC6725) and National Science Foundation (CHE0952337). We thank Minnesota Supercomputing Institute for CPU time and access to Gaussian03. We thank Donald Truhlar and Andreas Heyden for helpful discussion.

## ■ REFERENCES

- (1) Warshel, A.; Levitt, M. Theoretical studies of enzymic reactions: dielectric, electrostatic and steric stabilization of the carbonium ion in the reaction of lysozyme. *J. Mol. Biol.* **1976**, *103*, 227–249.
- (2) Singh, U. C.; Kollmann, P. A. A combined ab initio quantum mechanical and molecular mechanical method for carrying out simulations on complex molecular systems: Applications to the CH<sub>3</sub>Cl + Cl<sup>-</sup> exchange reaction and gas phase protonation of polyethers. *J. Comput. Chem.* **1986**, *7*, 718–730.
- (3) Field, M. J.; Bash, P. A.; Karplus, M. A combined quantum mechanical and molecular mechanical potential for molecular dynamics simulations. *J. Comput. Chem.* **1990**, *11*, 700–733.
- (4) Gao, J. Methods and applications of combined quantum mechanical and molecular mechanical potentials. *Rev. Comput. Chem.* **1996**, *7*, 119–185.
- (5) Friesner, R. A.; Beachy, M. D. Quantum mechanical calculations on biological systems. *Curr. Opin. Struct. Biol.* **1998**, *8*, 257–262.

(6) *Combined quantum mechanical and molecular mechanical methods*: ACS Symp. Ser. 712; Gao, J., Thompson, M. A.; Eds.; American Chemical Society: Washington, DC, 1998, 310 pp.

(7) Ruiz-López, M. F.; Rivail, J. L. Combined quantum mechanics and molecular mechanics approaches to chemical and biochemical reactivity. In *Encyclopedia of computational chemistry*; von Ragué Schleyer, P., Ed.; Wiley: Chichester, U.K., 1998; Vol. 1, pp 437–448.

(8) Monard, G.; Merz, K. M., Jr. Combined quantum mechanical/molecular mechanical methodologies applied to biomolecular systems. *Acc. Chem. Res.* **1999**, *32*, 904–911.

(9) Hillier, I. H. Chemical reactivity studied by hybrid QM/MM methods. *THEOCHEM* **1999**, *463*, 45–52.

(10) Hammes-Schiffer, S. Theoretical perspectives on proton-coupled electron transfer reactions. *Acc. Chem. Res.* **2000**, *34*, 273–281.

(11) Sherwood, P. Hybrid quantum mechanics/molecular mechanics approaches. In *Modern Methods and Algorithms of Quantum Chemistry*; Grotendorst, J., Ed.; John von Neumann Institute: Jülich, 2000; Vol. 3, pp 285–305.

(12) Gao, J.; Truhlar, D. G. Quantum mechanical methods for enzyme kinetics. *Annu. Rev. Phys. Chem.* **2002**, *53*, 467–505.

(13) Morokuma, K. New challenges in quantum chemistry: quests for accurate calculations for large molecular systems. *Philos. Trans. R. Soc. London, A* **2002**, *360*, 1149–1164.

(14) Lin, H.; Truhlar, D. G. QM/MM: What have we learned, where are we, and where do we go from here? *Theor. Chem. Acc.* **2007**, *117*, 185–199.

(15) Senn, H. M.; Thiel, W. QM/MM methods for biological systems. *Top. Curr. Chem.* **2007**, *268*, 173–290.

(16) Zhang, Y.; Lin, H. Flexible-boundary QM/MM calculations: II. Partial charge transfer across the QM/MM boundary that passes through a covalent bond *Theor. Chem. Acc.* **2010**, *126*, 315–322.

(17) Kerdcharoen, T.; Liedl, K. R.; Rode, B. M. A QM/MM simulation method applied to the solution of Li<sup>+</sup> in liquid ammonia. *Chem. Phys.* **1996**, *211*, 313–323.

(18) Kerdcharoen, T.; Morokuma, K. ONIOM-XS: an extension of the ONIOM method for molecular simulation in condensed phase. *Chem. Phys. Lett.* **2002**, *355*, 257–262.

(19) Heyden, A.; Lin, H.; Truhlar, D. G. Adaptive partitioning in combined quantum mechanical and molecular mechanical calculations of potential energy functions for multiscale simulations. *J. Phys. Chem. B* **2007**, *111*, 2231–2241.

(20) Buló, R. E.; Ensing, B.; Sikkema, J.; Visscher, L. Toward a practical method for adaptive QM/MM simulations. *J. Chem. Theory Comput.* **2009**, *9*, 2212–2221.

(21) Nielsen, S. O.; Buló, R. E.; Moore, P. B.; Ensing, B. Recent progress in adaptive multiscale molecular dynamics simulations of soft matter. *Phys. Chem. Chem. Phys.* **2010**, *12*, 12401–12414.

(22) Poma, A. B.; Delle Site, L. Classical to path-integral adaptive resolution in molecular simulation: Towards a smooth quantum-classical coupling. *Phys. Rev. Lett.* **2010**, *104*, 250201.

(23) Tongraar, A.; Liedl, K. R.; Rode, B. M. Solvation of Ca<sup>2+</sup> in water studied by Born-Oppenheimer ab initio QM/MM dynamics. *J. Phys. Chem. A* **1997**, *101*, 6299–6309.

(24) Tongraar, A.; Liedl, K. R.; Rode, B. M. Born-Oppenheimer ab initio QM/MM dynamics simulations of Na<sup>+</sup> and K<sup>+</sup> in water: From structure making to structure breaking effects. *J. Phys. Chem. A* **1998**, *102*, 10340–10347.

(25) Marini, G. W.; Liedl, K. R.; Rode, B. M. Investigation of Cu<sup>2+</sup> hydration and the Jahn-Teller effect in solution by QM/MM monte carlo simulations. *J. Phys. Chem. A* **1999**, *103*, 11387–11393.

(26) Schwenk, C. F.; Rode, B. M. New insights into the Jahn-Teller effect through ab initio quantum-mechanical/molecular-mechanical molecular dynamics simulations of Cu-II in water. *ChemPhysChem* **2003**, *4*, 931–943.

(27) Yagüe, J. I.; Mohammed, A. M.; Loeffler, H.; Rode, B. M. Classical and mixed quantum mechanical/molecular mechanical simulation of hydrated manganese ion. *J. Phys. Chem. A* **2001**, *105*, 7646–7650.

- (28) Inada, Y.; Loeffler, H. H.; Rode, B. M. Librational, vibrational, and exchange motions of water molecules in aqueous Ni(II) solution: classical and QM/MM molecular dynamics simulations. *Chem. Phys. Lett.* **2002**, *358*, 449–458.
- (29) Loeffler, H. H.; Rode, B. M. The hydration structure of the lithium ion. *J. Chem. Phys.* **2002**, *117*, 110–117.
- (30) Loeffler, H. H.; Yague, J. I.; Rode, B. M. QM/MM-MD simulation of hydrated vanadium(II) ion. *Chem. Phys. Lett.* **2002**, *363*, 367–371.
- (31) Remsungnen, T.; Rode, B. M. Dynamical properties of the water molecules in the hydration shells of Fe(II) and Fe(III) ions: ab initio QM/MM molecular dynamics simulations. *Chem. Phys. Lett.* **2003**, *367*, 586–592.
- (32) Tongraar, A.; Rode, B. M. The hydration structures of F<sup>-</sup> and Cl<sup>-</sup> investigated by ab initio QM/MM molecular dynamics simulations. *Phys. Chem. Chem. Phys.* **2003**, *5*, 357–362.
- (33) Tongraar, A.; Rode, B. M. Ab initio QM/MM dynamics of anion-water hydrogen bonds in aqueous solution. *Chem. Phys. Lett.* **2005**, *403*, 314–319.
- (34) Tongraar, A.; Tangkawanwanit, P.; Rode, B. M. A combined QM/MM molecular dynamics simulations study of nitrate anion (NO<sub>3</sub><sup>-</sup>) in aqueous solution. *J. Phys. Chem. A* **2006**, *110*, 12918–12926.
- (35) Payaka, A.; Tongraar, A.; Rode, B. M. Combined QM/MM MD study of HCOO<sup>-</sup>-water hydrogen bonds in aqueous solution. *J. Phys. Chem. A* **2009**, *113*, 3291–3298.
- (36) Kerdcharoen, T.; Morokuma, K. Combined quantum mechanics and molecular mechanics simulation of Ca<sup>2+</sup>/ammonia solution based on the ONIOM-XS method: Octahedral coordination and implication to biology. *J. Chem. Phys.* **2003**, *118*, 8856–8862.
- (37) Maseras, F.; Morokuma, K. IMOMM: a new integrated ab initio + molecular mechanics geometry optimization scheme of equilibrium structures and transition states. *J. Comput. Chem.* **1995**, *16*, 1170–1179.
- (38) Assfeld, X.; Rivail, J.-L. Quantum chemical computations on parts of large molecules: the ab initio local self consistent field method. *Chem. Phys. Lett.* **1996**, *263*, 100–106.
- (39) Bakowies, D.; Thiel, W. Hybrid models for combined quantum mechanical and molecular mechanical approaches. *J. Phys. Chem.* **1996**, *100*, 10580–10594.
- (40) Gao, J.; Amara, P.; Alhambra, C.; Field, M. J. A generalized hybrid orbital (GHO) method for the treatment of boundary atoms in combined QM/MM calculations. *J. Phys. Chem. A* **1998**, *102*, 4714–4721.
- (41) Antes, I.; Thiel, W. Adjusted connection atoms for combined quantum mechanical and molecular mechanical methods. *J. Phys. Chem. A* **1999**, *103*, 9290–9295.
- (42) de Vries, A. H.; Sherwood, P.; Collins, S. J.; Rigby, A. M.; Rigutto, M.; Kramer, G. J. Zeolite structure and reactivity by combined quantum-chemical-classical calculations. *J. Phys. Chem. B* **1999**, *103*, 6133–6141.
- (43) Zhang, Y.; Lee, T.-S.; Yang, W. A pseudobond approach to combining quantum mechanical and molecular mechanical methods. *J. Chem. Phys.* **1999**, *110*, 46–54.
- (44) Das, D.; Eurenus, K. P.; Billings, E. M.; Sherwood, P.; Chatfield, D. C.; Hodoscek, M.; Brooks, B. R. Optimization of quantum mechanical molecular mechanical partitioning schemes: Gaussian delocalization of molecular mechanical charges and the double link atom method. *J. Chem. Phys.* **2002**, *117*, 10534–10547.
- (45) DiLabio, G. A.; Hurley, M. M.; Christiansen, P. A. Simple one-electron quantum capping potentials for use in hybrid QM/MM studies of biological molecules. *J. Chem. Phys.* **2002**, *116*, 9578–9584.
- (46) Amara, P.; Field, M. J. Evaluation of an ab initio quantum mechanical/molecular mechanical hybrid-potential link-atom method. *Theor. Chem. Acc.* **2003**, *109*, 43–52.
- (47) Sherwood, P.; de Vries, A. H.; Guest, M. F.; Schreckenbach, G.; Catlow, C. R. A.; French, S. A.; Sokol, A. A.; Bromley, S. T.; Thiel, W.; Turner, A. J.; Billeter, S.; Terstegen, F.; Thiel, S.; Kendrick, J.; Rogers, S. C.; Casci, J.; Watson, M.; King, F.; Karlsen, E.; Sjøvoll, M.; Fahmi, A.; Schafer, A.; Lennartz, C. QUASI: A general purpose implementation of the QM/MM approach and its application to problems in catalysis. *THEOCHEM* **2003**, *632*, 1–28.
- (48) Pu, J.; Gao, J.; Truhlar, D. G. Generalized hybrid orbital (GHO) method for combining ab initio hartree-fock wave functions with molecular mechanics. *J. Phys. Chem. A* **2004**, *108*, 632–650.
- (49) Lin, H.; Truhlar, D. G. Redistributed charge and dipole schemes for combined quantum mechanical and molecular mechanical calculations. *J. Phys. Chem. A* **2005**, *109*, 3991–4004.
- (50) Shao, Y.; Kong, J. YinYang atom: A simple combined ab initio quantum mechanical molecular mechanical model. *J. Phys. Chem. A* **2007**, *111*, 3661–3671.
- (51) Zhang, Y.; Lin, H.; Truhlar, D. G. Self-consistent polarization of the boundary in the redistributed charge and dipole scheme for combined quantum-mechanical and molecular-mechanical calculations. *J. Chem. Theory Comput.* **2007**, *3*, 1378–1398.
- (52) Zhang, Y.; Lin, H. Flexible-boundary quantum-mechanical/molecular-mechanical calculations: Partial charge transfer between the quantum-mechanical and molecular-mechanical subsystems. *J. Chem. Theory Comput.* **2008**, *4*, 414–425.
- (53) Amara, P.; Field, M. J.; Alhambra, C.; Gao, J. The generalized hybrid orbital method for combined quantum mechanical/molecular mechanical calculations: formulation and tests of the analytical derivatives. *Theor. Chem. Acc.* **2000**, *104*, 336–343.
- (54) Lin, H.; Zhang, Y.; Pezeshki, S.; Truhlar, D. G. QMMM, version 1.4.0, University of Minnesota: Minneapolis, MN, 2011.
- (55) Jorgensen, W. L.; Maxwell, D. S.; Tirado-Rives, J. Development and testing of the OPLS all-atom force field on conformational energetics and properties of organic liquids. *J. Am. Chem. Soc.* **1996**, *118*, 11225–11236.
- (56) Jorgensen, W. L.; McDonald, N. A. Development of an all-atom force field for heterocycles. Properties of liquid pyridine and diazenes. *THEOCHEM* **1998**, *424*, 145–155.
- (57) McDonald, N. A.; Jorgensen, W. L. Development of an all-atom force field for heterocycles. Properties of liquid pyrrole, furan, diazoles, and oxazoles. *J. Phys. Chem. B* **1998**, *102*, 8049–8059.
- (58) Rizzo, R. C.; Jorgensen, W. L. OPLS all-atom model for amines: Resolution of the amine hydration problem. *J. Am. Chem. Soc.* **1999**, *121*, 4827–4836.
- (59) Kaminski, G. A.; Friesner, R. A.; Tirado-Rives, J.; Jorgensen, W. L. Evaluation and reparametrization of the OPLS-AA force field for proteins via comparison with accurate quantum chemical calculations on peptides. *J. Phys. Chem. B* **2001**, *105*, 6474–6487.
- (60) Kahn, K.; Bruice, T. C. Parameterization of OPLS-AA force field for the conformational analysis of macrocyclic polyketides. *J. Comput. Chem.* **2002**, *23*, 977–996.
- (61) Jorgensen, W. L.; Chandrasekhar, J.; Madura, J. D.; Impey, R. W.; Klein, M. L. Comparison of simple potential functions for simulating liquid water. *J. Chem. Phys.* **1983**, *79*, 926–935.
- (62) Dewar, M. J. S.; Zoebisch, E. G.; Healy, E. F.; Stewart, J. J. P. AM1: A new general purpose quantum mechanical molecular model. *J. Am. Chem. Soc.* **1985**, *107*, 3902–3909.
- (63) Roothaan, C. C. J. New developments in molecular orbital theory. *Rev. Mod. Phys.* **1951**, *23*, 69–89.
- (64) Becke, A. D. Density-functional exchange-energy approximation with correct asymptotic behavior. *Phys. Rev. A* **1988**, *38*, 3098–3100.
- (65) Becke, A. D. Density-functional thermochemistry. III. The role of exact exchange. *J. Chem. Phys.* **1993**, *98*, 5648–5652.
- (66) Lee, C.; Yang, W.; Parr, R. G. Development of the Colle-Salvetti correlation-energy formula into a functional of the electron density. *Phys. Rev. B: Condens. Matter* **1988**, *37*, 785–789.
- (67) Møller, C. M. S.; Plesset, M. S. Note on an approximation treatment for many-electron systems. *Phys. Rev.* **1934**, *46*, 618–622.
- (68) Ditchfield, R.; Hehre, W. J.; Pople, J. A. Self-consistent molecular-orbital methods. IX. An extended Gaussian-type basis for molecular-orbital studies of organic molecules. *J. Chem. Phys.* **1971**, *54*, 724–728.
- (69) Francl, M. M.; Pietro, W. J.; Hehre, W. J.; Binkley, J. S.; DeFrees, D. J.; Pople, J. A.; Gordon, M. S. Self-consistent molecular

orbital methods. XXIII. A polarization-type basis set for second-row elements. *J. Chem. Phys.* **1982**, *77*, 3654–3665.

(70) Clark, T.; Chandrasekhar, J.; Spitznagel, G. W.; Schleyer, P. v. R. Efficient diffuse function-augmented basis sets for anion calculations. III. The 3-21+G basis set for first-row elements, Li-F. *J. Comput. Chem.* **1983**, *4*, 294–301.

(71) Frisch, M. J.; Pople, J. A.; Binkley, J. S. Quadratic configuration interaction. A general technique for determining electron correlation energies. *J. Chem. Phys.* **1984**, *80*, 3265–3269.

(72) Hehre, W. J.; Ditchfield, R.; Pople, J. A. Self-consistent molecular orbital methods. XII. Further extensions of Gaussian-type basis sets for use in molecular orbital studies of organic molecules. *J. Chem. Phys.* **1972**, *56*, 2257–2261.

(73) Neese, F. ORCA, Version 2.8.0, University of Bonn: Bonn, 2011.

(74) Frisch, M. J.; Trucks, G. W.; Schlegel, H. B.; Scuseria, G. E.; Robb, M. A.; Cheeseman, J. R.; Montgomery, J., J. A.; Vreven, T.; Kudin, K. N.; Burant, J. C.; Millam, J. M.; Iyengar, S. S.; Tomasi, J.; Barone, V.; Mennucci, B.; Cossi, M.; Scalmani, G.; Rega, N.; Petersson, G. A.; Nakatsuji, H.; Hada, M.; Ehara, M.; Toyota, K.; Fukuda, R.; Hasegawa, J.; Ishida, M.; Nakajima, T.; Honda, Y.; Kitao, O.; Nakai, H.; Klene, M.; Li, X.; Knox, J. E.; Hratchian, H. P.; Cross, J. B.; Adamo, C.; Jaramillo, J.; Gomperts, R.; Stratmann, R. E.; Yazyev, O.; Austin, A. J.; Cammi, R.; Pomelli, C.; Ochterski, J. W.; Ayala, P. Y.; Morokuma, K.; Voth, G. A.; Salvador, P.; Dannenberg, J. J.; Zakrzewski, V. G.; Dapprich, S.; Daniels, A. D.; Strain, M. C.; Farkas, O.; Malick, D. K.; Rabuck, A. D.; Raghavachari, K.; Foresman, J. B.; Ortiz, J. V.; Cui, Q.; Baboul, A. G.; Clifford, S.; Cioslowski, J.; Stefanov, B. B.; Liu, G.; Liashenko, A.; Piskorz, P.; Komaromi, I.; Martin, R. L.; Fox, D. J.; Keith, T.; Al-Laham, M. A.; Peng, C. Y.; Nanayakkara, A.; Challacombe, M.; Gill, P. M. W.; Johnson, B.; Chen, W.; Wong, M. W.; Gonzalez, C.; Pople, J. A. *Gaussian03*, version E.01, Gaussian, Inc.: Pittsburgh, PA, 2003.

(75) Ponder, J. W. *TINKER*, version 5.1, Washington University: St. Louis, MO, 2010.

(76) Swope, W. C.; Andersen, H. C.; Berens, P. H.; Wilson, K. R. A computer simulation method for the calculation of equilibrium constants for the formation of physical clusters of molecules: Application to small water clusters. *J. Chem. Phys.* **1982**, *76*, 637–649.

(77) Berendsen, H. J. C.; Postma, J. P. M.; Gunsteren, W. F. v.; DiNola, A.; Haak, J. R. Molecular dynamics with coupling to an external bath. *J. Chem. Phys.* **1984**, *81*, 3684–3690.

(78) Shaw, K. E.; Woods, C. J.; Mulholland, A. J. Compatibility of quantum chemical methods and empirical (MM) water models in quantum mechanics/molecular mechanics liquid water simulations. *J. Phys. Chem. Lett.* **2009**, *1*, 219–223.

(79) 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. *J. Chem. Phys.* **1996**, *105*, 1968–1986.

(80) Riccardi, D.; Li, G. H.; Cui, Q. Importance of van der Waals interactions in QM/MM simulations. *J. Phys. Chem. B* **2004**, *108*, 6467–6478.

(81) Wang, B.; Truhlar, D. G. Combined quantum mechanical and molecular mechanical methods for calculating potential energy surfaces: Tuned and balanced redistributed-charge algorithm. *J. Chem. Theory Comput.* **6**, 359–369.

(82) Wang, B.; Truhlar, D. G. Geometry optimization using tuned and balanced redistributed charge schemes for combined quantum mechanical and molecular mechanical calculations. *Phys. Chem. Chem. Phys.* **13**, 10556–10564.

(83) Kästner, J.; Senn, H. M.; Thiel, S.; Otte, N.; Thiel, W. QM/MM free-energy perturbation compared to thermodynamic integration and umbrella sampling: Application to an enzymatic reaction. *J. Chem. Theory Comput.* **2006**, *2*, 452–461.

(84) Singh, U. C.; Kollman, P. A. An approach to computing electrostatic charges for molecules. *J. Comput. Chem.* **1984**, *5*, 129–145.

(85) Besler, B. H.; Merz, K. M., Jr.; Kollman, P. A. Atomic charges derived from semi-empirical methods. *J. Comput. Chem.* **1990**, *11*, 431–439.

(86) Zhang, Y.; Liu, H.; Yang, W. Free energy calculation on enzyme reactions with an efficient iterative procedure to determine minimum energy paths on a combined ab initio QM/MM potential energy surface. *J. Chem. Phys.* **2000**, *112*, 3483–3492.

(87) Donini, O.; Darden, T.; Kollman, P. A. QM-FE calculations of aliphatic hydrogen abstraction in citrate synthase and in solution: Reproduction of the effect of enzyme catalysis and demonstration that an enolate rather than an enol is formed. *J. Am. Chem. Soc.* **2000**, *122*, 12270–12280.

(88) Rod, T. H.; Ryde, U. Quantum mechanical free energy barrier for an enzymatic reaction. *Phys. Rev. Lett.* **2005**, *94*, 138302.

(89) Higashi, M.; Truhlar, D. G. Electrostatically embedded multi-configuration molecular mechanics based on the combined density functional and molecular mechanical method. *J. Chem. Theory Comput.* **2008**, *4*, 790–803.

# Parallel Generalized Born Implicit Solvent Calculations with NAMD

David E. Tanner,<sup>†,‡</sup> Kwok-Yan Chan,<sup>§,‡</sup> James C. Phillips,<sup>‡</sup> and Klaus Schulten<sup>\*,†,§</sup><sup>†</sup>Center for Biophysics and Computational Biology, University of Illinois at Urbana–Champaign, Urbana, Illinois<sup>‡</sup>Beckman Institute, University of Illinois at Urbana–Champaign, Urbana, Illinois<sup>§</sup>Department of Physics, University of Illinois at Urbana–Champaign, Urbana, Illinois

**ABSTRACT:** Accurate electrostatic descriptions of aqueous solvent are critical for simulation studies of biomolecules, but the computational cost of explicit treatment of solvent is very high. A computationally more feasible alternative is a generalized Born implicit solvent description which models polar solvent as a dielectric continuum. Unfortunately, the attainable simulation speedup does not transfer to the massive parallel computers often employed for simulation of large structures. Longer cutoff distances, spatially heterogeneous distribution of atoms, and the necessary 3-fold iteration over atom pairs in each timestep combine to challenge efficient parallel performance of generalized Born implicit solvent algorithms. Here, we report how NAMD, a parallel molecular dynamics program, meets the challenge through a unique parallelization strategy. NAMD now permits efficient simulation of large systems whose slow conformational motions benefit most from implicit solvent descriptions due to the inherent low viscosity. NAMD's implicit solvent performance is benchmarked and then illustrated in simulating the ratcheting *Escherichia coli* ribosome involving ~250 000 atoms.

## INTRODUCTION

Molecular dynamics (MD) is a computational method<sup>1</sup> employed for studying the dynamics of nanoscale biological systems on nanosecond to microsecond time scales.<sup>2</sup> Using MD, researchers can utilize experimental data from crystallography and cryo-electron microscopy (cryo-EM) to explore the functional dynamics of biological systems.<sup>3</sup>

Because biological processes take place in the aqueous environment of the cell, a critical component of any biological MD simulation is the solvent model employed.<sup>4,5</sup> An accurate solvent model must reproduce water's effect on solutes such as the free energy of solvation, dielectric screening of solute electrostatic interactions, hydrogen bonding, and van der Waals interactions with solute. For typical biological MD simulations, the solute is comprised of proteins, nucleic acids, lipids, or other small molecules.

Two main categories of solvent models are explicit and implicit solvents. Explicit solvents, such as SPC<sup>6</sup> and TIP3P,<sup>7</sup> represent water molecules explicitly as a collection of charged interacting atoms and calculate a simple potential function, such as Coulomb electrostatics, between solvent and solute atoms. Implicit solvent models, instead, ignore atomic details of the solvent and represent the presence of water indirectly through complex interatomic potentials between solute atoms only.<sup>8–10</sup> There are advantages and disadvantages of each solvent model.

Simulation of explicit water is both accurate and natural for MD but often computationally too demanding, not only since the inclusion of explicit water atoms increases a simulation's computational cost through the higher atom count but also because water slows down association and disassociation processes due to the relatively long relaxation times of interstitial water.<sup>11</sup> The viscous drag of explicit water also retards large conformational changes of macromolecules.<sup>12</sup>

An alternative representation of water is furnished by implicit solvent descriptions, which eliminate the need for explicit solvent

molecules. Implicit water remains always equilibrated to the solute. The absence of explicit water molecules also eliminates the viscosity imposed on simulated solutes, allowing faster equilibration of solute conformations and better conformational sampling. Examples of popular implicit solvent models are Poisson–Boltzmann electrostatics,<sup>13,14</sup> screened Coulomb potential,<sup>9,15</sup> analytical continuum electrostatics,<sup>16</sup> and generalized Born implicit solvent.<sup>17</sup>

The generalized Born implicit solvent (GBIS) model, used by MD programs CHARMM,<sup>18,19</sup> Gromacs,<sup>20,21</sup> Amber,<sup>22</sup> and NAMD,<sup>23,24</sup> furnishes a fast approximation for calculating the electrostatic interaction between atoms in a dielectric environment described by the Poisson–Boltzmann equation. The GBIS electrostatics calculation determines first the Born radius of each atom, which quantifies an atom's exposure to solvent, and, therefore, its dielectric screening from other atoms. The solvent exposure represented by Born radii can be calculated with varying speeds and accuracies<sup>25</sup> either by integration over the molecule's interior volume<sup>26,27</sup> or by pairwise overlap of atomic surface areas.<sup>17</sup> GBIS calculations then determine the electrostatic interaction between atoms based on their separation and Born radii.

GBIS has benefited MD simulations of small molecules.<sup>28</sup> For the case of large systems, whose large conformational motions<sup>29</sup> may benefit most from an implicit solvent description, but which must be simulated on large parallel computers,<sup>30</sup> the challenge to develop efficient parallel GBIS algorithms remains. In the following, we outline how NAMD addresses the computational challenges of parallel GBIS calculations and efficiently simulates large systems, demonstrated through benchmarks as well as simulations of the *Escherichia coli* ribosome, a RNA–protein complex involving ~250 000 atoms.

**Received:** August 22, 2011

**Published:** September 28, 2011

## METHODS

In order to characterize the challenges of parallel generalized Born implicit solvent (GBIS) simulations, we first introduce the key equations employed. We then outline the specific challenges that GBIS calculations pose for efficient parallel performance as well as how NAMD's implementation of GBIS achieves highly efficient parallel performance. GBIS benchmark simulations, which demonstrate NAMD's performance, as well as the ribosome simulations, are then described.

**Generalized Born Implicit Solvent Model.** The GBIS model<sup>8</sup> represents polar solvent as a dielectric continuum and, accordingly, screens electrostatic interactions between solute atoms. GBIS treats solute atoms as spheres of low protein dielectric ( $\epsilon_p = 1$ ), whose radius is the Bondi<sup>31</sup> van der Waals radius, in a continuum of high solvent dielectric ( $\epsilon_s = 80$ ).

The total electrostatic energy for atoms in a dielectric solvent is modeled as the sum of Coulomb and generalized Born (GB) energies:<sup>8</sup>

$$E_T^{\text{Elec}} = E_T^{\text{Coul}} + E_T^{\text{GB}} \quad (1)$$

The total Coulomb energy for the system of atoms is the sum over pairwise Coulomb energies:

$$E_T^{\text{Coul}} = \sum_i \sum_{j>i} E_{ij}^{\text{Coul}} \quad (2)$$

where the double summation represents all unique pairs of atoms within the interaction cutoff; the interaction cutoff for GBIS simulations is generally in the range 16–20 Å, i.e., longer than for explicit solvent simulations, where it is typically 8–12 Å. The reason for the wider cutoff is that particle-mesh Ewald summations, used to describe long-range Coulomb forces, cannot be employed for the treatment of long-range GBIS electrostatics.

The pairwise Coulomb energy,  $E_{ij}^{\text{Coul}}$  in eq 2, is

$$E_{ij}^{\text{Coul}} = (k_e/\epsilon_p)q_iq_j/r_{ij} \quad (3)$$

where  $k_e = 332$  (kcal/mol)Å/e<sup>2</sup> is the Coulomb constant,  $q_i$  is the charge on atom  $i$ , and  $r_{ij}$  is the distance between atoms  $i$  and  $j$ . The total GB energy for the system of atoms is the sum over pairwise GB energies and self-energies given by the expression

$$E_T^{\text{GB}} = \underbrace{\sum_i \sum_{j>i} E_{ij}^{\text{GB}}}_{\text{pair}} + \underbrace{\sum_i E_{ii}^{\text{GB}}}_{\text{self}} \quad (4)$$

where the pair-energies and self-energies are defined as<sup>8</sup>

$$E_{ij}^{\text{GB}} = -(k_e D_{ij})q_iq_j/f_{ij}^{\text{GB}} \quad (5)$$

Here,  $D_{ij}$  is the pairwise dielectric term,<sup>32</sup> which contains the contribution of an implicit ion concentration to the dielectric screening, and is expressed as

$$D_{ij} = (1/\epsilon_p) - \exp(-\kappa f_{ij}^{\text{GB}})/\epsilon_s \quad (6)$$

where  $\kappa^{-1}$  is the Debye screening length, which represents the length scale over which mobile solvent ions screen electrostatics. For an ion concentration of 0.2 M, room temperature water has a Debye screening length of  $\kappa^{-1} = \sim 7$  Å.  $f_{ij}^{\text{GB}}$  is<sup>8</sup>

$$f_{ij}^{\text{GB}} = \sqrt{r_{ij}^2 + \alpha_i \alpha_j \exp(-r_{ij}^2/4\alpha_i \alpha_j)} \quad (7)$$

The form of the pairwise GB energy in eq 5 is similar to the form of the pairwise Coulomb energy in eq 3 but is of opposite sign

and replaces the  $1/r_{ij}$  distance dependence by  $1/f_{ij}^{\text{GB}}$ . The GB energy bears a negative sign because the electrostatic screening counteracts the Coulomb interaction. The use of  $f_{ij}^{\text{GB}}$ , instead of  $r_{ij}$ , in eq 5 heavily screens the electrostatic interaction between atoms which are either far apart or highly exposed to solvent. The more exposed an atom is to high solvent dielectric, the more it is screened electrostatically, represented by a smaller Born radius,  $\alpha_i$ .

Accurately calculating the Born radius is central to a GBIS model as the use of perfect Born radii allows the GBIS model to reproduce, with high accuracy, the electrostatics and solvation energies described by the Poisson–Boltzmann equation<sup>33</sup> and does it much faster than a Poisson–Boltzmann or explicit solvent treatment.<sup>34</sup> Different GBIS models vary in how the Born radius is calculated; models seek to suggest computationally less expensive algorithms without undue sacrifice in accuracy. Many GBIS models<sup>35</sup> calculate the Born radius by assuming atoms are spheres whose radius is the Bondi<sup>31</sup> van der Waals radius and determine an atom's exposure to solute through the sum of overlapping surface areas with neighboring spheres.<sup>36</sup> The more recent GBIS model of Onufriev, Bashford, and Case (GB<sup>OBC</sup>), applied successfully to MD of macromolecules<sup>37,38</sup> and adopted in NAMD, calculates the Born radius as

$$\alpha_i = [(1/\rho_{i0}) - (1/\rho_i) \tanh(\delta\psi_i - \beta\psi_i^2 + \gamma\psi_i^3)]^{-1} \quad (8)$$

where  $\psi_i$ , the sum of surface area overlap with neighboring spheres, is calculated through

$$\psi_i = \rho_{i0} \sum_j H(r_{ij}, \rho_i, \rho_j) \quad (9)$$

As explained in prior studies,<sup>35,36,38</sup>  $H(r_{ij}, \rho_i, \rho_j)$  is the surface area overlap of two spheres based on their relative separation,  $r_{ij}$ , and radii,  $\rho_i$  and  $\rho_j$ ; the parameters  $\delta$ ,  $\beta$ , and  $\gamma$  in eq 8 have been calculated to maximize agreement between Born radii described by eq 8 and those derived from Poisson–Boltzmann electrostatics.<sup>38</sup>  $\rho_i$  and  $\rho_j$  are the Bondi<sup>31</sup> van der Waals radii of atoms  $i$  and  $j$ , respectively, while  $\rho_{i0}$  is the reduced radius,  $\rho_{i0} = \rho_i - 0.09$  Å, as required by GB<sup>OBC</sup>.<sup>38</sup>

The total electrostatic force acting on an atom is the sum of Coulomb and GB forces; the net Coulomb force on an atom is given by

$$\vec{F}_i^{\text{Coul}} = - \sum_j [dE_T^{\text{Coul}}/dr_{ij}] \hat{r}_{ij} \quad (10)$$

whose derivative ( $dE_T^{\text{Coul}}/dr_{ij}$ ) is inexpensive to calculate. The required derivatives ( $dE_T^{\text{GB}}/dr_{ij}$ ) for the GB force, however, are much more expensive to calculate because  $E_{ij}^{\text{GB}}$  depends on interatomic distances,  $r_{ij}$ , both directly (c.f., eqs 5 and 7) and indirectly through the Born radius (c.f., eqs 5, 7, 8, and 9). The net GB force on an atom is given by

$$\begin{aligned} \vec{F}_i^{\text{GB}} &= - \sum_j [dE_T^{\text{GB}}/dr_{ij}] \hat{r}_{ij} \\ &= - \sum_j \left[ \sum_k \sum_{l>k} (\partial E_T^{\text{GB}}/\partial r_{kl})(dr_{kl}/dr_{ij}) \right. \\ &\quad \left. + \sum_k (\partial E_T^{\text{GB}}/\partial \alpha_k)(d\alpha_k/dr_{ij}) \right] \hat{r}_{ij} \\ &= - \sum_j \left[ \partial E_T^{\text{GB}}/\partial r_{ij} + (\partial E_T^{\text{GB}}/\partial \alpha_i)(d\alpha_i/dr_{ij}) \right. \\ &\quad \left. + (\partial E_T^{\text{GB}}/\partial \alpha_j)(d\alpha_j/dr_{ij}) \right] \hat{r}_{ij} \end{aligned} \quad (11)$$



with  $\vec{r}_{ij} = \vec{r}_j - \vec{r}_i$ . The required partial derivative of  $E_T^{\text{GB}}$  with respect to a Born radius,  $\alpha_k$ , is

$$\begin{aligned} \partial E_T^{\text{GB}} / \partial \alpha_k &= \sum_i \sum_{j>i} [\partial E_{ik}^{\text{GB}} / \partial \alpha_k + \partial E_{kj}^{\text{GB}} / \partial \alpha_k] \\ &+ \sum_i \partial E_{ii}^{\text{GB}} / \partial \alpha_k \end{aligned} \quad (12)$$

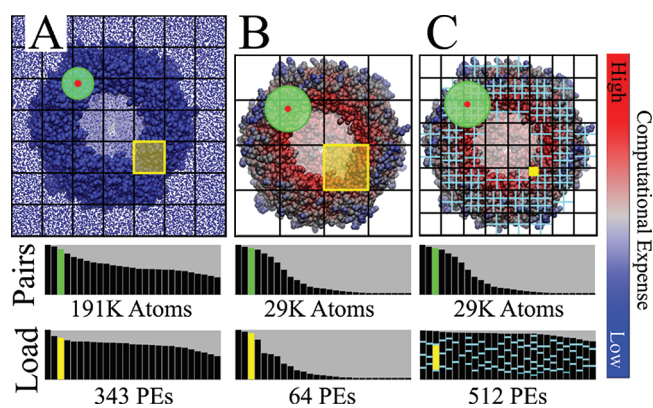
The summations in eqs 9, 11, and 12 require three successive iterations over all pairs of atoms for each GBIS force calculation, whereas calculating Coulomb forces for an explicit solvent simulation requires only one such iteration over atom pairs. Also, because of the computational complexity of the above GBIS equations, the total cost of calculating the pairwise GBIS force between pairs of atoms is  $\sim 7\times$  higher than the cost for the pairwise Coulomb force. For large systems and long cutoffs, the computational expense of implicit solvent simulations can exceed that of explicit solvent simulations; however, in this case, an effective speed-up over explicit solvent still arises due to faster conformational exploration, as will be illustrated below. The trade-off between the speed-up of implicit solvent models and the higher accuracy of explicit solvent models is still under investigation.<sup>34</sup> Differences between GBIS and Coulomb force calculations create challenges for parallel GBIS simulations that do not arise in explicit solvent simulations.

**Challenges in Parallel Calculation of GBIS Forces.** Running a MD simulation in parallel requires a scheme to decompose the simulation calculation into independent work units that can be executed simultaneously on parallel processors; the scheme employed for decomposition strongly determines how many processors the simulation can efficiently utilize and, therefore, how fast the simulation will be. For example, a common decomposition scheme, known as spatial or domain decomposition, divides the simulated system into a three-dimensional grid of spatial domains whose side length is the interaction cutoff distance. Because the atoms within each spatial domain are simulated on a single processor, the number of processors utilized equals the number of domains.

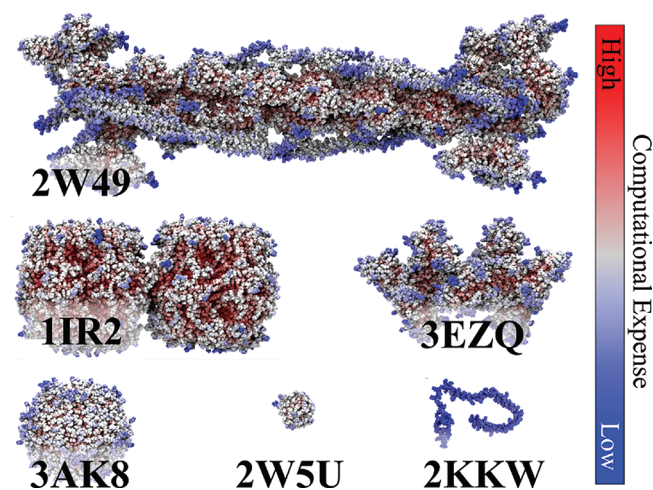
Although explicit solvent MD simulations perform efficiently in parallel, even for simple schemes such as naive domain decomposition, the GBIS model poses unique challenges for simulating large systems on parallel computers. We outline here the three challenges arising in parallel GBIS calculations and how NAMD addresses them. For the sake of concreteness, we use the SEp22 dodecamer (PDB ID: 3AK8) as an example, as shown in Figure 1.

*Challenge 1: Dividing Workload among Processors.* With a 12 Å cutoff, traditional domain decomposition divides the SEp22 dodecamer explicit solvent simulation (190 000 protein and solvent atoms) into  $7 \times 7 \times 7 = 343$  domains, which efficiently utilize 343 processors (see Figure 1A). Unfortunately, with a 16 Å cutoff for the implicit solvent treatment, the same decomposition scheme divides the system (30 000 protein atoms) into  $4 \times 4 \times 4 = 64$  domains, which can only utilize 64 processors (see Figure 1B). An efficient parallel GBIS implementation must employ a decomposition scheme which can divide the computational work among many (hundreds or thousands) processors (see Figure 1C).

*Challenge 2: Workload Imbalance from Spatially Heterogeneous Atom Distribution.* Due to the lack of explicit water atoms, the spatial distribution of atoms in a GBIS simulation (see Figure 2) is not uniform, as it is for an explicit solvent simulation



**Figure 1.** Work decomposition for implicit solvent and explicit solvent simulations. An SEp22 dodecamer is shown (front removed to show interior) with an overlaid black grid illustrating domain decomposition for explicit solvent (A), implicit solvent (B), and NAMD's highly parallel implicit solvent (C). Atoms are colored according to the relative work required to calculate the net force, with blue being the least expensive and red being the most expensive. The number of neighbor Pairs within the interaction cutoff (green circle, pairs) for an atom (red circle) varies more in implicit solvent than explicit solvent, as does the number of atoms within a spatial domain (yellow box, Load), each domain being assigned to a single processor (PE). (A) Because explicit solvent has a spatially homogeneous distribution of atoms, it has a balanced workload among processors using simple domain decomposition. (B) Domain decomposition with implicit solvent suffers from the spatially heterogeneous atom distribution; the workload on each processor varies highly. (C) NAMD's implicit solvent model (cyan grid representing force decomposition and partitioning), despite having a varying number of atoms per domain and varying computational cost per atom, still achieves a balanced workload among processors.



**Figure 2.** Biomolecular systems in benchmark. The performance of NAMD's parallel GBIS implementation was tested on six structures (Protein Data Bank IDs shown); the number of atoms and interactions are listed in Table 2. To illustrate the spatially heterogeneous distribution of work, each atom is colored by the relative time required to compute its net force, with blue being the fastest and red the being slowest.

(see Figure 1A). Some domains contain densely packed atoms, while others are empty (see Figure 1B). Because the number of atoms varies highly among implicit solvent domains, the workload assigned to each processor also varies highly. The highly

varying workload among processors for domain decomposition causes the naive decomposition scheme to be inefficient and, therefore, slow. An efficient parallel GBIS implementation must assign and maintain an equal workload on each processor (see Figure 1C).

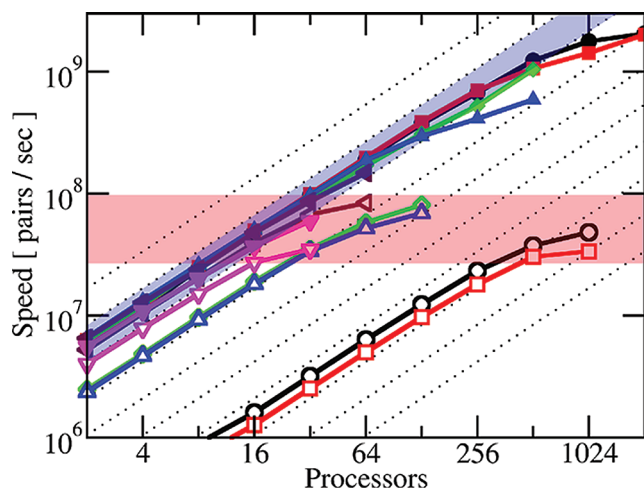
**Challenge 3: Three Iterations over Atom Pairs per Timestep.** Instead of requiring one iteration over atom pairs to calculate electrostatic forces, GBIS requires three independent iterations over atom pairs (c.f., eqs 9, 11, and 12). Because each of these iterations depends on the previous iteration, the cost associated with communication and synchronization is tripled. An efficient parallel GBIS implementation must schedule communication and computation on each processor as to maximize efficiency.

**Parallelization Strategy.** NAMD's unique strategy<sup>39</sup> for fast parallel MD simulations<sup>23</sup> enables it to overcome the three challenges of parallel GBIS simulations. NAMD divides GBIS calculations into many small work units using a three-tier decomposition scheme, assigns a balanced load of work units to processors, and schedules work units on each processor to maximize efficiency.

NAMD's three-tier work decomposition scheme<sup>40</sup> directly addresses challenge 1 of parallel GBIS calculations. NAMD first employs domain decomposition to initially divide the system into a three-dimensional grid of spatial domains. Second, NAMD assigns a force work unit to calculate pairwise forces within each domain and between each pair of adjacent domains. Third, each force work unit is further partitioned into up to 10 separate work units, where each partition calculates only one-tenth of the atom pairs associated with the force work unit. Dividing force work units into partitions based on computational expense avoids the unnecessary communication overhead arising from further partitioning already inexpensive force work units belonging to underpopulated domains. Adaptively partitioning the force work units based on computational expense improves NAMD's parallel performance even for non-implicit solvent simulations. NAMD's decomposition scheme is able to finely divide simulations into many (~40 000 for SEp22 dodecamer) small work units and, thereby, utilize thousands of processors.

NAMD's load balancer, a tool employed to ensure each processor carries an equivalent workload, initially distributes work units evenly across processors, thus partially overcoming challenge 2 of parallel GBIS calculations. However, as atoms move during a simulation, the number of atoms in each domain fluctuates (more so than for the explicit solvent case), which causes the computational workload on each processor to change. NAMD employs a measurement-based load balancer to maintain a uniform workload across processors during a simulation; periodically, NAMD measures the computational cost associated with each work unit and redistributes work units to new processors as required to maintain a balanced workload among processors. By continually balancing the workload, NAMD is capable of highly efficient simulations despite spatially heterogeneous atom distributions common to implicit solvent descriptions.

Though the three iterations over atom pairs hurt parallel efficiency by requiring additional (compared to the explicit solvent case) communication and synchronization during each timestep, NAMD's unique communication scheme is able to maintain parallel efficiency. Unlike most MD programs, NAMD is capable of scheduling work units on each processor in an order which overlaps communication and computation, thus maximizing efficiency. NAMD's overall parallel strategy of work decomposition, workload balancing, and work unit scheduling permits fast and efficient parallel GBIS simulations of even very large systems.



**Figure 3.** Parallel performance of NAMD and domain decomposition implicit solvent. Computational speed (pairwise interactions per second) for six biomolecular systems (see Figure 2): 11R2 (black circle), 2W49 (red square), 3AK8 (green diamond), 3EZQ (blue up triangle), 2WSU (maroon left triangle), and 2KKW (magenta down triangle). NAMD's parallel implicit solvent implementation (solid shapes) performs extremely well in parallel, as seen by the performance increasing linearly (blue highlight) with the number of processors, and is independent of system size. Performance of domain decomposition implicit solvent, however, suffers in parallel (empty shapes). Not only does there appear to be a maximum speed of  $10^8$  pairs/s (red highlight) regardless of processor count but the large systems (11R2, 2W49) also perform at significantly lower efficiency than the small systems (2ESU, 2KKW). Diagonal dotted lines represent perfect speedup. See also Table 2.

**Performance Benchmark.** To demonstrate the success of NAMD's parallel GBIS strategy, protein systems of varying sizes and configurations were simulated on 2–2048 processor cores using NAMD version 2.8. We also compare against an implementation of domain decomposition taken from Amber's PMEMD version 9,<sup>22</sup> which also contains the original implementation of the GB<sup>OB</sup>C implicit solvent model.<sup>38</sup> The benchmark consists of six systems, listed in Table 2 and displayed in Figure 2, chosen to represent small (2000 atoms), medium (30 000 atoms), and large (150 000 atoms) systems.

The following simulation parameters were employed for the benchmark simulations. A value of 16 Å was used for both non-bonded interaction cutoff and the Born radius calculation cutoff. An implicit ion concentration of 0.3 M was assumed. A timestep of 1 fs was employed with all forces being evaluated every step. System coordinates were not written to a trajectory file. For the explicit solvent simulation (Table 2: 3AK8-E), nonbonded interactions were cut off and smoothed between 10 and 12 Å, with PME<sup>41</sup> electrostatics, which require periodic boundary conditions, being evaluated every four steps.

Simulations were run on 2.3 GHz processors with 1 GB/s network interconnect for 600 steps. NAMD's speed is reported during simulation and was averaged over the last 100 steps (the first 500 steps are dedicated to initial load balancing). The speed of the domain decomposition implementation, in units seconds per timestep, was calculated as ("Master NonSetup CPU time")/(total steps); simulating up to 10 000 steps did not return noticeably faster speeds. Table 2 reports the simulation speeds in seconds/step for each benchmark simulation; Figure 3 presents simulation speeds scaled by system size

**Table 1. Total Electrostatic Energy of Benchmark Systems<sup>a</sup>**

	NAMD	Amber	error
2KKW	-5271.42	-5271.23	$3.6 \times 10^{-5}$
2WSU	-6246.10	-6245.88	$3.5 \times 10^{-5}$
3EZQ	-94 291.43	-94 288.17	$3.4 \times 10^{-5}$
3AK8	-89 322.23	-89 319.13	$3.4 \times 10^{-5}$
2W49	-396 955.31	-396 941.54	$3.4 \times 10^{-5}$
1IR2	-426 270.23	-426 255.45	$3.4 \times 10^{-5}$

<sup>a</sup>To validate NAMD's implicit solvent implementation, the total electrostatic energy (in units kcal/mol) of the six benchmark systems was calculated by NAMD and Amber implementations of the GB<sup>DBC</sup> implicit solvent;<sup>38</sup> error is calculated through eq 13.

in terms of the number of pairwise interactions per second (pips) calculated.

**Implementation Validation.** The correctness of our GBIS implementation was validated by comparison to the method's<sup>38</sup> original implementation in Amber.<sup>22</sup> Comparing the total electrostatic energy (see eq 1) of the six test systems as calculated by NAMD and Amber demonstrates their close agreement. Indeed, Table 1 shows that the relative error, defined through

$$\text{error} = |(E_T^{\text{Elec}}(\text{NAMD}) - E_T^{\text{Elec}}(\text{Amber})) / E_T^{\text{Elec}}(\text{Amber})| \quad (13)$$

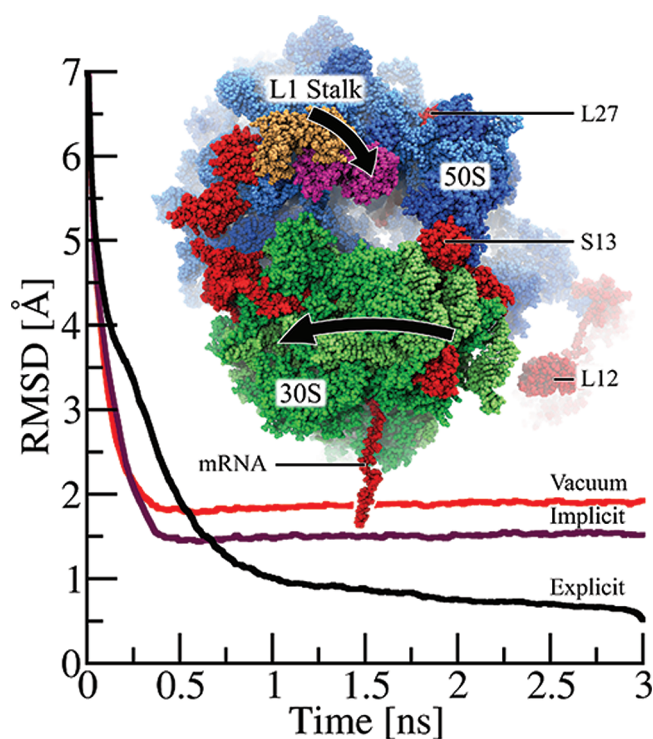
is less than  $4 \times 10^{-5}$  for all structures in Figure 2.

**Molecular Dynamics Flexible Fitting of Ribosome.** To illustrate the utility of NAMD's parallel GBIS implementation, we simulate the *Escherichia coli* ribosome. The ribosome is the cellular machine that translates genetic information on mRNA into protein chains.

During translation, tRNAs, with their anticodon loops to be matched to the genetic code on mRNA, carry amino acids to the ribosome. The synthesized protein chain is elongated by one amino acid each time a cognate tRNA (with its anticodon loop complementary to the next mRNA codon) brings an amino acid to the ribosome; a peptide bond is formed between the new amino acid and the existing protein chain. During protein synthesis, the ribosome complex fluctuates between two conformational states, namely, the so-called classical and ratcheted state.<sup>42</sup>

During transition from the classical to ratcheted state, the ribosome undergoes multiple, large conformational changes, including an intersubunit rotation between its 50S and 30S subunits<sup>42</sup> and the closing of its L1 stalk in the 50S subunit<sup>43</sup> (see Figure 4). The large conformational changes during the transition from the classical to ratcheted state are essential for translation,<sup>44</sup> as suggested by previous cryo-EM data.<sup>45</sup> To demonstrate the benefits of NAMD GBIS, we simulate the large conformational changes during ratcheting of the ~250 000-atom ribosome using molecular dynamics flexible fitting.

The molecular dynamics flexible fitting (MDFF) method<sup>3,46,47</sup> is a MD simulation method that matches crystallographic structures to an electron microscopy (EM) map. Crystallographic structures often correspond to nonphysiological states of biopolymers, while EM maps correspond often to functional intermediates of biopolymers. MDFF-derived models of the classical and ratcheted state ribosome provide atomic-level details crucial to understanding protein elongation in the ribosome. The MDFF method adds to a conventional MD simulation an EM map-derived potential, thereby driving a crystallographic structure



**Figure 4.** Molecular dynamics flexible fitting (MDFF) of the ribosome with NAMD's GBIS method. While matching the 250 000-atom classical ribosome structure into the EM map of a ratcheted ribosome, the 30S subunit (green) rotates relative to the 50S subunit (blue) and the L1 stalk moves 30 Å from its classical (tan) to its ratcheted (magenta) position. Highlighted (red) are regions where the implicit solvent structure agrees with the explicit solvent structure much more closely than does the in vacuo structure. The root-mean-squared deviation ( $\text{RMSD}_{\text{sol,exp}}(t)$ ) of the ribosome, defined in eq 14, with the final fitted explicit solvent structure as a reference, is plotted over time for explicit solvent ( $\text{RMSD}_{\text{exp,exp}}(t)$  in black), implicit solvent ( $\text{RMSD}_{\text{imp,exp}}(t)$  in purple), and in vacuo ( $\text{RMSD}_{\text{vac,exp}}(t)$  in red) MDFF. While the explicit solvent MDFF calculation requires ~1.5–2 ns to converge to its final structure, both implicit solvent and vacuum MDFF calculation require only 0.5 ns to converge. As seen by the lower RMSD values for  $t > 0.5$  ns, the structure derived from the implicit solvent fitting agrees more closely with the final explicit solvent structure than does the in vacuo structure. While this plot illustrates only the overall improvement of the implicit solvent structure over the in vacuo structure, the text discusses key examples of ribosomal proteins (L27, S13, and L12) whose structural quality is significantly improved by the use of implicit solvent.

toward the conformational state represented by an EM map. MDFF was previously applied to successfully match crystallographic structures of the ribosome to ribosome functional states as seen in EM.<sup>48–52</sup> Shortcomings of MDFF are largely due to the use of in vacuo simulations; such use was necessary hitherto as simulations in explicit solvent proved too cumbersome. Implicit solvent MDFF simulations promise a significant improvement of the MDFF method. We applied MDFF here, therefore, to fit an atomic model of a classical state ribosome into an EM map of a ratcheted state ribosome.<sup>45</sup>

The classical state in our simulations is an all-atom ribosome structure<sup>53</sup> with 50S and 30S subunits taken from PDB IDs 2J2V and 2J2U, respectively,<sup>54</sup> and the complex was fitted to an 8.9 Å resolution classical state EM map.<sup>45</sup> In the multistep protocol for fitting this classical state ribosome to a ratcheted state map,<sup>46</sup> the actual ribosome is fitted first, followed by fitting the tRNAs.

Since the fitting of the ribosome itself exhibits the largest conformational changes (intersubunit rotation and L1-stalk closing), we limit our MDFF calculation here to the ribosome and do not include tRNAs.

Three MDFF simulations were performed using NAMD<sup>23</sup> and analyzed using VMD.<sup>55</sup> The MDFF simulations are carried out in explicit TIP3P<sup>7</sup> solvent, in implicit solvent, and in vacuo. All simulations were performed in the NVT ensemble with the AMBER99 force field,<sup>56</sup> employing the SB<sup>57</sup> and BSC0<sup>58</sup> corrections and accounting for modified ribonucleosides.<sup>59</sup> The grid scaling parameter,<sup>3</sup> which controls the balance between the MD force field and the EM-map derived force field, was set to 0.3. Simulations were performed using a 1 fs timestep with non-bonded forces being evaluated every two steps. Born radii were calculated using a cutoff of 14 Å, while the nonbonded forces were smoothed and cut off between 15 and 16 Å. An implicit ion concentration of 0.1 M was assumed with protein and solvent dielectric set to 1 and 80, respectively. A Langevin thermostat with a damping coefficient of 5 ps<sup>-1</sup> was employed to hold the temperature to 300 K. In the explicit solvent simulation, the ribosome was simulated in a periodic box of TIP3P water<sup>7</sup> including an explicit ion concentration of 0.1 M, with nonbonded forces cut off at 10 Å and long-range electrostatics calculated by PME every four steps. The in vacuo simulation utilized the same parameters as explicit solvent, but without the inclusion of solvent or bulk ions, and neither PME nor periodicity were employed.

Each system was minimized for 5000 steps before performing MDFF for 3 ns. For the explicit solvent simulation, an additional 0.5 ns equilibration of water and ions was performed, with protein and nucleic acids restrained, before applying MDFF.

To compare the behavior of solvent models during the ribosome simulations, we calculate the root-mean-square deviation between models as

$$\text{RMSD}_{\text{sol,ref}}(t) = \sqrt{\frac{\sum_i^N [\vec{r}_{i,\text{sol}}(t) - \vec{r}_{i,\text{ref}}]^2}{N}} \quad (14)$$

where  $\vec{r}_{i,\text{sol}}(t)$  denotes the atomic coordinates at time  $t$  of the simulation corresponding to one of the three solvent models (exp, imp, or vac) and  $\vec{r}_{i,\text{ref}}$  denotes the atomic coordinates for the last timestep ( $t_f = 3$  ns) of the simulation using the reference solvent model (exp, imp, or vac) as specified below. Unless otherwise specified, the summation is over the  $N = 146\,000$  heavy atoms excluding the mRNA, L10, and L12 protein segments, which are too flexible to be resolved with the cry-EM method.

## RESULTS

**Performance Benchmarks.** The results of the GBIS benchmark simulations are listed in Table 2. Figure 3 illustrates the simulation speeds, scaled by system size, as the number of pairwise interactions per second (pips) calculated. For a perfectly efficient algorithm, pips would be independent of system size or configuration and would increase proportionally with processor count. NAMD's excellent parallel GBIS performance is demonstrated, as pips is nearly the same for all six systems and increases almost linearly with processor count, as highlighted (in blue) in Figure 3.

The domain decomposition algorithm<sup>22</sup> performs equally well for small systems, but its performance suffers significantly when system size and processor count are increased. The domain

**Table 2. NAMD and Domain Decomposition Benchmark Data<sup>a</sup>**

	PDB ID						
	2KKW	2W5U	3EZQ	3AK8	3AK8-E	2W49	1IR2
atoms	2016	2412	27 600	29 479	191 686	138 136	149 860
pairs	0.25 M	1 M	13.2 M	15.7 M	63.8 M	65.5 M	99.5 M
# procs	NAMD s/step						
2	0.0440	0.167	2.01	2.87	2.25	10.4	16.3
4	0.0225	0.0853	1.00	1.44	1.12	5.47	8.70
8	0.0122	0.0456	0.505	0.726	0.568	2.61	4.09
16	0.00664	0.0228	0.260	0.371	0.286	1.31	2.05
32	0.00412*	0.0126	0.136	0.191	0.146	0.661	1.03
64		0.00736*	0.0700	0.105	0.0868	0.333	0.520
128			0.0447	0.0575	0.0523	0.169	0.267
256			0.0321	0.0340	0.0288*	0.0935	0.148
512			0.0224*	0.0171*		0.0618	0.0806
1024						0.0461*	0.0563*
2048						0.0326	0.0486
# procs	domain decomposition s/step						
2	0.0613	0.208	5.57	7.15	306	373	
4	0.0312	0.104	2.83	3.67	169	203	
8	0.0163	0.0535	1.43	1.84	92.5	111	
16	0.0090*	0.0277	0.727	0.934	51.6	62.0	
32	0.0070	0.0162*	0.391	0.506	25.9	31.3	
64		0.013	0.254*	0.307*	13.1	15.7	
128			0.191	0.220	6.74	8.08	
256					3.63	4.28	
512					2.16*	2.66*	
1024					1.95	2.07	

<sup>a</sup>Speed, in units seconds/step, for both NAMD and the domain decomposition algorithm for the six test systems (see Figure 2) on 2-2048 processors (procs). Also listed are the number of atoms and pairs of atoms (in millions, M) within the cutoff (16 Å for implicit solvent and 12 Å for explicit solvent) in the initial structure. Data are not presented for higher processor counts with slower simulation speeds. An asterisk marks the highest processor count for which doubling the number of processors increased simulation speed by at least 50%. 3AK8-E uses explicit solvent. Each simulated system demonstrates that NAMD can efficiently utilize at least twice the number of processors as domain decomposition and, thereby, achieves simulation speeds much faster than for domain decomposition.

decomposition implementation also appears to be limited to a pips maximum of  $10^8$  pairs/s across all system sizes, no matter how many processors are used, as highlighted (in red) in Figure 3. NAMD runs efficiently on twice the number of processors compared to domain decomposition and greatly outperforms it for the large systems tested. The SEp22 dodecamer timings for both implicit (3AK8) and explicit (3AK8-E) solvent reported in Table 2 demonstrate that NAMD's parallel GBIS implementation is as efficient as its parallel explicit solvent capability. We note that the simulation speed for GBIS can be further increased, without a significant loss of accuracy, by shortening either the interaction or Born radius calculation cutoff distance.

**Ribosome Simulation.** To demonstrate the benefit of NAMD's GBIS capability for simulating large structures, a high-resolution classical state ribosome structure was fitted into a low-resolution

**Table 3. Root-Mean-Square Deviation (RMSD<sub>sol,ref</sub>(3 ns) in Å) between the Three Final Ribosome Structures Matched Using Explicit Solvent, GBIS, and in Vacuo MDFF<sup>a</sup>**

		reference		
		exp	imp	vac
solvent	exp	0	1.5	1.9
	imp	1.5	0	2.1
	vac	1.9	2.1	0

<sup>a</sup> The GBIS and explicit solvent MDFF structures closely agree, as seen by RMSD<sub>imp,exp</sub>(3 ns) = 1.5 Å, while the in vacuo MDFF structure deviates from the explicit solvent MDFF structure by RMSD<sub>vac,exp</sub>(3 ns) = 1.9 Å. See also Figure 4.

ratcheted state EM map in an in vacuo MDFF simulation as well as MDFF simulations employing explicit and implicit solvent. During the MDFF simulation, the ribosome undergoes two major conformational changes: closing of the L1 stalk and rotation of the 30S subunit relative to the 50S subunit (see Figure 4).

To compare the rate of convergence and relative accuracy of solvent models, the RMSD values characterizing the three MDFF simulations were calculated using eq 14. Figure 4 plots RMSD<sub>exp,exp</sub>(*t*), RMSD<sub>imp,exp</sub>(*t*), and RMSD<sub>vac,exp</sub>(*t*), which compare each MDFF simulation against the final structure reached in the explicit solvent case. We note that using the initial rather than final structure as the reference could yield a slightly different characterization of convergence,<sup>60</sup> e.g., a slightly different convergence time. As manifested by RMSD<sub>imp,exp</sub>(*t*) and RMSD<sub>vac,exp</sub>(*t*), the implicit solvent and vacuum MDFF calculations converge to their respective final structures in 0.5 ns compared to ~1.5–2 ns for the explicit solvent case, i.e., for RMSD<sub>exp,exp</sub>(*t*).

The final structures obtained from the MDFF simulations are compared in Table 3 through the RMSD<sub>sol,ref</sub>(*t*) values for *t* = 3 ns. The ribosome structure from GBIS MDFF closely agrees with the one from explicit solvent MDFF, as indicated by the value RMSD<sub>imp,exp</sub>(3 ns) = 1.5 Å; the in vacuo MDFF ribosome structure, however, compares less favorably with the explicit solvent MDFF structure, as suggested by the larger value RMSD<sub>vac,exp</sub>(3 ns) = 1.9 Å. While the 0.4 Å improvement in RMSD of the GBIS MDFF, over in vacuo MDFF, structure implies an overall enhanced quality, certain regions of the ribosome are particularly improved.

The regions with the highest structural improvement (highlighted red in Figure 4) belong to segments at the exterior of the ribosome and to segments not resolved by and, therefore, not coupled to the EM map, i.e., not being directly shaped by MDFF. For proteins at the exterior of the ribosome, GBIS MDFF produces higher quality structures than in vacuo MDFF, because these proteins are highly exposed to the solvent and, therefore, require a solvent description. The structural improvement for several exterior solvated proteins, calculated by RMSD<sub>vac,exp</sub>(3 ns) – RMSD<sub>imp,exp</sub>(3 ns), is 3.5, 2.4, and 1.6 Å for ribosomal proteins S6, L27, and S13 (highlighted red in Figure 4), respectively. Accurate modeling of these proteins is critical for studying the translation process of the ribosome. The L27 protein, for example, not only facilitates the assembly of the 50S subunit but it also ensures proper positioning of the new amino acid for peptide bond formation.<sup>61</sup> The S13 protein, located at the interface between subunits, is critical to the control of mRNA and tRNA translocation within the ribosome.<sup>62</sup>

The use of GBIS for MDFF also increases structural quality in regions where the EM map does not resolve the ribosome's structure, and therefore, MDFF does not directly influence conformation. Though it is most important that MDFF correctly models structural regions defined in the EM map, it is also desirable that it correctly describes regions of crystal structures not resolved by the EM map. The structural improvement, over in vacuo MDFF, of the unresolved segments is 8.3 Å for mRNA and 4.9 Å for L12 (highlighted red in Figure 4). The L12 segment is a highly mobile ribosomal protein in the 50S subunit that promotes binding of factors which stabilize the ratcheted conformation; L12 also promotes GTP hydrolysis, which leads to mRNA translocation.<sup>63</sup> As clearly demonstrated, the use of GBIS MDFF, instead of in vacuo MDFF, improves the MDFF method's accuracy for matching crystallographic structures to EM maps, particularly for highly solvated or unresolved proteins.

To compare the computational performance of the solvent models for MDFF, each ribosome simulation was benchmarked on 1020 processor cores (3.5 GHz processors with 5 GB/s network interconnect). The simulation speed for explicit solvent MDFF is 3.6 ns/day. For implicit solvent MDFF, it is 5.2 ns/day, and for vacuum MDFF, it is 37 ns/day. GBIS MDFF performs 50% faster than explicit solvent MDFF but 7 times slower than in vacuo MDFF. NAMD's GBIS implementation is clearly able to achieve a more accurate MDFF match of the ribosome structure (see Table 3) than does an in vacuo MDFF calculation and does so at a lower computational cost than explicit solvent MDFF.

## CONCLUSIONS

The generalized Born implicit solvent (GBIS) model has long been employed for molecular dynamics simulations of relatively small biomolecules. NAMD's unique GBIS implementation can also simulate very large systems, such as the entire ribosome, and does so efficiently on large parallel computers. The new GBIS capability of NAMD will be beneficial to accelerating in simulations the slow motions common to large systems by eliminating viscous drag from water.

## AUTHOR INFORMATION

### Corresponding Author

\*E-mail: kschulte@ks.uiuc.edu.

## ACKNOWLEDGMENT

The authors thank Chris Harrison for helpful discussion as well as Laxmikant Kalé's Parallel Programming Laboratory for parallel implementation advice. This work was supported by the National Science Foundation (NSF PHY0822613) and National Institutes of Health (NIH P41-RR005969) grants to K.S. and a Molecular Biophysics Training Grant fellowship to D.E.T. This research was supported in part by the National Science Foundation through TeraGrid HPC resources provided by the Texas Advanced Computing Center (TACC) at The University of Texas at Austin (<http://www.tacc.utexas.edu>) under grant number TG-MCA93S028. This work used the Extreme Science and Engineering Discovery Environment (XSEDE), which is supported by National Science Foundation grant number OCI-1053575.

## REFERENCES

- (1) Lee, E. H.; Hsin, J.; Sotomayor, M.; Comellas, G.; Schulten, K. *Structure* **2009**, *17*, 1295–1306.

- (2) Freddolino, P. L.; Schulten, K. *Biophys. J.* **2009**, *97*, 2338–2347.
- (3) Trabuco, L. G.; Villa, E.; Schreiner, E.; Harrison, C. B.; Schulten, K. *Methods* **2009**, *49*, 174–180.
- (4) Daura, X.; Mark, A. E.; van Gunsteren, W. F. *Comput. Phys. Commun.* **1999**, *123*, 97–102.
- (5) Daidone, I.; Ulmschneider, M. B.; Nola, A. D.; Amadei, A.; Smith, J. C. *Proc. Natl. Acad. Sci. U.S.A.* **2007**, *104*, 15230–15235.
- (6) Berendsen, H. J. C.; Postma, J. P. M.; van Gunsteren, W. F.; Hermans, J. Interaction models for water in relation to protein hydration. In *Intermolecular Forces*; Pullman, B., Ed.; D. Reidel Publishing Company: Dordrecht, The Netherlands, 1981; pp 331–342.
- (7) Jorgensen, W. L.; Chandrasekhar, J.; Madura, J. D.; Impey, R. W.; Klein, M. L. *J. Chem. Phys.* **1983**, *79*, 926–935.
- (8) Still, W. C.; Tempczyk, A.; Hawley, R. C.; Hendrickson, T. *J. Am. Chem. Soc.* **1990**, *112*, 6127–6129.
- (9) Hassan, S. A.; Mehler, E. L. *Proteins: Struct., Funct., Genet.* **2002**, *47*, 45–61.
- (10) Holst, M.; Baker, N.; Wang, F. *J. Comput. Chem.* **2000**, *21*, 1343–1352.
- (11) Nandi, N.; Bagchi, B. *J. Phys. Chem. B* **1997**, *101*, 10954–10961.
- (12) Rhee, Y. M.; Pande, V. S. *J. Phys. Chem. B* **2008**, *112*, 6221–6227.
- (13) Lu, B.; Cheng, X.; Huang, J.; McCammon, J. A. *Comput. Phys. Commun.* **2010**, *181*, 1150–1160.
- (14) Baker, N. A. *Curr. Opin. Struct. Biol.* **2005**, *15*, 137–143.
- (15) Hassan, S. A.; Mehler, E. L.; Zhang, D.; Weinstein, H. *Proteins: Struct., Funct., Genet.* **2003**, *51*, 109–125.
- (16) Schaefer, M.; Karplus, M. *J. Phys. Chem.* **1996**, *100*, 1578–1599.
- (17) Qiu, D.; Shenkin, P. S.; Hollinger, F. P.; Still, W. C. *J. Phys. Chem.* **1997**, *101*, 3005–3014.
- (18) Brooks, B. R.; Brucoleri, R. E.; Olafson, B. D.; States, D. J.; Swaminathan, S.; Karplus, M. *J. Comput. Chem.* **1983**, *4*, 187–217.
- (19) Brooks, B. R.; et al. *J. Comput. Chem.* **2009**, *30*, 1545–1614.
- (20) Hess, B.; Kutzner, C.; van der Spoel, D.; Lindahl, E. *J. Chem. Theory Comput.* **2008**, *4*, 435–447.
- (21) Larsson, P.; Lindahl, E. *J. Comput. Chem.* **2010**, *31*, 2593–2600.
- (22) Pearlman, D. A.; Case, D. A.; Caldwell, J. W.; Ross, W. S.; Cheatham, T. E.; DeBolt, S.; Ferguson, D.; Seibel, G.; Kollman, P. *Comput. Phys. Commun.* **1995**, *91*, 1–41.
- (23) Phillips, J. C.; Braun, R.; Wang, W.; Gumbart, J.; Tajkhorshid, E.; Villa, E.; Chipot, C.; Skeel, R. D.; Kale, L.; Schulten, K. *J. Comput. Chem.* **2005**, *26*, 1781–1802.
- (24) Tanner, D. E.; Ma, W.; Chen, Z.; Schulten, K. *Biophys. J.* **2011**, *100*, 2548–2556.
- (25) Feig, M.; Onufriev, A.; Lee, M. S.; Im, W.; Case, D. A.; Brooks, C. L. *J. Comput. Chem.* **2004**, *25*, 265–284.
- (26) Lee, M. S.; Salsbury, F. R.; Brooks, C. L. *J. Chem. Phys.* **2002**, *116*, 10606–10614.
- (27) Im, W.; Lee, M. S.; Brooks, C. L. *J. Comput. Chem.* **2003**, *24*, 1691–1702.
- (28) Shivakumar, D.; Deng, Y.; Roux, B. *J. Chem. Theory Comput.* **2009**, *5*, 919–930.
- (29) Grant, B. J.; Gorfie, A. A.; McCammon, J. A. *Theor. Chim. Acta* **2010**, *20*, 142–147.
- (30) Schulz, R.; Lindner, B.; Petridis, L.; Smith, J. C. *J. Chem. Theory Comput.* **2009**, *5*, 2798–2808.
- (31) Bondi, A. *J. Phys. Chem.* **1964**, *68*, 441–451.
- (32) Srinivasan, J.; Trevathan, M. W.; Beroza, P.; Case, D. A. *Theor. Chim. Acta* **1999**, *101*, 426–434.
- (33) Onufriev, L.; Case, D. A.; Bashford, D. *J. Comput. Chem.* **2002**, *23*, 1297–1304.
- (34) Feig, M.; Brooks, C. L. *Curr. Opin. Struct. Biol.* **2004**, *14*, 217–224.
- (35) Hawkins, G. D.; Cramer, C. J.; Truhlar, D. G. *J. Phys. Chem.* **1996**, *100*, 19824–19839.
- (36) Schaefer, M.; Froemmel, C. *J. Mol. Biol.* **1990**, *216*, 1045–1066.
- (37) Onufriev, A.; Bashford, D.; Case, D. A. *J. Phys. Chem.* **2000**, *104*, 3712–3720.
- (38) Onufriev, A.; Bashford, D.; Case, D. A. *Proteins: Struct., Funct., Bioinf.* **2004**, *55*, 383–394.
- (39) Kalé, L. V.; Krishnan, S. Charm++: Parallel Programming with Message-Driven Objects. In *Parallel Programming using C++*; Wilson, G. V., Lu, P., Eds.; MIT Press: Cambridge, MA, 1996; pp 175–213.
- (40) Kalé, L.; Skeel, R.; Bhandarkar, M.; Brunner, R.; Gursoy, A.; Krawetz, N.; Phillips, J.; Shinozaki, A.; Varadarajan, K.; Schulten, K. *J. Comput. Phys.* **1999**, *151*, 283–312.
- (41) Darden, T.; York, D.; Pedersen, L. *J. Chem. Phys.* **1993**, *98*, 10089–10092.
- (42) Cornish, P. V.; Ermolenko, D. N.; Noller, H. F.; Ha, T. *Mol. Cell* **2008**, *30*, 578–588.
- (43) Cornish, P.; Ermolenko, D. N.; Staple, D. W.; Hoang, L.; Hickerson, R. P.; Noller, H. F.; Ha, T. *Proc. Natl. Acad. Sci. U.S.A.* **2009**, *106*, 2571–2576.
- (44) Horan, L. H.; Noller, H. F. *Proc. Natl. Acad. Sci. U.S.A.* **2007**, *104*, 4881–4885.
- (45) Agirrezabala, X.; Lei, J.; Brunelle, J. L.; Ortiz-Meoz, R. F.; Green, R.; Frank, J. *Mol. Cell* **2008**, *32*, 190–197.
- (46) Trabuco, L. G.; Villa, E.; Mitra, K.; Frank, J.; Schulten, K. *Structure* **2008**, *16*, 673–683.
- (47) Wells, D. B.; Abramkina, V.; Aksimentiev, A. *J. Chem. Phys.* **2007**, *127*, 125101.
- (48) Villa, E.; Sengupta, J.; Trabuco, L. G.; LeBarron, J.; Baxter, W. T.; Shaikh, T. R.; Grassucci, R. A.; Nissen, P.; Ehrenberg, M.; Schulten, K.; Frank, J. *Proc. Natl. Acad. Sci. U.S.A.* **2009**, *106*, 1063–1068.
- (49) Seidelt, B.; Innis, C. A.; Wilson, D. N.; Gartmann, M.; Armache, J.-P.; Villa, E.; Trabuco, L. G.; Becker, T.; Mielke, T.; Schulten, K.; Steitz, T. A.; Beckmann, R. *Science* **2009**, *326*, 1412–1415.
- (50) Becker, T.; Bhushan, S.; Jarasch, A.; Armache, J.-P.; Funes, S.; Jossinet, F.; Gumbart, J.; Mielke, T.; Berninghausen, O.; Schulten, K.; Westhof, E.; Gilmore, R.; Mandon, E. C.; Beckmann, R. *Science* **2009**, *326*, 1369–1373.
- (51) Agirrezabala, X.; Scheiner, E.; Trabuco, L. G.; Lei, J.; Ortiz-Meoz, R. F.; Schulten, K.; Green, R.; Frank, J. *EMBO J.* **2011**, *30*, 1497–1507.
- (52) Frauenfeld, J.; Gumbart, J.; van der Sluis, E. O.; Funes, S.; Gartmann, M.; Beatrix, B.; Mielke, T.; Berninghausen, O.; Becker, T.; Schulten, K.; Beckmann, R. *Nat. Struct. Mol. Biol.* **2011**, *18*, 614–621.
- (53) Trabuco, L. G.; Schreiner, E.; Eargle, J.; Cornish, P.; Ha, T.; Luthey-Schulten, Z.; Schulten, K. *J. Mol. Biol.* **2010**, *402*, 741–760.
- (54) Berk, V.; Zhang, W.; Pai, R. D.; Cate, J. H. D. *Proc. Natl. Acad. Sci. U.S.A.* **2006**, *103*, 15830–15834.
- (55) Humphrey, W.; Dalke, A.; Schulten, K. *J. Mol. Graphics* **1996**, *14*, 33–38.
- (56) Cornell, W. D.; Cieplak, P.; Bayly, C. I.; Gould, I. R.; Merz, K. M., Jr.; Ferguson, D. M.; Spellmeyer, D. C.; Fox, T.; Caldwell, J. W.; Kollman, P. A. *J. Am. Chem. Soc.* **1995**, *117*, 5179–5197.
- (57) Hornak, V.; Abel, R.; Okur, A.; Strockbine, B.; Roitberg, A.; Simmerling, C. *Proteins* **2006**, *65*, 712–725.
- (58) Perez, A.; Marchan, I.; Svozil, D.; Sponer, J.; Cheatham, T. E.; Loughton, C. A.; Orozco, M. *Phys. J.* **2007**, *92*, 3817–3829.
- (59) Aduri, R.; Psciuk, B. T.; Saro, P.; Taniga, H.; Schlegel, H. B.; SantaLucia, J. *J. Chem. Theory Comput.* **2007**, *3*, 1464–1475.
- (60) Stella, L.; Melchionna, S. *J. Chem. Phys.* **1998**, *109*, 10115–10117.
- (61) Wower, I. K.; Wower, J.; Zimmermann, R. A. *J. Biol. Chem.* **1998**, *273*, 19847–19852.
- (62) Cukras, A. R.; Southworth, D. R.; Brunelle, J. L.; Culver, G. M.; Green, R. *Mol. Cell* **2003**, *12*, 321–328.
- (63) Diaconu, M.; Kothe, U.; Schlünzen, F.; Fischer, N.; Harms, J. M.; Tonevitsky, A. G.; Stark, H.; Rodnina, M. V.; Wahl, M. C. *Cell* **2005**, *121*, 991–1004.

# Linear-Scaling Time-Dependent Density Functional Theory Based on the Idea of “From Fragments to Molecule”

Fangqin Wu, Wenjian Liu,\* Yong Zhang, and Zhendong Li

Beijing National Laboratory for Molecular Sciences, Institute of Theoretical and Computational Chemistry, State Key Laboratory of Rare Earth Materials Chemistry and Applications, College of Chemistry and Molecular Engineering, and Center for Computational Science and Engineering, Peking University, Beijing 100871, People's Republic of China

**ABSTRACT:** To circumvent the cubic scaling and convergence difficulties encountered in the standard top-down localization of the global canonical molecular orbitals (CMOs), a bottom-up localization scheme is proposed based on the idea of “from fragments to molecule”. That is, the global localized MOs (LMOs), both occupied and unoccupied, are to be synthesized from the primitive fragment LMOs (pFLMOs) obtained from subsystem calculations. They are orthonormal but are still well localized on the parent fragments of the pFLMOs and can hence be termed as “fragment LMOs” (FLMOs). This has been achieved by making use of two important factors. Physically, it is the transferability of the locality of the fragments that serves as the basis. Mathematically, it is the special block-diagonalization of the Kohn–Sham matrix that allows retention of the locality: The occupied–occupied and virtual–virtual diagonal blocks are only minimally modified when the occupied–virtual off-diagonal blocks are annihilated. Such a bottom-up localization scheme is applicable to systems composed of all kinds of chemical bonds. It is then shown that, by a simple prescreening of the particle-hole pairs, the FLMO-based time-dependent density functional theory (TDDFT) can achieve linear scaling with respect to the system size, with a very small prefactor. As a proof of principle, representative model systems are taken as examples to demonstrate the accuracy and efficiency of the algorithms. As both the orbital picture and integral number of electrons are retained, the FLMO-TDDFT offers a clear characterization of the nature of the excited states in line with chemical/physical intuition.

## 1. INTRODUCTION

Time-dependent density functional theory (TDDFT)<sup>1,2</sup> has in the past two decades evolved into a powerful tool for investigating electronic excitations of small- to medium-sized systems. Yet, the formal cubic scaling precludes its applicability to large systems such as luminescent materials and biological molecules. For such real-life systems, some linear scaling TDDFT ought to be developed. The first attempt in this direction was made by Chen and co-workers.<sup>3,4</sup> They worked with the linearized time-dependent Kohn–Sham (KS) equation or equivalently the equation of motion for the one-electron reduced density matrix and made full use of the locality of the density matrix in an orthogonal atomic orbital (OAO) representation. The formalism was recently extended by Yang et al.<sup>5</sup> to nonorthogonal localized molecular orbitals, in terms of which both the metric and the molecular orbital coefficients become very sparse. Such time-domain approaches have the ability to capture all of the bright states in one run. However, this should be viewed as a disadvantage rather than an advantage, simply because the high-lying excited states are not really meaningful due to the approximate nature of the exchange–correlation potential and the incompleteness in the basis set. Therefore, the inclusion of such states merely increases the computational expenses. Moreover, time-domain simulations usually have to adopt very small time steps to ensure the accuracy and long simulation times to resolve energetically adjacent states. These would result in a large prefactor and hence a late crossover point, after which the linear scaling calculation becomes more efficient than the conventional one. As such, time-domain TDDFT often

has to be combined with approximate Hamiltonians for the time evolution of large systems.

Apart from the above linear scaling TDDFT in the time domain, efficient implementations of TDDFT in the frequency domain have also been available in the literature. As a straightforward generalization of the ground state fragment molecular orbital (FMO) method, the FMO-TDDFT due to Chiba and co-workers<sup>6</sup> decomposes the excitation energy into a sum of monomer excitation energy and many-body increments. Despite its high efficiency, the FMO-TDDFT has a severe limitation in that it relies on the local nature of the excitations. That is, the principle monomer should be large enough so as to fully accommodate the target excitations. Otherwise, the truncation of the many-body expansion at low order would not work.<sup>7</sup> The “density-fragment interaction” (DFI)-based TDDFT proposed by Fujimoto and Yang<sup>8</sup> employs instead the orbitals from a self-consistent treatment of the electron density interactions between fragments of fixed numbers of electrons. The main drawback of the DFI-TDDFT lies in the DFI treatment of the ground state, where interfragment exchange–correlation (XC) interactions and charge-transfer effects cannot be accounted for. As a result, the DFI-TDDFT energies are rather sensitive to the size of fragments. The AO-TDDFT by van Gisbergen et al.<sup>9</sup> and Coriani et al.<sup>10</sup> works directly with atomic orbitals (AO) and can achieve linear scaling by means of prescreening techniques and sparse-matrix algebra. At variance with both the FMO- and DFI-TDDFT,

Received: April 2, 2011

Published: September 14, 2011

the AO-TDDFT amounts to a uniform treatment of all kinds of excitations. However, it will become inefficient when the basis set consisting of a number of diffuse functions for the cutoff of the AO pairs is then ineffective.

The crucial issue is then whether it is possible to design a linear scaling TDDFT that is well balanced between accuracy and efficiency for all types of excitations of large systems composed of any kind of chemical bonds and meanwhile allows for easy interpretations of the excitations. To address this issue, we realize that there exist two types of localities, in energy and in space. In the CMO representation of TDDFT, the KS orbital energy difference is the leading term of the excitation energy and is usually dominant over the coupling term. As such, a low-lying TDDFT excited state often involves only a few particle–hole (p–h) pairs. This feature stems directly from the intrinsic nature of the KS orbitals. That is, both the occupied and virtual KS orbitals stem from the same local potential of  $N - 1$  electrons such that the latter are closely related to the excited states of an  $N$ -electron system. For comparison, one may recall that the Hartree–Fock (HF) virtual orbitals arise from the nonlocal potential of  $N$  electrons and are hence more related to the electron attachments of an  $(N + 1)$ -electron system rather than the excited states of an  $N$ -electron system. As a result, low-lying TDHF excited states are typically heavy mixtures of many p–h pairs. However, this particular feature of TDDFT cannot be employed *a priori* because the coupling term is a dense matrix: The CMOs are usually delocalized throughout the whole space such that neither the Coulomb nor the XC integrals in the coupling term can be truncated. The KS CMO can hence be characterized as “local in energy but delocalized in space”. On the contrary, the AOs are local in space but “delocalized in energy”, manifested by strong couplings among themselves. What are in between are the localized molecular orbitals (LMOs), which may have good localities both in energy and in space. Moreover, it is the p–h pairs that serve as the basis in TDDFT. Therefore, one should try to explore as much as possible the locality of the p–h pairs. Conceptually, the LMO-based TDDFT can be viewed as an intermediate between the CMO- and AO-based TDDFT.

An immediate question is then how to generate the desired LMO. As is well-known, the traditional top-down schemes<sup>11–13</sup> for generating the LMO from unitary transformations of the global CMO scale at least cubically with respect to the size of systems. Even worse is that the localization, especially that of the virtual CMO, may fail miserably for large systems. To avoid this problem, we propose here the idea of “from fragments to molecule”, i.e., a bottom-up approach for constructing the LMO. More specifically, the primitive fragment LMOs (pFLMOs) are first obtained from subsystem calculations via the standard schemes.<sup>11–13</sup> They are then taken as the basis for the global self-consistent field (SCF) calculations. To maintain the locality of the MO, only the off-diagonal blocks of the KS matrix between the occupied and virtual MOs are to be annihilated, whereas the diagonal blocks are to be modified as little as possible. It will be shown that the so-obtained MOs are only minor modifications of the pFLMOs. That is, they are still well localized on their parent fragments and can hence be termed as “fragment LMOs” (FLMOs). It is then shown that, by a simple prescreening of the FLMO p–h pairs, the FLMO-based TDDFT can achieve linear scaling with respect to the system size, with a very small prefactor (and hence no crossover point). As both the orbital picture and integral number of electrons are retained, the nature of the excitations, whether local, delocalized, or charge transfer,

can readily be deduced from analysis of the excited state eigenvectors in terms of transitions between fragments.

The remainder of the paper is organized as follows. Section 2.1 is devoted to the algorithm for the linear scaling LMO-TDDFT, which is followed by the construction of the FLMO in section 2.2. The validity and efficiency of the FLMO-based TDDFT are examined in detail in section 3. The account ends with conclusions and perspectives in section 4.

## 2. THEORY

**2.1. LMO-Based Linear Scaling TDDFT.** For the time being, we only consider closed-shell systems and further neglect spin–orbit couplings. The frequency-domain TDDFT then amounts to solving the following generalized eigenvalue equation:<sup>2</sup>

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B} & \mathbf{A} \end{bmatrix} \begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \end{bmatrix} = \omega \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & -\mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \end{bmatrix} \quad (1)$$

where the matrix elements of  $\mathbf{A}$  and  $\mathbf{B}$  are defined as

$$A_{ai,bj} = \Delta_{ai,bj} + K_{ai,bj} \quad (2)$$

$$B_{ai,bj} = K_{ai,jb} \quad (3)$$

$$\Delta_{ai,bj} = \delta_{ij}f_{ab} - \delta_{ab}f_{ji} \quad (4)$$

$$K_{pq,rs} = (pq|sr) + (pq|f_{xc}|sr) \quad (5)$$

Here and henceforth, the indices  $\{ij,k,l,\dots\}$ ,  $\{a,b,c,d,\dots\}$ , and  $\{p,q,r,s,\dots\}$  designate occupied, virtual, and unspecified (real-valued) orbitals, respectively. The Mulliken notation has been adopted for the two-electron integrals. Under the adiabatic approximation (ALDA), the exchange–correlation (XC) kernel  $f_{xc}$  reads

$$f_{xc}^{\sigma\sigma'}(\mathbf{r},\mathbf{r}') = \frac{\delta^2 E_{xc}}{\delta\rho_{\sigma}(\mathbf{r})\delta\rho_{\sigma'}(\mathbf{r}')} \delta(\mathbf{r}-\mathbf{r}') \quad (6)$$

which implies that the left-hand side of eq 1 is frequency independent. Therefore, eq 1 can in principle be solved by a single diagonalization. The cost would be of  $O(N_{\text{ph}}^3)$ , with  $N_{\text{ph}}$  being the product of the numbers of the occupied ( $N_o$ ) and virtual ( $N_v$ ) orbitals. In other words, such a one-step diagonalization would scale roughly as  $O(N^6)$ , with  $N$  characterizing the size of the system. In practice, only a few low-lying excited states are of interest, which can be obtained more efficiently by solving eq 1 in an iterative manner. Here, we adopt the modified Davidson iterative scheme<sup>14</sup> by rewriting eq 1 as

$$(\mathbf{A} - \mathbf{B})(\mathbf{A} + \mathbf{B})\mathbf{Z} = \omega^2\mathbf{Z} \quad (7)$$

with

$$\mathbf{Z} = \mathbf{X} + \mathbf{Y} \quad (8)$$

Specifically, for a given trial vector  $\mathbf{b}$ , the matrix-vector products  $(\mathbf{A} + \mathbf{B})\mathbf{b}$  and  $(\mathbf{A} - \mathbf{B})\mathbf{b}$  need to be formed. While the computational cost for the product  $\Delta\mathbf{b}$  is negligible, the contraction between the coupling matrix  $\mathbf{K}$  and the trial vector  $\mathbf{b}$  is very expensive. In the BDF package,<sup>15–18</sup> the contraction  $\mathbf{Kb}$  involves three steps in each iteration. The induced density  $\rho_{\text{ind}}(\mathbf{r})$

$$\rho_{\text{ind}}(\mathbf{r}) = \sum_{ai} \Omega_{ai}(\mathbf{r}) b_{ai}, \quad \Omega_{ai}(\mathbf{r}) = \phi_a(\mathbf{r}) \phi_i(\mathbf{r}) \quad (9)$$



is first evaluated and tabulated on the grids. The induced Coulomb potential  $V_{\text{ind}}(\mathbf{r})$

$$V_{\text{ind}}(\mathbf{r}) = \int \frac{\rho_{\text{ind}}(\mathbf{r}')}{|\mathbf{r} - \mathbf{r}'|} d\mathbf{r}' \quad (10)$$

is then evaluated by using the multipolar expansion technique.<sup>15</sup> The contraction  $\mathbf{Kb}$  is finally formed as

$$[\mathbf{Kb}]_{ai} = \int \Omega_{ai}(\mathbf{r}) [V_{\text{ind}}(\mathbf{r}) + f_{\text{xc}}^{\text{ALDA}}(\mathbf{r}) \rho_{\text{ind}}(\mathbf{r})] d\mathbf{r} \quad (11)$$

The computational costs of  $\rho_{\text{ind}}(\mathbf{r})$  and  $\mathbf{Kb}$  both scale formally as  $O(N_{\text{ph}}N_{\text{g}})$ , with  $N_{\text{g}}$  being the number of grid points. As  $N_{\text{g}}$  is proportional to the system size, both steps are  $O(N^3)$ . The cost for constructing the Coulomb potential  $V_{\text{ind}}(\mathbf{r})$  is  $O(N^2)$  with a small prefactor. In sum, without any truncation of the p–h space and the grid points, the scaling for the matrix-vector product  $\mathbf{Kb}$  is  $O(N^3)$  for each trial vector. The scaling can only be reduced if the locality of the p–h pairs is fully taken into account.

It is clear that it is the orbital overlap that dictates the significance of a p–h pair. A natural measure of the spatial overlap between a given occupied orbital  $\phi_i$  and a virtual orbital  $\phi_a$  is the inner product  $O_{ai}$  of their moduli:

$$O_{ai} = \int |\Omega_{ai}(\mathbf{r})| d\mathbf{r} \leq \sqrt{\int |\phi_a|^2 d\mathbf{r} \int |\phi_i|^2 d\mathbf{r}} = 1 \quad (12)$$

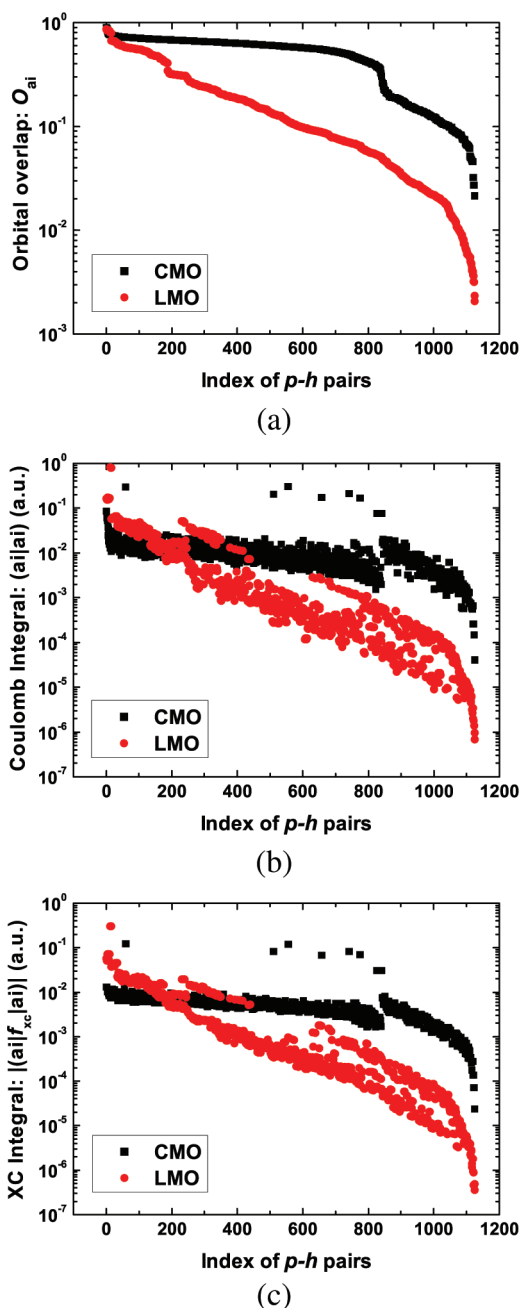
When the orbitals are orthonormalized, the quantity  $O_{ai}$  lies in the interval  $[0,1]$  due to the Schwarz inequality, as shown above. The products  $O_{ai}O_{bj}$  are closely correlated with the magnitudes of the Coulomb and XC integrals in the coupling matrix  $\mathbf{K}$ . To see this, the  $O_{ai}$  values for the p–h pairs in the CMO and LMO representations are compared for *trans*-1,3-butadiene ( $\text{C}_4\text{H}_6$ ) in Figure 1, alongside the absolute values of the diagonal Coulomb and XC integrals. A Slater-type double- $\zeta$  polarized basis set (DZP) was used for each atom, and the LMOs were obtained via the Boys localization<sup>11</sup> of the CMO calculated with the LDA. The p–h pairs are sorted in descendent order of  $O_{ai}$ . It is clearly seen that all of the quantities are significantly smaller in the LMO representation than in the CMO representation. In particular, both the Coulomb and XC integrals over the LMO decay quickly as  $O_{ai}$  decreases. As expected, they are roughly proportional to the square of  $O_{ai}$ . Taking a threshold of  $10^{-4}$ , about one-third of the diagonal Coulomb and XC integrals can be screened out in the LMO representation even for this small system, whereas almost all of the integrals are above this threshold in the CMO representation.

To facilitate the prescreening, the grid points are classified into several batches, each of which consists of a fixed number (e.g., 128) of grid points. The parameter  $\eta_{ai}^B$

$$\eta_{ai}^B = -\log O_{ai}^B, \quad \eta_{ai}^B \in [0, +\infty) \quad (13)$$

$$O_{ai}^B = \sum_{p \in B} w(\mathbf{r}_p) |\Omega_{ai}(\mathbf{r}_p)|, \quad O_{ai}^B \in [0, 1] \quad (14)$$

is then calculated for every p–h pair. Here,  $w(\mathbf{r}_p)$  are the weights of the grid points. The larger the distance between the p–h pair  $ai$  and the grid batch  $B$  or the distance between the orbitals  $\phi_a$  and  $\phi_i$ , the smaller the  $O_{ai}^B$  and, hence, the larger the  $\eta_{ai}^B$ . Therefore,  $\eta_{ai}^B$  is an effective measure of the significance of the pair  $ai$  on batch  $B$ . An appropriate threshold  $\eta$  can be introduced to screen out all of the pairs with  $\eta_{ai}^B > \eta$ . The larger the threshold



**Figure 1.** Comparison between the p–h pairs in the CMO and LMO representations for  $\text{C}_4\text{H}_6$  calculated with DZP basis set and LDA. (a) Orbital overlap  $O_{ai}$ . (b) Diagonal elements of the Coulomb kernel. (c) Diagonal elements of the ALDA kernel.

$\eta$  is, the more p–h pairs are retained. In the extreme case of  $\eta = +\infty$ , all of the p–h pairs are retained, implying no truncation to the coupling matrix  $\mathbf{K}$ . On the other hand, all of the p–h pairs are to be discarded with  $\eta = 0$ , leading to zero  $\mathbf{K}$ , i.e., the independent particle approximation (IPA) of TDDFT. Therefore, the balance between accuracy and efficiency can be monitored by the single threshold  $\eta$ . As the system size increases, the number of significant p–h pairs with  $\eta_{ai}^B < \eta$  will become constant for a given batch of grid points in the LMO representation. This is the key for achieving linear scaling in forming the contraction  $\mathbf{Kb}$ .

However, there is a price to pay for the LMO representation in solving eq 7 with the Davidson iterative diagonalization. Unlike the CMO representation, where the leading term  $\Delta^{\text{CMO}}$  (cf. eq 4) of  $\mathbf{A}$  is diagonal such that the preconditioner  $(\Delta^{\text{CMO}} - \omega_0 \mathbf{I})^{-1}$  for updating the trial vectors can trivially be evaluated, the  $\Delta^{\text{LMO}}$  matrix in the LMO representation is not diagonal, and the evaluation of  $(\Delta^{\text{LMO}} - \omega_0 \mathbf{I})^{-1}$  is a heavy task (for detailed discussions, see ref 10). As we have tested, taking only the diagonal elements of  $\Delta^{\text{LMO}}$  in the preconditioning may factually decelerate the convergence, particularly when  $\Delta^{\text{LMO}}$  is not diagonally dominant. To accelerate the convergence, Miura and Aoki<sup>19</sup> proposed in their LMO-based TDHF to use the “localized CMO” obtained by canonically orthogonalizing the LMO in each preset region. The transformed  $\Delta^{\text{LMO}}$  matrix is then block (regional)-diagonal. Further combined with a projector guess, the convergence is significantly improved compared with the original  $\Delta^{\text{LMO}}$ . However, it is still not as good as that in the CMO representation.

As a matter of fact, one can combine the good of the CMO (diagonal  $\Delta^{\text{CMO}}$ ) and the LMO (sparse  $\mathbf{K}^{\text{LMO}}$ ) representations by introducing a unitary transformation  $\mathbf{U}_{\text{VO}}$  between the CMO and LMO p–h bases:

$$\begin{aligned} \mathbf{U}_{\text{VO}} &= \mathbf{U}_{\text{vv}} \otimes \mathbf{U}_{\text{oo}}, \\ \mathbf{F}_{\text{oo}}^{\text{LMO}} \mathbf{U}_{\text{oo}} &= \mathbf{U}_{\text{oo}} \mathbf{F}_{\text{oo}}^{\text{CMO}}, \\ \mathbf{F}_{\text{vv}}^{\text{LMO}} \mathbf{U}_{\text{vv}} &= \mathbf{U}_{\text{vv}} \mathbf{F}_{\text{vv}}^{\text{CMO}} \end{aligned} \quad (15)$$

where  $\mathbf{U}_{\text{oo}}$  ( $\mathbf{U}_{\text{vv}}$ ) is the unitary transformation between the occupied (virtual) CMO and LMO (see Appendix ). More specifically, the matrix-vector product is to be carried out as

$$\mathbf{K}^{\text{CMO}} \mathbf{b}^{\text{CMO}} = \mathbf{U}_{\text{VO}}^{\dagger} \mathbf{K}^{\text{LMO}} \mathbf{b}^{\text{LMO}} \quad (16)$$

with the trial vector  $\mathbf{b}^{\text{LMO}}$  transformed as

$$\begin{aligned} b_{ai}^{\text{LMO}} &= [\mathbf{U}_{\text{VO}} \mathbf{b}^{\text{CMO}}]_{ai} = \sum_{bj} [\mathbf{U}_{\text{VO}}]_{ai,bj} b_{bj}^{\text{CMO}} \\ &= \sum_{bj} [\mathbf{U}_{\text{vv}}]_{ab} [\mathbf{U}_{\text{oo}}]_{ij} b_{bj}^{\text{CMO}} = [\mathbf{U}_{\text{vv}} \mathbf{b}_{\text{vo}}^{\text{CMO}} \mathbf{U}_{\text{oo}}^{\dagger}]_{ai} \end{aligned} \quad (17)$$

Note that in the above last equality the vector  $\mathbf{b}^{\text{CMO}}$  in the p–h space has been recast into an  $N_v \times N_o$  matrix in the orbital basis. The same procedure can also be applied to the right-hand side of eq 16. Both transformations scale computationally as  $O(N^3)$  but with a very small prefactor. For comparison, the AO-TDDFT<sup>9</sup> can be implemented in the same way, just with the following CMO to AO transformation:

$$b_{\mu\nu}^{\text{AO}} = [\mathbf{C}_v \mathbf{b}_{\text{vo}}^{\text{CMO}} \mathbf{C}_o^{\dagger}]_{\mu\nu} \quad (18)$$

where  $\mathbf{C}_o$  ( $\mathbf{C}_v$ ) is the coefficient matrix of the occupied (virtual) CMO. Yet, there exists a significant difference between the AO-TDDFT and LMO-TDDFT: The full dimension of  $[\mathbf{K}^{\text{AO}} \mathbf{b}^{\text{AO}}]$  in the AO-TDDFT is  $M^2$ , with  $M$  being the number of basis functions. This is to be compared with the dimension,  $N_{\text{ph}}$ , of  $[\mathbf{K}^{\text{LMO}} \mathbf{b}^{\text{LMO}}]$  in the LMO-TDDFT. The ratio  $N_{\text{ph}}/M^2$  is less than  $N_o/M$  and becomes smaller for a larger basis. In particular, when diffuse functions are employed, the cutoff of the AO pairs becomes ineffective whereas the LMO p–h pairs can still be significantly cut off because the locality of the occupied MO is not much affected by diffuse functions.

In sum, linear scaling TDDFT calculations can be achieved by utilizing the locality of the p–h basis. The salient feature of the above algorithm is evaluating all of the quantities in the LMO

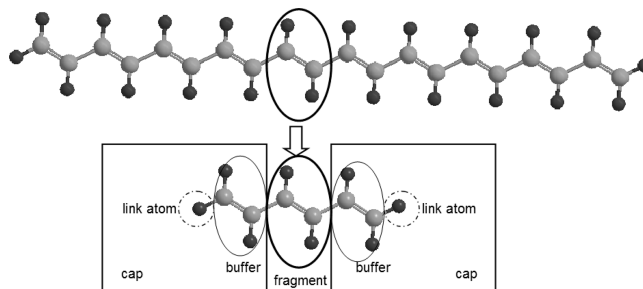


Figure 2. Fragmentation of  $\text{C}_{20}\text{H}_{22}$  ( $\text{C}_{2h}$ ).

representation but solving eq 7 in the CMO representation such that both the diagonality of  $\Delta^{\text{CMO}}$  and the sparsity of  $\mathbf{K}^{\text{LMO}}$  can be employed. As an additional benefit, full molecular symmetry of arbitrary order<sup>20</sup> can be incorporated by symmetrizing the CMO, which is necessary for the proper assignment of the calculated excited states.

**2.2. Efficient Construction of LMO.** The above algorithm for TDDFT can be combined with any kind of LMO as long as they can be generated efficiently. To circumvent the difficulties encountered in the top-down localization of the global CMO, we propose here a bottom-up localization scheme based on the idea “from fragments to molecule”. That is, the global orthonormal LMOs are to be synthesized from the pFLMOs obtained from subsystem calculations.

Like most fragment-based approaches (for a recent review, see ref 21), the whole molecule is first divided into several fragments (denoted as  $I, J, K$ , etc.) based on chemical intuition. Each fragment is then capped to form a subsystem (denoted as  $\bar{I}, \bar{J}, \bar{K}$ , etc.). To closely mimic the local chemical environment, the caps are just parts of the whole system directly bonded to the fragment, which are further saturated by link atoms. Taking polyacetylene  $\text{C}_{20}\text{H}_{22}$  as an illustration (Figure 2), the whole molecule is divided into 10 fragments, each of which is composed of two doubly bonded carbon atoms. Two carbon atoms on each side of every fragment are taken as the buffer. Every subsystem is finally formed by saturating the dangling bonds with hydrogen atoms.

With this fragmentation, the basis functions  $V_{\text{subsystem}}^{\bar{I}}$  for a subsystem  $\bar{I}$  contain three parts, viz.,

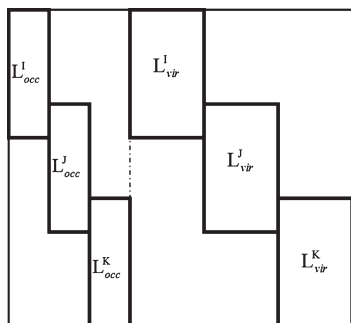
$$\begin{aligned} V_{\text{subsystem}}^{\bar{I}} &= V_{\text{fragment}}^I \oplus V_{\text{cap}}^I \\ &= V_{\text{fragment}}^I \oplus V_{\text{buffer}}^I \oplus V_{\text{link-atoms}}^I \end{aligned} \quad (19)$$

where  $V_{\text{fragment}}^I$ ,  $V_{\text{buffer}}^I$ , and  $V_{\text{link-atoms}}^I$  denote the respective basis functions for the fragment  $I$ , buffer, and link atoms. Conventional KS calculations are carried out for each subsystem in parallel

$$\mathbf{F}^{\bar{I}} \mathbf{C}^{\bar{I}} = \mathbf{S}^{\bar{I}} \mathbf{C}^{\bar{I}} \mathbf{E}^{\bar{I}} \quad (20)$$

where  $\mathbf{F}^{\bar{I}}$ ,  $\mathbf{C}^{\bar{I}}$ , and  $\mathbf{S}^{\bar{I}}$  are the KS, CMO coefficient, and overlap matrices of subsystem  $\bar{I}$ . The standard localization procedures<sup>11–13</sup> can be employed to localize the occupied and virtual CMO separately. As the size of the subsystems is very small, the computation cost of the SCF and localization steps is negligible. The so-obtained LMO coefficient matrix for subsystem  $\bar{I}$  is to be denoted as  $\mathbf{L}^{\bar{I}}$ . To identify the location of the LMO, a Löwdin population analysis is carried out for each LMO:

$$n_p^I = \sum_{\mu \in I} (\mathbf{L}_p^{\bar{I}} \mathbf{S}^{1/2})_{\mu} (\mathbf{S}^{1/2} \mathbf{L}_p^{\bar{I}})_{\mu} \quad (21)$$



**Figure 3.** Schematic illustration of  $\mathbf{L}^{\text{pFLMO}}$  in the AO basis.

where the summation is confined to the AO of fragment  $I$ . The LMOs with  $n_p^I \geq \theta^I$  are assigned to fragment  $I$ , while the remainder are to be discarded. For a fragment  $J$  directly connected with  $I$ , the selection threshold will be  $n_p^J \geq \theta^J = 1 - \theta^I$ . In practice, a value of 0.6 can be chosen for  $\theta^I$ . In this way, the assignment of the same LMO, especially the one of  $n_p^I = n_p^J = 0.5$ , to two fragments can be avoided. Since the link atoms do not belong to the system, their basis functions  $V_{\text{link-atoms}}^I$  have to be projected out from the fragment-centered LMO. Note in passing that this projection does not affect discernibly the norms of the LMO. The basis functions  $V_{\text{buffer}}^I$  of the buffer are instead retained, as they are necessary for describing the tails of the fragment LMO penetrating the system. To facilitate subsequent SCF calculations, the resultant LMOs are symmetrically orthonormalized to form the pFLMO. The coefficient matrix of the whole set of pFLMOs is in the AO representation of the following structure:

$$\mathbf{L}^{\text{pFLMO}} = \mathbf{L}^I \oplus \mathbf{L}^J \oplus \dots \oplus \mathbf{L}^K \quad (22)$$

which is further depicted in Figure 3. Here,  $\mathbf{L}^I$  collects the coefficients of the pFLMO of fragment  $I$ . The occupied and virtual parts of  $\mathbf{L}^{\text{pFLMO}}$  can be identified according to the occupations in the subsystem calculations.

If the caps are sufficiently large, further global SCF calculations may not be necessary. However, here, we decide to use the smallest caps possible. Therefore, the above pFLMO will be taken as the basis and the superposition of the fragment densities as the initial guess for further SCF calculations of the whole system. The KS matrix takes the following block structure:

$$\mathbf{F}^{\text{pFLMO}} = \begin{pmatrix} \mathbf{F}^{II} & \mathbf{F}^{IJ} & \dots & \mathbf{F}^{IK} \\ \mathbf{F}^{JI} & \mathbf{F}^{JJ} & \dots & \mathbf{F}^{JK} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{F}^{KI} & \mathbf{F}^{KJ} & \dots & \mathbf{F}^{KK} \end{pmatrix} \quad (23)$$

which is sparse since the elements  $F_{pq}^{\text{pFLMO}}$  are vanishingly small if the pFLMOs  $p$  and  $q$  are located on two distant fragments. As the pFLMOs represent well the chemical bondings of the whole molecule, only a few (macro) SCF iterations are needed to reach convergence. Because of this conquering, the dividing (fragmentation) is actually not very crucial, at variance with other schemes reported in the literature. The price to pay here is the cubic scaling in the diagonalization of the sparse KS matrix. However, it has a very small prefactor. Moreover, the energetic locality of the pFLMO can be employed to freeze the core-like ones and cut off the high-lying ones, both of which have nothing

to do with interfragment interactions. This can be facilitated by constructing fragment-centered CMOs. As the primary interest here is the efficiency of TDDFT rather than the SCF itself, such pruning of the pFLMO is not to be considered.

What is more crucial here is how to retain the locality of the MO during the SCF iterations or resume the locality in the very end of the conventional SCF calculation. In the former case, only the off-diagonal blocks of the KS matrix between the occupied and virtual MOs are to be annihilated in each SCF cycle. In principle, the Jacobi sweep of iterations can be applied here. However, the convergence is very slow (typically hundreds of iterations). To solve this issue, we propose a novel block-diagonalization approach in Appendix , where the decoupling condition can be solved either iteratively or noniteratively. For all of the cases encountered here, the convergence of the iterative block-diagonalization can be reached within just two to three (micro) cycles. In the latter case, the standard SCF procedure is invoked (i.e., a full-diagonalization of the KS matrix is carried out in each SCF iteration), and the block-diagonalization is done only when the convergence has been reached. At first glance, this is nothing but a top-down localization of the global CMO. However, this is possible only in the pFLMO basis. Therefore, such a one-step scheme is still within the spirit of “from fragments to molecule”. The MOs by the two types of calculations, i.e., one-step and multiple-step block-diagonalizations, are related by a unitary transformation and have been verified to have very much the same locality. Therefore, it is the latter one-step, noniterative block-diagonalization that is to be used here. It will be shown later on that the resultant MOs of the whole system are very close to the initial pFLMOs and can hence be dubbed as FLMOs. This feature stems from both the good transferability of the pFLMO and the particular block-diagonalization algorithm that leads by construction to minimal modifications of the diagonal blocks, the key for retaining the locality of both the occupied and virtual MOs. As such, the present approach provides a bottom-up construction of LMO, manifesting the idea of “from fragments to molecule”. It is an effective means for resolving the computational bottlenecks in both the SCF and localization procedures.

Apart from the computational savings in the TDDFT calculations, the nature of the excited states can readily be revealed by using the fact that every FLMO belongs to a specific parent fragment. That is, for a given excited state, the contribution of transitions from fragment  $I$  to  $J$  can be measured by the weight  $W_{IJ}$  defined as

$$W_{IJ} = \sum_{i \in I, a \in J} Z_{ai}^2 \quad (24)$$

with  $\mathbf{Z}$  being the normalized eigenvector of eq 7. The distribution of  $W_{IJ}$  versus the fragment indices or simply the distance  $R_{IJ}$  between the fragments can then clearly distinguish local, delocalized, and charge-transfer excitations.

### 3. RESULTS AND DISCUSSION

The above algorithms for the FLMO-based DFT and TDDFT have been implemented into the BDF package.<sup>15–18</sup> As a proof of concept, only the VWN5 form<sup>22</sup> of LDA is used for both the ground and excited state calculations. As far as computational efficiency is concerned, the use of other types of functionals merely changes the prefactor but not the scaling. A Slater-type DZP basis set and  $75 \times 302$  grid points are used for each atom. More specifically, each hydrogen atom has five functions (2s1p)

Table 1. Convergence of the AO- and FLMO-based KS-LDA/DZP SCF Calculations of C<sub>20</sub>H<sub>22</sub> and C<sub>60</sub>H<sub>62</sub><sup>a</sup>

	iteration	AO	FLMO				
			$\Delta E$ (1,1)	$\Delta E$ (1,2)	$\Delta E$ (1,3)	$\Delta E$ (5,1)	$\Delta E$ (5,2)
C <sub>20</sub> H <sub>22</sub>	0	-3.061596	0.242335	0.019939	0.004306	0.032149	0.002878
	1	0.937293	0.001311	0.000393	0.000091	0.000411	0.000122
	2	1.603761	0.001329	0.000089	0.000040	0.001890	0.000174
	3	0.006505	0.000502	0.000023	-0.000002	0.000049	0.000003
	4	0.009614	0.000028	0.000005	0.000000	-0.000007	0.000000
	5	0.000725	0.000008	0.000000		0.000004	
	6	0.000074	0.000000			0.000000	
	7	-0.000005					
	8	-0.000026					
	9	0.000003					
	10	-0.000001					
C <sub>60</sub> H <sub>62</sub>	0	-9.684877	0.858260	0.078572	0.019252	0.129093	0.011839
	1	2.625106	0.005090	0.001368	0.000521	0.001657	0.000674
	2	4.361038	0.000484	-0.000757	-0.000410	0.005922	-0.000090
	3	0.009323	0.001170	0.000174	-0.000006	-0.000086	-0.000013
	4	0.020549	0.000051	0.000028	0.000000	-0.000043	-0.000001
	5	0.000933	0.000003	-0.000008		0.000009	0.000000
	6	-0.000063	0.000000	0.000000		0.000000	
	7	-0.000038					
	8	-0.000033					
	9	-0.000001					
	10	-0.000001					
11	0.000000						

<sup>a</sup> Iteration 0 corresponds to the superposition of the fragment densities.  $\Delta E$  denotes the difference between the iterative and converged energies. The numbers in parentheses indicate the numbers of double bonds in each (fragment, cap).

while each non-hydrogen atom has 15 functions (4s2p1d). An ADZP basis set, obtained by augmenting the DZP set with 1s1p1d diffuse functions for carbon and 1s1p diffuse functions for hydrogen, will also be used for testing the basis set effects. In order to make a fair comparison between the LMO- and CMO-based TDDFT, no symmetry is employed. All of the calculations are carried out on an AMD Opteron Quad-core 8374 HE 2.2 GHz processor.

**3.1. Construction of FLMO.** The first point to be checked is the dependence of the pFLMO on the cap size. Taking linear polyacetylene C<sub>20</sub>H<sub>22</sub> as an example (cf. Figure 2), the cap on each side of a double-bond fragment may be composed of one, two, or three double bonds. The results are to be denoted as Cap-*n* (*n* = 1, 2, 3). After separate calculations of the subsystems, the Boys scheme is adopted to localize the CMO, which amounts to minimizing the orbital self-extension (OSE):

$$I = \sum_p O[\phi_p],$$

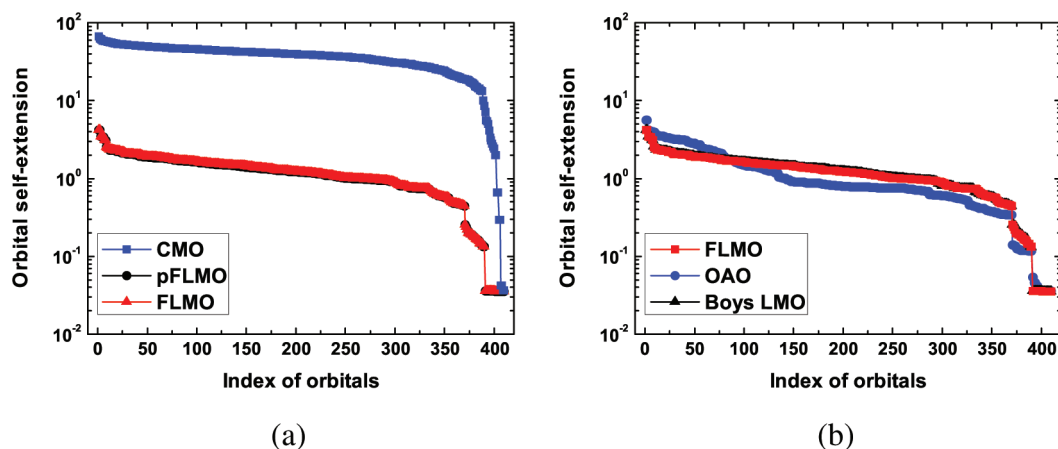
$$O[\phi_p] = \int \phi_p(\mathbf{r}) \phi_p(\mathbf{r}') r_{12}^2 \phi_p(\mathbf{r}) \phi_p(\mathbf{r}') \mathbf{dr} \mathbf{dr}'$$

$$= 2 \sum_{i=1,2,3} (\langle \phi_p | x_i^2 | \phi_p \rangle - \langle \phi_p | x_i | \phi_p \rangle^2) \quad (25)$$

After projecting out the basis functions of the link hydrogen atoms, the pFLMOs are set up via the Löwdin symmetric

orthogonalization. They are then taken as the basis and the superposition of the fragment densities,  $\sum_I \rho_I(\mathbf{r})$ , as the initial guess for further SCF calculations of the whole system. As for the AO-based SCF calculation, the superposition of the atomic densities is taken as the initial guess. The results are documented in Table 1. As expected, the pFLMOs provide a much better initial guess such that the number of SCF iterations is much reduced. It is noticeable that the error  $\Delta E_0$  in the energy of the zeroth iteration decreases as the cap size increases, indicating that the superimposed density  $\sum_I \rho_I(\mathbf{r})$  becomes increasingly close to the converged molecular density. It is also found that the error  $\Delta E_0$  remains a constant in the vicinity of the equilibrium, meaning that the energy estimated from  $\sum_I \rho_I(\mathbf{r})$  is accurate enough for geometry optimizations. The converged results with different caps agree with each other since the energetic locality of the pFLMO has not yet been employed for truncations. A second point to be checked here is the dependence of the pFLMO on the fragment size (denoted as as Frag-*n*, *n* = 1, 5). The results documented in Table 1 show that, for the same cap size, larger fragments also tend to decrease both the error of  $\Delta E_0$  and the number of iterations. All of these findings apply also to linear polyacetylene C<sub>60</sub>H<sub>62</sub>.

What is more interesting here is whether the global SCF calculation would spoil the locality of the pFLMO. To check this, the OSEs defined in eq 25 are compared in Figure 4 for each CMO, pFLMO, and FLMO of C<sub>20</sub>H<sub>22</sub>. The orbitals are sorted in descending order of the OSE. It is seen that most of the CMOs

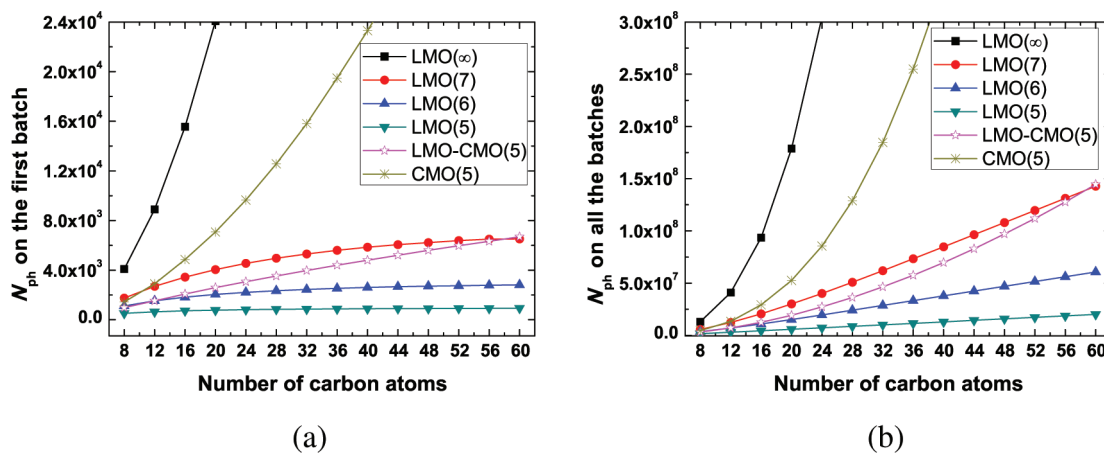


**Figure 4.** Comparison of the orbital self-extensions for different types of orbitals of  $C_{20}H_{22}$  at the LDA/DZP level. Each cap has one double bond. (a) CMO, pFLMO, and FLMO. (b) FLMO, OAO, and Boys LMO. The pFLMO/FLMO and FLMO/Boys LMO curves overlap each other in a and b, respectively.

**Table 2.** Wall Times (in seconds) of the AO- and FLMO-based SCF Calculations<sup>a</sup>

molecule	step	AO-SCF	FLMO-SCF				
			(1,1)	(1,2)	(1,3)	(5,1)	(5,2)
$C_{20}H_{22}$	subsystem		530[53]	1350[135]	2610[261]	620[310]	806[403]
	subsystem <sup>b</sup>		73	202	403	310	403
	SCF	565	415	381	347	415	347
	localization		0.1	0.1	0.1	0.1	0.1
	total	565	945	1731	2957	1035	1153
	total <sup>b</sup>		488	583	750	725	750
$C_{60}H_{62}$	subsystem		1590[53]	4050[135]	7830[261]	2142[357]	3360[560]
	subsystem <sup>b</sup>		73	202	403	403	716
	SCF	2784	2209	2209	1947	2209	2078
	localization		2	2	2	2	2
	total	2784	3799	6259	9777	4351	5438
	total <sup>b</sup>		2284	2413	2352	2614	2796

<sup>a</sup> Averaged wall times per subsystem are in brackets. The numbers of double bonds in each (fragment, cap) are in parentheses. The one-step, non-iterative block-diagonalization is adopted to construct the FLMO. For comparison, the wall times for the Boys localization scheme are 455 s for  $C_{20}H_{22}$  and 16 920 s for  $C_{60}H_{62}$ . <sup>b</sup> Only the time for the largest subsystem is counted, as the subsystems are calculated on parallel nodes.



**Figure 5.** Dependence of the number of p-h pairs of  $C_nH_{n+2}$  on the threshold  $\eta$ . (a) Number of p-h pairs on the first batch of grid points. (b) Number of p-h pairs on all of the batches.

are delocalized with large OSEs, whereas the locality of the pFLMOs and FLMOs is virtually identical. This reveals that the pFLMOs constructed from subsystem calculations describe the chemical bondings very well, such that only their tails are subject to minor changes when brought together to form the whole molecule. At this point, it is instructive to compare the present FLMO with the LMO from the Boys localization of the global CMO as well as the OAO. As judged from the OSE in Figure 4b, the locality of the FLMO and Boys LMO is practically the same, whereas the OAOs are marginally more local. A similar situation is found also in calculations with larger basis sets such as ADZP and TZ2P. As the LDA potential falls off too quickly, the

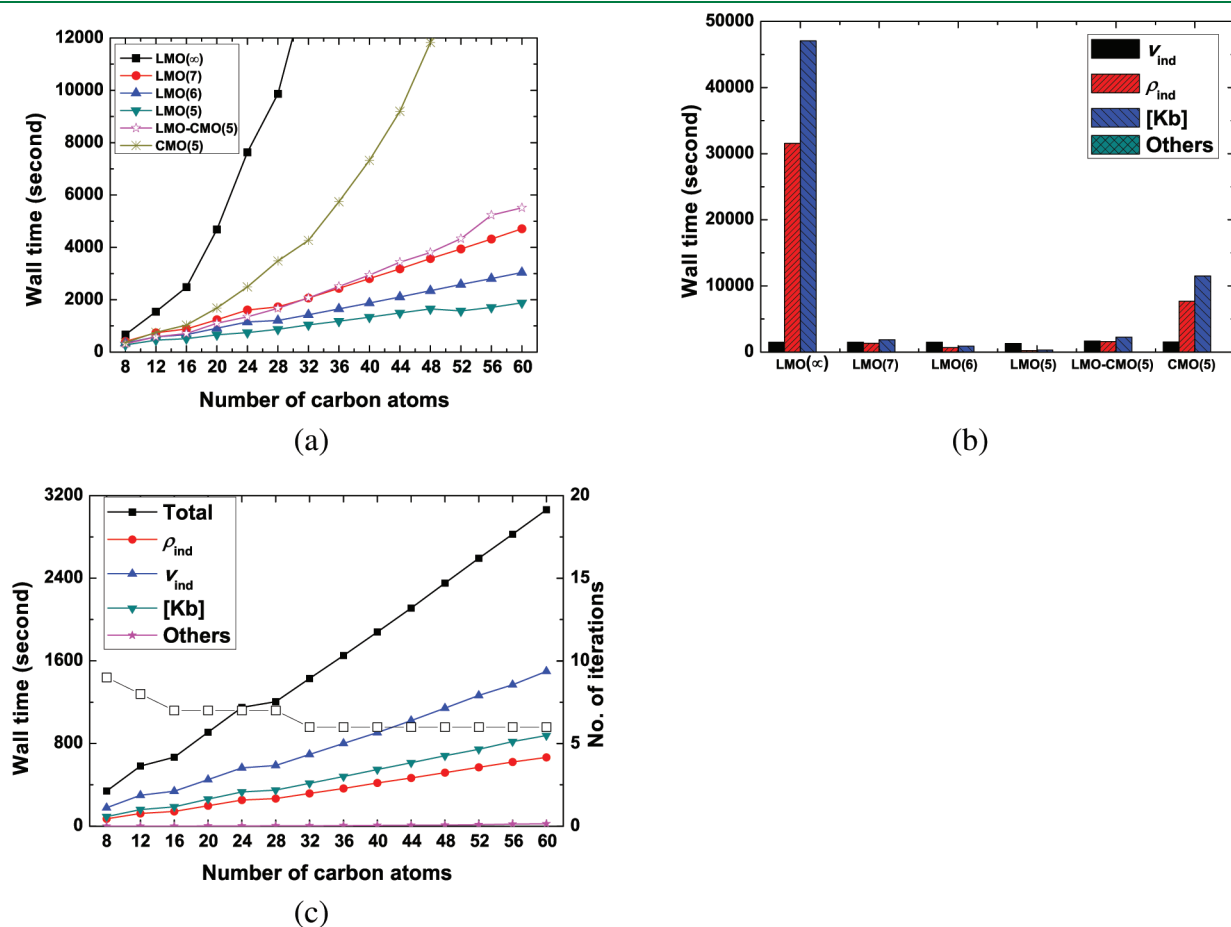
**Table 3. Polynomial Fittings of the  $N_{\text{ph}}-N_{\text{C}}$  Curves in Figure 5b and the  $T-N_{\text{C}}$  Curves (in seconds) in Figure 6a<sup>a</sup>**

scheme	$N_{\text{ph}}$		$T$	
	expression	$R^2$	expression	$R^2$
LMO( $\infty$ )	$2.9 \times 10^4 N_{\text{C}}^{2.9}$	0.9999	$3.6 N_{\text{C}}^{2.5}$	0.9955
LMO(7)	$3.0 \times 10^6 N_{\text{C}} - 2.0 \times 10^7$	0.9969	$82.8 N_{\text{C}} - 415.2$	0.9929
LMO(6)	$1.0 \times 10^6 N_{\text{C}} - 7.0 \times 10^7$	0.9993	$51.6 N_{\text{C}} - 134.5$	0.9942
LMO(5)	$3.6 \times 10^5 N_{\text{C}} - 1.0 \times 10^6$	0.9999	$30.9 N_{\text{C}} - 45.1$	0.9877
LMO-CMO(5)	$6.4 \times 10^4 N_{\text{C}}^{0.9}$	0.9995	$15.2 N_{\text{C}}^{1.5}$	0.9954
CMO(5)	$1.7 \times 10^3 N_{\text{C}}^{2.7}$	1.0000	$4.3 N_{\text{C}}^{2.1}$	0.9903

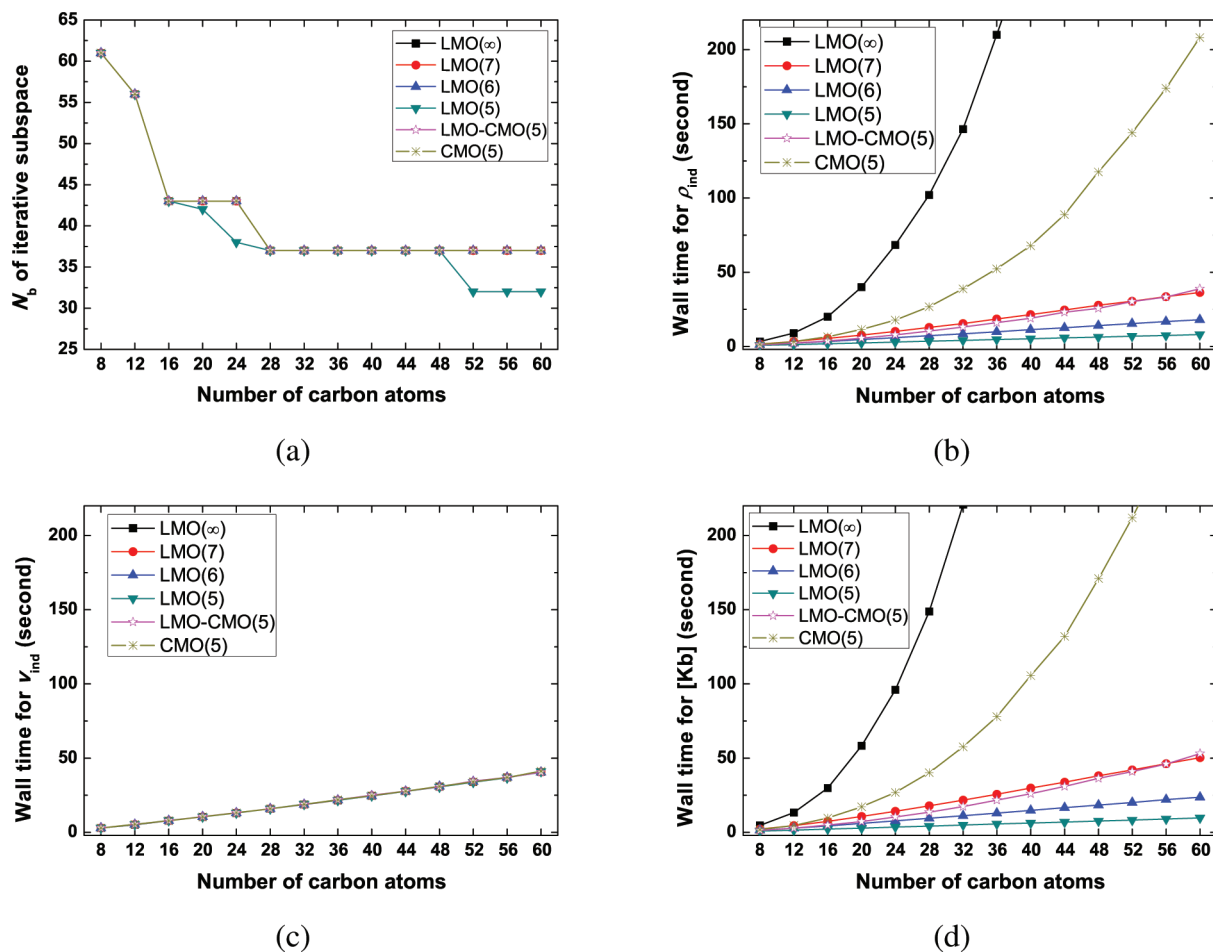
<sup>a</sup>  $R^2$ : Correlation coefficient.

corresponding virtual MOs are artificially too diffuse. Therefore, the present findings can safely be extended to other functionals of correct long-range behavior that tend to yield more compact MOs. The close proximity between the FLMOs, Boys LMOs, and OAOs indicates that they perform equally well in any linear scaling calculations. Yet, the FLMOs are clearly advantageous over the Boys LMOs in the construction and advantageous over the OAOs in the physical interpretation of the results.

To compare the relative efficiency of the FLMO- and AO-based SCF calculations, the wall times for the key steps are collected in Table 2. Noticeably, the subsystem calculations do cause a significant overhead (roughly a factor of 2) compared with the AO-SCF. However, this is only true if the subsystems, including the identical ones, are all calculated one after another on a single node. As a matter of fact, only the time for the largest subsystem is relevant here, since the subsystems have actually been calculated on parallel nodes. The overhead is then well compensated by the reduced number of SCF iterations. More importantly, compared with the very expensive Boys localization of the global CMO in the AO basis, the present one-step, noniterative block-diagonalization of the KS matrix in the pFLMO basis costs essentially nothing. Note in particular that the Boys localization fails even for  $\text{C}_{20}\text{H}_{22}$  when the ADZP basis set is used. Therefore, even in the sequential treatment of the subsystems, the overhead is not really an issue, as it will be overcompensated by the subsequent benefit from the LMO.



**Figure 6.** (a) Total wall times of different calculations of the five lowest excited states of  $\text{C}_n\text{H}_{n+2}$ . (b) Wall times of the key steps in calculations of  $\text{C}_{60}\text{H}_{62}$ . (c) Wall times of the key steps in scheme LMO(6). Empty squares indicate the number of iterations on the right vertical axis.



**Figure 7.** Comparison of different calculations of the five lowest excited states of  $C_nH_{n+2}$ . (a) Final dimension  $N_b$  of the iterative subspace in the Davidson diagonalization. (b) Averaged wall time per trial vector for  $\rho_{ind}$ . (c) Averaged wall time per trial vector for  $V_{ind}$ . (d) Averaged wall time per trial vector for [Kb].

For massive calculations of real-life systems, the coefficients of the pFLMO can be preconstructed and stored such that they can directly be used as building blocks. Making full use of the molecular symmetry can also reduce the costs of the subsystems as they usually have higher symmetry than the whole system. Overall, the FLMO-SCF provides an effective means for handling large and complex systems in the spirit of “from fragments to molecule”. On one hand, it can avoid the convergence difficulties encountered in the conventional SCF calculations. On the other hand, it can produce the desired LMO so as to accelerate subsequent post-SCF calculations such as TDDFT (vide post) and wave-function-based local correlation methods.

**3.2. FLMO-Based TDDFT.** **3.2.1. Efficiency and Accuracy.** Having discussed the efficient construction of the FLMO, we now come to the main point of the present work, i.e., linear scaling TDDFT. Linear polyacetylenes ( $C_nH_{n+2}$ ) of different chain lengths are first chosen to examine the efficiency and accuracy of the FLMO-TDDFT, which can be monitored by a single threshold  $\eta$  (cf. eq 13). Four values of  $\eta$  ( $\infty$ , 7, 6, 5) are considered for cutting off the FLMO pairs. For comparison,  $\eta = 5$  is also used for the CMO pairs as well as the occupied LMO–virtual CMO pairs. The so-defined schemes are to be denoted as LMO( $\eta$ ), LMO-CMO( $\eta$ ), or CMO( $\eta$ ). Note that both LMO( $\infty$ ) and LMO-CMO( $\infty$ ) are equivalent to CMO( $\infty$ ).

As the computational costs of the two expensive steps (??) in forming  $Kb$  are proportional to the number  $N_{ph}$  of p–h pairs, we first examine how the threshold  $\eta$  is correlated with the  $N_{ph}$  on the batches of grid points. In Figure 5, the  $N_{ph}$  with  $\eta_{ai}^B \leq \eta$  on the first batch (centered on the leftmost carbon atom) and that on all of the batches are depicted for different values of  $\eta$ . Note that, for brevity, the number  $N_C$  of carbon atoms has been taken to characterize the system size. Since no pairs are cut off in scheme LMO( $\infty$ ), the  $N_{ph}$  on the first batch scales as  $O(N_C^2)$  and that on all of the batches scales as  $O(N_C^3)$ . It is noticeable that, even with  $\eta = 5$ , the  $N_{ph}$  in the CMO representation cannot significantly be reduced. The scaling of CMO(5) is almost the same as LMO( $\infty$ ), just with a smaller prefactor. By contrast, in the LMO representation with  $\eta = 7, 6$ , and 5, the  $N_{ph}$ 's on the first batch all approach constants, and the  $N_{ph}$ 's on all of the batches exhibit a perfect linear scaling with the system size; see the linearly fitted functions in Table 3. For  $C_{60}H_{62}$ , only 3%, 1%, and 0.4% of the total  $N_{ph}$  in LMO( $\infty$ ) are kept in LMO(7), LMO(6), and LMO(5), respectively. Therefore, a significant reduction of the computational costs can be achieved in forming  $Kb$  by prescreening the p–h pairs in the LMO presentation. As for the LMO–CMO(5) case, the  $N_{ph}$  scales roughly as  $O(N_C^2)$  due to the strong locality of the occupied LMO. The  $N_{ph} - N_C$  curves lie between the LMO(7) and LMO(6) ones for  $N_C$ 's smaller than 60 and then exceed the LMO(7) ones for larger  $N_C$ 's.

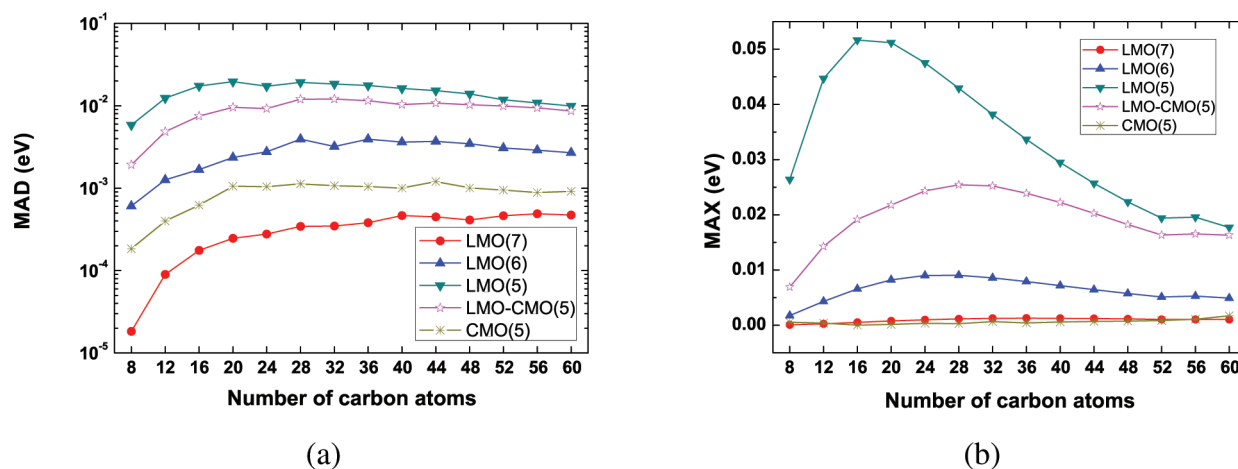


Figure 8. (a) Mean absolute deviations (MAD) and (b) maximum absolute deviations (MAX) of  $LMO(\eta)$  from  $LMO(\infty)$  in the energies of the five lowest excited states of  $C_nH_{n+2}$ .

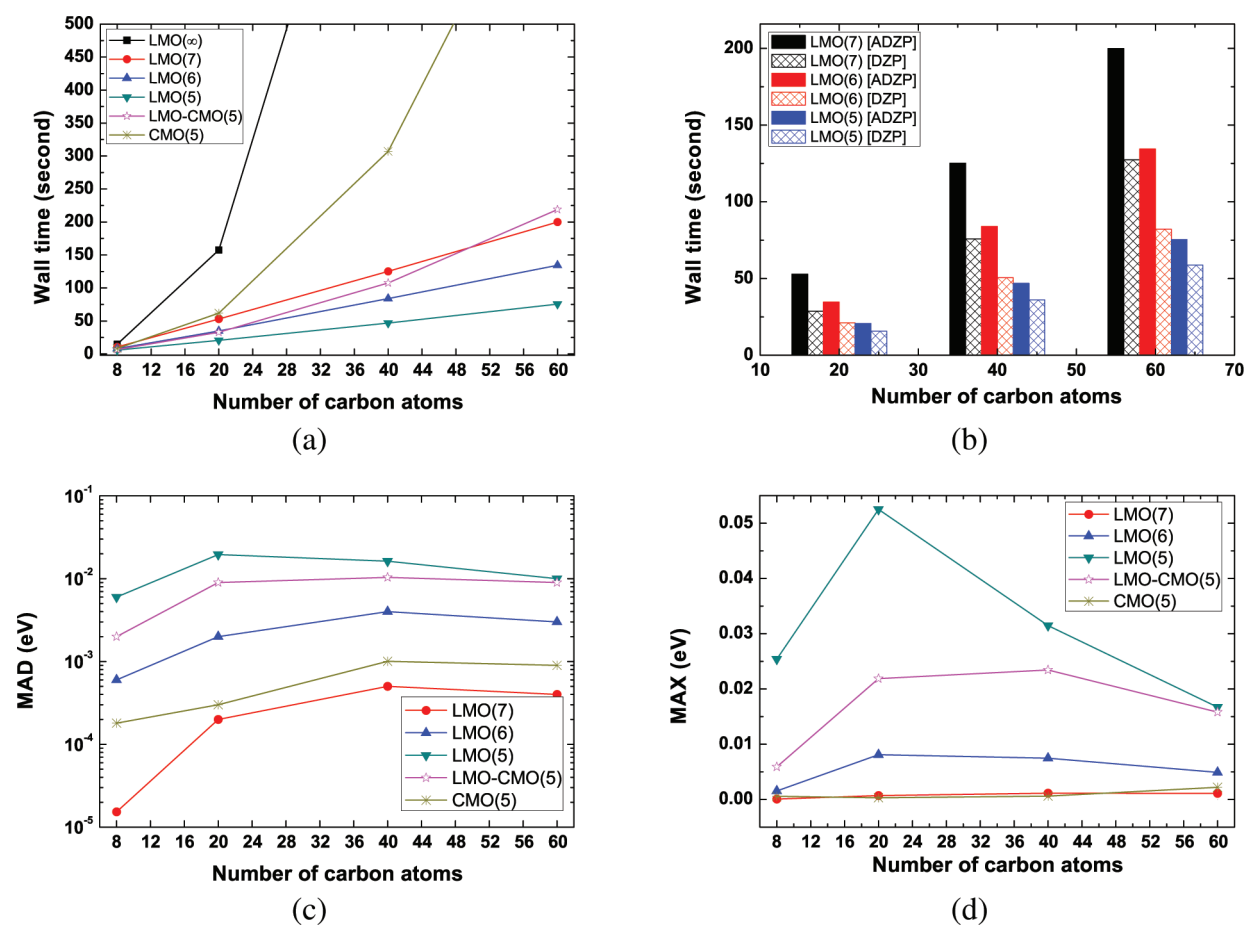


Figure 9. Calculations of the five lowest excited states of  $C_nH_{n+2}$  at the ALDA/ADZP level. (a) Averaged total wall times per trial vector for  $\rho_{\text{ind}}$ ,  $V_{\text{ind}}$  and [Kb]. (b) Comparison between the averaged total wall times per trial vector with the ADZP and DZP basis sets. (c) Mean absolute deviations (MAD) and (d) maximum absolute deviations (MAX) of  $LMO(\eta)$  from  $LMO(\infty)$  in the excitation energies.

The overall computational costs of different schemes are compared in Figure 6a for linear polyacetylenes  $C_nH_{n+2}$ . The five lowest excited states are calculated for each molecule. It is seen from Table 3 that both  $LMO(\infty)$  (equivalent to  $CMO(\infty)$  and  $LMO-CMO(\infty)$ ) and  $CMO(5)$  scale sharper than

quadratic, whereas  $LMO(5)$  to  $LMO(7)$  exhibit linear scaling and are all cheaper than  $CMO(5)$ , even for molecules as small as  $C_8H_{10}$ . Although it is not strictly comparable, it deserves to be mentioned that the  $LMO(6)$  ( $LMO(5)$ ) calculations of the five states of  $C_nH_{n+2}$  are only marginally more expensive (cheaper)



Table 4. LDA/DZP-TDDFT Calculations of Selected Molecules<sup>a</sup>

molecule	scheme	$N_{\text{ph}}$	$N_{\text{it}}$	$N_{\text{b}}$	$T_1$	$T_2$	$T_3$	$T_{\text{tot}}$
$\text{C}_{20}\text{H}_{22}$	LMO( $\infty$ )	$1.8 \times 10^8$ (100%)	7	43	1720	449	2510	4700
	LMO(7)	$3.0 \times 10^7$ (17%)	7	43	329	449	459	1244
	LMO(6)	$1.5 \times 10^7$ (8%)	7	43	198	449	261	915
	LMO(5)	$5.8 \times 10^6$ (3%)	7	42	100	439	120	664
	LMO-CMO(5)	$1.9 \times 10^7$ (11%)	7	47	259	491	343	1100
	CMO(5)	$5.3 \times 10^7$ (29%)	7	43	490	449	740	1690
$\text{C}_{20}\text{H}_{22}$ (ADZP)	LMO( $\infty$ )	$3.2 \times 10^8$ (100%)	8	52	2964	546	4680	8195
	LMO(7)	$7.6 \times 10^7$ (24%)	8	51	887	536	1264	2697
	LMO(6)	$3.6 \times 10^7$ (11%)	8	51	479	536	653	1678
	LMO(5)	$1.3 \times 10^7$ (4%)	8	48	198	504	241	948
	LMO-CMO(5)	$3.4 \times 10^7$ (11%)	8	52	494	546	666	1711
	CMO(5)	$9.5 \times 10^7$ (30%)	8	57	1191	599	1739	3533
$\text{C}_{20}\text{H}_{42}$	LMO( $\infty$ )	$3.8 \times 10^8$ (100%)	4	30	1980	690	3240	5912
	LMO(7)	$4.1 \times 10^7$ (11%)	4	30	339	690	456	1487
	LMO(6)	$1.9 \times 10^7$ (5%)	4	30	189	690	240	1121
	LMO(5)	$6.6 \times 10^6$ (2%)	4	30	99	690	117	908
	LMO-CMO(5)	$2.8 \times 10^7$ (7%)	4	30	270	690	348	1310
	CMO(5)	$1.1 \times 10^8$ (28%)	4	30	684	690	1029	2405
$\text{C}_{25}\text{H}_{44}$	LMO( $\infty$ )	$5.9 \times 10^8$ (100%)	6	50	5400	1315	8100	14820
	LMO(7)	$5.1 \times 10^7$ (9%)	6	50	670	1315	925	2915
	LMO(6)	$2.3 \times 10^7$ (4%)	6	50	385	1315	495	2200
	LMO(5)	$8.2 \times 10^6$ (1%)	6	50	190	1315	225	1735
	LMO-CMO(5)	$3.6 \times 10^7$ (6%)	6	50	565	1315	735	2620
	CMO(5)	$1.4 \times 10^8$ (23%)	6	50	1495	1315	2165	4980
$\text{C}_{21}\text{H}_{34}\text{N}_{10}\text{O}_{10}$	LMO( $\infty$ )	$1.3 \times 10^9$ (100%)	1	5	900	105	1740	2747
	LMO(7)	$6.6 \times 10^7$ (5%)	1	5	76	105	120	303
	LMO(6)	$3.1 \times 10^7$ (2%)	1	5	45	105	65	216
	LMO(5)	$1.1 \times 10^7$ (1%)	1	5	23	105	28	158
	LMO-CMO(5)	$5.6 \times 10^7$ (4%)	1	5	72	105	110	289
	CMO(5)	$2.3 \times 10^8$ (18%)	1	5	175	105	305	587
$\text{C}_{38}\text{H}_{58}\text{N}_{19}\text{O}_{19}$ ( $\alpha$ -helix)	LMO( $\infty$ )	$7.9 \times 10^9$ (100%)	1	5	7150	437	11300	18900
	LMO(7)	$2.9 \times 10^8$ (4%)	1	5	329	437	467	1243
	LMO(6)	$1.1 \times 10^8$ (1%)	1	5	150	437	200	797
	LMO(5)	$3.1 \times 10^7$ (0.4%)	1	5	59	437	70	576
	LMO-CMO(5)	$1.9 \times 10^8$ (2%)	1	5	215	437	293	957
	CMO(5)	$9.5 \times 10^8$ (12%)	1	5	1100	437	1515	3067

<sup>a</sup>  $N_{\text{ph}}$ : Number of effective p–h pairs;  $N_{\text{it}}$ : Number of iterations of the Davidson diagonalization;  $N_{\text{b}}$ : Final dimension of the Davidson iterative subspace;  $T_1$ : Wall time (in second) for  $\rho_{\text{ind}}$ ;  $T_2$ : Wall time for  $V_{\text{ind}}$ ;  $T_3$ : Wall time for  $\mathbf{Kb}$ ;  $T_{\text{tot}}$ : Total wall time.

than the corresponding SCF calculations. Different portions of the wall time in the calculations of  $\text{C}_{60}\text{H}_{62}$  are further displayed in Figure 6b. It is clear that the evaluations of  $\mathbf{Kb}$  and  $\rho_{\text{ind}}$  are significantly more expensive than that of  $V_{\text{ind}}$  in the CMO representation but become even cheaper in the LMO representation. It is also noticeable from Figure 6a and Table 3 that the time-size ( $T-N_{\text{C}}$ ) curves in the LMO representation are not really straight lines. Such deviations from perfect linearity arise from the fact that the number of iterations in the Davidson diagonalization varies with the system size. As the cost of the Davidson diagonalization is proportional to the final dimension  $N_{\text{b}}$  of the iterative subspace, i.e., the total number of  $\mathbf{Kb}$  products to reach convergence, it is more appropriate to take  $N_{\text{b}}$  to characterize the scaling. As shown in Figure 7a, the  $N_{\text{b}}$  with different thresholds  $\eta$  are roughly the same for a given system and decrease as the system gets larger. When averaged over  $N_{\text{b}}$ , the

$T-N_{\text{C}}$  curves for  $\rho_{\text{ind}}$  and  $\mathbf{Kb}$  become perfectly linear for all of the LMO schemes, as shown in Figure 7b and d. They are well correlated with the corresponding curves for the number of effective p–h pairs shown in Figure 5b. One particular remark should be made on the  $V_{\text{ind}}$  curves in Figure 7c, where all of the curves coincide and scale linearly. This is because the same treatment of  $V_{\text{ind}}$  has been made in all of the schemes. That is, the monopole approximation for  $V_{\text{ind}}$  is invoked for grid points outside a radius of 20 au from the position of a given atom. Without the monopole approximation, the evaluation of  $V_{\text{ind}}$  would scale quadratically, albeit with a small prefactor.

To demonstrate the accuracy of the LMO( $\eta$ ) schemes, the mean absolute deviations (MAD) and maximum absolute deviations (MAX) from LMO( $\infty$ ) are displayed in Figure 8 for the five lowest excited states of  $\text{C}_n\text{H}_{n+2}$ . It is seen that the MAD of LMO( $\eta$ ) is roughly  $10^{-\eta+3}$  eV. That is, the accuracy

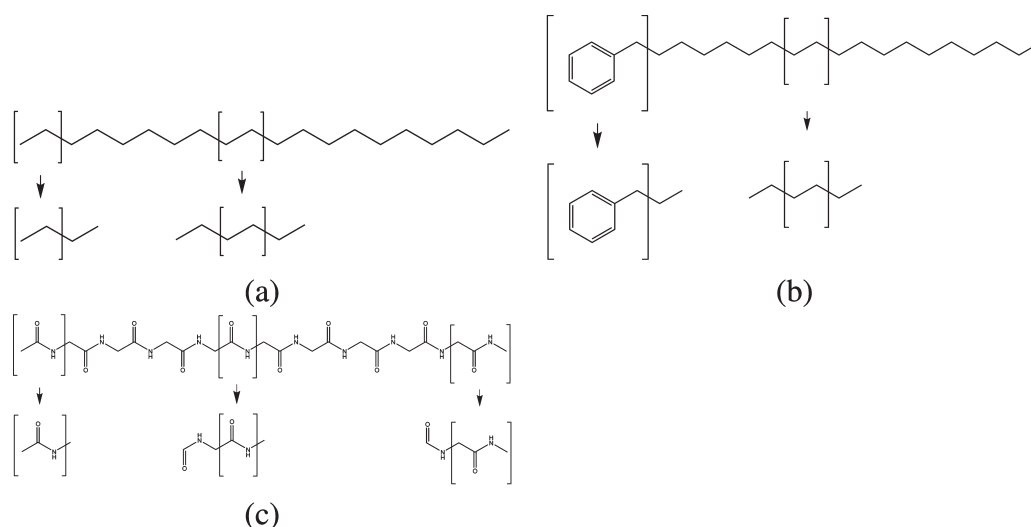


Figure 10. Fragmentation of model linear systems. (a) Polyethylene  $C_{20}H_{42}$ , (b) n-nonadecyl benzene  $C_{25}H_{44}$ , and (c) polypeptide  $C_{21}H_{34}N_{10}O_{10}$ .

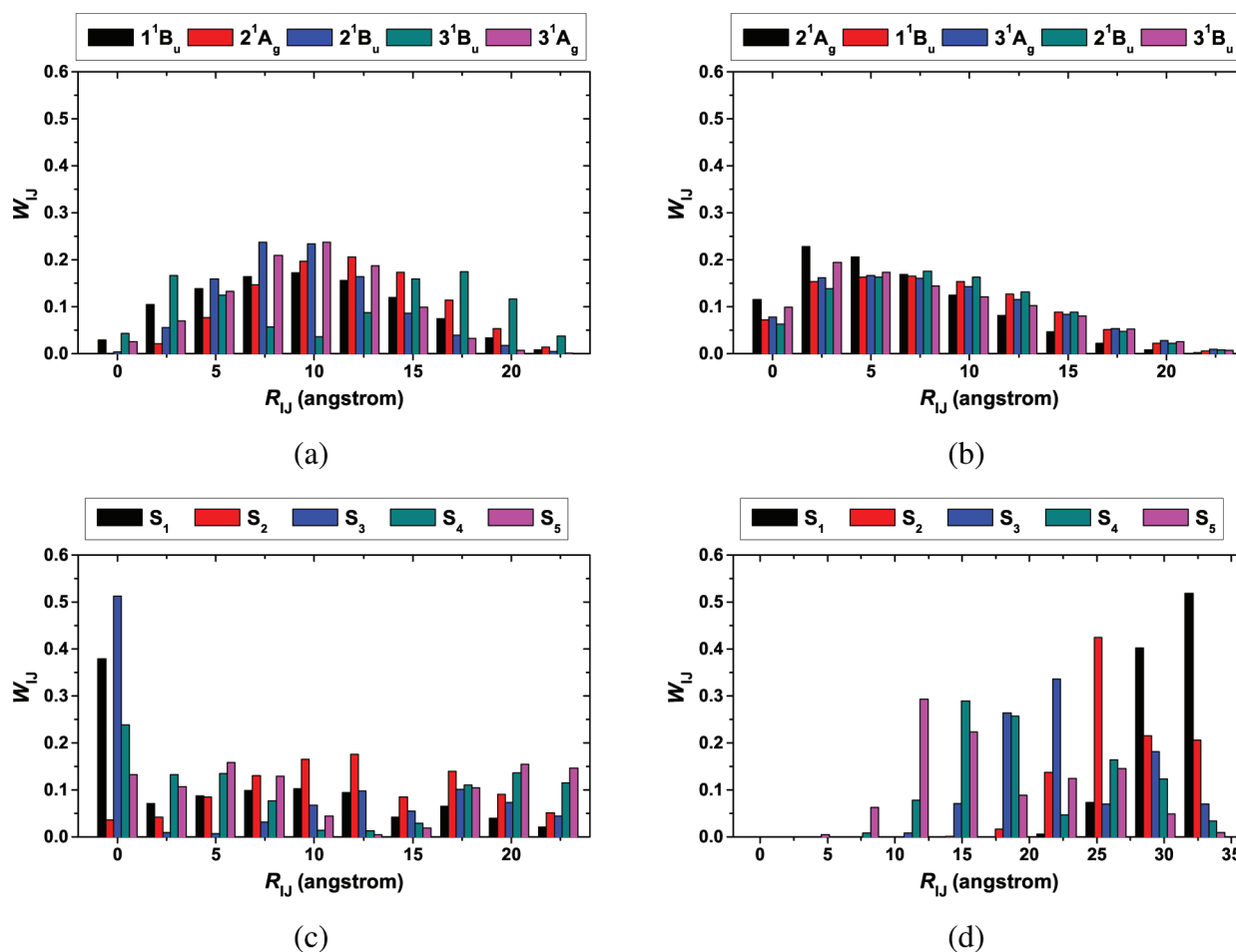


Figure 11. The distributions  $W_{IJ}$  of the excited states. (a)  $C_{20}H_{22}$ , (b)  $C_{20}H_{40}$ , (c)  $C_{25}H_{44}$ , (d)  $C_{21}H_{34}N_{10}O_{10}$ .

of  $LMO(\eta)$  increases monotonically as the threshold  $\eta$  increases to include more p-h pairs. In view of this relation, it is at first glance surprising that the accuracy of  $CMO(5)$  is lower than  $LMO(7)$ , although significantly more p-h pairs are included in the former (cf. Figure 5). However, this can be understood by

realizing that the spatial distributions of the p-h pairs are dramatically different in the CMO and LMO representations. Most of the CMO p-h pairs are delocalized with small  $\eta_{ai}^B$  values. That is, essentially every CMO p-h pair makes a small contribution, the truncation of which is only possible by using a small

**Table 5.** LDA/DZP-TDDFT Excitation Energies (eV) and Oscillator Strengths (au, in parentheses) for the Five Lowest Excited States of Selected Molecules<sup>a</sup>

molecule	state	IPA	LMO( $\infty$ )	LMO(7)	LMO(6)	LMO(5)	LMO-CMO(5)	CMO(5)
C <sub>20</sub> H <sub>22</sub>	1 <sup>1</sup> B <sub>u</sub>	1.9640	2.4841 (2.2709)	2.4833 (2.2647)	2.4754 (2.2113)	2.4302 (1.9288)	2.4623 (2.0958)	2.4839 (2.2625)
	2 <sup>1</sup> A <sub>g</sub>	2.5270	2.5619 (0.0000)	2.5619 (0.0000)	2.5618 (0.0000)	2.5612 (0.0000)	2.5630 (0.0000)	2.5619 (0.0000)
	2 <sup>1</sup> B <sub>u</sub>	2.5758	3.2419 (0.0832)	3.2418 (0.0833)	3.2416 (0.0849)	3.2402 (0.0956)	3.2402 (0.0637)	3.2417 (0.0832)
	3 <sup>1</sup> B <sub>u</sub>	3.1388	3.3752 (1.0148)	3.3750 (1.0126)	3.3733 (0.9995)	3.3637 (0.9256)	3.3798 (1.0670)	3.3742 (0.9969)
	3 <sup>1</sup> A <sub>g</sub>	3.1775	3.4765 (0.0000)	3.4767 (0.0000)	3.4749 (0.0000)	3.4439 (0.0000)	3.4586 (0.0000)	3.4721 (0.0000)
	MAD			0.0003 (0.0017)	0.0025 (0.0153)	0.0201 (0.0888)	0.0094 (0.0494)	0.0011 (0.0053)
C <sub>20</sub> H <sub>22</sub> (ADZP)	1 <sup>1</sup> B <sub>u</sub>	1.9468	2.4570 (2.3056)	2.4563 (2.3004)	2.4489 (2.2510)	2.4045 (1.9725)	2.4351 (2.1204)	2.4573 (2.2987)
	2 <sup>1</sup> A <sub>g</sub>	2.5164	2.5425 (0.0000)	2.5425 (0.0000)	2.5425 (0.0000)	2.5419 (0.0000)	2.5437 (0.0000)	2.5425 (0.0000)
	2 <sup>1</sup> B <sub>u</sub>	2.5469	3.2214 (0.0639)	3.2214 (0.0641)	3.2212 (0.0654)	3.2200 (0.0766)	3.2201 (0.0477)	3.2213 (0.0635)
	3 <sup>1</sup> B <sub>u</sub>	3.1165	3.3418 (1.0197)	3.3416 (1.0174)	3.3401 (1.0038)	3.3301 (0.9279)	3.3467 (1.0788)	3.3412 (1.0079)
	3 <sup>1</sup> A <sub>g</sub>	3.1714	3.4365 (0.0000)	3.4368 (0.0000)	3.4363 (0.0000)	3.4052 (0.0000)	3.4208 (0.0000)	3.4359 (0.0000)
	MAD			0.0002 (0.0015)	0.0020 (0.0144)	0.0195 (0.0875)	0.0090 (0.0521)	0.0003 (0.0038)
C <sub>20</sub> H <sub>42</sub>	2 <sup>1</sup> A <sub>g</sub>	6.5067	6.5190 (0.0000)	6.5190 (0.0000)	6.5190 (0.0000)	6.5188 (0.0000)	6.5187 (0.0000)	6.5190 (0.0000)
	1 <sup>1</sup> B <sub>u</sub>	6.6067	6.6159 (0.0005)	6.6159 (0.0005)	6.6159 (0.0004)	6.6158 (0.0004)	6.6157 (0.0002)	6.6159 (0.0003)
	3 <sup>1</sup> A <sub>g</sub>	6.7576	6.7691 (0.0000)	6.7691 (0.0000)	6.7691 (0.0000)	6.7689 (0.0000)	6.7688 (0.0000)	6.7690 (0.0000)
	2 <sup>1</sup> B <sub>u</sub>	6.9483	6.9623 (0.0021)	6.9623 (0.0020)	6.9623 (0.0020)	6.9621 (0.0017)	6.9619 (0.0011)	6.9622 (0.0018)
	3 <sup>1</sup> B <sub>u</sub>	6.9752	6.9845 (0.0002)	6.9845 (0.0002)	6.9845 (0.0002)	6.9844 (0.0002)	6.9844 (0.0002)	6.9845 (0.0001)
	MAD			0.0000 (0.0000)	0.0000 (0.0000)	0.0001 (0.0001)	0.0002 (0.0002)	0.0000 (0.0001)
C <sub>25</sub> H <sub>44</sub>	S <sub>1</sub>	4.7960	5.0138 (0.0019)	5.0138 (0.0019)	5.0137 (0.0019)	5.0132 (0.0019)	5.0135 (0.0019)	5.0139 (0.0019)
	S <sub>2</sub>	4.8845	5.1524 (0.0296)	5.1523 (0.0293)	5.1520 (0.0275)	5.1505 (0.0204)	5.1516 (0.0196)	5.1523 (0.0283)
	S <sub>3</sub>	5.1037	5.1760 (0.0012)	5.1760 (0.0012)	5.1760 (0.0012)	5.1760 (0.0014)	5.1760 (0.0015)	5.1759 (0.0011)
	S <sub>4</sub>	5.1382	5.6355 (0.0980)	5.6355 (0.0982)	5.6350 (0.0978)	5.6299 (0.0876)	5.6338 (0.0883)	5.6354 (0.1012)
	S <sub>5</sub>	5.1923	5.6464 (0.0073)	5.6464 (0.0073)	5.6463 (0.0070)	5.6460 (0.0055)	5.6462 (0.0061)	5.6464 (0.0074)
	MAD			0.0000 (0.0001)	0.0002 (0.0005)	0.0017 (0.0043)	0.0006 (0.0042)	0.0001 (0.0009)
C <sub>21</sub> H <sub>34</sub> N <sub>10</sub> O <sub>10</sub>	S <sub>1</sub>	4.0930	4.0930 (0.0000)	4.0930 (0.0000)	4.0930 (0.0000)	4.0930 (0.0000)	4.0930 (0.0000)	4.0930 (0.0000)
	S <sub>2</sub>	4.1973	4.1973 (0.0000)	4.1973 (0.0000)	4.1973 (0.0000)	4.1973 (0.0000)	4.1973 (0.0000)	4.1973 (0.0000)
	S <sub>3</sub>	4.2651	4.2651 (0.0000)	4.2651 (0.0000)	4.2651 (0.0000)	4.2651 (0.0000)	4.2651 (0.0000)	4.2651 (0.0000)
	S <sub>4</sub>	4.3283	4.3283 (0.0000)	4.3283 (0.0000)	4.3283 (0.0000)	4.3283 (0.0000)	4.3283 (0.0000)	4.3283 (0.0000)
	S <sub>5</sub>	4.3933	4.3933 (0.0000)	4.3933 (0.0000)	4.3933 (0.0000)	4.3933 (0.0000)	4.3933 (0.0000)	4.3933 (0.0000)
	MAD			0.0000 (0.0000)	0.0000 (0.0000)	0.0000 (0.0000)	0.0000 (0.0000)	0.0000 (0.0000)
C <sub>38</sub> H <sub>58</sub> N <sub>19</sub> O <sub>19</sub> ( $\alpha$ -helix)	S <sub>1</sub>	0.7467	0.7467 (0.0000)	0.7467 (0.0000)	0.7467 (0.0000)	0.7467 (0.0000)	0.7467 (0.0000)	0.7467 (0.0000)
	S <sub>2</sub>	0.9320	0.9320 (0.0000)	0.9320 (0.0000)	0.9320 (0.0000)	0.9320 (0.0000)	0.9320 (0.0000)	0.9320 (0.0000)
	S <sub>3</sub>	0.9929	0.9929 (0.0000)	0.9929 (0.0000)	0.9929 (0.0000)	0.9929 (0.0000)	0.9929 (0.0000)	0.9929 (0.0000)
	S <sub>4</sub>	1.1782	1.1782 (0.0000)	1.1782 (0.0000)	1.1782 (0.0000)	1.1782 (0.0000)	1.1782 (0.0000)	1.1782 (0.0000)
	S <sub>5</sub>	1.2153	1.2153 (0.0000)	1.2153 (0.0000)	1.2153 (0.0000)	1.2153 (0.0000)	1.2153 (0.0000)	1.2153 (0.0000)
	MAD			0.0000 (0.0000)	0.0000 (0.0000)	0.0000 (0.0000)	0.0000 (0.0000)	0.0000 (0.0000)

<sup>a</sup>IPA: Independent particle approximation. MAD: Mean absolute deviations from LMO( $\infty$ ).

threshold  $\eta$  at the expense of introducing sizable errors. By contrast, the LMO p–h pairs are mostly local with  $\eta_{ai}^B$  spanning a very large range (cf. Figure 1a). That is, only a small number of LMO p–h pairs make significant contributions.

At this stage, one may wonder what happens if the chosen basis set consists also of diffuse functions. To address this, we consider the aforementioned ADZP basis set. It is seen from Figure 9a that the FLMO-TDDFT still scales linearly with respect to the system size. The computational costs by the ADZP and DZP basis sets are further compared in Figure 9b for C<sub>n</sub>H<sub>n+2</sub> with  $n = 20, 40,$  and  $60$ . When going from DZP to ADZP, the costs increase by 65% for both LMO(7) and LMO(6) and by 30% for LMO(5), in line with the increase of 65% in the number of basis functions. More detailed comparisons between the ADZP and DZP basis sets can be found from Table 4 in the case of C<sub>20</sub>H<sub>22</sub>. Note that the use of diffuse functions has a much smaller effect on the FLMO pairs

than on the AO pairs. Again take C<sub>20</sub>H<sub>22</sub> as an example. To achieve the desired accuracy (i.e., 0.01 eV in excitation energy and 0.05 a.u. in oscillator strength),  $\eta = 5$  is already sufficient for the DZP-AO-TDDFT but which is only marginally cheaper (by 7%) than the DZP-FLMO-TDDFT with  $\eta = 6$ . However, the threshold has to be increased to 7 for the ADZP-AO-TDDFT to avoid numerical instabilities (imaginary energies), which is then more expensive than the ADZP-FLMO-TDDFT by a factor of 2.2 for  $\eta = 6$  and a factor of 1.5 for  $\eta = 7$ . Similar findings are also observed for larger molecules like C<sub>40</sub>H<sub>42</sub>. Therefore, the FLMO-TDDFT is clearly more efficient than the AO-TDDFT in the presence of diffuse functions. It is also seen from eqs 10 and 9 that the LMO( $\eta$ ) results of the two basis sets have very much the same accuracy.

3.2.2. Application to Different Types of Systems. Having examined the efficiency and accuracy, we now apply the FLMO-TDDFT

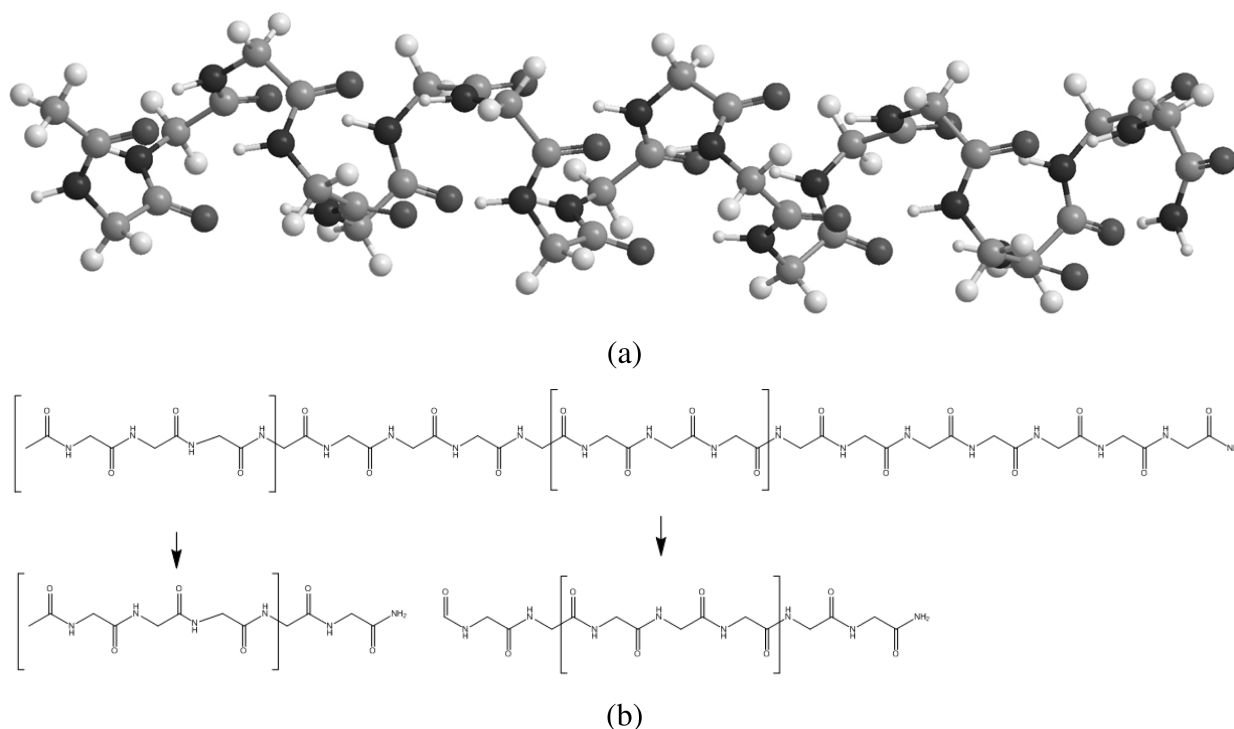


Figure 12. Fragmentation of  $\alpha$ -helix polypeptides. (a) 3D geometry. (b) Fragmentation.

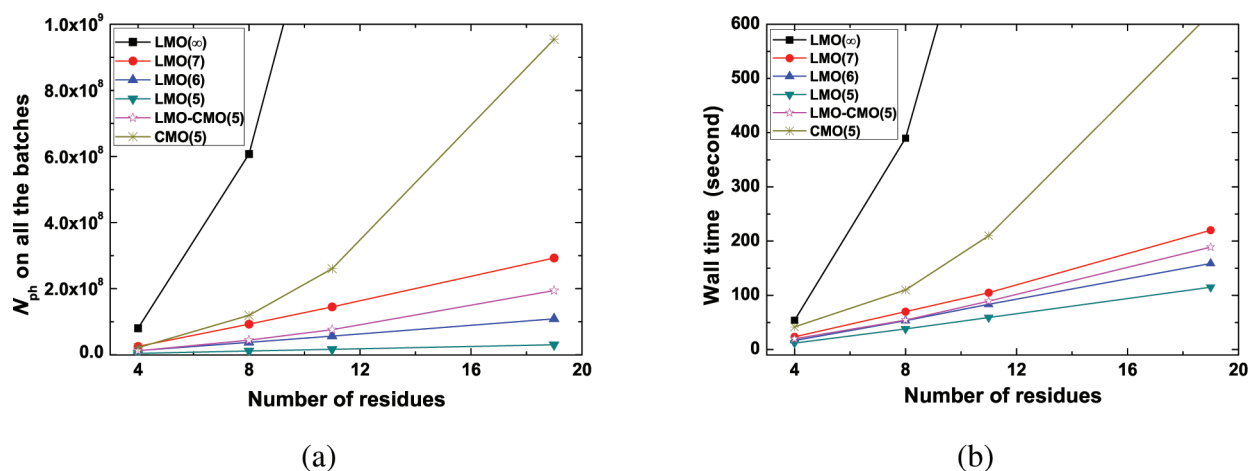


Figure 13. Calculations of the five lowest excited states of  $\alpha$ -helix polypeptides with 4, 8, 11, and 19 residues. (a) Number of p-h pairs. (b) Averaged total wall time per trial vector for  $\rho_{ind}$ ,  $V_{ind}$ , and  $[Kb]$ .

to different kinds of model systems. First, consider linear systems. In addition to the linear polyacetylene  $C_{20}H_{22}$  ( $C_{2h}$ , cf. Figure 2), the linear polyethylene  $C_{20}H_{42}$  ( $C_{2h}$ ), n-nonadecyl benzene  $C_{25}H_{44}$  ( $C_1$ ), and linear polypeptide with 10 residues  $C_{21}H_{34}N_{10}O_{10}$  ( $C_1$ ) are also included. The structures are depicted in Figure 10, all of which are divided into 10 minimal fragments.

The number  $N_{ph}$  of effective p-h pairs, the number  $N_{it}$  of iterations of the Davidson diagonalization, the final dimension  $N_b$  of the Davidson iterative subspace, the wall time for the key steps in constructing the contraction  $Kb$ , and the total wall time are listed in Table 4. It is seen that, compared with the CMO representation, significant reductions of the LMO pairs are possible for all of the systems. As a result, the respective speedups of LMO(7), LMO(6), and LMO(5) amount to 4, 5, and 7 times

for the first three linear systems and are enhanced to 9, 13, and 17 times for the linear polypeptide. The latter is due to stronger locality in the virtual FLMO on one hand and a reduced number  $N_b$  of constructing  $V_{ind}$  on the other.

The excitation energies and oscillator strengths for the five lowest excited states of the four molecules are documented in Table 5. It deserves to be mentioned that the DZP excitation energies for  $C_{20}H_{22}$  differ from the corresponding TZ2P and ADZP values by only 0.03 eV or less, confirming the quality of the calculations. For the other three linear molecules, the truncated LMO schemes essentially reproduce the LMO( $\infty$ ) results. This is a direct consequence of the very small couplings, as can be verified from the proximity between the IPA and TDDFT energies. In particular, identical results have been

**Table 6.** Excitation Energies (in eV) and Oscillator Strengths (au, in parentheses) of the Five Lowest Excited States of Linear Polyacetylenes<sup>a</sup>

molecule	state	LMO( $\infty$ )	(1,1)	(1,2)	(1,3)	(5,1)	(5,2)
C <sub>20</sub> H <sub>22</sub>	1 <sup>1</sup> B <sub>u</sub>	2.4841 (2.2709)	2.4634 (2.2422)	2.4715 (2.2241)	2.4744 (2.2164)	2.4719 (2.2216)	2.4743 (2.2157)
	2 <sup>1</sup> A <sub>g</sub>	2.5619 (0.0000)	2.5489 (0.0004)	2.5578 (0.0000)	2.5608 (0.0001)	2.5596 (0.0000)	2.5609 (0.0002)
	2 <sup>1</sup> B <sub>u</sub>	3.2419 (0.0832)	3.2322 (0.0657)	3.2384 (0.0741)	3.2402 (0.0799)	3.2362 (0.0764)	3.2397 (0.0821)
	3 <sup>1</sup> B <sub>u</sub>	3.3752 (1.0148)	3.3628 (1.0082)	3.3702 (1.0058)	3.3730 (1.0032)	3.3745 (1.0000)	3.3735 (1.0011)
	3 <sup>1</sup> A <sub>g</sub>	3.4765 (0.0000)	3.4692 (0.0000)	3.4731 (0.0000)	3.4742 (0.0000)	3.4717 (0.0000)	3.4760 (0.0000)
	MAD		0.0126 (0.0106)	0.0057 (0.0130)	0.0034 (0.0139)	0.0051 (0.0142)	0.0030 (0.0140)
	MAX		0.0207 (0.0288)	0.0126 (0.0469)	0.0097 (0.0545)	0.0122 (0.0493)	0.0098 (0.0552)
	C <sub>40</sub> H <sub>42</sub>	1 <sup>1</sup> B <sub>u</sub>	1.8407 (2.6041)	1.8153 (2.4929)	1.8266 (2.4637)	1.8311 (2.4528)	1.8269 (2.4620)
2 <sup>1</sup> A <sub>g</sub>		1.8916 (0.0000)	1.8730 (0.0001)	1.8846 (0.0000)	1.8893 (0.0000)	1.8852 (0.0000)	1.8891 (0.0000)
2 <sup>1</sup> B <sub>u</sub>		2.2114 (0.1252)	2.1957 (0.0965)	2.2057 (0.1113)	2.2095 (0.1137)	2.2067 (0.0898)	2.2094 (0.1158)
3 <sup>1</sup> B <sub>u</sub>		2.2672 (2.5045)	2.2484 (2.4612)	2.2585 (2.4331)	2.2625 (2.4268)	2.2587 (2.4513)	2.2624 (2.4246)
3 <sup>1</sup> A <sub>g</sub>		2.2860 (0.0000)	2.2640 (0.0000)	2.2744 (0.0000)	2.2782 (0.0001)	2.2761 (0.0001)	2.2764 (0.0002)
MAD			0.0201 (0.0367)	0.0094 (0.0451)	0.0053 (0.0481)	0.0087 (0.0462)	0.0058 (0.0485)
MAX			0.0254 (0.1112)	0.0141 (0.1404)	0.0096 (0.1513)	0.0138 (0.1421)	0.0099 (0.1529)
C <sub>60</sub> H <sub>62</sub>		1 <sup>1</sup> B <sub>u</sub>	1.6572 (2.5425)	1.6323 (2.3884)	1.6451 (2.3543)	1.6503 (2.3456)	1.6488 (2.3490)
	2 <sup>1</sup> A <sub>g</sub>	1.6940 (0.0000)	1.6732 (0.0002)	1.6860 (0.0000)	1.6911 (0.0000)	1.6898 (0.0000)	1.6914 (0.0000)
	2 <sup>1</sup> B <sub>u</sub>	1.8795 (0.1260)	1.8609 (0.1055)	1.8723 (0.1231)	1.8771 (0.1223)	1.8762 (0.0979)	1.8775 (0.1074)
	3 <sup>1</sup> A <sub>g</sub>	1.8979 (0.0000)	1.8794 (0.0162)	1.8847 (0.0000)	1.8915 (0.0000)	1.8900 (0.0000)	1.8896 (0.0000)
	3 <sup>1</sup> B <sub>u</sub>	1.9094 (3.1896)	1.8889 (3.0736)	1.9001 (3.0376)	1.9047 (3.0293)	1.9034 (3.0589)	1.9048 (3.0393)
	MAD		0.0207 (0.0614)	0.0099 (0.0686)	0.0047 (0.0722)	0.0060 (0.0705)	0.0048 (0.0743)
	MAX		0.0249 (0.1540)	0.0121 (0.1882)	0.0070 (0.1969)	0.0084 (0.1935)	0.0067 (0.2027)

<sup>a</sup> Column LMO( $\infty$ ) employs the fully converged FLMO, while the remaining columns refer to LMO(6) employing the FLMO from the block-diagonalization of the KS matrix due to the superposition of the fragment densities (cf. Table 1). The numbers in parentheses indicate the numbers of double bonds in each (fragment, cap). MAD/MAX: Mean/maximum absolute deviations from LMO( $\infty$ ).

obtained by all of the schemes for the linear polypeptide C<sub>21</sub>H<sub>34</sub>N<sub>10</sub>O<sub>10</sub> where the lowest excited states are dominated by long-range charge-transfers<sup>23</sup> for which the ALDA couplings vanish.

As stated before, the LMO representation offers also a clear interpretation of the excited states in line with chemical/physical intuition. To see this, the distributions  $W_{IJ}$  (cf. 24) of the excited state eigenvectors are depicted in Figure 11 as function of the distances  $R_{IJ}$  between the centers of mass of the fragments. The  $W_{IJ}$ 's corresponding to the same  $R_{IJ}$ 's are summed together. It is clearly seen that the low-lying excited states of C<sub>20</sub>H<sub>22</sub> and C<sub>20</sub>H<sub>42</sub> are combinations of interfragment transitions and are hence completely delocalized. By contrast, the first and third excited states of C<sub>25</sub>H<sub>44</sub> are dominated by local excitations within one fragment ( $R_{IJ} = 0$ ). The excited states of the linear polypeptide C<sub>21</sub>H<sub>34</sub>N<sub>10</sub>O<sub>10</sub> are instead dominated by charge-transfer transitions between well separated fragments. In particular, the first excited state passes through a distance as large as 25 Å. Of course, this might not be very realistic due to inherent problems of the ALDA kernel.

To further reveal the performance of the FLMO-TDDFT for 3D systems, we now consider  $\alpha$ -helix polypeptides, the largest of which consists of 19 residues (C<sub>38</sub>H<sub>58</sub>N<sub>19</sub>O<sub>19</sub>). Its fragmentation is shown in Figure 12. Each fragment is composed of four residues (one  $\alpha$ -helix turn) and is capped with two residues on each side. Note that different fragmentations affect merely the number of SCF iterations but not the FLMO-TDDFT calculations. The wall times and the excitation energies of C<sub>38</sub>H<sub>58</sub>N<sub>19</sub>O<sub>19</sub> are documented in Tables 4 and 5, respectively. It is clearly seen from Figure 13 that the FLMO-TDDFT scales still linearly with respect to the system size characterized by the number of residues.

**3.3. Possible Approximations.** Since it is the orbital overlaps rather than the whole orbitals themselves that determine the contributions of the p–h pairs to the excitations, it should be possible to use approximate instead of fully converged orbitals in the TDDFT calculations. To confirm this viewpoint, we just take the FLMO from the block-diagonalization of the KS matrix due to the superposition of the fragment densities (iteration 0 in Table 1). The so-calculated excitation energies and oscillator strengths are given in Table 6 for three linear polyacetylenes, i.e., C<sub>20</sub>H<sub>22</sub>, C<sub>40</sub>H<sub>42</sub>, and C<sub>60</sub>H<sub>62</sub>. It is seen that the results are very good, with the mean absolute errors in the energies only of a few hundredths of an electronvolt. This again results from the high transferability of the pFLMO.

At variance with the FMO-TDDFT,<sup>6,7</sup> where the same many-body expansion is applied to both the ground and excited state energies, the ground and excited state calculations can in the present case be approximated separately with well controlled accuracy. For instance, caps varying in size can be chosen for different fragments in the subsystem calculations, and different thresholds can be applied to prescreen the interfragment matrix elements in the global SCF calculation. Likewise, in the FLMO-TDDFT calculation, different thresholds  $\eta$  can be applied to p–h pairs from different fragments. Such hierarchical approximations allow one to treat different portions of the whole system with different accuracies and even with different Hamiltonians.<sup>24–28</sup>

Moreover, like the KS matrix (see 23), the **A** and **B** matrices in eq 7 can be partitioned into fragment contributions, viz.:

$$\begin{aligned} \mathbf{M} &= \sum_I \mathbf{M}_{II,II} + \sum_{I \neq J} \mathbf{M}_{II,JJ} + \sum_I \sum_{K \neq L} (\mathbf{M}_{II,KL} + \mathbf{M}_{KL,II}) \\ &+ \sum_{I \neq J} \sum_{K \neq L} \mathbf{M}_{IJ,KL}, (\mathbf{M} = \mathbf{A}, \mathbf{B}) \end{aligned} \quad (26)$$

where the first term on the right-hand side represents pure local excitations, the second term the couplings between local excitations, the third term the couplings between local and change-transfer excitations, and the last term pure change-transfer excitations. A many-body expansion is also possible, in the spirit of the FMO-TDDFT.<sup>6,7</sup> However, there are at most four-body terms here, whereas up to  $N_m$ -body terms are required in the FMO-TDDFT, with  $N_m$  being the total number of monomers. In this way, for a given system, some central blocks can be selected on the basis of chemical intuition, and the accuracy can be systematically improved by adding more blocks.

#### 4. CONCLUSIONS AND OUTLOOK

A very efficient linear-scaling TDDFT has been developed for uniform treatments of all kinds of excitations of large systems composed of arbitrary chemical bonds. This has been achieved by making full use of the locality of the p-h basis in the LMO representation. The required FLMO can readily be generated by synthesizing the pFLMO obtained from subsystem calculations, in the spirit of “from fragments to molecule”. The novel block-diagonalization scheme presented here plays an essential role in retaining/resuming the locality both of the occupied and virtual MO. Very interestingly, this particular construction of the FLMO has much in common with the formulation of exact two-component relativistic theories (see Appendix A). Another salient feature of the FLMO-TDDFT is to combine the good of the LMO and CMO representations. That is, all matrix elements are evaluated in the LMO representation, but the eigenvalue problem is solved in the CMO representation via efficient unitary transformations. In this way, not only the convergence of the iterative diagonalization can rapidly be achieved but also the molecular symmetry of arbitrary order can fully be employed. Moreover, the orbital picture and number of electrons are retained so as to allow a clear interpretation of the nature of the excited states, whether local, delocalized, or charge transfer. Also because of this, the algorithm is fully compatible with the spin-adapted TDDFT for open-shell systems<sup>29,30</sup> as well as the relativistic counterpart of TDDFT.<sup>31–36</sup>

Further enhancement of the efficiency of the algorithm is still possible. Apart from those schemes already outlined in section 3.3, the “energetic locality” of both the LMO and CMO can further be explored. The “primitive fragment CMO” from canonical orthogonalization of the pFLMO can be classified into inner core, outer core, valence, and high-lying ones, based on which valence type of pFLMO can be identified. Only this subset rather than the whole set of the pFLMO needs to be employed for accounting for the interfragment interactions. The SCF calculation of the whole system can then be greatly simplified. Likewise, the concept of active space can be introduced in the TDDFT calculation as is done in the CMO representation. The algorithms apply to all kinds of XC functionals/kernels. That is, the use of more refined functionals such as hybrids with some portion of the HF exchange affects only the prefactor but not the scaling. Last but not least, the present FLMO can also be employed in wave-function-based local correlation methods. Progress along these directions is being made in our laboratory.

#### APPENDIX A. MATRIX BLOCK-DIAGONALIZATION

The purpose here is to find a unitary transformation matrix  $U$  that can block-diagonalize a given Hermitian matrix  $F$

partitioned as

$$F = \begin{pmatrix} F_{11} & F_{12} \\ F_{21} & F_{22} \end{pmatrix}, F_{11} = F_{11}^\dagger, F_{21} = F_{12}^\dagger, F_{22} = F_{22}^\dagger \quad (27)$$

Noticeably, there are infinitely many such unitary transformations. However, it is the one leading to least modifications of the diagonal blocks of  $F$  that is to be sought. Without a loss of generality, we can write the  $U$  matrix as

$$U = DPQ \quad (28)$$

$$D = \begin{pmatrix} I & -X^\dagger \\ X & I \end{pmatrix} \quad (29)$$

$$P = \begin{pmatrix} P_+ & 0 \\ 0 & P_- \end{pmatrix}, P_+ = (I + X^\dagger X)^{-1/2}, P_- = (I + XX^\dagger)^{-1/2} \quad (30)$$

$$Q = \begin{pmatrix} Q_+ & 0 \\ 0 & Q_- \end{pmatrix} = Q^\dagger \quad (31)$$

It is readily shown that, under the condition

$$\tilde{F}_{21} = F_{21} - XF_{11} + F_{22}X - XF_{12}X = 0 \quad (32)$$

the  $F$  matrix will be block-diagonalized as

$$\tilde{F} = U^\dagger F U = \begin{pmatrix} \tilde{F}_{11} & 0 \\ 0 & \tilde{F}_{22} \end{pmatrix} \quad (33)$$

where

$$\tilde{F}_{11} = Q_+^\dagger P_+^\dagger (F_{11} + X^\dagger F_{21} + F_{12}X + X^\dagger F_{22}X) P_+ Q_+ \quad (34)$$

$$\tilde{F}_{22} = Q_-^\dagger P_-^\dagger (F_{22} - XF_{12} - F_{21}X^\dagger + XF_{11}X^\dagger) P_- Q_- \quad (35)$$

It is hence clear that  $D$  does the decoupling through condition 32.  $P$  does the renormalization, while  $Q$  cannot be uniquely determined. However, according to Lemma 2 in ref 37, for a given  $X$  and hence positive-definite  $DP$ , the minimization problem

$$\min \|I - DPQ\|_F, Q^\dagger Q = I \quad (36)$$

has the unique solution  $Q = I$ . That is, if  $Q = I$  in eq 28, the resultant  $U$  matrix will lead to least modifications of the diagonal blocks of  $F$  in the sense of a Frobenius norm  $\|\cdot\|_F$ . This is precisely the result we want.

Instead of solving the nonlinear algebraic Riccati eq 32, we propose to solve the simpler linear Sylvester equation:

$$F_{21}^{(i)} - X^{(i)} F_{11}^{(i)} + F_{22}^{(i)} X^{(i)} = 0 \quad (37)$$

iteratively by diagonalizing both  $F_{11}^{(i)}$  and  $F_{22}^{(i)}$ . The starting point is  $F^{(0)} = F$ . In terms of the so-obtained  $X^{(i)}$ ,  $U^{(i)} = D^{(i)} P^{(i)}$ , and hence  $F^{(i+1)} = U^{(i)\dagger} F^{(i)} U^{(i)}$  can be constructed. The iterations continue until the off-diagonal blocks of  $F^{(i+1)}$  vanish. Typically, only 2–3 cycles are needed to achieve convergence. Actually, it can rigorously be proven that the convergence is globally

monotonic and locally cubic as long as there exists a gap between the largest eigenvalue of  $F_{11}$  and the smallest eigenvalue of  $F_{22}$ . The mathematical aspects of the scheme will be presented elsewhere.

Alternatively, the  $X$  matrix can be constructed noniteratively by first solving the standard eigenvalue problem

$$FC = CE \quad (38)$$

which can be rewritten in block form

$$\begin{pmatrix} F_{11} & F_{12} \\ F_{21} & F_{22} \end{pmatrix} \begin{pmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{pmatrix} = \begin{pmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{pmatrix} \begin{pmatrix} E_1 & 0 \\ 0 & E_2 \end{pmatrix} \quad (39)$$

It can readily be shown that

$$X = C_{21}C_{11}^{-1}, \quad -X^\dagger = C_{12}C_{22}^{-1} \quad (40)$$

both satisfy condition 32 by inserting  $C_{21} = XC_{11}$  or  $C_{12} = -X^\dagger C_{22}$  into eq 39.

In view of eqs 33, 39, and 40, we have

$$\tilde{F}\tilde{C} = \tilde{C}E \quad (41)$$

where

$$\begin{aligned} \tilde{C} &= U^\dagger C = \begin{pmatrix} \sqrt{I + X^\dagger X} C_{11} & 0 \\ 0 & \sqrt{I + X X^\dagger} C_{22} \end{pmatrix} \\ &= \begin{pmatrix} \tilde{C}_{11} & 0 \\ 0 & \tilde{C}_{22} \end{pmatrix} \end{aligned} \quad (42)$$

After the full diagonalization of  $F$ , the constructions of  $X$  by the first equality of eq 40, as well as the remaining matrices, including  $\tilde{F}_{11}$ ,  $\tilde{F}_{22}$ ,  $\tilde{C}_{11}$ , and  $\tilde{C}_{22}$ , etc., are very cheap. Therefore, this noniterative block-diagonalization is much favored over the above iterative scheme, which is more expensive than the full diagonalization by a factor roughly equal to the number of iterations.

If the input  $F$  is the KS matrix in the orthonormal pFLMO basis and partitioned according to the occupation of the CMO,  $\tilde{F}_{11}$  in eq 34 and  $\tilde{F}_{22}$  in 35 would be the  $F_{oo}^{LMO}$  and  $F_{vv}^{LMO}$ , respectively. Likewise, the  $\tilde{C}_{11}$  and  $\tilde{C}_{22}$  blocks of eq 42 would be the respective  $U_{oo}$  and  $U_{vv}$  matrices required in eq 15 as  $E = F^{CMO}$ . The relationship between the FLMOs, CMOs, pFLMOs, and AOs reads

$$\phi^{LMO} = \phi^{CMO} \tilde{C}^\dagger = \phi^{pFLMO} \tilde{C} \tilde{C}^\dagger = \phi^{pFLMO} U = \phi^{AO} L^{pFLMO} U \quad (43)$$

Finally, it deserves to be mentioned that the above iterative and noniterative matrix block-diagonalization schemes have actually been employed to reduce the four-component matrix Dirac equation down to the exact two-component ones.<sup>24–28</sup> However, they are used here for a completely different purpose, viz., localization of the CMO represented in the pFLMO basis. The present iterative scheme is more robust than the previous ones,<sup>24,25</sup> where it is the Riccati eq 32 rather than the Sylvester eq 37 that has been solved, with the elements of  $F$  fixed throughout the iterations rather than updated as done here. Moreover, the present idea of “from fragments to molecule” for synthesizing the molecular wave function can also be regarded as an extension of the previous idea of “from atoms/fragments to

molecule” for synthesizing the relativistic molecular Hamiltonian.<sup>26,27</sup> As such, the two very different fields, linear-scaling TDDFT and exact two-component relativistic theories, share the same mathematics and just differ in making use of different aspects of physical locality.

## AUTHOR INFORMATION

### Corresponding Author

\*E-mail: liuwj@pku.edu.cn.

## ACKNOWLEDGMENT

The principal author (W.L.) thanks Mr. W. Zhang for testing and implementing the block-diagonalization algorithms presented in Appendix during his visit to Peking University. The research of this work was supported by grants from the National Natural Science Foundation of China (Project No. 21033001).

## REFERENCES

- (1) Runge, E.; Gross, E. K. U. *Phys. Rev. Lett.* **1984**, *52*, 997.
- (2) Casida, M. E. *Recent Advances in Density Functional Methods*, 1st ed.; World Scientific: Singapore, 1995.
- (3) Yam, C.; Yokojima, S.; Chen, G. *Phys. Rev. B* **2003**, *68*, 153105.
- (4) Wang, F.; Yam, C.; Chen, G.; Fan, K. *J. Chem. Phys.* **2007**, *126*, 134104.
- (5) Cui, G.; Fang, W.; Yang, W. *Phys. Chem. Chem. Phys.* **2010**, *12*, 416.
- (6) Chiba, M.; Fedorov, D. G.; Kitaura, K. *J. Chem. Phys.* **2007**, *127*, 104108.
- (7) Chiba, M.; Koido, T. *J. Chem. Phys.* **2010**, *133*, 044113.
- (8) Fujimoto, K.; Yang, W. *J. Chem. Phys.* **2008**, *129*, 054102.
- (9) van Gisbergen, S. J. A.; Guerra, C. F.; Baerends, E. J. *J. Comput. Chem.* **2000**, *21*, 1511.
- (10) Coriani, S.; Høst, S.; Jansik, B.; Thøgersen, L.; Olsen, J.; Jørgensen, P.; Reine, S.; Pawłowski, F.; Helgaker, T.; Sałek, P. *J. Chem. Phys.* **2007**, *126*, 154108.
- (11) Boys, S. F. *Quantum Theory of Atoms, Molecules, and the Solid State*; Academic Press: New York, 1966; p 253.
- (12) Edmiston, C.; Ruedenberg, K. *Rev. Mod. Phys.* **1963**, *35*, 457.
- (13) Pipek, J.; Mezey, P. G. *J. Chem. Phys.* **1989**, *90*, 4916.
- (14) Stratmann, R. E.; Scuseria, G. E.; Frisch, M. J. *J. Chem. Phys.* **1998**, *109*, 8218.
- (15) Liu, W.; Hong, G.; Dai, D.; Li, L.; Dolg, M. *Theor. Chem. Acc.* **1997**, *96*, 75.
- (16) Liu, W.; Wang, F.; Li, L. *J. Theor. Comput. Chem.* **2003**, *2*, 257.
- (17) Liu, W.; Wang, F.; Li, L. *Recent Advances in Computational Chemistry*; World Scientific: Singapore, 2004; Vol. 5, p 257.
- (18) Liu, W.; Wang, F.; Li, L. *Encyclopedia of Computational Chemistry*; Chichester, U.K., 2004.
- (19) Miura, M.; Aoki, Y. *Mol. Phys.* **2010**, *108*, 205.
- (20) Peng, D.; Ma, J.; Liu, W. *Int. J. Quantum Chem.* **2009**, *109*, 2149.
- (21) Li, S.; Li, W. *Annu. Rep. Prog. Chem., Sect. C* **2008**, *104*, 256.
- (22) Vosko, S. J.; Wilk, L.; Nusair, M. *Can. J. Phys.* **1980**, *58*, 1200.
- (23) Peach, M. J. G.; Benfield, P.; Helgaker, T.; Tozer, D. J. *J. Chem. Phys.* **2008**, *128*, 044118.
- (24) Kutzelnigg, W.; Liu, W. *J. Chem. Phys.* **2005**, *123*, 241102.
- (25) Liu, W.; Kutzelnigg, W. *J. Chem. Phys.* **2007**, *126*, 114107.
- (26) Liu, W.; Peng, D. *J. Chem. Phys.* **2006**, *125*, 044102.
- (27) Peng, D.; Liu, W.; Xiao, Y.; Cheng, L. *J. Chem. Phys.* **2007**, *127*, 104106.
- (28) Liu, W. *Mol. Phys.* **2010**, *108*, 1679.
- (29) Li, Z.; Liu, W. *J. Chem. Phys.* **2010**, *133*, 064106.
- (30) Li, Z.; Liu, W.; Zhang, Y.; Suo, B. *J. Chem. Phys.* **2011**, *134*, 134101.
- (31) Gao, J.; Liu, W.; Song, B.; Liu, C. *J. Chem. Phys.* **2004**, *121*, 6658.

- (32) Gao, J.; Zou, W.; Liu, W.; Xiao, Y.; Peng, D.; Song, B.; Liu, C. *J. Chem. Phys.* **2005**, *123*, 054102.
- (33) Peng, D.; Zou, W.; Liu, W. *J. Chem. Phys.* **2005**, *123*, 144101.
- (34) Xu, W.; Ma, J.; Peng, D.; Zou, W.; Liu, W.; Staemmler, V. *Chem. Phys.* **2009**, *356*, 219.
- (35) Xu, W.; Zhang, Y.; Liu, W. *Sci. China Ser. B Chem.* **2009**, *1945*, 52.
- (36) Zhang, Y.; Xu, W.; Sun, Q.; Zou, W.; Liu, W. *J. Comput. Chem.* **2010**, *532*, 31.
- (37) Dieci, S.; Friedman, M. J. *Numer. Linear Algebra Appl.* **2001**, *8*, 317.



# Some Comments on Topological Approaches to the $\pi$ -Electron Currents in Conjugated Systems

Timothy K. Dickens,<sup>\*,†</sup> José A. N. F. Gomes,<sup>‡</sup> and Roger B. Mallion<sup>§</sup>

<sup>†</sup>University Chemical Laboratory, University of Cambridge, Lensfield Road, Cambridge, CB2 1EW, England, United Kingdom

<sup>‡</sup>REQUIMTE, Departamento de Química, Faculdade de Ciências, Universidade do Porto, Rua do Campo Alegre 687, 4169-007 Porto, Portugal

<sup>§</sup>School of Physical Sciences, University of Kent, Canterbury, CT2 7NH, England, United Kingdom

**ABSTRACT:** Within the past two years, three sets of independent authors (Mandado, Ciesielski et al., and Randić) have proposed methods in which  $\pi$ -electron currents in conjugated systems are estimated by invoking the concept of circuits of conjugation. These methods are here compared with ostensibly similar approaches published more than 30 years ago by two of the present authors (Gomes and Mallion) and (likewise independently) by Gayoso. Patterns of bond currents and ring currents computed by these methods for the nonalternant isomer of coronene that was studied by Randić are also systematically compared with those calculated by the Hückel–London–Pople–McWeeny (HLPMP) “topological” approach and with the *ab initio*, “ipso-centric” current-density maps of Balaban et al. These all agree that a substantial diamagnetic  $\pi$ -electron current flows around the periphery of the selected structure (which could be thought of as a “perturbed” [18]-annulene), and consideration is given to the differing trends predicted by these several methods for the  $\pi$ -electron currents around its central six-membered ring and in its internal bonds. It is observed that, for any method in which calculated  $\pi$ -electron currents respect Kirchhoff’s Laws of current conservation at a junction, consideration of *bond* currents—as an alternative to the more-traditional *ring* currents—can give a different insight into the magnetic properties of conjugated systems. However, provided that charge/current conservation is guaranteed—or Kirchhoff’s First Law holds for bond currents instead of the more-general current-densities—then ring currents represent a more efficient way of describing the molecular reaction to the external magnetic field: ring currents are independent quantities, while bond currents are not.

## 1. INTRODUCTION

Thirty-five years ago, Randić proposed<sup>1</sup> the approach for calculating resonance energies of conjugated systems that has become known as the method of conjugated circuits. Very recently, the same author<sup>2</sup> and—independently and almost simultaneously—Mandado<sup>3</sup> and Ciesielski et al.<sup>4</sup> have adapted this formalism in order to estimate the relative intensities of the  $\pi$ -electron currents that (classically) are considered to flow along the individual bonds of such conjugated systems when excited to do so by the presence of an external magnetic field; this magnetic field may be assumed, without a loss of generality, to be oriented in a direction at right angles to the molecular plane of the conjugated system in question (taken to be geometrically planar). This phenomenon is generally known as the “ring-current effect” (see refs 5–7 for reviews)—but Randić specifically, in his recent note,<sup>2</sup> has been especially careful not to invoke any explicit assumptions about “rings” *per se*.

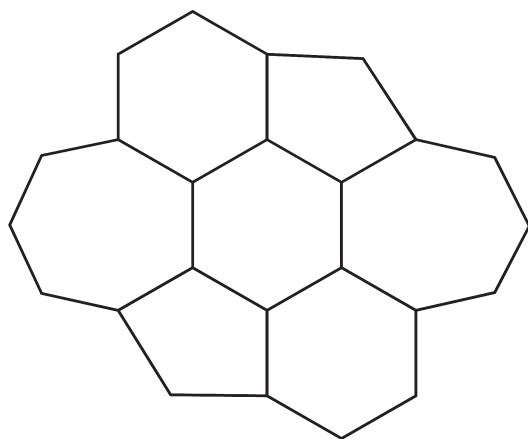
Randić’s work,<sup>2</sup> which follows from his earlier, preliminary thoughts on the matter,<sup>8</sup> has developed from Balaban et al.’s recently published<sup>9</sup> maps pictorially indicating the patterns of  $\pi$ -electron current densities in 18 isomers that include, and are all related to, coronene. Ciesielski et al.<sup>4</sup> have also compared predictions (on the carcinogen 3,4-benzopyrene) arising from their own method<sup>4</sup> with a current-density map provided by Fowler.<sup>10</sup> The computations that give rise to Fowler’s maps<sup>9,10</sup>—like the calculations of Randić<sup>2</sup>—also make no presuppositions about the existence of rings. The so-called “ipso-centric” method that the Fowler school routinely adopts<sup>9,10</sup> invokes *ab initio* Gaussian computations in order to produce what Randić<sup>2</sup> describes as “well converged current maps

with rather modest basis-sets.” The origin of the contributions to the overall  $\pi$ -electron current density is also closely tracked and traced by this approach:<sup>9</sup> diamagnetic contributions come from  $[4n + 2]$  cycles and paramagnetic ones from  $[4n]$  cycles. For the purpose of comparing the predictions of Balaban et al.<sup>9</sup> with the results of his own calculations based on the method described in ref 2, Randić<sup>2</sup> singled out a particular one of the 18 isomers studied by Balaban et al.<sup>9</sup> (though he has since extended his investigation to include all of them<sup>11</sup>); this was the conjugated system that Balaban et al. labeled<sup>9</sup> “13” and called—on a systematic notation that they introduced and defined<sup>9</sup>—“[567567]”. It has 180° rotational symmetry about a perpendicular axis through the center of the middle six-membered ring; if the structure is planar, its symmetry is  $C_{2h}$ . The carbon-atom skeleton of this system is shown in Figure 1. Randić<sup>2</sup> obtained encouraging agreement between his results and the pictorial, qualitative current-density map of Balaban et al.<sup>9</sup> and thereby drew some conclusions of a philosophical nature about the virtues of quantum mechanical vs graph-theoretical approaches to conjugated systems of this type—conclusions that were somewhat akin to those that had been expressed by Gayoso, in a similar context, in the 1979 *Comptes Rendus*.<sup>12,13</sup>

As is to be expected of a theory that is concerned with the microscopic equivalent of classical electrical networks,<sup>6</sup> many previous treatments<sup>14</sup> of the “ring-current” effect—for example, those due to London, Coulson, Pople, McWeeny (as generalized by Veillard and by Gayoso and Boucekkinne), Aihara, and

Received: April 14, 2011

Published: September 19, 2011



**Figure 1.** The carbon-atom skeleton of the conjugated system named “[567567]” by Balaban et al.<sup>9</sup> It is an isomer of coronene and, on the assumption that it is planar, it is of symmetry  $C_{2h}$ , having 2-fold rotational symmetry about an axis through the center of the middle six-membered ring, at right angles to the plane of the paper.

Mizoguchi<sup>14</sup>—have naturally invoked the concept of *circuits*,<sup>14</sup> *per se*, but these have not, in general, been circuits, specifically, of *conjugation*. The sudden and (to us) unexpected recent resurgence of interest<sup>2–4,8</sup> in approaches to magnetic properties that involve, specifically, *circuits of conjugation*—a topic that was originally studied independently some 30 years ago by two of us (Gomes and Mallion)<sup>15–19</sup> and by Gayoso<sup>12</sup>—has therefore motivated us to draw attention to three areas, developed in this paper:

- Previous Similar Work.** We point out that, in the latter part of the 1970s, Gomes and Mallion,<sup>15–19</sup> and (independently) Gayoso,<sup>12</sup> applied what the former authors called “conjugation circuits”<sup>20</sup> when calculating the magnetic properties of conjugated systems. We point out that the approaches of Gomes and Mallion,<sup>15–19</sup> and that of Gayoso,<sup>12</sup> have many similarities to those recently proposed by Randić,<sup>2,8,11,21</sup> by Mandado,<sup>3</sup> and by Ciesielski et al.,<sup>4</sup> and we apply these (and other) methods to the particular conjugated system [567567] (Figure 1) that was selected for study by Randić in ref 2.
- Bond Currents.** We draw attention to the fact that when what have been called “topological ring currents”<sup>22–25</sup>—computed for [567567] (Figure 1) by the recently defined Hückel–London–Pople–McWeeny (hereafter HLPMP) approach<sup>23,24</sup>—are expressed (entirely equivalently) as *bond* currents, qualitative and even semiquantitative agreement is frequently seen between these “topological” bond currents, Randić’s bond currents,<sup>2</sup> the bond currents calculated by the methods of Mandado<sup>3</sup> and of Ciesielski et al.,<sup>4</sup> and the qualitative  $\pi$ -electron current-density maps of Balaban et al.<sup>9</sup>
- Kirchhoff’s Law of Current Conservation.** We observe that, in the context of any method in which calculated currents strictly obey Kirchhoff’s Law on conservation of current at a junction,<sup>26–29</sup> consideration of *bond* currents—as distinct from (entirely equivalent, but more traditional) *ring* currents—can give a different conceptual insight into the reaction of the molecule to the external magnetic field in the case of conjugated systems like [567567] (Figure 1).

## 2. THE “CONJUGATION CIRCUITS” METHOD OF GOMES AND MALLION<sup>15,16</sup> (1976 AND 1979)

**Details of the Method.** In order clearly to recount the method that two of us (Gomes and Mallion) proposed in the second half of the 1970s,<sup>15,16</sup> we here quote *verbatim* what we at the time described as our “prescription”—directly and *in extenso*—from the 1979 *Revista Portuguesa de Química*.<sup>16</sup> This paper was itself a distillation of the method first proposed by one of us (Gomes) in a thesis,<sup>15</sup> written three years earlier; some theoretical justification for the “prescription” in terms of valence-bond theory was attempted in the early 1980s.<sup>17–19</sup> We adopt this procedure because, although the details of the method were openly published more than 30 years ago,<sup>16</sup> the Gomes–Mallion formalism is evidently *not* well-known—two of the sets of very recent authors,<sup>2,4</sup> for example, were, it seems, not aware of it, and the third<sup>3</sup> cited it only in passing, as being a general reference relating merely to the concept of conjugation circuits and not, specifically, to the calculation of magnetic properties *per se*.<sup>30</sup> We therefore quote from ref 16, as follows:<sup>31</sup>

It is assumed in this prescription that the effect of the magnetic field on a molecule is felt independently by every one of the various “conjugation circuits” which are extant in each Kekulé-structure; as far as magnetic properties are concerned, an individual Kekulé-structure may be regarded as a superposition of its constituent “conjugation circuits”, the effects of which are simply additive. The system of “ring currents” in the actual molecule is then obtained by finally averaging the contributions from individual Kekulé-structures over all possible Kekulé-structures which can be devised for the molecule as a whole. Accordingly, the method proposed here for estimating the relative “ring-current” intensities in a given molecule is based on the following postulates:

- The method of Baer et al.<sup>32</sup> gives reliable estimates of the relative “ring-current” intensities in regular annulenes, when an amplitude of 3.60 eV (*ca.* 348 kJ mol<sup>-1</sup>) is taken<sup>33</sup> for the harmonic potential that occurs in their calculations.<sup>15–19,34</sup>
- A “conjugation circuit” within a given Kekulé-structure of an arbitrary, planar, polycyclic, conjugated hydrocarbon is a circuit that consists entirely of alternating single- and double bonds.<sup>1,15,16,35</sup>
- If a particular ring lies entirely within a given “conjugation circuit” — even if no bond of that ring actually lies on the “conjugation circuit” itself — this ring shall be said to “participate” in that “conjugation circuit”.<sup>36</sup>
- The “ring-current” intensity in any particular ring of such a polycyclic hydrocarbon receives a nonzero contribution from each “conjugation circuit” that occurs in all the various Kekulé-structures that can be devised for the molecule as a whole, *provided* that the ring in question *participates* in that “conjugation circuit”. These contributions are strictly additive. If the ring in question *does not participate* in a specific “conjugation circuit”, that particular “conjugation circuit” makes no contribution to the “ring-current” intensity in the ring under discussion.
- The nonzero contribution to the “ring-current” intensity in a given ring from an individual “conjugation circuit” comprising  $N$  bonds is equal to the “ring-current” intensity calculated (*via* (i), above) to be associated with a

**Table 1.** Data Needed for a Series of Idealized Annulenes (With Ring Sizes from [4] up to [22]) When Applying the Method Described in Refs 15 and 16 (Adapted with permission from ref 16. Copyright 1979 by The Portuguese Chemical Society.)

number of bonds ( $N$ )	ring area <sup>a</sup> ( $A_N$ )	ring-current intensity <sup>b,c</sup> ( $J_N$ )
4	0.385	-2.19
5	0.662	
6	1	1
7	1.399 <sup>d</sup>	
8	1.858	-1.27
10	2.962	+0.72
12	4.309	-0.69
14	5.902	+0.38
16	7.740	-0.38
18	9.823	+0.17
20	12.151	-0.2 <sup>e,f</sup>
22	14.724	+0.0 <sup>e,g</sup>

<sup>a</sup> Expressed (to three decimal places) as a ratio to the area of a standard benzene hexagon. For idealized regular [ $N$ ]-gons, all of uniform side length, it can be shown by elementary trigonometry that:

$$\left(\frac{\text{area of regular } [N]\text{-gon}}{\text{area of regular hexagon}}\right) = \left(\frac{N \cot(\pi/N)}{6 \cot(\pi/6)}\right)$$

This formula is the source of the figures listed in the middle column, above.

<sup>b</sup> Expressed (to two decimal places) as a ratio to the benzene ring-current intensity calculated, by the same method,<sup>32</sup> for benzene. Extrapolated values (see footnotes *e*, *f*, and *g* to this table) are given to fewer decimal places.

<sup>c</sup> Calculated by the free-electron, one-dimensional model of Baer et al.<sup>32,34</sup> with a periodic potential as described in rule (i) of the method of Gomes and Mallion,<sup>15,16</sup> presented in the text. <sup>d</sup> The value of 1.339 given in ref 16 is a misprint. <sup>e</sup> (Extrapolated) Baer et al.<sup>32</sup> did not report ring currents for [20]- and [22]-annulenes; because the present calculations require ring-current data for the [22]-annulene, the (virtually zero) value for it was estimated<sup>15,16</sup> by extrapolation. <sup>f</sup> In ref 16, this extrapolation was estimated to be -0.1, rather than the -0.2 estimated here. Such small differences are insignificant<sup>34</sup> as far as the ring-current estimates reported here are concerned because of the rare occurrence—indeed, nonoccurrence, in the present case—of circuits of size 20, and the relatively large number of Kekulé structures (9) and sets of conjugation circuits ( $9(9-1) = 72$ ) that are extant in the conjugated system under study ([567567]—see Figure 1). <sup>g</sup> This estimated, extrapolated ring-current contribution is virtually zero.<sup>34</sup> Again, our final results will not be sensitive to any errors that there might be in this extrapolation because only four of the 72 conjugation circuits that arise among the relatively large number (9) of Kekulé structures involved in the present calculation involve annulenic circuits of length [22].

model [ $N$ ]-annulene, except for a correction which takes into account the difference between the area of the model [ $N$ ]-annulene and the actual area of the “conjugation circuit” in question; in applying this correction it is assumed that the “ring currents” are proportional to the ring areas. Any one, specified, “conjugation circuit” contributes equally in this manner to the intensities of the “ring currents” in all the rings that participate in it.

- (vi) The relative “ring-current” intensity in a given ring may finally be obtained by averaging all such contributions (including the zero ones) over the total number of Kekulé-structures possessed by the complete molecule.

Consistent with the rule (v), the “ring-current” contribution due to the  $n^{\text{th}}$  conjugated circuit, of  $N$  sides and area  $A^{(n)}$ , is taken

to be proportional to the quantity  $J^{(n)}$ , where

$$J^{(n)} = J_N \left(\frac{A^{(n)}}{A_N}\right) \quad (1)$$

in which  $J_N$  and  $A_N$  are, respectively, the “ring-current” intensity and the ring area associated with an idealized, regular, planar [ $N$ ]-annulene [given in a table, reproduced, and slightly modified, in Table 1]. By rules (iv) and (vi), the relative “ring-current” intensity,  $J_r$  in a given ring,  $r$ , is then

$$J_r = \frac{1}{K} \sum_{\substack{\text{All "conjugation circuits"} \\ \text{in which ring participates}}} J^{(n)} \quad (2)$$

where  $K$  is the total number of Kekulé structures that may be devised for the molecule as a whole; the summation runs over all “conjugation circuits”,  $n$ , in which the ring  $r$  participates and all Kekulé structures are to be considered, one at a time.

**Discussion of the Gomes–Mallion Formalism and the More Recent Approaches.** Gomes and Mallion<sup>16</sup> concluded by giving a worked example of their “prescription” (to calculate the ring currents in naphthalene), and they then proceeded to apply it to a total of 15 structures and to compare the ring currents so calculated with those evaluated by what we nowadays refer to as the HLP<sup>23,24</sup> “topological” approach—with, overall, encouraging results.<sup>15,16</sup> If the above, 30-year-old description<sup>16</sup> is compared with the formulations recently presented by Randić,<sup>2</sup> by Mandado,<sup>3</sup> and by Ciesielski et al.,<sup>4</sup> it will be seen that the older theory ostensibly has the following features:

- (a) The method presented in ref 16 does make some attempt to take into account (by its rule (v) and its eq 1, above) the effect of differing ring areas—as, also, do the methods of Mandado<sup>3</sup> and Ciesielski et al.,<sup>4</sup> but the formalism of ref 2 does not. In the Mandado method,<sup>3</sup> the proper dependence of the current (as well as that of the resonance energy) on the size of the circuit is obtained by numerical fitting (to arrive at the parameter  $b = 2$ ), while Gomes and Mallion<sup>15,16</sup> separately use the circuit area (as it defines the magnetic flux) and the number of alternating single and double bonds (as this defines the quantum mechanical response of the electronic system).
- (b) As Gomes and Mallion stated:<sup>16</sup> “...the magnetic effect is taken to be proportional to the true area of the circuit — as indeed it is, both classically and in simple quantum-mechanical calculations.”<sup>5–7,14,22–25</sup> Since the external magnetic-field manifests itself in this phenomenon by means of magnetic fluxes through rings,<sup>5–7,14,22–25,29</sup> it is clear that any satisfactory account of it must recognize the influence of the areas of the different rings. For example, in structure 13 of ref 9 ([567567], shown in our Figure 1), five-membered, six-membered, and seven-membered rings lie side-by-side in the same molecule, and if they were isolated regular polygons of the same side length, their areas would vary<sup>24</sup> between about 66% (in the case of the five-membered rings) and about 140% (for the seven-membered rings) of the area of a standard benzene hexagon (see footnote *a* of Table 1). It should, however, be emphasized that, especially for the larger conjugation circuits, the actual areas are always smaller than the areas of the idealized, regular annulenes of the same perimeter.<sup>16</sup>

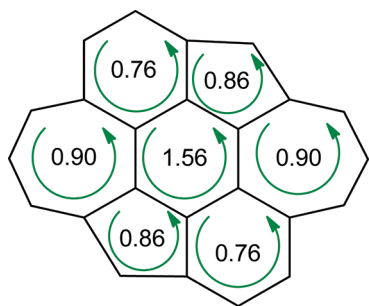
- (c) The method proposed in refs 15 and 16 attempts to differentiate contributions not just from  $[4n]$  and  $[4n + 2]$  circuits *per se*, but also for  $[4n]$  circuits with different values of  $n$  and for  $[4n + 2]$  circuits with different values of  $n$ . It does this partly by acknowledging the effect of differing ring areas and partly by incorporating into its founding tenants—by means of rule (i) of ref 16—the numerical values of annulenic ring-current intensities reported by Baer et al.<sup>32</sup> The methods described in refs 3 and 4 take the size of the conjugation circuit into account by recognizing the effect of the different *areas* of the several conjugation circuits, but the method of ref 2 does not.
- (d) By virtue of its rule (vi) and its eq 2, above, the method of Gomes and Mallion<sup>16</sup> effectively “normalizes” the final ring-current intensity by a division, at the very end of the arithmetical process of calculation, by the total number of Kekulé structures—as was also done by Gayoso<sup>12</sup> who, like Gomes and Mallion,<sup>16</sup> published independently on this subject in 1979; consequently, the resulting ring currents can easily be compared between one molecule and another—as, indeed, was done in refs 16–19, over a wide range of different conjugated systems. Furthermore—as is conventional and usually convenient, and also as in ref 16—ring currents so-calculated can easily be presented as *dimensionless quantities* (and hence as *pure numbers*) by the simple device of expressing them as a *ratio* to the ring-current intensity calculated, by the same method, for benzene (the ring-current intensity in which is, therefore, by definition, precisely 1). This conventional procedure<sup>5–7</sup> of expressing quantities as a ratio to benzene is also adopted by Gayoso<sup>12</sup> and by Mandado<sup>3</sup> but not by Randić in ref 2 nor—at least in the case of the bond currents—by Ciesielski et al.<sup>4</sup> (even though certain “local” and “global” quantities immediately calculated by Ciesielski et al.,<sup>4</sup> once they have obtained the computed bond currents, *are* themselves “normalized” by an appropriate division (in this case,  $(1/2)K(K - 1)$ , the number<sup>37</sup> of sets of distinct conjugation circuits<sup>38</sup>). In the approach of ref 2, however—and in the case of the bond currents calculated by the algorithm presented in ref 4—there appears to be no such averaging over all Kekulé structures or sets of conjugation circuits, and so not only are the units of measurement of the calculated  $\pi$ -electron currents not obvious but comparisons from one molecule to another would appear to be difficult.<sup>39</sup> We believe that this “averaging factor” is very important because of the physics that it conveys. In the old work of Gomes,<sup>17,19</sup> and in the recent work of Mandado,<sup>3</sup> it is not an averaging but a quantum-mechanical normalization factor; in our original work in ref 16, it was essentially an averaging factor. Accordingly, although it would be possible to compare a series of similar molecules without such “averaging” or “normalization”—for example,<sup>11</sup> the family of coronene and its 17 isomers that were studied by Balaban et al.<sup>9</sup>—it is not clear how a comparison would be made between the  $\pi$ -electron currents outside such a closely related series without doing so. As an illustration of this claim, we note that the method described in ref 2 would appear to give a  $\pi$ -electron current of size 2 for benzene, whereas, in ref 2,  $\pi$ -electron currents as high as 36, in these units, are reported for structure [567567] (Figure 1). Likewise, if

no “normalization” is done on the bond currents, the numbers presented by Ciesielski et al. in Figure 6 of ref 4 would seem to imply that a bond current more than 20 times the benzene ring current is extant in 3,4-benzopyrene.<sup>4</sup> It is appropriate at this stage to note that both Gomes and Mallion<sup>16</sup> (1979) and Gayoso<sup>12</sup> (also 1979) independently chose to divide, at the end of the calculation, by the number ( $K$ ) of Kekulé structures—this procedure later being justified by an application of valence-bond theory.<sup>17–19</sup> An essentially similar procedure was invoked by Mandado,<sup>3</sup> while Ciesielski et al.,<sup>4</sup> in their approach, have (as stated) opted (like Randić<sup>2</sup>) *not* to normalize their bond currents at all—though, as mentioned, they do immediately divide some “local” and “global” quantities calculated from the bond currents by the number<sup>37</sup>  $((1/2)K(K - 1))$  of sets of distinct conjugation circuits.<sup>38</sup> Finally, it should be noted that, in subsequent versions<sup>11,21</sup> of his basic method,<sup>2</sup> Randić and his collaborators *have* invoked “normalizations”—by dividing by  $K$  (for example, in ref 21), as did Gomes and Mallion, Gayoso,<sup>11</sup> and Mandado,<sup>3</sup> or by dividing by  $K(K - 1)$  (as, for example, in ref 11). For more, very recent, discussions on normalizing bond currents in molecules of different sizes, see refs 11, 21, and 40.

The recent methods of Randić<sup>2</sup> and of Ciesielski et al.<sup>4</sup> do have the aesthetic virtue of being entirely graph-theoretical in nature,<sup>41</sup> while the Mandado approach,<sup>3</sup> relying, as it does, on some parametrization, may be considered to be not purely graph-theoretical. By virtue of its rule (i), above, the approach of Gomes and Mallion<sup>16</sup> is likewise not purely graph-theoretical, either, for it borrows from quantum mechanics,<sup>32</sup> and, because of its rule (v) and its eq 1, above, it invokes what, on the face of it, is a non-graph-theoretical procedure<sup>29</sup> in an attempt to take ring areas<sup>41</sup> into account (but see ref 24 for an argument that ring areas should be treated as “topological”). This is why, at the time, the method presented in ref 16 was described as a “quasi-topological” one. Randić’s method<sup>2</sup> and that of Ciesielski et al.<sup>4</sup> do, though, stop at the graph-theoretical analysis; Gomes and Mallion<sup>16</sup> go further:

- (i) to use estimates of the conjugation increment that depend on its area—as also do the recent approaches of Mandado<sup>3</sup> and Ciesielski et al.<sup>4</sup>—as well as on a simple quantum-mechanical model calculation for the ring current in an annulene of appropriate size<sup>32</sup> (a feature that is not explicitly adopted by any of the modern authors<sup>2–4</sup> but the parametrization applied by Mandado<sup>3</sup> effectively serves the same purpose<sup>42</sup>) and
- (ii) to average over all Kekulé structures

The prescription described in refs 15 and 16 was later given a theoretical foundation by one of us (Gomes);<sup>17–19</sup> this was based on a simple valence-bond formalism with a nonempirical parametrization that was valid for resonance energies and magnetic ring currents. The same author has suggested<sup>43</sup> that this approach may be generalized in order to provide a more realistic description of currents outside the conventional “bond lines”. Mandado’s approach<sup>3</sup> is based on first-order response theory and, at the level of the formalism, is essentially equivalent to that of Gomes.<sup>19</sup> To ensure that the method was theoretically well grounded, Gomes<sup>19</sup> avoided parametrizations by fitting and using a simple quantum-mechanical model



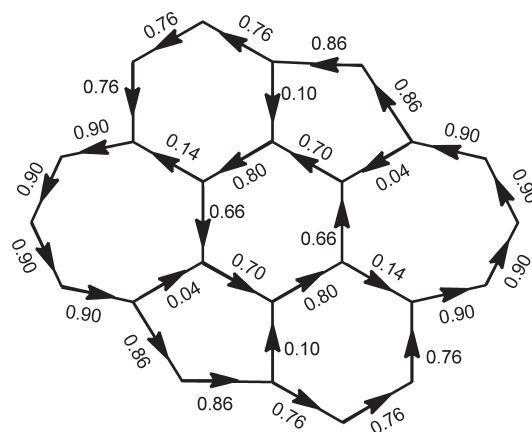
**Figure 2.** Ring currents in the conjugated system [567567], calculated by the method of Gomes and Mallion.<sup>15,16</sup>

calculation of the energies and magnetically induced currents of the annulenes.

**Application of the Gomes–Mallion Method<sup>16</sup> to [567567].** Ring currents calculated for the structure [567567] (Figure 1) by an application of the method of Gomes and Mallion<sup>16</sup> (in conjunction with the data given in Table 1, which is modified from refs 15 and 16) are presented in Figure 2. These diamagnetic currents are, by convention, considered to circulate in an anticlockwise sense around their respective rings (as indicated by the arrows in Figure 2). Because the annulene ring-current intensities of Baer et al.<sup>32</sup> are available to only two decimal places (Figure 1, right-hand column), ring-current (and, later, bond-current) intensities calculated using the Gomes–Mallion method<sup>16</sup> are quoted only to that accuracy, whereas such currents predicted by all the other methods dealt with in this study are quoted to three or more places of decimals. (Reporting data to apparently higher accuracy is in any case perhaps not entirely justified as we are here dealing with rather crude approaches to experimental observables.)

### 3. RING CURRENTS AND BOND CURRENTS

**General Considerations.** Randić's recent calculation<sup>2</sup> has presented  $\pi$ -electron currents not as ring currents but as bond currents. Likewise, the method of Ciesielski et al.,<sup>4</sup> for example, is also initially aimed at calculating bond currents (from which other quantities—both “local” and “global”—are subsequently computed). Therefore, in order conveniently to compare the  $\pi$ -electron currents calculated by these methods with the predictions of other approaches, we shall deal here with bond currents, as well as ring currents. In the literature, over the course of many years, consideration has been given overwhelmingly to ring currents rather than bond currents.<sup>5–7,14,22–25</sup> It does not, however, seem to be widely emphasized that, for any method of calculation that guarantees the applicability of Kirchhoff's Law<sup>26–29</sup> for conservation of currents at a junction (in a classical, macroscopic electrical network), *the two representations are entirely equivalent.*<sup>6</sup> By analogy with the theory of such macroscopic classical networks, the “ring current” in a conjugated molecule is the microscopic analog of the “loop current” (e.g., refs 6, 27, and 28) in a macroscopic “Kirchhoff” network,<sup>26</sup> while the “bond current”, considered as a “line current”<sup>43</sup> along the bond, is the microscopic analog of the current in a wire that constitutes a single branch of the (macroscopic) Kirchhoff network<sup>26</sup> in question.<sup>29</sup> This idea has been evaluated by two of us (Gomes and Mallion) in a review.<sup>6</sup> However, in a molecular context, it was originally discussed—with bond currents being regarded classically as “line currents”—by Longuet-Higgins and



**Figure 3.** Bond currents in the conjugated system [567567], calculated by the method of Gomes and Mallion.<sup>15,16</sup> These bond currents—which are entirely consistent with the ring currents presented in Figure 2—obey Kirchhoff's Law<sup>26–29</sup> of conservation of currents at a junction.

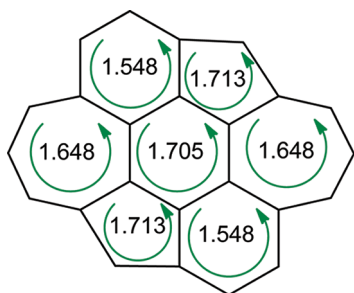
Salem,<sup>44</sup> some 50 years ago, and since then, it has been capitalized upon from time to time by several authors (e.g., refs 45–47). Adopting ring currents may be said to have the advantage of allowing a description by means of a set of independent numbers (or variables), while bond currents are related among themselves by Kirchhoff's Law.<sup>26–29</sup> The sum of currents coming out of a given junction is the (algebraic) sum of all those going into it; one of the currents involved in that junction is, therefore, not independent of the others. However, if—as here and in refs 2, 4, and 9—there is interest in the currents flowing in a particular bond, or around a certain region of the molecule, then bond currents may be more informative; in order to obtain the bond currents for shared bonds, ring currents in adjacent rings have (algebraically) to be added.

By no means do all theories of the ring-current effect, however, give rise to calculated bond currents that respect Kirchhoff's Law for conservation of currents at a junction; this law is violated, for example, in some quantum-mechanical approaches such as the “uncoupled Hartree–Fock” SCF one of Amos and Roberts<sup>48,49</sup> where, unlike in, for example, the HLP method<sup>23,24</sup> (in which “Hückel”-type assumptions<sup>50</sup> are made about neglect of non-neighboring interactions in the Hamiltonian), matrix elements between non-neighboring centers are, in general, nonzero.<sup>51</sup>

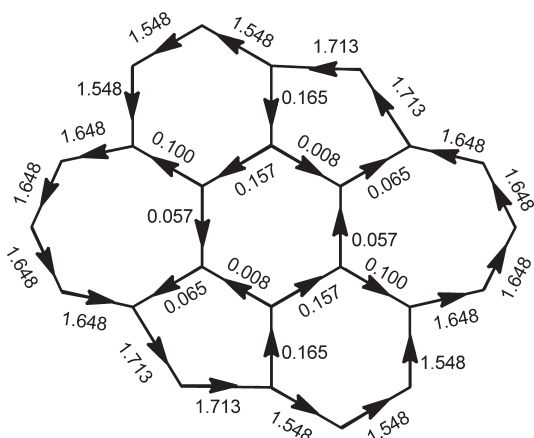
**Bond Currents by the Method of Gomes and Mallion<sup>16</sup> in the Structure [567567] (Figure 1).** Accordingly, as a result of our discussion above, the ring-current data presented in Figure 2 are shown in Figure 3—entirely equivalently—as bond currents. These will be discussed later (in section 4).

**“Topological” Ring Currents and (Entirely Equivalent) “Topological” Bond Currents in the Structure [567567] (Figure 1).** The idea of what one of us (Mallion<sup>23</sup>) has called “topological” ring currents (originally discussed informally, 35 years ago<sup>22</sup>) was only recently well-defined, initially for benzenoid hydrocarbons,<sup>23</sup> and, later, its definition was formally extended<sup>24</sup> to encompass conjugated systems containing rings of more than one size. We believe that the HLP method<sup>23,24</sup> is self-evidently the most appealing of the so-called “topological” approaches to the calculation of  $\pi$ -electron currents in conjugated systems because

- it is based on the well-established Hückel–London–Pople–McWeeny formalism,<sup>14,23,24</sup> and it is, thereby, legitimately founded on sound physics and quantum mechanics, and yet



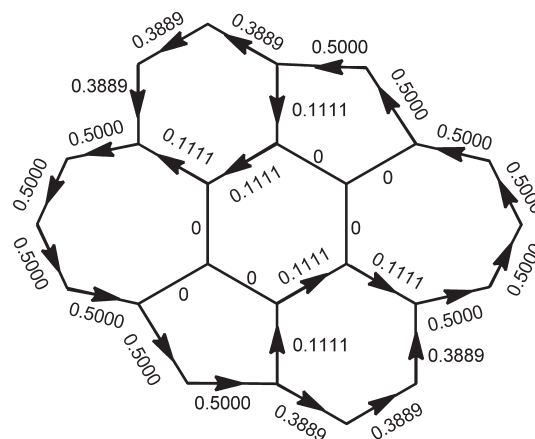
**Figure 4.** “Topological” ring currents in the conjugated system [567567], calculated by the HLP method.<sup>23,24</sup>



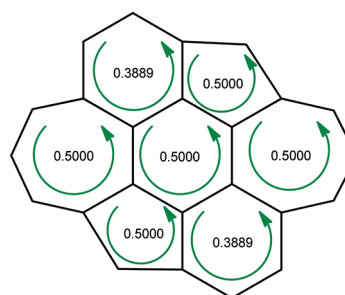
**Figure 5.** “Topological” bond currents in the conjugated system [567567], calculated by the HLP method.<sup>23,24</sup> These currents—which are entirely consistent with the ring currents presented in Figure 4—obey Kirchhoff’s Law<sup>26–29</sup> of conservation of currents at a junction.

- (b) it has all the advantages of a graph-theoretical approach because, once the carbon–carbon connectivity of the conjugated system under study has been written down, and values are agreed for its ring areas, then ring-current and bond-current intensities calculated by the HLP method<sup>23,24</sup> do not depend on any empirical (or, indeed, on any other) parameters, provided (as their definition requires<sup>23,24</sup>) that such ring- and bond-current intensities are expressed as a ratio to the corresponding quantity, calculated by the same method, for benzene.

The isomer of coronene that is under study, the structure [567567] illustrated in Figure 1, is a system of the appropriate type containing, as it does, five-, six-, and seven-membered rings. In Figure 4, therefore, we present topological ring currents (“loop currents”<sup>27,28</sup> on the macroscopic classical-network analogy described earlier) for [567567]. As is conventional (and as was done in Figure 2), the (diamagnetic) ring currents are presented as circulating anticlockwise around the rings that are their respective domains. Their relative intensities have been calculated by the HLP method on the detailed assumptions carefully specified in refs 23 and 24; the same assumptions were also recently adopted for calculations on a family of benzo-annelated perylenes.<sup>25</sup> In accordance with the required definition of what constitutes “topological” ring currents,<sup>23,24</sup> ring areas were calculated according to the formula given in footnote *a* of Table 1.



**Figure 6.** “Normalized”  $\pi$ -electron bond currents in the structure [567567], calculated using the method of Randić.<sup>2</sup> These currents obey Kirchhoff’s Law<sup>26–29</sup> of conservation of currents at a junction.



**Figure 7.** “Normalized”  $\pi$ -electron ring currents in the structure [567567], deduced from the normalized bond currents shown in Figure 6, calculated using the method of Randić.<sup>2</sup> These currents, which are conventionally defined in the anticlockwise direction around each ring, obey Kirchhoff’s Law<sup>26–29</sup> of conservation of currents at a junction.

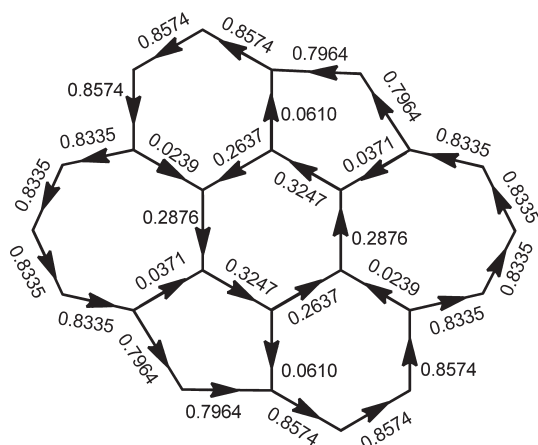
In Figure 5, we present these same “topological” ring-current data but, this time, broken down into individual bond currents, as previously described.

These “topological” bond currents will be discussed later (in section 4).

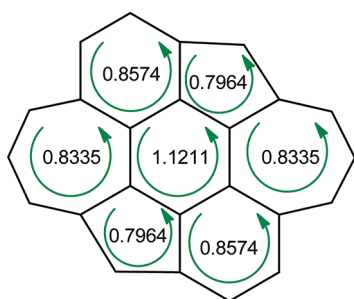
**Randić’s Bond Currents and Ring Currents in Structure [567567] (Figure 1).** In Figure 6, we display Randić’s bond currents<sup>2</sup> for later discussion and comparison with the bond currents already encountered and with those calculated by the methods of Mandado<sup>3</sup> and of Ciesielski et al.,<sup>4</sup> to be reported later in the paper. In order to facilitate such comparisons, we have taken the liberty of “normalizing” the “raw” (integral) bond currents depicted by Randić in Figure 5 of ref 2 by dividing them by  $K(K - 1)$  ( $= 72$ , in this case, as, here,  $K = 9$ )—as Randić and co-workers<sup>11</sup> themselves did in later work—and rounding the results to four decimal places.

Figure 7 shows Randić’s data equivalently presented as ring currents.

**Bond Currents and Ring Currents in Structure [567567] (Figure 1) by the Method of Mandado<sup>3</sup>.** In the initially submitted version of this paper, we did not report a calculation using Mandado’s method<sup>3</sup> and an anonymous reviewer very kindly supplied us with one for the structure [567567]. In order to have an independent check on the data that the reviewer had provided, we asked Professor P. W. Fowler and his colleagues

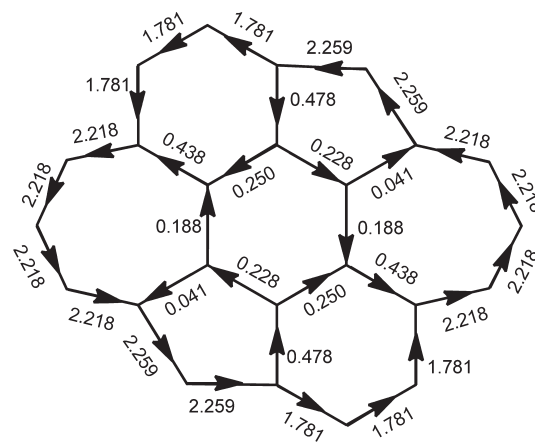


**Figure 8.**  $\pi$ -electron bond currents in the structure [567567], calculated using the method of Mandado<sup>3</sup> (reproduced here by the kind permission of Professor P. W. Fowler<sup>52</sup>). These currents obey Kirchhoff's Law<sup>26–29</sup> of conservation of currents at a junction.

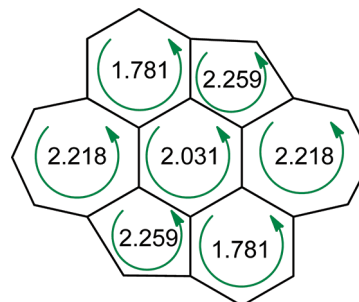


**Figure 9.**  $\pi$ -electron ring currents in the structure [567567], deduced from the bond currents shown in Figure 8, calculated using the method of Mandado<sup>3</sup> (reproduced here by the kind permission of Professor P. W. Fowler<sup>52</sup>). These currents, which are conventionally defined in the anticlockwise direction around each ring, obey Kirchhoff's Law<sup>26–29</sup> of conservation of currents at a junction.

W. Myrvold, W. Bird, and S. Cotton at the Universities of Sheffield (England) and Victoria (Alberta) to perform a calculation on [567567] using the Mandado<sup>3</sup> method. The results that they obtained are somewhat different from the bond currents provided by the reviewer. What we present in Figure 8 are the results of Professor Fowler et al. (used with his kind permission<sup>52</sup>)—effected by Mandado's model<sup>3</sup> (with Mandado's parameter “a” = 1, as is appropriate for this structure—see ref 3) rather than the reviewer's data. We adopt this policy because we are sure of the provenance of Professor Fowler's calculations (which, furthermore, we know have been effected automatically, by application of a computer algorithm, rather than by hand). It should be pointed out that Mandado's parametrization, designed for benzenoid structures,<sup>3</sup> might not be entirely appropriate for structures (like [567567]—Figure 1) that contain rings of other sizes. (We may observe in passing here that, although he does not stress it in ref 3, we feel that one of the strengths of Mandado's approach<sup>3</sup> is that magnetic susceptibilities and currents are produced in parallel, without the need for independent parametrizations.) The calculations presented in Figures 8 and 9 were effected adopting ring areas calculated according to the formula quoted in footnote *a* of Table 1. Bond currents (relative to benzene) are presented in Figure 8, and the ring currents that have been



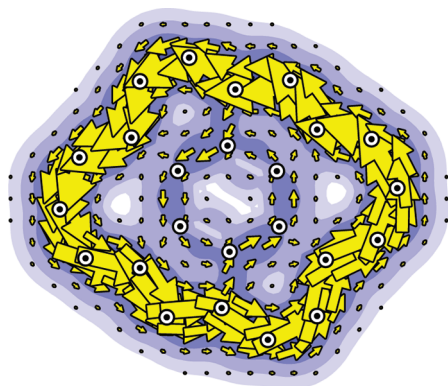
**Figure 10.** Bond currents (additionally—and against the prescription proposed in ref 4—divided by  $K(K - 1)$ ) in the conjugated system [567567], computed from the “raw” (unnormalized) bond currents calculated using the method of Ciesielski et al.<sup>4</sup> These currents—which are entirely consistent with the ring currents presented in Figure 11—obey Kirchhoff's Law<sup>26–29</sup> of conservation of currents at a junction.



**Figure 11.** “Normalized”  $\pi$ -electron ring currents in [567567], using the method of Ciesielski et al.<sup>4</sup> These ring currents, which are conventionally defined in the anticlockwise direction around each ring, have been deduced from the calculated bond currents depicted in Figure 10; they obey Kirchhoff's Law<sup>26–29</sup> of conservation of currents at a junction (within the round-off error displayed).

deduced from them are depicted in Figure 9. The numerical values quoted in these Figures differ only slightly—and, for the purposes of our discussion in this paper, not significantly—from the computations offered by the anonymous reviewer. The data in Figures 8 and 9 will be discussed later (in section 4).

**Bond Currents and Ring Currents in Structure [567567] (Figure 1) Using the Method of Ciesielski et al.<sup>4</sup>** We have also applied the third of the recent methods, that of Ciesielski et al.,<sup>4</sup> to the structure [567567]. Bond currents are presented in Figure 10, and the ring currents that have been deduced from them are depicted in Figure 11. It should be noted that these bond currents have here been “normalized” by division by  $K(K - 1)$  ( $= 72$ , in this case, as  $K = 9$ ) and—as a reviewer has pointed out—should not, therefore, strictly be called “Ciesielski et al.” currents at all; however, we take this small liberty for reasons of comparability between different molecules and different methods of calculation, as explained elsewhere in this paper. We emphasize that each bond current was separately and independently calculated by application of the formalism of Ciesielski et al.,<sup>4</sup> thereby enabling many independent checks to be made on the computations and verifying, in actual practice, that Kirchhoff's



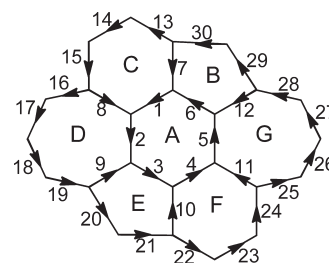
**Figure 12.** Pseudo- $\pi$  current-density map for [567567], calculated by use of the maximum-symmetry B3LYP/6-31G\*\* geometry,<sup>9</sup> reported in ref 9 (and reproduced here by the kind permission of Professor P. W. Fowler and the Slovenian Chemical Society). Maps of this kind show the current densities  $1a_0$  above the molecular plane after they have been projected *into* the molecular plane; as such, the projected current densities do not necessarily respect Kirchhoff's Law<sup>26–29</sup> of conservation of currents at a junction<sup>54</sup> because they are just one component of the real current. Therefore, although they are visually appealing, these current-density maps are not directly comparable with the “bond currents”—which are regarded strictly as classical “line-currents”<sup>26–29,44–47</sup>—that are being considered elsewhere in this paper.

First Law<sup>26–29</sup> does indeed hold within the context of Ciesielski et al.'s method.<sup>4,53</sup> The required ring areas were again calculated by means of the expression in footnote *a* of Table 1. The data in Figures 10 and 11 will be discussed later (in section 4).

#### 4. COMPARISON OF $\pi$ -ELECTRON CURRENTS FROM THE FIVE METHODS STUDIED WITH THE CURRENT DENSITY MAP OF BALABAN ET AL.<sup>9</sup>

**Overall Approach.** In this section, we compare the calculated current densities in the *ab initio* (“ipso-centric”<sup>9</sup>) pictorial current-density map for the structure [567567] (Figure 1), due to Balaban et al.<sup>9</sup> (Figure 12), with the following five sets of quantities calculated for the same structure:

- Ring currents (Figure 2), and bond currents deduced from them (Figure 3), computed using the 1976/1979 “quasi-topological” method of Gomes and Mallion,<sup>15,16</sup> which invokes the concept of “conjugation circuits”<sup>1,15–19</sup>
- “Topological” ring currents (Figure 4) evaluated by the HLPMP approach<sup>23,24</sup> and the bond currents that have been deduced from them (Figure 5)
- Bond currents (Figure 6) and the ring currents (Figure 7) that are consistent with them, calculated using Randić's recent, purely graph-theoretical, method<sup>2</sup>—which, like the methods of Gomes and Mallion,<sup>16</sup> Gayoso,<sup>12</sup> Mandado,<sup>3</sup> and Ciesielski et al.,<sup>4</sup> is also based on the idea of “conjugated circuits”<sup>1,8,11,15–19,21</sup>
- Bond currents calculated (Figure 8) using the method of Mandado<sup>3</sup> and ring currents deduced from them (Figure 9), on the assumption of Kirchhoff's Laws of current conservation at a junction<sup>26–29</sup>
- Bond currents calculated (Figure 10) using the method of Ciesielski et al.<sup>4</sup> and ring currents deduced from them (Figure 11), on the assumption of Kirchhoff's Laws. The method of Ciesielski et al.<sup>4</sup> is also based on conjugation



**Figure 13.** Labeling of rings (A–G), labeling of bonds (1–30), and (arbitrary) definition of bond directions (indicated by the direction of the arrow on each bond) in the structure [567567].

circuits, and may likewise be considered to be graph-theoretical even though (unlike in the method of Randić,<sup>2</sup> but as in the HLPMP “topological” approach and that of Mandado<sup>3</sup>) it does take into account the effect of ring areas.<sup>41,42</sup>

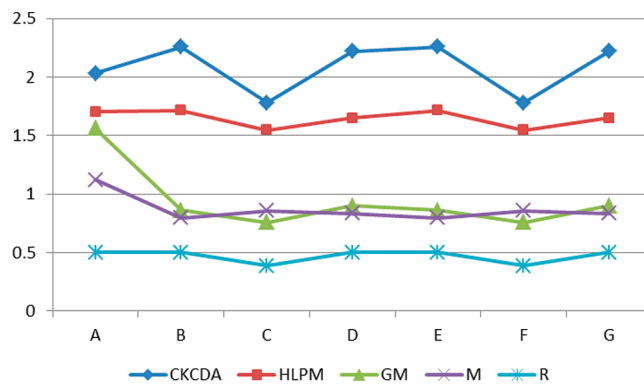
Qualitatively, the current-density map of Balaban et al.<sup>9</sup> (Figure 12) indicates a strong diamagnetic current around the periphery, with minimal activity, so far as currents are concerned, in the region of the central six-membered ring, A (which is completely surrounded by the six other rings—two five-membered (B and E), two six-membered (C and F), and two seven-membered (D and G)—see Figure 13). Very weak currents are also predicted in what Balaban et al.<sup>9</sup> refer to as the “spokes” bonds, connecting the central six-membered ring to the periphery of the structure (these are the bonds labeled 7, 8, 9, 10, 11, and 12 in Figure 13). It should be emphasized that the current-density maps under discussion<sup>9</sup> do not directly allow a quantitative estimate of the relative size of the bond currents in the molecule. In fact, Balaban et al. calculate<sup>9</sup> a current-density field—as described in the caption to Figure 12—and do not purport to be evaluating a bond current. The latter could in principle be effected by, for example, use of the integration technique of Atkins and Gomes<sup>55</sup>—though, to do this, knowledge would be needed of the current field (vector) over the surface that was used for the integration.

The current-density pattern of Figure 12 is well reproduced by the  $\pi$ -electron bond currents (Figure 6) that Randić<sup>2</sup> has calculated for this structure. It is also faithfully reflected in the patterns of the HLPMP “topological” bond currents displayed in Figure 5. Comparison of Figures 5, 6, and 12 shows that all three of these models predict that, within the central six-membered ring (A) itself, although the overall circulation around this ring is weak, a (relatively) stronger (diamagnetic) current is apparent in the bonds to the “northwest” (bond 1 of Figure 13) and to the “southeast” (bond 4 of Figure 13) in that ring (when the structure is depicted in the orientation shown in Figures 1–13) than in the other four bonds (2, 3, 5, and 6) in the central six-membered ring (ring A of Figure 13). The methods of Ciesielski et al.<sup>4</sup> (Figure 10) and of Gomes and Mallion<sup>4</sup> (Figure 3) likewise concur—though much less markedly—about the “northwest” and “southeast” bonds in the central ring, while, by contrast, the Mandado<sup>3</sup> approach suggests that the “northeast” (bond 6 in Figure 13) and the “southwest” bonds (bond 3) have the largest bond currents in the central ring (A).

Overall, the Gomes–Mallion approach<sup>15,16</sup> and that of Mandado<sup>3</sup>—contradicting the other methods<sup>2,4,23,24</sup>—predict a ring-current intensity in the central ring, A (Figure 13), that is considerably stronger than that in the rings around the perimeter.



Because of this, only partial cancellation takes place in the bonds that are shared by the central six-membered ring (A) and the outer rings (B–G). As a result, reasonably substantial diamagnetic  $\pi$ -electron currents are still predicted, by these methods, to be extant in all of the bonds of that central six-membered ring (Figures 3 and 8). In the case of the Gomes–Mallion method,<sup>15,16</sup> this observation is undoubtedly explained almost entirely by the fact that this latter method incorporates into its foundations the quantum-mechanically calculated annulene ring currents of Baer et al.<sup>32</sup> Inevitably, therefore, this method, as it stands at the moment, does intrinsically have built into it the phenomenon that contributions from both  $[4n]$  and  $[4n + 2]$  circuits decrease rapidly as  $n$  becomes larger—to the extent that, by the time circuits of size  $[20]$  and  $[22]$  are encountered, the ring-current contribution, according to the calculations of Baer et al.,<sup>32</sup> is virtually zero (see Table 1). The ring areas, of course, become larger as the conjugation circuits increase in length—it is the quantum-mechanically estimated ring current in the  $[4n]$  or  $[4n + 2]$  conjugation circuit that, according to Baer et al.,<sup>32</sup> shrinks rapidly as  $n$  becomes larger (see Table 1, right-hand column). Furthermore, the conjugation circuits with the longer lengths (lengths of  $[16]$ ,  $[18]$ ,  $[20]$ , and  $[22]$ ) that arise in the course of a calculation on  $[567567]$  invariably “contain” (that is, enclose within them—in the sense of rule (iii) of the Gomes–Mallion method described in section 2—and thus make contributions to the ring currents in) some or all of the peripheral rings (rings B–G in Figure 13). Inevitably, therefore, those rings eventually accumulate smaller calculated ring-current intensities than they would if the larger circuits all contributed equally (as they do in Randić’s method<sup>2</sup>) or if the contributions from conjugation circuits increased solely in proportion to their area (as they do, for example, in the method of Ciesielski et al.<sup>4</sup>). Although the central ring (A) is also “contained within” (see rule (iii) of section 2) these larger conjugation circuits—and thus, in the Gomes–Mallion approach, this ring A likewise receives diminished contributions from these larger circuits—actually carrying out such a calculation by hand (as we have)<sup>56</sup> shows that by far the greatest effect on the ring current in the central ring, A, arises from many conjugation circuits that include that ring as a circuit of length  $[6]$ . In fact, eight of the 72 sets of conjugation circuits (including the disjoint ones<sup>2,4</sup>) for  $[567567]$  involve a contribution to that central ring from a  $[6]$ -membered circuit; each of these (as can be seen from Table 1) makes a relatively large contribution (of 1, in these units) to the calculated ring current in ring A. Hence, this rapid diminution of annulenic ring current for the conjugation circuits of larger size would exaggerate the ring current associated with the central ring (A) and underestimate the ring currents for the peripheral rings (B–G). This would give rise to the prediction of a substantially greater current circulation in the six bonds of ring A, as a result of only partial cancellation of the (smaller) ring currents in adjacent rings. This finding could thus possibly be merely an artifact of the Gomes–Mallion method, as originally formulated.<sup>15,16</sup> To investigate this point, we have carried out some simple “topological” calculations, using the HLPM method,<sup>23,24</sup> on the ring currents in the family of  $[4n+2]$ -annulenes and have found—to our surprise—that, far from decreasing with annulene size, they actually increased quite dramatically. This matter will be the subject of future study. Meanwhile, it should be observed in passing that the Mandado approach also predicts a relatively large ring-current intensity in the central ring. We have not investigated why this should be so, but we speculate that it might be connected (a) with the fact that the Mandado method<sup>3</sup> does not take into account contributions from disjoint conjugation circuits and (b) by virtue of its initially assigning

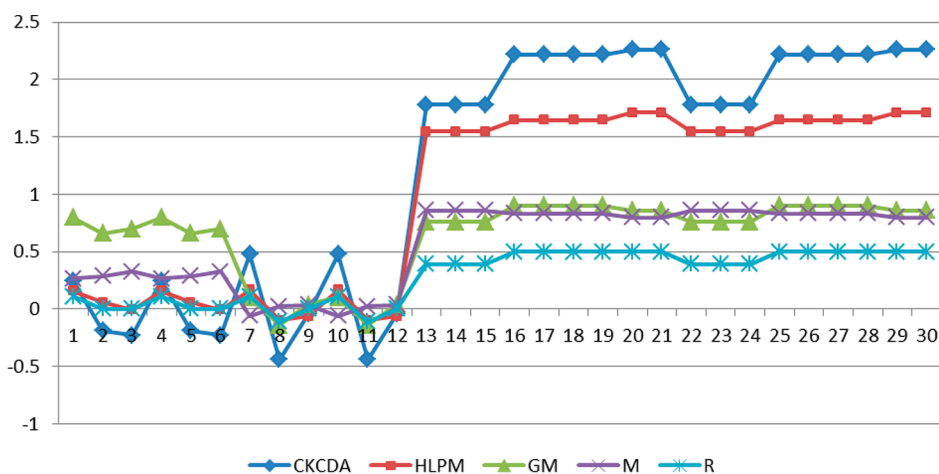


**Figure 14.** Comparator diagram for ring currents. Key to methods of calculation: CKCDA = Ciesielski et al.<sup>4</sup> (but “normalized”, contrary to the prescription of Ciesielski et al.,<sup>4</sup> by division by  $(1/2)(K(K - 1))$ ); HLPM = Hückel–London–Pople–McWeeny (“topological”);<sup>23,24</sup> GM = Gomes–Mallion;<sup>16</sup> M = Mandado;<sup>3</sup> R = Randić<sup>2</sup> (but “normalized” by division<sup>55</sup> by  $K(K - 1)$ ). The horizontal axis refers to the seven rings, A–G, labeled as in Figure 13. The vertical axis gives ring-current intensities expressed, effectively, as a ratio to the benzene ring current, calculated by the corresponding method; the ring currents may thus be regarded as dimensionless quantities.

a weighting to conjugation-circuit contributions that varies as the reciprocal of the ring area (rather than being proportional to conjugation-circuit ring areas, as in the formalisms of Gomes and Mallion<sup>15,16</sup> and of Ciesielski et al.<sup>4</sup>); the consequence might be that the Mandado method<sup>3</sup> likewise contrives to minimize the importance of the contributions from the larger conjugation circuits.<sup>42</sup> Mandado<sup>3</sup> takes  $H_{LL} = H_0 / (af_L)^b$ . For  $b = 2$ , the energy becomes independent of the area ( $f$ ), and the current becomes proportional to  $1/f$ . In the Gomes–Mallion method, the dependence is different, but the current also decreases when the area increases.

**Analysis of “Comparator” Diagrams for Ring Currents and Bond Currents Calculated by the Five Methods Studied.** In this section, we use “comparator diagrams” in order visually to compare trends in ring currents and in bond currents using the five methods that we have studied. We first define (by means of Figure 13) the ring labelings, the bond labelings, and the bond directions adopted in the comparator diagrams themselves (Figures 14 and 15). The seven rings of the structure  $[567567]$  are labeled A–G, and the 30 bonds are labeled 1–30 (as in Figure 13). Diamagnetic (that is, on our conventions, *positive*) ring currents are defined as running anticlockwise around the ring in question; bonds are defined in the directions of the arrows depicted in Figure 13. These directions are arbitrary. If the net current calculated for a given bond, by any of the five methods applied, is in the direction of the arrow shown in Figure 13, it is counted positive; if against the direction of the arrow, it is negative. With these conventions, we present the comparator diagrams for the ring currents (Figure 14) and the bond currents (Figure 15) calculated for the structure  $[567567]$  using the five methods<sup>2–4,16,23,24</sup> that we have studied. The reader is directed to the captions of Figures 14 and 15 for explanations about the axes, scales, and units that feature in these comparator diagrams, and for the key to the five methods of calculation that have been considered.

The feature discussed in the last paragraph of the previous subsection is immediately and strikingly illustrated by the green curve in the ring-current comparator diagram (Figure 12), which



**Figure 15.** Comparator diagram for bond currents. Key to methods of calculation: CKCDA = Ciesielski et al.<sup>4</sup> (but “normalized”, contrary to the prescription of Ciesielski et al.,<sup>4</sup> by division by  $(1/2)(K(K - 1))$ ); HLP = Hückel–London–Pople–McWeeny (“topological”),<sup>23,24</sup> GM = Gomes–Mallion;<sup>16</sup> M = Mandado;<sup>3</sup> R = Randić<sup>2</sup> (but “normalized” by division<sup>55</sup> by  $K(K - 1)$ ). The horizontal axis refers to the 30 bonds, labeled 1–30 as in Figure 13, with their directions defined as in that figure. The vertical axis gives bond-current intensities, as dimensionless quantities, effectively expressed as a ratio to the bond-current intensity calculated, by the corresponding method, for benzene.

concerns the Gomes–Mallion (GM) method, and by the purple curve representing the Mandado<sup>3</sup> (M) method. The central ring, A, bears a ring-current intensity materially greater than those in the peripheral rings, B–G, whose ring currents are themselves noticeably smaller in size than the corresponding ones calculated either by the HLP “topological” method<sup>23,24</sup> (the brown curve) or by the method of Ciesielski et al.<sup>4</sup> (CKCDA, the dark-blue curve). The (light-blue) curve from the Randić<sup>2</sup> method is more attenuated and less variable, but it generally follows the pattern of the HLP and CKCDA curves. However, with the exception of that central ring, A, the relative pattern of variation of (smaller) ring currents in the peripheral rings (B–G) that is observed along the green (GM) curve in Figure 14 does follow quite closely those of at least three of the other methods studied (HLP, R, and CKCDA).

In turning to a consideration of the comparator diagram for bond currents (Figure 15), we recall (Figure 13) that the bonds labeled 1–6 are those comprising the central six-membered ring, A; those labeled 7–12 are what Balaban et al.<sup>9</sup> call the “spoke” bonds, connecting the perimeter to the entirely internal, central ring, A (Figure 13); and the bonds labeled 13–30 are those that lie around the periphery of the structure [567567]. When assessing the bond currents in the bonds (1–6) situated in the central ring (A) and those in the so-called<sup>9</sup> “spoke” bonds (7–12), it should be borne in mind that these bond currents (a) are small and (b) are the result of “cancellation” (by subtraction) of two much larger, but approximately equal, quantities—the ring currents in the two adjacent rings that flank any of these bonds that are labeled 1–12. Sometimes this process of cancellation results in a small positive current in the (arbitrary) direction in which the bond in question has been defined in Figure 13; sometimes it results in a small negative one—and, furthermore, this is the situation for *each* of these 12 bonds *and* for *each* of the five methods (HLP, GM, R, M, and CKDA) that we have applied. Any correspondences among the five methods in these regions (i.e., those involving bonds 1–12) of the bond-current comparator diagram (Figure 15) are, therefore, difficult to discern visually—though it can be seen that the correspondence in variation between curve R (light blue) and curve CKCDA (dark blue)

is in fact close, even in this region. However, we arbitrarily chose to define all of the (unshared) peripheral bonds, 13–30 (Figure 13), in the *same* direction as the diamagnetic (*i. e.*, anticlockwise) ring currents in the rings of which these bonds form a part. Examination of *this* area of the bond-currents comparator diagram (Figure 15)—that for bonds 13–30—reveals an entirely consistent pattern of trends among the five methods. (Recall that, in comparator diagrams, the *pattern* of variation is what counts.) Once again, for the reasons discussed in the last paragraph of the previous subsection, the GM method and method M (as well as method R) predict much lower bond-current intensities in these peripheral rings than do the HLP and CKCDA approaches.

## 5. UNITS IN BOND-CURRENT AND RING-CURRENT CALCULATIONS: “NORMALIZATION”

Before concluding, we draw attention to the *units* to be adopted when various types of  $\pi$ -electron currents are presented. The bond currents evident in Figures 3, 5, 6, 8, and 10 and the ring currents depicted in Figures 2, 4, 7, 9, and 11 are, effectively, all expressed as a *ratio* to the ring-current/bond-current intensity calculated, by the corresponding method, for benzene; accordingly, currents calculated in this way are, as already noted, *dimensionless* quantities, with the benzene value being identically 1, by definition. This conventional approach was followed by Gomes and Mallion,<sup>14</sup> Gayoso,<sup>12</sup> and Mandado<sup>3</sup> and was followed in the definition of “topological ring current” in the context of the HLP formalism.<sup>22–25</sup> This, however, does not appear to be the case with Randić’s  $\pi$ -electron bond currents presented in Figure 5 of ref 2. These are not expressed as a ratio to the corresponding value for benzene, and neither are they (unlike in the Gomes–Mallion “conjugation-circuit” method<sup>16</sup> and in that of Gayoso<sup>11</sup>) “normalized” by averaging the contributions to each  $\pi$ -electron bond-current over all Kekulé structures ( $K$ ) that can be devised for the conjugated system as a whole. Because of this, it is not clear in what *units* such  $\pi$ -electron bond currents are actually expressed. The same criticism could strictly be leveled at the “raw” bond currents calculated by

the method of Ciesielski et al.<sup>4</sup> and displayed, for example, in Figure 6 of ref 4. In their calculation, although they do at a later stage “normalize” by dividing by  $(1/2)(K(K - 1))$ , this (as a reviewer pointed out to us) is *not* in fact done at the stage of the calculation when the actual *bond currents* themselves are calculated. Consequently, if the method of Ciesielski et al. were taken at its face value, as presented in ref 4, the initial, “raw” bond currents (like those arising from Randić’s initial formulation<sup>2</sup>) would likewise be difficult to interpret. Another consequence of this lack of division by the number of Kekulé structures (or, if preferred, the number of sets of conjugation circuits) is that it would appear, as a general rule, that the greater the number of Kekulé structures a system has, the larger its calculated  $\pi$ -electron bond currents are likely to be—by a process of sheer accumulation—when estimated by the methods described in refs 2 and 4. Because of this, it is not obvious how diverse types of conjugated systems (such as, for example, the range of structures treated in ref 16) can be compared, one with the other, on this model. For example, in [567567], a  $\pi$ -electron bond current of “36” is encountered (Figure 6). Now, as already noted, if the approach of ref 2 were applied to benzene, the  $\pi$ -electron current would be calculated to be “2”, on these units. It would clearly be unreasonable to deduce from this—and, indeed, we emphasize that Randić<sup>2</sup> does not claim to do so—that there is, in some of the peripheral bonds in [567567], a  $\pi$ -electron current that is 18 times the size of the  $\pi$ -electron current in benzene (which, of the face of it, seems unlikely). An analogous comment could be made about the “raw” bond currents in 3,4-benzopyrene reported in Figure 6 of ref 4. Nevertheless, these examples do illustrate the difficulties, in the context of the method described in refs 2 and 4, that can arise when comparisons between different types of conjugated structures are sought.

In any case, these problems may easily be averted, in the case of the methods of refs 2 and 4, if bond currents are simply “normalized”, by a suitable division. This is why we have taken the liberty of normalizing our bond currents and ring currents calculated by the methods of Randić<sup>2</sup> (in Figures 6 and 7) and of Ciesielski et al.<sup>4</sup> (in Figures 10 and 11). In fairness, it ought to be noted that, in subsequent versions of his method,<sup>11,21,40</sup> Randić (and his co-workers) have performed a “normalization” process, either by dividing<sup>21</sup> by the number ( $K$ ) of Kekulé structures (as Gomes and Mallion,<sup>15,16</sup> Gayoso,<sup>12</sup> and Mandado<sup>3</sup> do) or by dividing,<sup>11</sup> not by  $K$ , but by  $K(K - 1)$ , the total number<sup>37</sup> of sets of conjugation circuits—or, if preferred  $(1/2)(K(K - 1))$ , the number of *distinct* sets of conjugation-circuits.<sup>40</sup>

## 6. CONCLUSIONS

- (a) We have personally repeated (and successfully reproduced) the calculations on [567567] presented by Randić<sup>2,58</sup> and the calculations on 3,4-benzopyrene reported by Ciesielski et al.<sup>4</sup> and can thereby verify that these methods do have considerable elegance and aesthetic appeal. As Randić points out,<sup>2</sup> and as we have noted previously, his approach has the philosophical virtue of being entirely graph-theoretical in nature. The method of Ciesielski et al.<sup>4</sup> is also purely graph-theoretical—if it is accepted (as, indeed, we do propose<sup>24</sup>) that ring areas can legitimately be considered as part of a graph-theoretical “prescription”.<sup>22–24,41,42</sup> We suggest, however, that, for maximum efficacy in practical applications, the conjugated-circuit approach outlined in ref 2—and, to a lesser extent, that presented in ref 4—would benefit from

- (i) Averaging the final computed  $\pi$ -electron currents by dividing at the end by the total number of Kekulé structures (as was done in refs 16, 12, and 21) or by the number of conjugation circuits (as was done in ref 11). This would aid comparability between diverse molecules (such as, for example, [567567]—Figure 1—and benzene), and it would go some way toward solving the vexed problem of *units*, discussed earlier in the context of Randić’s method<sup>2</sup> and that of Ciesielski et al.<sup>4</sup> Randić has dealt now with this point.<sup>11,21,40</sup>
- (ii) Weighting the contributions of individual conjugation circuits according to the actual (or, failing this, the idealized<sup>22–25,41</sup>) areas of the rings that lie within them
- (iii) Taking account of the fact that, quite apart from the area factor,<sup>41,42</sup> just referred to in (ii), above, not all conjugated circuits should be considered to contribute to bond currents to an equal extent. Randić’s approach does allow distinction between the diamagnetic contributions from  $[4n + 2]$  circuits and the paramagnetic ones arising from  $[4n]$  circuits; however, there would appear to be no provision in the method of ref 2 for specifying that, for example, a  $[4n + 2]$  circuit with, say,  $n = 3$ , should contribute differently from one with, say,  $n = 4$ . Likewise, the (paramagnetic) contribution appropriate for a  $[4n]$  circuit with (say)  $n = 2$  is different from that properly due to a  $[4n]$  circuit with (say)  $n = 4$ —but there is no mechanism for taking this into account in the method described in ref 2. This problem is partially considered in ref 4 by a consideration of ring areas,<sup>41</sup> but in the recent methods,<sup>2–4</sup> no further account is taken of the variation in annulenic ring currents<sup>32</sup> for annulenes of different sizes<sup>32</sup>—though, in ref 3 this effect does seem successfully to be mimicked by an appropriate parametrization.
- (b) If the provisions suggested in (a), above, were to be adjoined to the methods of Randić<sup>2</sup> and of Ciesielski et al.,<sup>4</sup> the end result would be something very similar to the old prescription of Gomes and Mallion.<sup>13,14</sup> Furthermore, the method of ref 2 is equivalent to that proposed in ref 16 if we
- (i) assume that the contributions from all conjugation circuits are equal and
- (ii) omit the last stage of averaging (or normalizing) over all Kekulé structures.
- The elegance of purely topological methods<sup>2,4</sup> and the simplicity of their calculation should thus be evaluated against the advantages of bringing in physical considerations through
- (i) the dependence of a magnetic effect on the circuit area
- (ii) the dependence—much discussed in the classical literature<sup>59</sup>—of the *size* (and not just the *sign*) of the annulenic ring current on the number of carbon atoms forming the ring.
- The method of Ciesielski et al.<sup>4</sup> does the former (i) but not the latter (ii).
- (c) Randić<sup>2</sup> pointed out that his  $\pi$ -electron (bond) currents for [567567] (Figure 1) compare favorably with the qualitative current-density maps presented by Balaban et al.<sup>9</sup> The same can be said for the HLP “topological” bond currents that we have illustrated in Figure 5. The calculations (presented in Figure 3) that were obtained *via* Gomes and Mallion’s 1979 method<sup>16</sup> based on “conjugation circuits” do,

however, only partially support this view, though they do concur with the other approaches—HLPM bond currents<sup>23,24</sup> (Figure 5), Randić bond currents<sup>2</sup> (Figure 6), Mandado's bond currents<sup>3</sup> (Figure 8), Ciesielski et al.'s bond currents<sup>4</sup> (Figure 10) and the current-density map of Balaban et al.<sup>9</sup> (Figure 12)—that, in the case of the structure [567567] (Figure 1) which could, perhaps, be thought of as a “perturbed [18]-annulene”),<sup>60,22</sup> the strongest current does flow around its perimeter (see also Figure 15).

- (d) Regarding the methods of Randić,<sup>2</sup> Mandado,<sup>3</sup> Ciesielski et al.,<sup>4</sup> and Gomes and Mallion,<sup>16</sup> we note the following similarities, differences, and comparisons:
- All four methods rely on knowledge of the conjugation circuits<sup>1,12,15,16,20,35</sup> in the structure under study.
  - The approaches of Gomes and Mallion,<sup>16</sup> Randić,<sup>2</sup> and Ciesielski et al.<sup>4</sup> explicitly consider *disjoint* conjugation circuits (as illustrated, for example in Figure 2 of ref 2 and Figure 4 of ref 4), and the method of Mandado<sup>3</sup> excludes<sup>35</sup> them (see, for example, ref 3 and the Supporting Information of that reference).
  - The methods of Mandado,<sup>3</sup> Ciesielski et al.,<sup>4</sup> and Gomes and Mallion<sup>16</sup> rely on knowledge of the various *ring areas* of the structure—but Randić's method<sup>2</sup> does not. Refs 4 and 16 do, however, incorporate consideration of the areas of conjugation circuits very differently from ref 3—see, for example, ref 42.
  - The method of Gomes and Mallion<sup>16</sup> requires “external” knowledge of the ring-current intensities in the family of [*N*]-annulenes, calculated using a quantum-mechanical method,<sup>32</sup> based on a one-dimensional cyclic model with a periodic potential, in order to mimic the nuclear positions by the troughs of the potential; Mandado's approach<sup>3</sup> requires a suitable parametrization.<sup>42</sup> (It should be noted the HLPM formalism<sup>23,24</sup>—the approach that we favor as being the least subjective of all of the methods considered here—requires for its application knowledge only of the carbon–carbon connectivity of the structure in question *and* the areas of its constituent rings.)
- (e) Finally, we remark that this study has demonstrated that consideration of bond currents, as distinct from the more traditional ring currents, can give an extra conceptual insight into the magnetic properties of conjugated systems like [567567] (Figure 1). This is evident from the detailed, semiquantitative deductions that we have been able to make from these computations when the information about the calculated  $\pi$ -electron currents is displayed in the form of individual bond currents, as it is in Figures 3, 5, 6, 8, and 10, rather than as ring currents (as in Figures 2, 4, 7, 9, and 11). As is self-evident—though the point is not often emphasized—both are rigorously equivalent representations in the case of any method that respects Kirchhoff's Law of current conservation at a junction.<sup>26–29</sup> Nevertheless, provided that charge/current conservation is guaranteed—or Kirchhoff's Law is valid for bond currents instead of the more-general current densities—ring currents do represent a more efficient way of describing the molecular reaction to the external magnetic field: ring currents are independent, while bond currents are not.<sup>61</sup>

## AUTHOR INFORMATION

### Corresponding Author

\*E-mail: tkd25@cam.ac.uk.

## ACKNOWLEDGMENT

We are extremely grateful to two anonymous reviewers who very kindly provided some calculations of their own using the methods of Mandado<sup>3</sup> and Ciesielski et al.<sup>4</sup> Numerical disagreements were kindly adjudicated and resolved by Professor P. W. Fowler and his colleagues W. Myrvold, W. Bird, and S. Cotton at the Universities of Sheffield (England) and Victoria (Alberta), who generously and independently checked those calculations. They also helpfully checked by computer the calculations (reported in Figures 3 and 10, respectively) that we ourselves had effected (by hand) using application of the Gomes–Mallion method and of the method due to Ciesielski et al. We also thank Professor Fowler, and the Slovenian Chemical Society, for generously allowing us to reproduce (in Figure 12) the current-density map of the structure [567567] (originally published in ref 9), and The Portuguese Chemical Society for kind permission to reproduce and adapt material from ref 16; we also acknowledge Professor M. Randić for his kindness in providing several preprints before their actual publication.

## REFERENCES

- (1) (a) Randić, M. *Chem. Phys. Lett.* **1976**, *38*, 68–70. (b) Randić, M. *Tetrahedron* **1977**, *33*, 1905–1920. (c) Randić, M. *J. Am. Chem. Soc.* **1977**, *99*, 444–450. (d) Randić, M. *Pure Appl. Chem.* **1980**, *52*, 1587–1596.
- (2) Randić, M. *Chem. Phys. Lett.* **2010**, *500*, 123–127.
- (3) Mandado, M. *J. Chem. Theory Comput.* **2009**, *5*, 2694–2701.
- (4) Ciesielski, A.; Krygowski, T. M.; Cyrański, M. K.; Dobrowolski, M. A.; Aihara, J.-I. *Phys. Chem. Chem. Phys.* **2009**, *11*, 11447–11455.
- (5) Haigh, C. W.; Mallion, R. B. Ring current theories in nuclear magnetic resonance. In *Progress in Nuclear Magnetic Resonance Spectroscopy*; Emsley, J. W., Feeney, J., Sutcliffe, L. H., Eds.; Pergamon Press: Oxford, England, United Kingdom, 1979/1980; Vol. 13, pp 303–344.
- (6) Gomes, J. A. N. F.; Mallion, R. B. The concept of ring currents. In *Concepts in Chemistry*; Rouvray, D. H., Ed.; Research Studies Press Limited: Taunton, Somerset, England, United Kingdom, 1997; John Wiley & Sons, Inc.: New York, 1997; Chapter 7; pp 205–253.
- (7) (a) Lazzarotti, P. Ring currents. In *Progress in Nuclear Magnetic Resonance Spectroscopy*; Emsley, J. W., Feeney, J., Sutcliffe, L. H., Eds.; Elsevier: Amsterdam, 2000; Vol. 39, pp 1–88. (b) Gomes, J. A. N. F.; Mallion, R. B. *Chem. Rev.* **2001**, *101*, 1349–1383.
- (8) Randić, M. *Chem. Rev.* **2003**, *103*, 3449–3605; especially pp 3539–3543.
- (9) Balaban, A. T.; Bean, D. E.; Fowler, P. W. *Acta Chim. Slov.* **2010**, *57*, 507–512.
- (10) See Figure 7 of Ciesielski et al.,<sup>4</sup> which depicts a  $\pi$ -electron current-density map of 3,4-benzopyrene that had been provided by Professor Fowler.
- (11) Randić, M.; Nović, M.; Vračko, M.; Vukičević, D.; Plavšić, D. *Int. J. Quantum Chem.* In press. DOI: 10.1002/qua.23081.
- (12) Gayoso, J. C. R. *Hebd. Séances Acad. Sci.* **1979**, *C 288*, 327–330.
- (13) A reviewer has suggested that we should explicitly quote Gayoso's “philosophical” conclusion in his 1979 paper.<sup>12</sup> Gayoso wrote: “*En relevant, sans passer la mécanique quantique, toute la richesse descriptive de la formule développée, la théorie chimique des graphes est selon nous à revivifier l'activité théorique en chimie.*” After translation into English, this reads approximately as follows: *Observing the wealth of descriptive formulation developed without quantum mechanics, we think that chemical graph theory is set to revive activity in theoretical chemistry.*

(14) (a) Pople, J. A. *Mol. Phys.* **1958**, *1*, 175–180. (b) McWeeny, R. *Mol. Phys.* **1958**, *1*, 311–321. (c) Veillard, A. *J. Chim. Phys. Phys. Chim.—Biol.* **1962**, *59*, 1056–1066. (d) Gayoso, J.; Boucekine, A. C. R. *Hebd. Séances Acad. Sci.* **1971**, *C272*, 184–187. (e) Aihara, J.-I. *J. Am. Chem. Soc.* **1979**, *101*, 558–560. (f) Mizoguchi, N. *Bull. Chem. Soc. Jpn.* **1987**, *60*, 2005–2010. (g) O’Leary, B.; Mallion, R. B. *J. Math. Chem.* **1989**, *3*, 323–342. (h) It should be noted that, despite its authorship, ref 14g is actually a historical paper reporting previously unpublished work (dating from 1953) by the late Professor C. A. Coulson, F.R.S., in which the London theory<sup>14i</sup> of diamagnetic susceptibility in conjugated systems is couched in terms of Coulson’s own contour-integral formalism.<sup>14j</sup> (i) London, F. *J. Phys. Radium (7<sup>e</sup> Sér.)* **1937**, *8*, 397–409. (j) Coulson, C. A. *Proc. Cambridge Philos. Soc.* **1940**, *36*, 201–203.

(15) Gomes, J. A. N. F. *Some Magnetic Effects in Molecules*. D. Phil. Thesis, University of Oxford (Linacre College): England, United Kingdom, 1976; pp 69–73.

(16) Gomes, J. A. N. F.; Mallion, R. B. *Rev. Port. Quim.* **1979**, *21*, 82–89.

(17) Gomes, J. A. N. F. *Mol. Phys.* **1980**, *40*, 765–769.

(18) Gomes, J. A. N. F. *Croat. Chem. Acta* **1980**, *53*, 561–569.

(19) Gomes, J. A. N. F. *Theoret. Chim. Acta* **1981**, *59*, 333–356.

(20) The term “conjugation circuit” was adopted<sup>15,16</sup> in order to emphasize the implication that the said circuits are considered as agents for the “transmission of conjugation”, rather than the fact that the circuits themselves are actually “conjugated”, as such—as is implied by the name (“conjugated circuit”) that Randić (independently) coined<sup>1</sup> for the same idea.

(21) Randić, M.; Plavšić, D.; Vukičević, D. *J. Indian Chem. Soc.* **2011**, *88*, 13–23.

(22) Coulson, C. A.; Mallion, R. B. *J. Am. Chem. Soc.* **1976**, *98*, 592–598.

(23) Mallion, R. B. *Croat. Chem. Acta* **2008**, *81*, 227–246.

(24) Balaban, A. T.; Dickens, T. K.; Gutman, I.; Mallion, R. B. *Croatica Chem. Acta* **2010**, *83*, 209–215.

(25) Dickens, T. K.; Mallion, R. B. *J. Phys. Chem. A* **2011**, *115*, 351–356.

(26) Kirchhoff, G. *Anal. Phys. Chem.* **1845**, *64*, 497–514.

(27) Cundy, H. M. *S.M.P. Advanced Mathematics Book 3 [Metric]*; Cambridge University Press: London, 1970; p 911.

(28) Smart, D. *Linear Algebra & Geometry*; Cambridge University Press: Cambridge, England, United Kingdom, 1988; S. M. P. Further Mathematics Series, pp 307–308.

(29) Mallion, R. B. *Proc. R. Soc. London, Ser. A* **1974/1975**, *341*, 429–449.

(30) Even more neglected is the pioneering work of Gayoso<sup>12</sup>—published in 1979, independently of the paper by Gomes and Mallion<sup>16</sup> (also 1979). In an attempt to apply Randić’s theory of conjugation-circuit increments<sup>1</sup> to estimate magnetic-susceptibility exaltations in conjugated hydrocarbons, Gayoso<sup>12</sup> used eight fitted parameters and applied his formalism to 19 test molecules; he obtained semiquantitative agreement with the experimental data. Gayoso’s contribution<sup>12</sup> has, however, been overlooked entirely by all three schools<sup>2–4</sup> reporting recent work on the relevance of conjugation circuits to the magnetic properties of conjugated systems.

(31) Note that the wording of rule (iii)—here printed in italics—has been modified; please see also ref 36.

(32) Baer, F.; Kuhn, H.; Regel, W. *Z. Naturforsch. A* **1967**, *22*, 103–112.

(33) The reasons for adopting this approach are explained in the Introduction to ref 16.

(34) We have frequently emphasized in the text that, as part of its intrinsic canon, the method of Gomes and Mallion<sup>15,16</sup> adopts the relative ring-current intensities calculated by Baer et al.<sup>32</sup> for the [N]-annulenes of a size up to [18] (and extrapolations of them beyond that size—please see footnotes *e* and *f* to Table 1). The magnitudes of these decrease very rapidly to zero as ring size increases (Table 1, right-hand column). The idea that contributions should decrease with the size of the conjugation circuit was also evident in Randić’s original study<sup>1</sup> on

resonance energies. Such a decline of contributions with the size of the conjugation circuit is also a feature of Mandado’s method<sup>3</sup> because, in that approach, contributions from conjugation circuits are initially weighted according to the reciprocals of their areas. Gayoso’s parametrization<sup>12</sup> in the context of magnetic susceptibilities did not, however, support a monotonic decrease with the size of the conjugation circuit. Furthermore, we have carried out some simple “topological” calculations, using the HLPM method,<sup>23,24</sup> on the ring currents in the family of [4n+2] annulenes and have found that, far from decreasing with annulene size, they actually *increased*—and quite dramatically. This aspect will be the subject of future investigation.

(35) As Randić<sup>2</sup> and Ciesielski et al.<sup>4</sup> explicitly point out—and as follows directly from the theory proposed by one of us (Gomes)<sup>17–19</sup>—full account should be taken of any *disjoint* conjugation circuits that are extant within the Kekulé structures that are under consideration—see, for example, Figure 2 of ref 2 and Figure 4 of ref 4. According to ref 3, and the Supporting Information to it, such disjoint conjugation circuits appear to be omitted from consideration in the method of Mandado.<sup>3</sup> Furthermore, the inclusion of disjoint conjugation circuits was neither explicitly prescribed nor yet expressly proscribed in the initial formulation of Gomes and Mallion’s 1979 approach<sup>16</sup> (this did not, in any case, affect the final numerical results for the majority of the structures initially studied in ref 16, in which no such disjoint conjugation circuits actually arise). For small molecules, simple inspection of the set of Kekulé structures (“perfect matchings”) allows identification of all of the required conjugated circuits.<sup>16</sup> In general, however, we recommend examining all of the overlaps between all pairs of Kekulé structures;<sup>19</sup> application of this algorithm will guarantee that sets of disjoint conjugation circuits are automatically identified and retained for consideration in the bond-current/ring-current calculations.

(36) (a) In view of the theoretical formulation of the method proposed by Gomes,<sup>17–19</sup> this rule—printed in the text in italics—is here given a wording that is more general than that adopted for the corresponding rule (iii) in ref 16. (The original wording in ref 16 was “(iii) If at least one bond of a particular ring forms a part of a given “conjugation circuit”, this ring shall be said to “participate” in that “conjugation circuit.”). Professor P. W. Fowler<sup>36b,c</sup> has recently proposed an equivalent version of rule (iii) that is couched entirely in terms of *bond* currents, rather than *ring* currents: (b) Fowler, P. W. Personal Communication, June 19th, 2011. (c) Fowler, P. W.; Myrvold, W. *J. Phys. Chem.* DOI: jp-2011-06548t.R1.

(37) Gutman, I.; Randić, M. *Chem. Phys.* **1979**, *41*, 265–270.

(38) (a) An anonymous reviewer has pointed out that a coefficient,  $c_L$ , is defined in Mandado’s approach<sup>3</sup> as the ratio of the number of conjugated circuits of a certain type,  $L$ , and the number of Kekulé structures ( $n_K$  in Mandado’s paper)—please see ref 3 for a definition of these quantities. This ratio emerges from the valence-bond treatment of the resonance energy proposed in ref 3, which starts with the valence-bond wave function proposed by Herndon,<sup>38b</sup> in which the root of the number of Kekulé structures is introduced in order to normalize the wavefunction (equation 1 in Mandado’s paper—see also ref 19). It seems that Mandado effectively “normalizes” by dividing by  $K/2$  (2 being the number of Kekulé structures in benzene). Regarding the approach of Ciesielski et al.,<sup>4</sup> we are unclear why the so-called “local” and “global” quantities calculated from the bond currents are normalized (by division by  $(1/2)K(K-1)$ , the number of sets of distinct conjugation circuits) but *not* the actual bond currents themselves—which, of course, are what we ourselves are primarily interested in in this paper: (b) Herndon, W. C. *J. Am. Chem. Soc.* **1973**, *95*, 2404–2406.

(39) If, in the context of ref 3, one were to parametrize or simulate a quantum-mechanical method, one would more naturally fit the quantity  $H_{LL}$  to the number of bonds in the circuit, rather than to the number of rings.

(40) (a) Even more elaborate “normalizations” have been discussed and argued for by Randić<sup>40b</sup> in the context, specifically, of structures related to perylene: (b) Randić, M. Personal communication to R.B.M., May 3, 2011.

(41) The method of Ciesielski et al.<sup>4</sup> does involve ring *areas*, but it has been claimed<sup>22</sup> that consideration of such areas, if properly effected

according to certain well-defined rules<sup>24</sup> (as specified in footnote *a* of Table 1), need not necessarily deprive an approach of being aptly described as “topological” or “graph-theoretical”.<sup>22–25</sup>

(42) (a) The method of Gomes and Mallion<sup>16</sup> takes into account the effect of the number of centers on the ring-current intensity by incorporating the numerical values of the annulenic ring-current intensities reported by Baer et al.<sup>32</sup> An anonymous reviewer has pointed out that this is also taken into consideration, in an equivalent way, in Mandado’s method.<sup>3</sup> In fact, Mandado introduces a variable  $f_L$  defined as the number of benzene rings enclosed by the circuit, in order to represent the size of the circuit. In the method of Gomes and Mallion,<sup>15–19</sup> the size of the circuits is measured in idealized molecular geometries with equal bond lengths; the energy and the current associated with a annulene are assumed to depend on the number of centers of that annulene. These energy terms depend mainly on the number of alternating double and single bonds, which, in turn, depend on the number of centers. Gomes and Mallion<sup>15,16</sup> adopted the annulenic ring currents of Baer et al.<sup>32</sup> as an alternative to a possibly arbitrary parametrization. Mandado<sup>3</sup> has instead opted for a parametrization. The following is a summary<sup>42b</sup> of how all four “conjugation circuits” approaches<sup>15,16,2–4</sup> count the contributions of conjugation circuits as they increase in size. In both the Gomes–Mallion<sup>15,16</sup> and Mandado<sup>3</sup> models, contributions taper off<sup>42b</sup> with ring size—they start off large for small conjugation circuits, settle down, and then virtually disappear altogether at  $[N] = 22$ . In the Randić<sup>2</sup> model, all conjugation circuits, of whatever size, give an equal contribution per occurrence. In the approach of Ciesielski et al.,<sup>4</sup> larger conjugation circuits contribute more per occurrence, as they have larger area. In Mandado’s model,<sup>3</sup> the larger conjugation circuits contribute less per occurrence, because, in this formulation, the area is in the denominator. In the approach of Gomes and Mallion,<sup>15,16</sup> although circuit area is in the numerator, contributions from the larger conjugation circuits fall off rapidly with the size of the conjugation circuit because of the fact that the Baer et al.<sup>32</sup> “annulenic” currents decrease virtually to zero at annulene size [22]. This phenomenon of the decrease in importance of the contributions from the larger conjugation circuits is thus essentially the same in both the Mandado<sup>3</sup> and Gomes–Mallion<sup>15,16</sup> formalisms: (b) Fowler, P. W. Personal communication, July 3, 2011.

(43) Gomes, J. A. N. F. *THEOCHEM* **1990**, *210*, 111–119.

(44) (a) Longuet-Higgins, H. C.; Salem, L. *Proc. R. Soc. London, Ser. A* **1960**, *257*, 445–456. (b) Salem, L. *The Molecular Orbital Theory of Conjugated Systems*; W. A. Benjamin: Reading, MA, 1966; Chapter 4.

(45) (a) Haddon, R. C. *Tetrahedron* **1972**, *28*, 3613–3633. (b) *idem ibid.* 3635–3655.

(46) Mallion, R. B. *Mol. Phys.* **1973**, *25*, 1415–1432.

(47) (a) Mallion, R. B. *Empirical Appraisal and Graph Theoretical Aspects of Simple Theories of the “Ring-Current” Effect in Conjugated Systems*. D. Phil. Thesis: University of Oxford (Christ Church): England, United Kingdom, 1979; pp 124–131. (b) Haigh, C. W.; Mallion, R. B. *Croat. Chim. Acta* **1989**, *62*, 1–26.

(48) Amos, A. T.; Roberts, H. G. F. *Mol. Phys.* **1971**, *20*, 1073–1080.

(49) (a) Kirchhoff’s Law of current conservation<sup>26–29</sup> is, by contrast, not violated in, for example, the “coupled Hartree–Fock” procedure described in: (b) Coulson, C. A.; Gomes, J. A. N. F.; Mallion, R. B. *Mol. Phys.* **1975**, *30*, 713–732.

(50) (a) Coulson, C. A.; O’Leary, B.; Mallion, R. B. *Hückel Theory for Organic Chemists*; Academic Press: London, 1978. (b) Yates, K. *Hückel Molecular Orbital Theory*; Academic Press: New York, 1978.

(51) The reasons why such methods lose current conservation at junctions are gone into in some detail by one of the present authors (R.B.M.) in a long footnote on page 1420 of ref 46.

(52) Fowler, P. W. Personal communication to R. B. M., June 8, 2011.

(53) An anonymous reviewer has, very conscientiously, repeated our calculations using the method of Ciesielski et al.<sup>4</sup> and claimed that some of our bond currents reported in Figure 10 are incorrect. We therefore invoked the help of Professor P. W. Fowler and his colleagues W. Myrvold, W. Bird, and S. Cotton at the Universities of Sheffield (England) and Victoria (Alberta), in order to effect an independent check.<sup>51</sup> When

“normalized” by division by the factor  $K(K - 1)$ —against, however, the recommendations of Ciesielski et al.,<sup>4</sup> at this stage of the calculation—the bond currents in [567567] calculated by Fowler et al. using the method of Ciesielski et al.<sup>4</sup> agreed entirely with ours (displayed in Figure 10) to four significant figures and disagreed with those that had been provided by the reviewer. We therefore persist with our original data, confident that they have been independently confirmed by Professor Fowler and his above-named colleagues.

(54) Fowler, P. W.; Bean, D. E. Personal Communication to R. B. M. at the Fifth Conference on Computers in Scientific Discovery, Sheffield, England, United Kingdom, July 2010.

(55) Atkins, P. W.; Gomes, J. A. N. F. *Mol. Phys.* **1976**, *32*, 1063–1074.

(56) Although our calculations using the Gomes–Mallion method<sup>15,16</sup> were effected by hand, they were independently checked by means a computer algorithm written and run by Professor P. W. Fowler and his colleagues W. Myrvold, W. Bird, and S. Cotton at the Universities of Sheffield (England) and Victoria (Alberta). Complete agreement was found (to the number of decimal places quoted in Figures 2 and 3).

(57) Ring currents by method R are determined by the count of conjugation circuits contributing to them, arising from the  $K$  perfect matchings/Kekulé structures and, in Randić’s later work, are “normalized” by dividing either<sup>11</sup> (i) by the total number,<sup>37</sup>  $K(K - 1)$ , of pairwise overlaps or<sup>21</sup> (ii) merely by  $K$  itself. In more recent work, on perylenes, even more delicate and elaborate “normalizations” have been argued for.<sup>40</sup>

(58) It should be noted that, in Figure 1 of ref 2, the Kekulé structure labeled “E” is in fact—in error—merely a repetition of the one labeled “B”. Kekulé structure “E” should actually have been the one that is displayed there *but* with the single/double bonds in the *central* ring running in the alternative way (the patterns in the other six rings remaining undisturbed). Professor Randić has evidently used the *correct* Kekulé structure “E” in his calculations, for, on repetition of them, we independently agree with his final “bond currents”, displayed in Figure 5 of ref 2. (and—after “normalization”<sup>11</sup> by division by 72—in our Figure 6).

(59) (a) Pople, J. A.; Untch, K. G. *J. Am. Chem. Soc.* **1966**, *88*, 4811–4815. (b) Haddon, R. C.; Haddon, V. R.; Jackman, L. M. *Top. Curr. Chem.* **1971**, *16*, 103–220. (c) Sondheimer, F. *Acc. Chem. Res.* **1972**, *5*, 81–91.

(60) Trost, B. M.; Bright, G. M.; Frihart, C.; Brittelli, D. *J. Am. Chem. Soc.* **1971**, *93*, 737–745.

(61) (a) This observation is re-enforced by Figure 3 of ref 61b, in which the bond currents in the 60 bonds of kekulene are expressed in terms of just three independent parameters,  $A$ ,  $B$ , and  $\delta$ . These are not actually the ring currents in the three symmetrically non-equivalent rings of kekulene, but they are related to them in the following way: the ring current in the six-membered ring that Steiner et al.<sup>61b</sup> label “I” is equal to  $(A + \delta)$ ; the ring current in ring II is  $(A - \delta)$ , and that in the internal, 18-membered, ring III is  $(A - B)$ . Once again, therefore, in this example, a bond-current description is entirely equivalent to a ring-current one, with as many independent variables being needed to specify the several bond currents as there are symmetrically non-equivalent rings (and, hence, distinct ring currents) in the structure under consideration: (b) Steiner, E.; Fowler, P. W.; Jenneskens, L. W.; Acocella, A. *Chem. Commun.* **2001**, 659–660.

# MSINDO-sCIS: A New Method for the Calculation of Excited States of Large Molecules

Immanuel Gadaczek,\* Katharina Krause, Kim Julia Hintze, and Thomas Bredow

Institut für Physikalische und Theoretische Chemie, Universität Bonn, Wegelerstr. 12, 53115 Bonn, Germany

**S** Supporting Information

**ABSTRACT:** Theoretical background, parametrization, and performance of the semiempirical configuration interaction singles (CIS) method MSINDO-sCIS designed for the calculation of optical spectra of large organic molecules are presented. The CIS Hamiltonian is modified by scaling of the Coulomb and exchange integrals and a semiempirical correction. For a recently proposed benchmark set of 28 medium-sized organic molecules, vertical excitation energies for singlet and triplet states are calculated and statistically evaluated. A full reparameterization of the MSINDO method for both ground and excited state properties was necessary. The results of the reparameterized MSINDO-sCIS method are compared to the currently best semiempirical method for excited states, OM3-CISDTQ, and to other standard methods, such as MNDO and INDO/S. The mean absolute deviation with respect to the theoretical best estimates (TBEs) for MSINDO-sCIS is 0.44 eV, comparable to the OM3 method but significantly smaller than for INDO/S. The computational effort is strongly reduced compared to OM3-CISDTQ and OM3-MRCISD, since only single excitations are taken into account. Higher excitations are implicitly included by parametrization and an empirical correction term. By application of the Davidson–Liu block diagonalization method, high computational efficiency is achieved. Furthermore, it is demonstrated that the MSINDO-sCIS method correctly describes charge-transfer (CT) states that represent a problem for time-dependent density functional theory (TD-DFT) methods.

## 1. INTRODUCTION

During recent decades, there has been substantial progress in the experimental and theoretical characterization of electronically excited states. New experimental techniques provide better insight into photophysical processes in molecules.<sup>2,3</sup> On the other hand, more and more efficient and accurate theoretical methods for the description of excited states have been developed.<sup>4,5</sup> Multireference configuration interaction (MRCI),<sup>6</sup> multistate complete active-space second order perturbation theory (MS-CASPT2),<sup>7</sup> and coupled cluster methods (CC2,<sup>8</sup> CCSDT,<sup>9</sup> CC3<sup>10</sup>) are well-established and highly accurate but, at the same time, extremely costly. Recently published benchmark calculations demonstrate the advance of CCx<sup>11</sup> methods. Even with relatively small basis sets, they provide reasonable agreement with more elaborate ab initio approaches. Since these highly accurate methods are only applicable for small molecules, time-dependent density functional theory (TD-DFT)<sup>12</sup> has become one of the most popular methods for spectra prediction and the description of excited states due to its reliability at low computational cost in most cases. Yet the application of TD-DFT is at present limited to systems with a few hundred atoms. Although TD-DFT is an attractive choice for the description of excited states, there are well-documented problems,<sup>13,14</sup> in particular for charge-transfer (CT) states. Even hybrid methods and the perturbative corrected double-hybrid methods<sup>15</sup> do not give a quantitatively correct description. Moreover, the accuracy of TD-DFT is limited in general. Recently, Thiel et al.<sup>16</sup> have shown that the vertical excitation energies obtained from TD-DFT have mean absolute deviations of 0.3–0.5 eV from the theoretical best estimates (TBEs). This raises the question of how, e.g., novel organic solar cells, which normally consist of

large donor and acceptor species with many hundred atoms, can be treated theoretically, when no applicable method exists. For such problems, a reliable method which is also efficient is desirable. In this context, semiempirical methods have come back into focus for the description of excited states in larger molecules.<sup>17</sup> Due to the integral approximations made in all semiempirical methods, the computational effort is decreased by orders of magnitude<sup>18,19</sup> compared to first-principles methods. Therefore systems with up to thousands of atoms can be treated. The error introduced by the vast integral approximations is compensated by calibration of the method against experimental reference data. Most semiempirical methods are parametrized to reproduce ground-state properties, which does not guarantee that the excited-state properties are also reliable. The INDO/S (intermediate neglect of differential overlap for spectroscopy) method is well-known to provide accurate results for vertical excitation energies.<sup>20</sup> Therefore, this method is still frequently used for the calculation of optical absorption spectra for large organic molecules and transition-metal compounds.<sup>21</sup> Due to the well-balanced parameters, as developed for the INDO/S and the reparameterized INDO/S2<sup>22</sup> methods, the calculated excitation energies are in good agreement with experimental results. Although these methods seem to be superior over standard quantum-chemical methods, they have certain limitations. INDO/S makes use of the configuration interaction singles (CIS) method<sup>23</sup> for the treatment of excited states. The lack of higher excitations restricts the INDO/S-CIS method to states which are dominated by single excitations. This is problematic if a

Received: June 22, 2011

Published: September 21, 2011

mixing of double excitations is dominating the electronic transition. Due to the use of a minimal basis set, Rydberg states cannot be treated. The main reason why INDO/S is not used in photochemistry, however, is the design of the method. INDO/S—CI targets spectroscopy but is not parametrized for ground-state properties, and potential energy surfaces (PES) may not be reliable. Therefore, the application to the calculation of fluorescence spectra and photoreactions is limited. This has led to some more elaborate concepts, where higher excitations are included. It has been recently shown by Silva-Junior and Thiel that semiempirical methods, especially the OMx<sup>24–26</sup> CISDTQ and MRCISD<sup>27</sup> approaches, give promising results for the calculation of excited states.<sup>17</sup> Accurate results for the exploration of PES for both ground and excited states have been reported. Therefore, this is from our point of view the most reliable semiempirical method for the calculation of excited states at the moment. However, due to the calculation of NDDO-type two-electron–two-center integrals<sup>28,29</sup> and the inclusion of double, triple, and quadruple excitations, the computational effort is dramatically increased with respect to INDO/S, even though the implementation of the GUGA leads to a speedup in the calculation of the CI matrix elements.<sup>27</sup>

This leads to our starting point. In our present implementation, we do not go beyond the explicit calculation of singles excitations. This gives a method similar to the INDO/S case, but at the full CIS level, where all possible single excitations are taken into account. In our implementation, we take advantage of the orthogonalization correction in the MSINDO method,<sup>30</sup> which is similar to the OMx methods. In the parametrization, we include both ground and excited states properties. In order to account for cases where higher excitations play an important role, we introduced a semiempirical correction term.

In this article, we present the basic equations of our MSINDO-sCIS method and then show the results of our method for the benchmark set of Thiel. We will show that our method shows similar accuracy compared to the OMx-CISDTQ methods for most compounds.

## 2. METHOD

In this section, we describe the implementation of the CIS equations into the MSINDO method. The first part reviews the CIS theory and the implementation of our empirical correction. In the second part, we discuss the implementation of the method with an efficient matrix diagonalization, where all semiempirical considerations are introduced. In the last subsections, we will describe the basic ideas of our parametrization for ground and excited states, including charge-transfer states.

**2.1. Theory.** Within the CIS approach, the excited-state wave function is formed in terms of the Hartree–Fock orbitals. With  $N$  occupied and  $M$  virtual orbitals,  $N \times M$  determinants are created by interchanging all pairs of occupied and virtual orbitals. Linear combinations of these determinants form the CIS wave function of the excited states. For closed-shell systems, the CIS wave function can be formed in terms of a singlet or triplet configuration state function (CSF):

$$|{}^{1,3}\Psi_I\rangle = \sum_i^{\text{occ}} \sum_a^{\text{vir}} t_i^a |{}^{1,3}\Phi_a^i\rangle \quad (1)$$

The unique set of the expansion factors  $\{t_i^a\}$ —the CIS amplitudes—defines the wave function of the given excited state  $I$ . Variational solution of the Schrödinger equation with this wave

function leads to an eigenvalue problem, where the CIS matrix  $H$  has to be diagonalized:

$${}^M H_{ia,jb} = \langle {}^M \Phi_a^i | \hat{H} | {}^M \Phi_b^j \rangle \quad M = 1, 3 \quad (2)$$

Since the original MSINDO parameters have been optimized for the ground-state properties,<sup>30</sup> we introduced two additional parameters for an improved parametrization of the excited states. These parameters are scaling factors for the two-electron integrals in the CIS matrices ( $c_1$  and  $c_2$  for singlets,  $c_T$  for triplets):

$${}^1 H_{ia,jb} = \delta_{ij} \delta_{ab} [\varepsilon_a - \varepsilon_i - d_{ia}^{\text{corr}}] + 2c_1 (ia|jb) - c_2 (ij|ab) \quad (3)$$

Considering that higher excitations have a large influence on certain totally symmetric excited states, we improved the CIS energy using a semiempirical correction. The correction term  $d_{ia}^{\text{corr}}$  (eq 4) is an empirical correction for totally symmetric single excited states that may have a strong mixing with double excitations  $ii \rightarrow aa$  from the ground state.

$$d_{ia}^{\text{corr}} = |(aa|ia) + (ii|ia)| \quad (4)$$

It is well-known that perturbative doubles corrections to the excited CIS states—the so-called CIS(D) method<sup>31,32</sup>—give a significant improvement in the excitation energies. In principle, the ground state and all excited states should be corrected by electron correlation due to mixing with higher excitations. In CIS(D), this is approximated by second-order perturbation theory that takes into account double excitations with regard to ground and excited state determinants. The full implementation of doubles correction is an  $O(N^5)$  process, which would counteract the philosophy of semiempirical methods. While the ground state is implicitly correlated in semiempirical methods, it is not clear how much this is the case for excited states.<sup>33</sup> Only in totally symmetric excited states can double excitations of the type  $ii \rightarrow aa$  and also higher excitations based hereupon mix with the singly excited determinants, which leads to a lowering of the corresponding excitation energy. Since the corresponding coupling integrals are in general larger in absolute value than those of other excitations  $ij \rightarrow aa$ ,  $ii \rightarrow ab$ , and  $ij \rightarrow ab$ , which may belong to a different irreducible representation, a correction is more important for totally symmetric excited states. This is in line with the deviations of calculated excitation energies at the ab initio CIS(D) level compared to CC2, CCSD, and CC3. For this purpose, we performed ab initio CIS(D) calculations for *E*-butadiene as a test molecule with the optimized MSINDO structures and aug-cc-TZVP basis sets and compared them to previous results obtained at CCn level with the same basis set.<sup>34</sup> The largest deviation (0.9 eV for CIS(D)) was observed for the singlet  $A_g$  state, whereas singlet  $B_u$  and triplet  $A_g$  and  $B_u$  states coincide within 0.2–0.3 eV. In preliminary test calculations, it was found that this approach led to the best overall agreement with the TBES. Therefore, we decided to empirically correct only these kinds of excited states. Without the correction term (eq 4) we observed in some cases, that also the ordering of the excited states was incorrect.  $d_{ia}^{\text{corr}}$  is exactly zero if orbitals  $i$  and  $a$  do not belong to the same irreducible representation. It is not necessary to make a similar correction for triplet states due to spin selection rules.

$${}^3 H_{ia,jb} = \delta_{ij} \delta_{ab} [\varepsilon_a - \varepsilon_i] - c_T (ij|ab) \quad (5)$$



Solving eq 2 gives then the excitation energy of a given state (eqs 6 and 7):

$${}^1E_{\text{CIS}} - E_{\text{RHF}} = \sum_i^{\text{occ}} \sum_a^{\text{vir}} t_i^a \left\{ t_i^a (\varepsilon_a - \varepsilon_i - d_{ia}^{\text{corr}}) + \sum_j^{\text{occ}} \sum_b^{\text{vir}} t_j^b [2c_1(ia|jb) - c_2(ij|ab)] \right\} \quad (6)$$

$${}^3E_{\text{CIS}} - E_{\text{RHF}} = \sum_i^{\text{occ}} \sum_a^{\text{vir}} t_i^a \left\{ t_i^a (\varepsilon_a - \varepsilon_i) - \sum_j^{\text{occ}} \sum_b^{\text{vir}} t_j^b c_T(ij|ab) \right\} \quad (7)$$

This approach follows the original ideas of Grimme<sup>35</sup> in his DFT/SCI method, where a scaling of the integrals gave a significant improvement of the vertical excitation energies. Different from the expressions in ref 35 where only a single scale factor  $c = c_2 = c_T$  is applied and  $c_1 = 1$  is used as a constant, we treat all parameters  $c_1$ ,  $c_2$ , and  $c_T$  as adjustable parameters in order to improve the calculated excitation energies.

The empirical shift introduced for Rydberg states and core excited states<sup>35</sup> is not applied in our implementation, since our semiempirical method is not intended to be applied to these kinds of problems. The scaling in eqs 3 and 5 does not affect the Fock matrix, and therefore these parameters are decoupled from the ground state.

**2.2. Matrix Diagonalization.** For the matrix diagonalization necessary to solve the CIS equations, we use the Davidson–Liu block diagonalization method,<sup>36</sup> which has been shown to be very efficient, especially for large sparse matrices.<sup>37</sup> This iterative procedure starts with the guess vectors  $|\mathbf{b}_i\rangle$ , which are expanded to form the best approximation of the CIS wave function. All eigenvectors of the CIS matrix are expanded in an  $L$ -dimensional orthonormal subspace of the eigenvectors:

$$|\mathbf{x}^k\rangle = \sum_i^L \alpha_i^k |\mathbf{b}_i\rangle \quad (8)$$

where  $|\mathbf{x}^k\rangle$  is the exact eigenvector and  $\alpha_i^k$  represents the expansion coefficients. With these new basis vectors, a projected CIS matrix is formed:

$$\langle \mathbf{b}_i | \hat{H} | \mathbf{b}_j \rangle \quad 1 \leq i, j \leq L \quad (9)$$

The expansion coefficients  $\alpha^k$  are given by the eigenvectors of this matrix with the eigenvalues  $\rho_k$  and the eigenvectors  $|\mathbf{c}^k\rangle$ . The approximated eigenvectors  $|\mathbf{c}^k\rangle$  are corrected by  $|\delta^k\rangle$ :

$$|\mathbf{c}^k\rangle - |\delta^k\rangle = |\mathbf{x}^k\rangle \quad (10)$$

The set of  $|\delta^k\rangle$  is directly related to the residual vectors  $|\mathbf{r}^k\rangle$ :

$$(\mathbf{H} - \lambda^k) |\delta^k\rangle = -(\mathbf{H} - \lambda^k) |\mathbf{x}^k\rangle = -|\mathbf{r}^k\rangle \quad (11)$$

The residuals are calculated by connecting eqs 8 and 11. Using the definition of the  $\sigma$  (eq 9) vectors, one obtains

$$|\mathbf{r}^k\rangle = \sum_i^L \alpha_i^k (|\sigma_i\rangle - \rho_k |\mathbf{b}_i\rangle) \quad (12)$$

Since the CIS matrix is sparse and dominated by the diagonal elements, the correction vector is approximated by

$$|\delta^k\rangle \approx -(\mathbf{D} - \rho^k \mathbf{E})^{-1} |\mathbf{r}^k\rangle \quad (13)$$

where  $\mathbf{D}$  is an arbitrary diagonal matrix, which is connected to  $\mathbf{H}$ . In our implementation, we do not use the diagonal elements of  $\mathbf{H}$ . Since these are dominated by the orbital energy differences, we simply set

$$D_{ia,ia} = \varepsilon_a - \varepsilon_i \quad (14)$$

This leads to a significantly faster calculation of the diagonal matrix  $\mathbf{D}$ , which improves the overall performance of the Davidson algorithm despite an increase in the number of iterations. The obtained correction vectors  $|\delta^k\rangle$  are then normalized with respect to the existing set of expansion vectors  $|\mathbf{b}_i\rangle$  by a modified Gram–Schmidt orthogonalization.<sup>38</sup> This leads to an increased size of the expansion space in every step of the algorithm up to a preselected threshold. As a standard criterion for convergence, we have chosen  $10^{-8}$  au for the eigenvalues and  $10^{-6}$  au for the norm of the residuals. In each step, the  $\sigma$  vectors are calculated according to eq 9. With a restricted Hartree–Fock reference, the CSFs are used, and two independent sets of  $\sigma$  vectors are defined:

$${}^1\sigma_a^i = (E_{\text{RHF}} - \varepsilon_i + \varepsilon_a - d_{ia}^{\text{corr}}) t_i^a + \sum_j^{\text{occ}} \sum_b^{\text{vir}} t_j^b [2c_1(ia|jb) - c_2(ij|ab)] \quad (15)$$

$${}^3\sigma_a^i = (E_{\text{RHF}} - \varepsilon_i + \varepsilon_a) t_i^a - \sum_j^{\text{occ}} \sum_b^{\text{vir}} t_j^b c_T(ij|ab) \quad (16)$$

By defining pseudo-Fock matrices of the form

$${}^1\tilde{F}_{ia} = \sum_j^{\text{occ}} \sum_b^{\text{vir}} t_j^b [2c_1(ia|jb) - c_2(ij|ab)] \quad (17)$$

$${}^3\tilde{F}_{ia} = -\sum_j^{\text{occ}} \sum_b^{\text{vir}} t_j^b c_T(ij|ab) \quad (18)$$

the  $\sigma$  vectors are much more efficiently calculated in the atomic orbital (AO) basis. This leads to the AO transformed pseudo-Fock matrix, which is given in the general case as

$${}^1\tilde{F}_{\mu\nu} = \sum_{\lambda\sigma} T_{\lambda\sigma} [2c_1(\mu\nu|\lambda\sigma) - c_2(\mu\lambda|\nu\sigma)] \quad (19)$$

$${}^3\tilde{F}_{\mu\nu} = \sum_{\lambda\sigma} T_{\lambda\sigma} c_T(\mu\lambda|\nu\sigma) \quad (20)$$

where  $T_{\lambda\sigma}$  are AO-transformed CIS amplitudes. Within the INDO approximation,<sup>39</sup> three kinds of elements of the pseudo-Fock matrix have to be distinguished, the diagonal elements, the intra-atomic off-diagonal blocks, and the two-center terms:

$${}^1F_{\mu\nu} = \begin{cases} \sum_{\lambda\sigma \in A} T_{\lambda\sigma} [2c_1(\mu\mu|\lambda\sigma) - c_2(\mu\lambda|\mu\sigma)] & \text{for } \mu = \nu \in A \\ + \sum_{\lambda \in B \neq A} T_{\lambda\lambda} c_1(\mu\mu|\lambda\lambda) & \\ \sum_{\lambda\sigma \in A} T_{\lambda\sigma} [2c_1(\mu\nu|\lambda\sigma) - c_2(\mu\lambda|\nu\sigma)] & \text{for } \mu, \nu \in A \\ -T_{\mu\nu} c_2(\mu\mu|\nu\nu) & \text{for } \mu \in A, \nu \in B \end{cases} \quad (21)$$

$${}^3F_{\mu\nu} = \begin{cases} -\sum_{\lambda\sigma\in A} T_{\lambda\sigma} c_T(\mu\lambda|\mu\sigma) & \text{for } \mu = \nu \\ -\sum_{\lambda\sigma\in A} T_{\lambda\sigma} c_T(\mu\lambda|\nu\sigma) & \text{for } \mu, \nu \in A \\ -T_{\mu\nu} c_T(\mu\mu|\nu\nu) & \text{for } \mu \in A, \nu \in B \end{cases} \quad (22)$$

Necessary matrix operations of the Davidson–Liu approach are efficiently performed by BLAS<sup>40,41</sup> routines in our implementation.

**2.3. Parameterization.** In order to obtain accurate values for the excitation energies, the whole MSINDO parameter set, which has repeatedly demonstrated its reliability,<sup>42–45</sup> has been reoptimized. This affects not only the excited states but also the ground-state properties. After augmenting the original MSINDO reference set,<sup>46</sup> the ground-state properties and some additional vertical excitation energies were optimized with respect to accurate reference data by varying the parameters in a nonlinear minimization algorithm, where the sums of least-squares errors are minimized.<sup>47</sup> The one-center two-electron integrals are evaluated in terms of the Slater–Condon parameters  $F^n$  and  $G^n$  ( $n = 0–3$ ). Following the ideas of Zerner and Ridley,<sup>20</sup> these are treated as adjustable parameters in the present MSINDO-sCIS implementation at variance with the agreed MSINDO method<sup>46</sup> (see Table 6, Supporting Information). This was necessary, since Slater–Condon factors play an important role in spectroscopy.<sup>48</sup> In the original MSINDO version,<sup>46</sup> the Slater–Condon factors were calculated analytically with a special set of orbital exponents (for a better comparison, their values are given in Table 6 in the Supporting Information). In the present approach, it was necessary to treat the  $F^n$  and  $G^n$  as independent empirical parameters in order to obtain reasonable agreement with experimental results for both ground and excited state properties. We had, to some extent, to give up the concept of a physical interpretation of the  $F$  and  $G$  factors. The  $F_{ss}^0$  values are expected to increase within the second row, but the optimized value for C is larger than that for N, O, and F. The obtained values of the Slater–Condon factors  $F_{ss}^0$ ,  $F_{sp}^0$ ,  $F_{pp}^0$ ,  $F_{pp}^2$ , and  $G_{sp}^1$  for the second-row elements and additionally  $F_{sd}^0$ ,  $F_{pd}^0$ ,  $G_{pd}^1$ ,  $G_{sd}^2$ , and  $G_{pd}^3$  for the third-row element sulfur are presented in Table 6 in the Supporting Information. Furthermore, we allowed the correction factors for the orthogonalization  $f^{B,orth}$  (see ref 46) to change. In the original version, these were formally preparameterized and then treated as constants, depending on the angular momentum of the orbital (1, 0.75, and 0.5). Additionally, the ionization potentials  $I$  and the shielding parameters, which are introduced by a distance-dependent exponential function with the exponent  $\kappa_{EP}$ , were also included in the parameterization. The last parameters, which were reoptimized, are the resonance integral parameters  $K_i$ , which depend on the local symmetry ( $s\sigma$ ,  $p\sigma$ ,  $p\pi$ ,  $d\sigma$ ,  $d\pi$ , or  $d\delta$ ) in a diatomic coordinate system. The complete new set of these parameters compared to the old parameters for common elements in organic molecules can be found in Table 7 in the Supporting Information. The excited-state parameters, which are introduced in the calculation of the excited states (eqs 3 and 5), have the following values:

$$c_i = \begin{cases} 0.9225 & \text{for } i = 1 \\ 0.9993 & \text{for } i = 2 \\ 0.8944 & \text{for } i = T \end{cases} \quad (23)$$

In the original reference,<sup>35</sup> the parameters  $c_1$ ,  $c_2$ , and  $c_T$  had the values 1.0, 0.317, and 0.317. However, these parameters have been obtained for density functional theory, which is completely

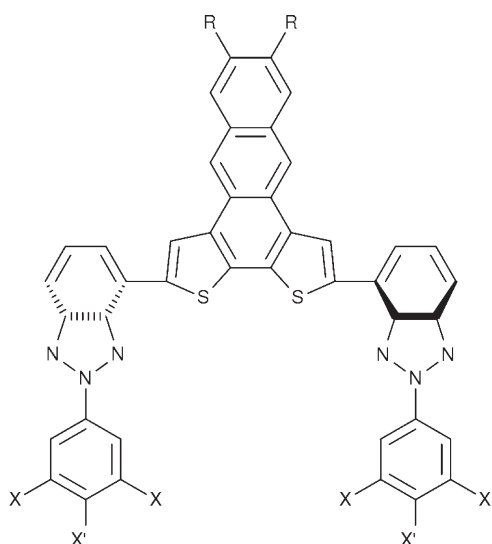
different from our present approach. Since both  $c_1$  and  $c_2$  are close to 1, the artificial problem of self-interaction (incomplete cancellation of Coulomb and exchange integrals) is nearly avoided for singlet states. These excited state parameters are in the spirit of DFT, where only a few global parameters are used to obtain accurate results, rather than in the spirit of semiempirical methods. But our statistical evaluation shows (see section 5.2) that it is not necessary to introduce bond-specific parameters for excited states. In summary, we have taken four steps in our parameterization. The first one is the inclusion of Slater–Condon factors in the parameterization. The second is the introduction of three global parameters for the excited states, which scale the Coulomb and exchange integrals. Third, an element-specific scaling of the orthogonalization correction is introduced. And the last step is the empirical correction  $d^{corr}$  (eq 4), which improves the excitation energies of totally symmetric excitations. Those ideas eventually define a new method, which is intended to be efficient and accurate. In a first step, only selected elements of the first two periods H, C, N, O, F, and, additionally, S have been reparameterized. Further work, the reparameterization of the remaining second and third row elements including the third-row transition metals, is in progress.

### 3. GENERAL CONSIDERATIONS

Ground-state geometry optimizations and vertical excitation energy calculations are carried out with the new version MSINDO-sCIS. Different from Thiel's statistical evaluation,<sup>17</sup> we did not take the geometries from MP2/6-31G\* calculations, because we want to demonstrate that the MSINDO-sCIS method is capable of providing reliable results for both kinds of properties. In preliminary tests, vertical excitation energies with MP2 geometries were calculated for selected molecules of the test set. The small change in the geometries did not affect the vertical excitation energies for both singlet and triplet states significantly (changes were below 0.1 eV). Since semiempirical methods employ a minimal valence basis set, Rydberg states or states with a valence-Rydberg mixing cannot be described properly; accordingly, Thiel's benchmark set<sup>1</sup> contains only valence excited states. Within the MSINDO-sCIS approach, all possible single excitations of the valence orbitals are taken into account in the active space. For the excited-state assignment, the electronically excited states are first classified according to their point group symmetry. Afterward, the most significant orbitals of the excitation are used for a further inspection of the kind of transition. With this information, it is possible to assign the states and distinguish between  $\pi-\pi^*$ ,  $n-\pi^*$ , and  $\sigma-\pi^*$  transitions.

### 4. COMPUTATIONAL PERFORMANCE

As a benchmark for the computational performance of our approach, we have calculated the first eight excited states of selected molecules of the class of bis(phenyl-benzotriazole) dithiophenes (see Figure 1) and compared the MSINDO-sCIS CPU timings to TD-DFT calculations (see Table 1). These molecules are expected to have charge-transfer states in the range of visible light, so that they should work in a solar cell device.<sup>49</sup> The TD-DFT calculations (B3LYP/TZV(P)) were carried out using the ORCA 2.7.0 program package.<sup>50</sup> The computational performance for both methods strongly depends on the number of Davidson cycles. Comparing the timings of MSINDO-sCIS and TD-B3LYP (shown in Table 1) gives a ratio of up to 1:7500 on a single Quad-Core AMD 2378 Opteron 2.4 GHz processor.



**Figure 1.** Class of molecules used for performance test in comparison with TDDFT (B3LYP/TZVP).

**Table 1.** MSINDO-sCIS and TD-DFT Calculation Timings for the Different Types of Bis(phenyl-benzotriazole) dithiophenes (see Figure 1)<sup>a</sup>

molecule	MSINDO-sCIS	B3LYP/TZVP
X' = H, X = NO <sub>2</sub> , R = OC(C <sub>2</sub> H <sub>5</sub> )C <sub>5</sub> H <sub>11</sub>	315s (N = 20)	11952s (N = 05)
X' = H, X = NO <sub>2</sub> , R = isoprop	195s (N = 20)	74986s (N = 05)
X' = H, X = NO <sub>2</sub> , R = CH <sub>3</sub>	150s (N = 19)	61303s (N = 06)
X' = H, X = F, R = OC(C <sub>2</sub> H <sub>5</sub> )C <sub>5</sub> H <sub>11</sub>	238s (N = 20)	259435s (N = 06)
X' = F, X = H, R = OC(C <sub>2</sub> H <sub>5</sub> )C <sub>5</sub> H <sub>11</sub>	190s (N = 20)	244928s (N = 12)
X' = X = H, R = NO <sub>2</sub>	89s (N = 17)	665596s (N = 52)
X' = X = H, R = phenyl	155s (N = 20)	230268s (N = 10)
X' = X = H, R = COOH	104s (N = 19)	166497s (N = 14)

<sup>a</sup>The values in parentheses are the number *N* of Davidson iterations.

This superior speed of MSINDO-sCIS allows applications of the excited-state description of really large molecules containing more than 1000 atoms, which are not feasible with TD-DFT or even more accurate methods.

## 5. RESULTS

In this section, we will compare ground- and excited-state properties calculated with the new method MSINDO-sCIS with experimental and other theoretical approaches. We will start with the ground state properties, where we will show that the new parameter set is as reliable as the original set developed by Ahlswede and Jug.<sup>46</sup> This is followed by a detailed description of the vertical excitation energies, where all excitations from the benchmark set of Silva-Junior et al.<sup>1,11</sup> were taken into account. The MSINDO results are then compared to the theoretical best estimates,<sup>1</sup> TD-B3LYP,<sup>16</sup> INDO and OMx results.<sup>17</sup> It should be mentioned that the comparison with the OM3-CISDTQ results

**Table 2.** MSINDO Mean Absolute Errors for Ground- and Excited-State Properties<sup>a</sup>

property	unit	no. references	parameter set	
			old <sup>b</sup>	new <sup>c</sup>
$\Delta_f H$	kcal/mol	89	5.25	4.97
<i>R</i>	Å	210	0.013	0.016
$\varphi$	deg	94	1.69	1.76
IP	eV	83	0.43	0.53
$\mu$	D	49	0.34	0.35
$\Delta E_{0-n}$	eV	36		0.34

<sup>a</sup>Heats of formation  $\Delta_f H$ , bond lengths *R*, bond angles  $\varphi$ , ionization potentials IP, dipole moments  $\mu$ , and vertical excitation energies  $\Delta E_{0-n}$ . In the old parameterization, no excited states were included. <sup>b</sup>Reference 46. <sup>c</sup>Present work.

is just for benchmarking purposes, since the common used method for excited states is OM2/MRCISD, which has the same accuracy.<sup>17</sup>

**5.1. Ground State Properties.** To avoid the well-known problems of INDO/S for ground-state geometries, while being parametrized for excited states, we included ground-state properties in our parametrization. A comparison of the results obtained with new MSINDO-sCIS parameters with those of the previous version is given in Table 2. Distances *R* and angles  $\varphi$  are only slightly less accurate than in the 1999 version of MSINDO.<sup>46,42</sup> On the other hand, it can be seen from Table 2 that the new set of parameters even slightly improves the heats of formation, while the ionization potentials are slightly less accurate than before. But in the overall view of these results, it is obvious that we now have reliable parameters for both ground and excited states.

**5.2. Accuracy of Vertical Excitation Energies.** **5.2.1. Vertical Singlet Excitations.** The benchmark of the calculated singlet excitation energies with respect to the TBES<sup>1</sup> is presented in Table 3. We will discuss the results of all groups of compounds in the following paragraphs.

**Unsaturated Aliphatic Hydrocarbons. Ethene.** The energy of the singlet  $\pi-\pi^*$  state of ethene (TBE-2 is 7.80 eV) is underestimated by 0.31 eV with MSINDO-sCIS. This error is larger than that of the OMx-CISDTQ suite of methods (0.02–0.05 eV) but much less compared to standard MNDO methods, where the energy is underestimated by more than 1 eV, and also compared to the INDO/S method, which overestimates the energy by more than 0.5 eV.

**E-Butadiene, E-Hexatriene, and E-Octatetraene.** The most interesting point in this series of C<sub>2n</sub>H<sub>2n+2</sub> polyenes is the gap between the bright 1B<sub>u</sub> and the 2A<sub>g</sub> state. In the benchmark set, this gap is reduced with increasing *n* until the ordering is reversed. Due to the large contributions from double excitations (33%)<sup>51</sup> and higher excited configurations (20% within CISDTQ)<sup>17</sup> for the A<sub>g</sub> state, this problem cannot be solved by parametrization. The MSINDO-sCIS errors for the 2A<sub>g</sub> states of E-butadiene, E-hexatriene, and E-octatetraene are 0.29 eV, 1.42 eV, and 1.31 eV, respectively. These relatively large errors indicate that our approach—including the empirical correction term (eq 4)—is not accurate for excited states with significant double excitation contributions. Even compared to the INDO/S method, where the errors are 0.32 eV, 0.81 eV, and 0.76 eV for the 2A<sub>g</sub> states, MSINDO is even slightly inferior.

**Cyclopropene.** The excited states of strained ring systems such as cyclopropene are problematic for standard semiempirical

**Table 3. Vertical Excitation Energies  $\Delta E$  [eV] for Singlet States: MSINDO Results Compared with the Theoretical Best Estimates TBE-2 from Ref 1**

molecule	state	type	TBE-2	MSINDO
ethene	$1^1B_{1u}$	$\pi-\pi^*$	7.80	7.49
E-butadiene	$1^1B_u$	$\pi-\pi^*$	6.18	6.09
	$2^1A_g$	$\pi-\pi^*$	6.55	6.84
E-hexatriene	$1^1B_u$	$\pi-\pi^*$	5.10	5.33
	$2^1A_g$	$\pi-\pi^*$	5.09	6.51
E-octatetraene	$1^1B_u$	$\pi-\pi^*$	4.66	4.88
	$2^1A_g$	$\pi-\pi^*$	4.47	5.78
cyclopropene	$1^1B_1$	$\pi-\pi^*$	6.67	6.07
	$1^1B_2$	$\pi-\pi^*$	6.68	5.97
cyclopentadiene	$1^1B_2$	$\pi-\pi^*$	5.55	5.33
	$2^1A_1$	$\pi-\pi^*$	6.28	6.57
norbornadiene	$1^1A_2$	$\pi-\pi^*$	5.37	5.42
	$1^1B_2$	$\pi-\pi^*$	6.21	5.69
benzene	$1^1B_{2u}$	$\pi-\pi^*$	5.08	5.34
	$1^1B_{1u}$	$\pi-\pi^*$	6.54	5.93
	$1^1E_{1u}$	$\pi-\pi^*$	7.13	7.02
	$1^1E_{2g}$	$\pi-\pi^*$	8.15	8.11
naphthalene	$1^1B_{3u}$	$\pi-\pi^*$	4.25	4.59
	$1^1B_{2u}$	$\pi-\pi^*$	4.82	5.40
	$2^1A_g$	$\pi-\pi^*$	5.90	6.03
	$1^1B_{1g}$	$\pi-\pi^*$	5.75	6.26
	$2^1B_{3u}$	$\pi-\pi^*$	6.11	6.42
	$2^1B_{2u}$	$\pi-\pi^*$	6.36	6.26
	$2^1B_{1g}$	$\pi-\pi^*$	6.46	7.73
	$3^1A_g$	$\pi-\pi^*$	6.49	6.77
furan	$1^1B_2$	$\pi-\pi^*$	6.32	5.59
	$2^1A_1$	$\pi-\pi^*$	6.57	6.27
	$3^1A_1$	$\pi-\pi^*$	8.13	6.63
pyrrole	$2^1A_1$	$\pi-\pi^*$	6.37	5.92
	$1^1B_2$	$\pi-\pi^*$	6.57	5.69
	$3^1A_1$	$\pi-\pi^*$	7.91	6.45
imidazole	$2^1A'$	$\pi-\pi^*$	6.25	5.45
	$1^1A''$	$n-\pi^*$	6.65	6.29
	$3^1A'$	$\pi-\pi^*$	6.73	6.13
pyridine	$1^1B_2$	$\pi-\pi^*$	4.85	5.89
	$1^1B_1$	$n-\pi^*$	4.59	4.66
	$1^1A_2$	$n-\pi^*$	5.11	5.29
	$2^1A_1$	$\pi-\pi^*$	6.26	5.54
	$2^1B_2$	$\pi-\pi^*$	7.27	6.99
	$3^1A_1$	$\pi-\pi^*$	7.18	7.10
pyrazine	$1^1B_{3u}$	$n-\pi^*$	4.13	3.44
	$1^1A_u$	$n-\pi^*$	4.98	4.66
	$1^1B_{2u}$	$\pi-\pi^*$	4.97	5.45
	$1^1B_{2g}$	$n-\pi^*$	5.65	5.46
	$1^1B_{1g}$	$n-\pi^*$	6.69	7.05
	$1^1B_{1u}$	$\pi-\pi^*$	6.83	5.58
	$2^1B_{2u}$	$\pi-\pi^*$	7.81	8.39
	$2^1B_{1u}$	$\pi-\pi^*$	7.86	5.98
pyrimidine	$1^1B_1$	$n-\pi^*$	4.43	4.41
	$1^1A_2$	$n-\pi^*$	4.85	4.74
	$1^1B_2$	$\pi-\pi^*$	5.34	5.95
	$2^1A_1$	$\pi-\pi^*$	6.82	6.03

**Table 3. Continued**

molecule	state	type	TBE-2	MSINDO
pyridazine	$1^1B_1$	$n-\pi^*$	3.85	3.85
	$1^1A_2$	$n-\pi^*$	4.44	4.38
	$2^1A_1$	$\pi-\pi^*$	5.20	5.07
	$2^1A_2$	$n-\pi^*$	5.66	5.11
s-triazine	$1^1A_1''$	$n-\pi^*$	4.70	4.79
	$1^1A_2''$	$n-\pi^*$	4.71	5.56
	$1^1E''$	$n-\pi^*$	4.75	4.93
	$1^1A_2'$	$\pi-\pi^*$	5.71	7.07
s-tetrazine	$1^1B_{3u}$	$n-\pi^*$	2.46	2.65
	$1^1A_u$	$n-\pi^*$	3.78	3.50
	$1^1B_{1g}$	$n-\pi^*$	4.87	5.40
	$1^1B_{2u}$	$\pi-\pi^*$	5.08	5.21
	$1^1B_{2g}$	$n-\pi^*$	5.28	5.26
	$2^1A_u$	$n-\pi^*$	5.39	5.19
formaldehyde	$1^1A_2$	$n-\pi^*$	3.88	3.74
	$1^1B_1$	$\sigma-\pi^*$	9.04	9.38
	$2^1A_1$	$\pi-\pi^*$	9.29	8.96
acetone	$1^1A_2$	$n-\pi^*$	4.38	4.37
	$1^1B_1$	$\sigma-\pi^*$	9.04	8.42
	$2^1A_1$	$\pi-\pi^*$	8.90	6.71
p-benzoquinone	$1^1B_{1g}$	$n-\pi^*$	2.74	2.62
	$1^1A_u$	$n-\pi^*$	2.86	3.17
	$1^1B_{3g}$	$\pi-\pi^*$	4.44	5.02
	$1^1B_{1u}$	$\pi-\pi^*$	5.47	5.92
	$1^1B_{3u}$	$n-\pi^*$	5.55	5.47
	$2^1B_{3g}$	$\pi-\pi^*$	7.16	6.20
formamide	$1^1A''$	$n-\pi^*$	5.55	5.40
	$2^1A'$	$\pi-\pi^*$	7.35	7.23
acetamide	$1^1A''$	$n-\pi^*$	5.62	5.59
	$2^1A'$	$\pi-\pi^*$	7.14	7.47
propanamide	$1^1A''$	$n-\pi^*$	5.65	5.36
	$2^1A'$	$\pi-\pi^*$	7.09	7.30
cytosine	$2^1A'$	$\pi-\pi^*$	4.66	4.50
	$1^1A''$	$n-\pi^*$	4.87	4.98
	$2^1A''$	$n-\pi^*$	5.26	5.18
	$3^1A'$	$\pi-\pi^*$	5.62	5.73
thymine	$1^1A''$	$n-\pi^*$	4.82	4.40
	$2^1A'$	$\pi-\pi^*$	5.20	4.69
	$3^1A'$	$\pi-\pi^*$	6.27	5.44
	$2^1A''$	$n-\pi^*$	6.16	5.53
	$4^1A'$	$\pi-\pi^*$	6.53	6.70
uracil	$1^1A''$	$n-\pi^*$	5.00	4.88
	$2^1A'$	$\pi-\pi^*$	5.25	5.01
	$3^1A'$	$\pi-\pi^*$	6.26	5.80
	$2^1A''$	$n-\pi^*$	6.10	5.84
	$4^1A'$	$\pi-\pi^*$	6.70	6.08
	$3^1A''$	$n-\pi^*$	6.56	6.21
adenine	$1^1A''$	$n-\pi^*$	5.12	5.12
	$2^1A'$	$\pi-\pi^*$	5.25	5.32
	$3^1A'$	$\pi-\pi^*$	5.25	5.60
	$2^1A''$	$n-\pi^*$	5.75	5.82

methods. Here, MSINDO-sCIS provides the best results within the considered semiempirical methods. While the OMx methods

still underestimate the  ${}^1B_1$  state by 0.43–0.83 eV, the MSINDO-sCIS error is not larger than 0.6 eV. For the problematic  ${}^1B_2$  state (TBE: 6.68 eV), the MSINDO-sCIS error of –0.61 eV is still high. But still, MSINDO-sCIS performs better than standard MNDO methods.<sup>17</sup> It has been pointed out by Thiel and Silva-Junior<sup>17</sup> that this state has significant contributions from diffuse orbitals<sup>52</sup> and therefore may be problematic for methods using minimal basis sets.

**Cyclopentadiene.** For the first two excited singlet states of cyclopentadiene, MSINDO-sCIS slightly underestimates the  ${}^1B_2$  state (TBE 5.55 eV) by 0.22 eV and overestimates the  ${}^1A_1$  (TBE 6.28 eV) state by 0.29 eV. This is a slightly better result than that obtained by OMx, where the error is between 0.41 and 0.79 eV.<sup>17</sup>

**Norbornadiene.** Comparing the results for norbornadiene shows that MSINDO reproduces the first excited  $A_2$  state well with an error of 0.05 eV. For the second excitation with  $B_2$  symmetry, MSINDO-sCIS underestimates the TBE by 0.52 eV. The absolute error is comparable to the OMx methods that overestimate the excitation energy for both states by more than 0.5 eV. It can be seen from the INDO/S results<sup>17</sup> that the underestimation of the excitation energies in this molecules seems to be a typical INDO problem. But the underestimation in the case of MSINDO is less than for INDO/S ( $A_2$  error, –0.87 eV;  $B_2$  error, –0.67 eV<sup>17</sup>).

**Aromatic Hydrocarbons and Heterocycles. Benzene.** For benzene, four excited singlet states in the range of 5.0–8.2 eV are found in the benchmark set.<sup>1</sup> For three of the states, MSINDO-sCIS calculations show good agreement with TBE-2. The first excited singlet state  $B_{2u}$  as well as the two higher excited states ( $E_{1u}$  and  $E_{2g}$ ) are well reproduced by MSINDO-sCIS with errors of 0.26 eV, 0.11 eV, and 0.04 eV. Only the  $B_{1u}$  state differs more from the TBE-2 by –0.49 eV. The severe problem of the  $\sigma-\sigma^*$  contamination in the INDO/S method<sup>17</sup> is apparently diminished with our ansatz. The OMx methods show similar performance with errors between 0.1 and 0.67 eV. Therefore, it can be concluded that the  $B_{1u}$  state is in general problematic to describe within semiempirical methods.

**Naphthalene.** The lower excited states of naphthalene are reasonably reproduced, while higher states are less accurate. The doubles mixing for the  $A_g$  states, which has been discussed before,<sup>17</sup> is well described with our empirical correction method. Comparing the errors of the  $A_g$  states of MSINDO-sCIS (0.13–0.28 eV) with the OMx errors (0.55–0.74 eV) shows that our simple correction in connection with a good parameterization can even outperform the explicit calculation of higher excited determinants (CISDTQ). But the main difference from the OMx method is that MSINDO-sCIS overestimates the excitation energy, while OMx excitation energies are too low. The MSINDO-sCIS errors are unfortunately much larger for higher excited states of a given irreducible representation. For example, the difference for the  $2B_{1g}$  state is 1.27 eV, which is considerably larger than for OMx (0.15–0.22 eV) and also INDO/S (0.07 eV). This seems to be a general trend.

**Furan.** The singlet excitation energies of furan obtained with MSINDO-sCIS show the same errors as the other semiempirical methods. The error range of 0.3–1.5 eV is similar to the OMx methods. It is an improvement over the INDO/S and INDO/S2 methods, where errors up to 2.0 eV are obtained.

**Pyrrrole.** For pyrrrole, which is isoelectronic with furan, the results for the excitation energies are similar. The error range is somewhat larger (0.65–1.46 eV), and again all energies are underestimated. This is similar to other semiempirical methods.<sup>17</sup>

**Imidazole.** The imidazole spectrum consists of two  $\pi-\pi^*$  transitions and one  $n-\pi^*$  transition. The  $n-\pi^*$  transition is well reproduced by MSINDO-sCIS with an error of 0.26 eV, but the  $\pi-\pi^*$  transitions are underestimated by approximately 0.8 eV. This leads to a wrong ordering of the excitations, but the MSINDO-sCIS errors are still smaller than those of other semiempirical methods.

**Pyridine.** Due to the break of symmetry by substituting one carbon atom by nitrogen in benzene, the four  $\pi-\pi^*$  transitions split into six transitions with  $A_1$  and  $B_2$  symmetry in pyridine. The lowest four of them are included in the benchmark set. Comparing the TBEs with the MSINDO-sCIS results, a large scattering of the excitation energies is observed. The errors are between 0.05 and 1.04 eV. The scattering is higher than with OMx or INDO/S, and even the  $1B_2 - 2A_1$  ordering is wrong. But here, the higher excitations are in better agreement with the TBE-2 benchmark results. In particular, the two  $n-\pi^*$  transitions are described rather well. Both values are overestimated (0.07 and 0.18 eV) but still close to the benchmark results.

**Pyrazine, Pyrimidine, and Pyridazine.** The performance of MSINDO-sCIS with the azabenzenes with two nitrogen atoms is similar to that for pyridine. The exception is pyridazine, where the excitation energies are underestimated by less than 0.2 eV. The pyrazine results, where eight reference energies are in the benchmark set, are very inhomogeneous. Except for the  ${}^1B_{2g}$  state, all energies are underestimated by more than 0.5 eV. The state ordering is in general parallel to the ab initio data, except for the change of the  ${}^1B_{1g}$  and  ${}^1B_{1u}$  states. The same effect can be observed in the OMx benchmark calculations of Thiel and Silva-Junior.<sup>17</sup>

**s-Triazine.** The three lowest excitations for s-triazine are nearly degenerate  $n-\pi^*$  transitions (4.70–4.75 eV). Here, MSINDO-sCIS overestimates the second excitation by 0.85 eV. Since the other two states are overestimated by only 0.09–0.18 eV, this also results in a loss of the near-degeneracy, similar to the results of AM1 and PM3.<sup>17</sup> The  $\pi-\pi^*$  excitation energy is overestimated by more than 1.0 eV. This is the opposite of all other semiempirical methods, where all energies are strongly underestimated, except for OM3, where the error is only –0.02 eV.

**s-Tetrazine.** The excited states of s-tetrazine are well described by MSINDO-sCIS, except for the third  $n-\pi^*$  excitation (error: 0.53 eV). The errors for the other states are all below 0.3 eV, which is superior to all other semiempirical methods, where the errors are typically over 0.5 eV. Even the OMx methods have problems with the  $n-\pi^*$  excitation, where the maximum error is over 1.0 eV.

**Aldehydes, Ketones, and Amides. Formaldehyde and Acetone.** Formaldehyde and acetone have nearly the same electronic spectrum. The results obtained by MSINDO-sCIS are in good agreement with those of TBEs-2 (errors of 0.01–0.34 eV), except for the  $A_1$  states (errors of 0.37–2.19 eV), which is problematic, since the energy is strongly underestimated in the case of acetone. Even though this is the maximum error in the complete benchmark set, this is an improvement compared with the INDO/S and also the INDO/S2 results, where the energy of these states is overestimated by 3–4 eV. Here, the OMx methods perform much better with errors between 0.39 and 0.82 eV.

**p-Benzoquinone.** The two lowest singlet excited states are dark states of the  $n-\pi^*$  type. MSINDO-sCIS reproduces the values of these states well. The errors are below 0.25 eV. The higher lying  $n-\pi^*$  state with  $B_{3u}$  symmetry is also correctly described with an error of 0.08 eV. The  $\pi-\pi^*$  states show typical

**Table 4. Vertical Excitation Energies  $\Delta E$  [eV] for Triplet States: MSINDO Results Compared with the Theoretical Best Estimates TBE-2 from Ref 1**

molecule	state	type	TBE-2	MSINDO
ethene	$1^3B_{1u}$	$\pi-\pi^*$	4.50	3.90
E-butadiene	$1^3B_u$	$\pi-\pi^*$	3.20	3.04
	$1^3A_g$	$\pi-\pi^*$	5.08	4.31
E-hexatriene	$1^3B_u$	$\pi-\pi^*$	2.40	2.62
	$1^3A_g$	$\pi-\pi^*$	4.15	3.69
E-octatetraene	$1^3B_u$	$\pi-\pi^*$	2.20	2.40
	$1^3A_g$	$\pi-\pi^*$	3.55	3.23
cyclopropene	$1^3B_2$	$\pi-\pi^*$	4.28	3.48
	$1^3B_1$	$\sigma-\pi^*$	6.40	5.93
cyclopentadiene	$1^3B_2$	$\pi-\pi^*$	3.26	2.82
	$1^3A_1$	$\pi-\pi^*$	5.09	4.09
norbornadiene	$1^3A_2$	$\pi-\pi^*$	3.68	3.17
	$1^3B_2$	$\pi-\pi^*$	4.16	3.14
benzene	$1^3B_{1u}$	$\pi-\pi^*$	4.15	3.30
	$1^3E_{1u}$	$\pi-\pi^*$	4.86	5.36
	$1^3B_{2u}$	$\pi-\pi^*$	5.88	6.32
	$1^3E_{2g}$	$\pi-\pi^*$	7.51	7.26
naphthalene	$1^3B_{2u}$	$\pi-\pi^*$	3.09	2.82
	$1^3B_{3u}$	$\pi-\pi^*$	4.09	4.56
	$1^3B_{1g}$	$\pi-\pi^*$	4.42	3.93
	$2^3B_{2u}$	$\pi-\pi^*$	4.56	4.90
	$2^3B_{3u}$	$\pi-\pi^*$	4.92	5.75
	$1^3A_g$	$\pi-\pi^*$	5.42	5.10
	$2^3B_{1g}$	$\pi-\pi^*$	6.12	6.89
	$2^3A_g$	$\pi-\pi^*$	6.17	7.15
	$3^3A_g$	$\pi-\pi^*$	6.65	7.40
	$3^3B_{1g}$	$\pi-\pi^*$	6.67	7.16
furan	$1^3B_2$	$\pi-\pi^*$	4.11	3.18
	$1^3A_1$	$\pi-\pi^*$	5.43	4.26
pyrrole	$1^3B_2$	$\pi-\pi^*$	4.44	3.30
	$1^3A_1$	$\pi-\pi^*$	5.42	6.13
imidazole	$1^3A'$	$\pi-\pi^*$	4.65	3.59
	$2^3A'$	$\pi-\pi^*$	5.64	4.94
	$1^3A''$	$n-\pi^*$	6.25	5.81
	$3^3A'$	$\pi-\pi^*$	6.38	6.15
pyridine	$1^3A_1$	$\pi-\pi^*$	4.06	3.71
	$1^3B_1$	$n-\pi^*$	4.25	4.33
	$1^3B_2$	$\pi-\pi^*$	4.64	5.16
	$2^3A_1$	$\pi-\pi^*$	4.91	5.62
	$1^3A_2$	$n-\pi^*$	5.28	5.40
	$2^3B_2$	$\pi-\pi^*$	6.08	6.35
s-tetrazine	$1^3B_{3u}$	$n-\pi^*$	1.87	2.05
	$1^3A_u$	$n-\pi^*$	3.49	3.31
	$1^3B_{1g}$	$n-\pi^*$	4.18	4.48
	$1^3B_{1u}$	$\pi-\pi^*$	4.36	3.86
	$1^3B_{2u}$	$\pi-\pi^*$	4.39	4.45
	$1^3B_{2g}$	$n-\pi^*$	4.89	4.77
	$2^3A_u$	$n-\pi^*$	4.96	4.92
	$2^3B_{1u}$	$\pi-\pi^*$	5.32	6.49
formaldehyde	$1^3A_2$	$\pi-\pi^*$	3.50	3.81
	$1^3A_1$	$\pi-\pi^*$	5.87	6.57
acetone	$1^3A_2$	$n-\pi^*$	4.05	4.51

**Table 4. Continued**

molecule	state	type	TBE-2	MSINDO
	$1^3A_1$	$\pi-\pi^*$	6.07	6.64
p-benzoquinone	$1^3B_{1g}$	$n-\pi^*$	2.50	2.70
	$1^3A_u$	$n-\pi^*$	2.61	3.22
	$1^3B_{1u}$	$\pi-\pi^*$	3.02	3.30
	$1^3B_{3g}$	$\pi-\pi^*$	3.37	3.44
formamide	$1^3A''$	$n-\pi^*$	5.28	5.60
	$1^3A'$	$\pi-\pi^*$	5.69	6.40
acetamide	$1^3A''$	$n-\pi^*$	5.35	5.80
	$1^3A'$	$\pi-\pi^*$	5.71	6.63
propanamide	$1^3A''$	$n-\pi^*$	5.38	5.54
	$1^3A'$	$\pi-\pi^*$	6.08	6.43

errors of round about 0.5 eV, which is similar to all other semiempirical methods.

*Formamide, Acetamide, and Propanamide.* The series of the three smallest amides—formamide, acetamide, and propanamide—is included in the benchmark set. All of these molecules have a similar spectrum, with a low lying  $n-\pi^*$  transition as the first excited singlet state. The errors of MSINDO-sCIS are between 0.03 and 0.29 eV, different from all other semiempirical methods, which strongly underestimate these states (errors of 0.4–1.1 eV). The errors for the  $\pi-\pi^*$  transitions (0.12–0.34 eV) are slightly higher than for the  $n-\pi^*$  state but still acceptable.

*Nucleobases. Cytosine, Thymine, Adenine, and Uracil.* The four nucleobases cytosine, thymine, adenine, and uracil play an important role in biochemistry. They are suitable representatives for the whole class of building blocks in larger biomolecules. All four benchmark states of cytosine are well reproduced by MSINDO-sCIS; their errors are in the range of 0.08–0.16 eV. For this molecule, the OMx methods have errors of 0.27–0.68 eV. Again, it can be observed that the error for  $n-\pi^*$  transitions is lower than for the  $\pi-\pi^*$  transitions. Thymine is a more problematic case; here, MSINDO-sCIS errors are 0.17–0.83 eV. This is comparable to the OMx results. The excited-state energies of uracil are underestimated by 0.12–0.65 eV. Again, we observe that the lower states are better described than the higher states. The vertical excitation energies of adenine are overestimated by 0.08–0.35 eV. Here, we find a correct state ordering compared to the TBE-2.

*5.2.2. Vertical Triplet Excitations.* The benchmark set of the calculated triplet excitations can be found in Table 4. We will discuss all of the results in the following paragraphs.

*Unsaturated Aliphatic Hydrocarbons.* The triplet state energies of the unsaturated hydrocarbons are underestimated for the smaller chains and the rings. For the larger chains (hexatriene and octatetraene), the energies of the first excited triplet states are overestimated, while the second excited triplet state energies are underestimated. The errors are in the range of 0.16–1.1 eV. This is similar to that for the OMx methods. Compared to common NDDO methods (MNDO, AM1, and PM3) and to INDO/S, where the excitation energies are underestimated by up to 3 eV, MSINDO-sCIS is an improvement.

*Aromatic Hydrocarbons and Heterocycles.* The vertical triplet excitations of benzene are described within an acceptable error range (0.25–0.85 eV) by MSINDO-sCIS. The  $^3E_{1u}$  and the  $^3B_{1u}$  state energies are underestimated, while the other states are slightly overestimated. The state ordering is correct and the errors decrease when the excitation energy is raised. For naphthalene,

**Table 5.** Deviations of Vertical Excitation Energies in eV for Singlet and Triplet Excited States from TBE for TD-B3LYP, INDO/S, and OM3-CISDTQ<sup>a</sup>

singlet states (count = 103) <sup>b</sup>	TD-B3LYP/TZVP <sup>c</sup>	INDO/S <sup>d</sup>	OM3-CISDTQ <sup>d</sup>	MSINDO-sCIS
mean error [eV]	-0.07	-0.23	-0.22	-0.10
mean abs. error [eV]	0.27	0.51	0.45	0.44
std. dev. [eV]	0.33	0.70	0.54	0.59
max.(+) dev. [eV]	1.02	2.79	1.76	2.19
max.(-) dev. [eV]	0.75	1.45	1.19	1.42
triplet states (count = 63)	TD-B3LYP/TZVP <sup>c</sup>	INDO/S <sup>d</sup>	OM3-CISDTQ <sup>d</sup>	MSINDO-sCIS
mean error [eV]	-0.45	-0.31	-0.26	-0.01
mean abs. error [eV]	0.45	0.65	0.45	0.50
std. dev. [eV]	0.49	0.86	0.54	0.59
max.(+) dev. [eV]		2.49	1.08	1.17
max.(-) dev. [eV]	0.93	2.01	1.17	1.17

<sup>a</sup> MSINDO-sCIS errors are given with respect to TBE-2. <sup>b</sup> Count for TBE-1 is 104 for OM3 and TD-B3LYP/TZVP and 103 for INDO/S. <sup>c</sup> TBE-1 values taken from ref 16. <sup>d</sup> TBE-1 values taken from ref 17.

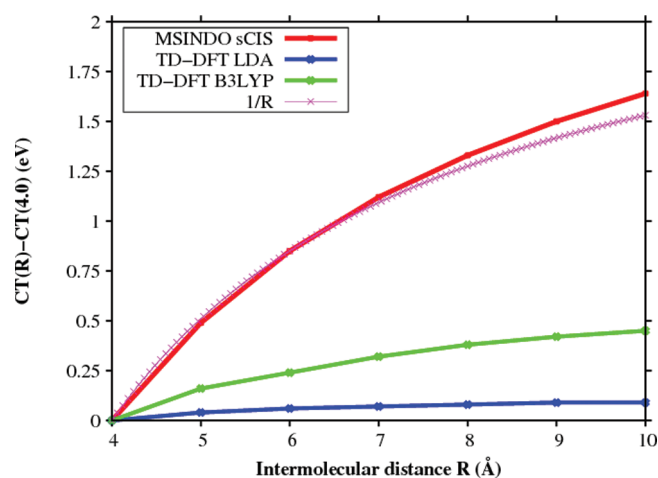
the errors scatter in the range of 0.27–0.88 eV. This is within the OMx error range and slightly better than the INDO/S results. The excited triplet states of furan are strongly underestimated by 0.93–1.17 eV but still less than in the case of the INDO/S method. Even INDO/S2 underestimates the second excited triplet state by more than 1.5 eV.

The lowest triplet states of pyrrole and imidazole are strongly underestimated by up to 1.0 eV. At the same time, the  $1^3A_1$  state of pyrrole is underestimated by 0.67 eV, which leads to a  $1^3B_2-1^3A_1$  splitting of 1.45 eV compared to 0.98 eV in TBE-2. The OMx methods result in a more correct  $1^3B_2-1^3A_1$  splitting between 0.75 and 0.89 eV. The errors for pyridine (0.08–0.52 eV) are comparably low and can compete with the OMx results, which are obtained on the CISDTQ level. Here, we see again that the totally symmetric excitation has the largest difference from the TBE value. The same holds for s-tetrazine (errors are 0.04–1.17 eV), where the MSINDO-sCIS results are closer to the TBE-2 values than any other semiempirical method, except for the highest vertical triplet excitation of s-tetrazine in the benchmark set. Here, we have a large error of 1.17 eV.

**Aldehydes, Ketones, and Amides.** The vertical triplet excitation energies of all aldehydes, ketones, and amides are overestimated by 0.07–0.88 eV. The qualitative state ordering is correct for all cases. Formaldehyde and acetone triplet excited states are described with acceptable errors comparable to OMx. Within the MSINDO-sCIS approach, we observe the general trend of an overestimation and no scattering as for all other semiempirical methods. Compared to OMx methods, the same behavior is observed for p-benzoquinone within the MSINDO-sCIS method. All states of this molecule are within a maximum overestimation of 0.2 eV except for the  $A_u$  state, where an error of 0.61 eV is obtained with MSINDO-sCIS. This is similar to that for the OMx methods, where the error for this state is in the range of 0.41–0.61 eV. In the case of amides, we see again the effect that the  $n-\pi^*$  transitions are described with a much smaller error than the  $\pi-\pi^*$  transitions. The errors are nearly twice as large for the  $\pi-\pi^*$  transitions. This is an effect which is not observed in the OMx or INDO/S methods. In the OMx methods, the  $n-\pi^*$  transitions are underestimated by nearly the same amount by which the  $\pi-\pi^*$  transitions are overestimated. This makes a simple shift of the values impossible for the OMx methods. Within the MSINDO-sCIS theory for amides, aldehydes, and

ketones, the values could in principle be shifted by a small amount to match it with experimental data.

**5.3. Statistical Evaluation.** To classify the MSINDO-sCIS method within the mainframe of computationally feasible methods, we compared the statistical evaluation to three methods. We have chosen INDO/S because of the conceptual equivalence, OM3-CISDTQ, because it has been demonstrated that this semiempirical method is currently the most accurate for the description of excited states, and TD-B3LYP/TZVP, because of its popularity in excited-state calculations. A statistical overview is given in Table 5. Since the benchmark results for TBE-1<sup>11</sup> and the more recent TBE-2<sup>1</sup> do not differ that much, the OM3,<sup>17</sup> INDO/S,<sup>17</sup> and TD-B3LYP<sup>16</sup> statistics available in the literature may be compared to the MSINDO-sCIS results benchmarked to TBE-2. It can be seen that for singlet states the MSINDO-sCIS method provides the same accuracy as the OM3-CISDTQ method. The mean errors are half those of the OM3 method, while the mean absolute errors are nearly the same. The standard deviation for singlets is comparable to those of the OM3 and lower compared to those of the INDO/S method. The maximum +/– deviations of MSINDO-sCIS for singlets are higher than in the OM3 method but still smaller than for the INDO/S method. The outlier in the MSINDO-sCIS method is the  $2^1A_1$  state of acetone, which is a  $\pi-\pi^*$  transition. The totally symmetric excitations are in general the problematic states in our CIS method, indicating that the simple correction term (eq 4) is not very accurate. A comparison with the INDO/S method shows that MSINDO-sCIS is slightly better in all statistical aspects, while TD-B3LYP/TZVP is in turn better by a factor of 1.5 to 2 in all values. But in the overall view on the important statistical points, it can be seen that the MSINDO-sCIS accuracy for singlets is sufficient. For triplet excitations, MSINDO-sCIS is well behaved. Although no higher excitations than singles are explicitly included, MSINDO-sCIS is on a similar level as OM3-CISDTQ. Even the maximum +/– deviations are comparable to those of OM3-CISDTQ. Comparing the MSINDO-sCIS statistics with those of INDO/S shows that, although the methods are conceptually similar, MSINDO-sCIS is superior. The comparison with TD-B3LYP shows that MSINDO-sCIS can compete with TD-DFT for triplet states. Here, MSINDO-sCIS has—similar to OM3-CISDTQ—only a slightly larger standard



**Figure 2.** Lowest excitation energy of a charge transfer (CT) state in the  $C_2H_4-C_2F_4$  complex.

deviation, which results in a comparable reliability to that of TD-B3LYP, although the maximum  $+/-$  deviations are higher.

## 6. CHARGE-TRANSFER STATES

The standard approach for excited state calculations, TD-DFT, has well-known problems with charge transfer (CT) states.<sup>13–15</sup> In order to compare our present approach with TD-DFT, we studied the common benchmark system, the  $C_2H_4-C_2F_4$  complex.<sup>53</sup> There is a high-lying CT state at around 13 eV, where one electron of ethene is transferred to the tetra-fluoro-ethene. To visualize the results, we plotted the excitation energy against the distance between both molecules (see Figure 2). For the TD-DFT calculations, we used the ORCA program package.<sup>50</sup> Starting with optimized structures of ethene and tetra-fluoro-ethene, the excitation energies were calculated starting with a distance of 4 Å. This distance was increased up to 10 Å with a step value of 1 Å. The energy of the CT state with respect to its value at  $R = 4$  Å was then plotted against the distance. It can be seen from Figure 2 that MSINDO-sCIS gives the correct  $1/R$  behavior, while the TD-DFT methods fail to give the correct description. This is a typical problem in TD-DFT and cannot even be solved by using double hybrid methods.<sup>15</sup> Within the ab initio CIS theory, on the other hand, this problem is totally absent. Therefore, it should not appear in semiempirical methods that are based on Hartree–Fock theory. Since we have introduced scaling parameters in the description of the excited state (eq 6), it was necessary to ensure that the CT error does not occur in our method. But according to the present results, the  $1/R$  behavior is still correctly reproduced with MSINDO-sCIS.

## 7. SUMMARY AND CONCLUSIONS

We have introduced a novel method for the calculation of excited states at a semiempirical level. Since the parametrization included ground state properties, the MSINDO-sCIS method yields reliable results for both ground and excited states. This is an improvement over the common INDO/S methods, which focus on excited states. We have demonstrated that the vertical excitation energies obtained with MSINDO-sCIS are in reasonable agreement with the TBEs, comparable to the OM3 methods. Most errors are in the range of 0.1–0.6 eV, with a trend toward larger errors for higher excitations. Compared to TD-B3LYP, the

MSINDO-sCIS method has a larger mean error, but it has two major advantages. First, the calculation times are orders of magnitude smaller, and second, the charge-transfer error is not present. A comparison of the accuracy with TD-B3LYP shows that MSINDO-sCIS can even compete for triplets. This is quite surprising, but it shows that the parameters are well balanced for the exchange part. For the triplet case, MSINDO-sCIS is superior to all other semiempirical methods except the OMx methods. Furthermore, due to the use of the Davidson–Liu algorithm, MSINDO-sCIS is computationally more efficient. Conceptually MSINDO-sCIS is below the OMx level, because OMx methods include NDDO integrals and also higher excitations. But we showed that a careful parametrization yields comparable results, although some of the parameters lost their physical significance. Therefore, it cannot be excluded that the new parametrization gives unbalanced results for systems that are quite different from those not included in the reference set. A couple of outliers have been observed in the present study, e.g., for the  $2^1A_1$  state of acetone. This may be an indicator for an unbalanced treatment of the excited state energies. There are other outliers where MSINDO-sCIS gives large errors for excited states with large doubles contributions (for example, in the polyenes). However, the overall performance is quite satisfactory.

Although analytical gradients for all sorts of semiempirical wave functions have been available for a number of years,<sup>54</sup> it is another advantage of the present approach that analytical gradients of a CIS wave function are much easier to implement<sup>55</sup> and faster to calculate than for higher excited determinants. This opens an efficient way for the calculation of electronic spectra, including vibrational coupling and excited-state geometries. Therefore, future applications to technically important systems, e.g., organic solar cells, are planned, where CT states play an important role and the molecules consist of several hundred atoms. Here, the larger errors for high-lying states do not play an important role since usually only the few lowest states are of interest.

Due to the implementation of the cyclic cluster model in MSINDO,<sup>56</sup> the calculation of excited state properties of solids and surfaces is another subject of future research with MSINDO-sCIS.<sup>57</sup>

## ■ ASSOCIATED CONTENT

**S Supporting Information.** Tables 6 and 7. This material is available free of charge via the Internet at <http://pubs.acs.org/>.

## ■ AUTHOR INFORMATION

### Corresponding Author

\*E-mail: [gadaczek@thch.uni-bonn.de](mailto:gadaczek@thch.uni-bonn.de).

## ■ ACKNOWLEDGMENT

The financial support by the Deutsche Forschungsgemeinschaft (SFB 813, project A2) is gratefully acknowledged.

## ■ REFERENCES

- (1) Silva-Junior, M.; Schreiber, M.; Sauer, S.; Thiel, W. *J. Chem. Phys.* **2010**, *133*, 174318.
- (2) Erman, P. In *Molecular Spectroscopy*; Barrow, R. F., Long, D. A., Sheridan, J., Eds.; The Royal Society of Chemistry: London, 1979; Vol. 6; pp 174–231.



- (3) Alerstam, E.; Andersson-Engels, S.; Svensson, T. *Optics Express* **2008**, *16*, 10440–10454.
- (4) Chattopadhyay, S.; Pahari, D.; Mahapatra, U.; Mukherjee, D. *Comput. Chem.: Rev. Curr. Trends* **2005**, *9*, 121.
- (5) Grimme, S. *Reviews in Computational Chemistry*; John Wiley & Sons, Inc.: New York, 2004; pp 153–218.
- (6) Buenker, R. J.; Peyerimhoff, S. D.; Butscher, W. *Mol. Phys.* **1978**, *35*, 771–791.
- (7) Finley, J.; Malmqvist, P.-A.; Roos, B. O.; Serrano-Andrés, L. *Chem. Phys. Lett.* **1998**, *288*, 299–306.
- (8) Christiansen, O.; Koch, H.; Jorgensen, P. *Chem. Phys. Lett.* **1995**, *243*, 409–418.
- (9) Koch, H.; Christiansen, O.; Jorgensen, P.; de Meras, A. M. S.; Helgaker, T. *J. Chem. Phys.* **1997**, *106*, 1808–1818.
- (10) Christiansen, O.; Koch, H.; Jorgensen, P. *J. Chem. Phys.* **1995**, *103*, 7429–7441.
- (11) Schreiber, M.; Silva-Junior, M.; Sauer, S.; Thiel, W. *J. Chem. Phys.* **2008**, *128*, 134110.
- (12) Runge, E.; Gross, E. K. U. *Phys. Rev. Lett.* **1984**, *52*, 997.
- (13) Dreuw, A.; Head-Gordon, M. *Chem. Rev.* **2005**, *105*, 4009–4037.
- (14) Casida, M. E. *THEOCHEM* **2009**, *914*, 3–18.
- (15) Grimme, S.; Neese, F. *J. Chem. Phys.* **2007**, *127*, 154116.
- (16) Silva-Junior, M. R.; Schreiber, M.; Sauer, S. P. A.; Thiel, W. *J. Chem. Phys.* **2008**, *129*, 104103.
- (17) Silva-Junior, M.; Thiel, W. *J. Chem. Theory Comput.* **2010**, *6*, 1546–1564.
- (18) Jug, K. *Theor. Chim. Acta* **1969**, *14*, 91–135.
- (19) Thiel, W. In *Theory and Applications of Computational Chemistry*; Dykstra, C. E., Frenking, G., Kim, K. S., Scuseria, G. E., Eds.; Elsevier: Amsterdam, 2005; pp 559–580.
- (20) Ridley, J.; Zerner, M. C. *Theor. Chim. Acta* **1973**, *32*, 111–134.
- (21) Zerner, M. C.; Loew, G. H.; Kirchner, R. F.; Mueller-Westerhoff, U. T. *J. Am. Chem. Soc.* **1980**, *102*, 589–599.
- (22) Li, J.; Williams, B.; Cramer, C. J.; Truhlar, D. G. *J. Chem. Phys.* **1999**, *110*, 724–733.
- (23) Foresman, J. B.; Head-Gordon, M.; Pople, J. A.; Frisch, M. J. *J. Phys. Chem.* **1992**, *96*, 135–149.
- (24) Kolb, M.; Thiel, W. *J. Comput. Chem.* **1993**, *14*, 775–789.
- (25) Weber, W.; Thiel, W. *Theor. Chem. Acc.* **2000**, *103*, 495–506.
- (26) Otte, N.; Scholten, M.; Thiel, W. *J. Phys. Chem. A* **2007**, *111*, 5751–5755.
- (27) Koslowski, A.; Beck, M. E.; Thiel, W. *J. Comput. Chem.* **2003**, *24*, 714–726.
- (28) Dewar, M. J. S.; Thiel, W. *Theor. Chim. Acta* **1977**, *46*, 89.
- (29) Dewar, M. J. S.; Thiel, W. *J. Am. Chem. Soc.* **1977**, *99*, 4899–4907.
- (30) Bredow, T.; Jug, K. In *MSINDO, Electronic Encyclopedia of Computational Chemistry*; v. Ragué Schleyer, P., Allinger, N. L., Clark, T., Gasteiger, J., Kollman, P. A., Schaefer, H. F., III, Schreiner, P. R., Eds.; Wiley: New York, 2004.
- (31) Head-Gordon, M.; Rico, R.; Oumi, M.; Lee, T. *Chem. Phys. Lett.* **1994**, *219*, 21–29.
- (32) Head-Gordon, M.; Oumi, M.; Maurice, D. *Mol. Phys.* **1999**, *96*, 593–602.
- (33) Baker, J. D.; Zerner, M. C. *Chem. Phys. Lett.* **1990**, *175*, 192–196.
- (34) Silva-Junior, M. R.; Sauer, S. P.; Schreiber, M.; Thiel, W. *Mol. Phys.* **2010**, *108*, 453–465.
- (35) Grimme, S. *Chem. Phys. Lett.* **1996**, *259*, 128–137.
- (36) Davidson, E. R. *J. Comput. Phys.* **1975**, *17*, 87–94.
- (37) Leininger, M. L.; Sherrill, C. D.; Allen, W. D.; Schaefer, H. F. S., III. *J. Comput. Chem.* **2001**, *22*, 1574–1589.
- (38) Hoffmann, W. *Computing* **1989**, *41*, 335–348.
- (39) Bacon, A. D.; Zerner, M. C. *Theor. Chim. Acta* **1979**, *53*, 21–54.
- (40) Blackford, L. S.; et al. *ACM Trans. Mater. Software* **2002**, *28*, 135–151.
- (41) Dongarra, J. J. *Int. J. High Perform. Appl. Supercomput.* **2002**, *16*, 1–199.
- (42) Ahlswede, B.; Jug, K. *J. Comput. Chem.* **1999**, *20*, 572–578.
- (43) Jug, K.; Kunert, L.; Köster, A. *Theor. Chem. Acc.* **2000**, *104*, 417–425.
- (44) Adachi, M.; Bredow, T.; Jug, K. *Dyes Pigments* **2004**, *63*, 225–230.
- (45) Ferro, N.; Bredow, T. *J. Comput. Chem.* **2010**, *31*, 1063–1079.
- (46) Ahlswede, B.; Jug, K. *J. Comput. Chem.* **1999**, *20*, 563–571.
- (47) Bartels, R. H. University of Texas and Center for Numerical Analysis and Report CNA-44. University of Texas: Austin, TX, 1972.
- (48) Slater, J. *Quantum Theory of Atomic Structure*; McGraw-Hill: New York, 1960; Vol. 1; pp 339–442.
- (49) Höger, S. Personal communications, 2010.
- (50) Petrenko, T.; Krylova, O.; Neese, F.; Sokolowski, M. *New J. Phys.* **2009**, *11*, 015001.
- (51) Head, J. *Int. J. Quantum Chem.* **2003**, *95*, 580–592.
- (52) González-Luque, R.; Merchán, M.; Roos, B. Z. *Phys. D* **1996**, *36*, 311–316.
- (53) Dreuw, A.; Weisman, J. L.; Head-Gordon, M. *J. Chem. Phys.* **2003**, *119*, 2943–2946.
- (54) Patchkovskii, S.; Koslowski, A.; Thiel, W. *Theor. Chem. Acc.* **2005**, *114*, 84–89.
- (55) Gadaczek, I.; Krause, K.; Hintze, K.; Bredow, T. Manuscript in preparation, 2011.
- (56) Bredow, T.; Geudtner, G.; Jug, K. *J. Comput. Chem.* **2001**, *22*, 89–101.
- (57) Gadaczek, I.; Hintze, K.; Bredow, T. Manuscript in preparation, 2011.

# Excited-State Studies of Polyacenes: A Comparative Picture Using EOMCCSD, CR-EOMCCSD(T), Range-Separated (LR/RT)-TDDFT, TD-PM3, and TD-ZINDO

K. Lopata,<sup>\*,†</sup> R. Reslan,<sup>‡</sup> M. Kowalska,<sup>§</sup> D. Neuhauser,<sup>\*,‡</sup> N. Govind,<sup>\*,†</sup> and K. Kowalski<sup>\*,†</sup>

<sup>†</sup>William R. Wiley Environmental Molecular Sciences Laboratory, Pacific Northwest National Laboratory, Richland, Washington 99352, United States

<sup>‡</sup>Department of Chemistry and Biochemistry, University of California, Los Angeles, California 90095-1569, United States

<sup>§</sup>Department of Chemistry, Washington State University Tri-Cities, Richland, Washington 99354, United States

**ABSTRACT:** The low-lying excited states ( $L_a$  and  $L_b$ ) of polyacenes from naphthalene to heptacene ( $N = 2-7$ ) are studied using various time-dependent computational approaches. We perform high-level excited-state calculations using equation of motion coupled cluster with singles and doubles (EOMCCSD) and completely renormalized equation of motion coupled cluster with singles, doubles, and perturbative triples (CR-EOMCCSD(T)) and use these results to evaluate the performance of various range-separated exchange-correlation functionals within linear-response (LR) and real-time (RT) time-dependent density functional theories (TDDFT). As has been reported recently, we find that the range-separated family of functionals addresses the well-documented TDDFT failures in describing these low-lying singlet excited states to a large extent and are as accurate as results from EOMCCSD on average. Real-time TDDFT visualization shows that the excited state charged densities are consistent with the predictions of the perimeter free electron orbital (PFE0) model. This corresponds to particle-on-a-ring confinement, which leads to the well-known red-shift of the excitations with acene length. We also use time-dependent semiempirical methods like TD-PM3 and TD-ZINDO, which are capable of handling very large systems. Once reparametrized to match the CR-EOMCCSD(T) results, TD-ZINDO becomes roughly as accurate as range-separated TDDFT, which opens the door to modeling systems such as large molecular assemblies.

## 1. INTRODUCTION

Polyacenes or acenes constitute a class of polycyclic organic compounds consisting of linearly fused benzene rings. These compounds, and their derivatives, have been studied extensively, and over the last several years, the larger representatives in this class have been used in a plethora of applications such as light-emitting diodes,<sup>1-4</sup> photovoltaic cells,<sup>5-7</sup> liquid crystal displays,<sup>8</sup> and organic field-effect transistors<sup>9,10</sup> to name a few. Pentacene, in particular, has received much attention because of its high charge-carrier (hole) mobility in films and molecular crystals.<sup>11-13</sup> For an overview of the electronic applications of acenes, see the reviews by Anthony.<sup>14,15</sup>

In a nutshell, the electronic properties of these materials are dictated by the  $\pi$  electrons which occupy the highest occupied and lowest unoccupied states; the  $\pi$  interactions between adjacent acene molecules, for example, give rise to the high hole mobility through molecular films. In a single molecule, the lowest valence excitations have  $\pi-\pi^*$  character, and the two lowest singlet excitations are commonly assigned as the  $L_a$  ( $B_{2u}$  symmetry) and  $L_b$  ( $B_{3u}$  symmetry) states, respectively. The former represents the polarization along the short axis, while the latter represents the polarization along the long axis. The  $L_b$  is the lowest excited state in naphthalene but switches positions with the  $L_a$  state for larger acenes, with the crossing happening around anthracene. It has long been suggested, from a valence-bond point of view, that the  $L_a$  state is mostly ionic in character involving significant rearrangement of the excited-state density, whereas the  $L_b$  state is mostly covalent where the excited-state density is similar to the ground state.

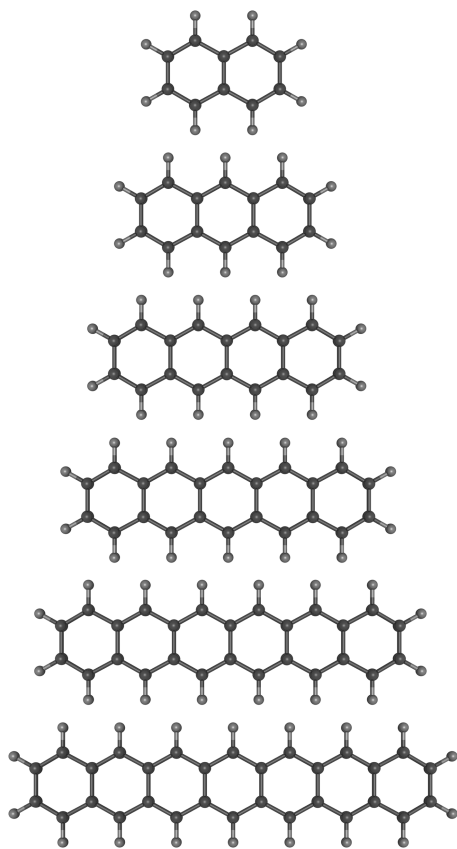
There has been significant progress in describing these excitations theoretically,<sup>16-23</sup> within which time-dependent

density functional theory (TDDFT)<sup>24-26</sup> has been the predominant method. It is now well-known, however, that for TDDFT traditional and global hybrid functionals fail to describe these lowest excitations. Grimme and Parac demonstrated that the ordering switches earlier than expected with both classes of functionals, and the excitation energy of the  $L_a$  state is severely underestimated and progressively worsens with system size.<sup>20</sup> Increasing the Hartree-Fock (HF) content in the exchange-correlation improves the picture, but  $L_a$  worsens the excitation energy of the  $L_b$  state. They concluded that it was impossible to capture both states accurately just by adjusting the HF content.

Very recently, range-separated hybrid (RSH) functionals have been applied to the  $L_a$  state in acenes.<sup>17-19,22</sup> RSHs correct the incorrect asymptotic behavior of the exchange by splitting the exchange into a short-range part and a long-range part. For many optically active charge transfer states, RSHs rival the accuracy of the equation of motion coupled cluster singles doubles (EOMCCSD) method on average. The success of RSHs in this case, however, is in many ways quite surprising, as the  $L_a$  state is an intramolecular transverse excitation (along short-axis of molecule) and clearly not long-range at all. Richard and Herbert labeled this a charge-transfer-like state in disguise,<sup>18</sup> which Kuritz et al. subsequently rationalized as arising from minimal overlap in auxiliary orbitals,<sup>19</sup> akin to minimal overlap of the hole/charge orbitals in a typical charge transfer excitation.

Received: July 26, 2011

Published: September 14, 2011



**Figure 1.** Structures of the acenes studied. From top to bottom: naphthalene ( $N = 2$ ), anthracene ( $N = 3$ ), tetracene ( $N = 4$ ), pentacene ( $N = 5$ ), hexacene ( $N = 6$ ), heptacene ( $N = 7$ ).

In some sense, acenes serve as a rough prototype for a more complicated light harvesting system and also as the fundamental building block for many molecular electronic devices. Careful analysis of the excitations in these deceptively simple molecules serves as a crucial test for the accuracy and predictive power of a theoretical technique, as indicated by the intense interest in benchmarking TDDFT results in these systems. In this light, our main goal in this paper is to examine the low-lying excited states of polyacenes from naphthalene to heptacene (Figure 1) using a wide selection of time-dependent approaches. We first perform a systematic analysis based on high-level coupled cluster (EOMCCSD and CR-EOMCCSD(T)) calculations. These calculations are used to benchmark the performance of various range-separated exchange-correlation functionals implemented within linear response and real-time TDDFT. Additionally, we explore the use of semiempirical time-dependent PM3 and ZINDO for describing these excitations and reparametrize their Hamiltonians to better match the results of high level theory. All structures were obtained using cc-pVTZ/B3LYP.

The rest of the paper is organized as follows: In section 2, we briefly review the various time-dependent approaches used in this study and provide the necessary computational details. The results are presented and discussed in section 3 and the concluding remarks in section 4.

## 2. METHODOLOGIES AND COMPUTATIONAL DETAILS

Below, we briefly review the formalisms for equation of motion coupled cluster (EOMCC), real-time time-dependent density

functional theory (RT-TDDFT), and real-time time-dependent PM3 and ZINDO. All results except the PM3 and ZINDO ones were obtained using NWChem.<sup>27</sup> The TD-PM3 results were obtained by a modification of the PM3 module from MOPAC 6.0,<sup>28,29</sup> to perform iterative time-dependent calculation of the TD-PM3 excitation energies.<sup>30</sup> The TD-ZINDO results were obtained by an analogous modification of ZINDO from the ZINDO-MN package.<sup>31</sup> The linear response TDDFT results were calculated using the module in NWChem; since the approach is widely used (e.g., refs 26 and 32), we omit the details.

**2.1. Equation of Motion Coupled Cluster.** The EOMCC formalism<sup>33</sup> can be viewed as an excited-state extension of the single-reference coupled cluster method, where the wave function corresponding to the  $K$ th state is represented as

$$|\Psi_K\rangle = R_K e^T |\Phi\rangle \quad (1)$$

where  $T$  and state-specific  $R_K$  operators are the cluster and excitation operators, respectively, and  $|\Phi\rangle$  is the so-called reference function usually chosen as a Hartree–Fock determinant. Various approximate schemes range from the basic EOMCCSD approximation where the cluster and correlation operators are represented as sums of scalar ( $R_{K,0}$  for excitation operator only), single ( $T_1, R_{K,1}$ ), and double ( $T_2, R_{K,2}$ ) excitations

$$|\Psi_K^{\text{EOMCCSD}}\rangle = (R_{K,0} + R_{K,1} + R_{K,2}) e^{T_1 + T_2} |\Phi\rangle \quad (2)$$

to the more advanced EOMCCSDT and EOMCCSDTQ approach, accounting for the effect of triple and/or quadruple excitations. It has been demonstrated that the progression of methods, EOMCCSD  $\rightarrow$  EOMCCSDT  $\rightarrow$  EOMCCSDTQ..., in the limit converges to the exact (full configuration interaction) energies. However, the rapid growth in the numerical complexity of the EOMCC methods makes calculations with the EOMCCSDT or EOMCCSDTQ methods very expensive, even for relatively small systems. Unfortunately, the EOMCCSD method is capable of providing reliable results only for singly excited states. However, as has recently been demonstrated,<sup>34</sup> errors in the range of 0.25–0.30 eV with respect to the experimental vertical excitation energies (VEE) persist with increasing system size.

In order to narrow the gap between the EOMCCSD and EOMCCSDT VEEs, several noniterative  $N^7$ -scaling methods that mimic the effect of triples in a perturbative fashion have been proposed in the past.<sup>35–40</sup> The completely renormalized EOMCCSD(T) approach, denoted CR-EOMCCSD(T),<sup>41</sup> falls into this class (see also refs 42 and 43–45 for the most recent developments). In this approach, the energy correction  $\delta_K^{\text{CR-EOMCCSD(T)}}$  is added to the EOMCCSD VEE ( $\omega_K^{\text{EOMCCSD}}$ )

$$\omega_K^{\text{CR-EOMCCSD(T)}} = \omega_K^{\text{EOMCCSD}} + \delta_K^{\text{CR-EOMCCSD(T)}} \quad (3)$$

where  $\delta_K^{\text{CR-EOMCCSD(T)}}$  is expressed through the trial wave function  $\langle\Psi_K|$  and the triply excited EOMCCSD moment operator  $M_{K,3}^{\text{EOMCCSD}}$  (see ref 41 for details):

$$\delta_K^{\text{CR-EOMCCSD(T)}} = \frac{\langle\Psi_K|M_{K,3}^{\text{EOMCCSD}}|\Phi\rangle}{\langle\Psi_K|(R_{K,0} + R_{K,1} + R_{K,2}) e^{T_1 + T_2}|\Phi\rangle} \quad (4)$$

Although the CR-EOMCCSD(T) method is characterized by the same  $N^7$  scaling as the ground-state CCSD(T) method,<sup>46</sup> the fact that triply excited EOMCCSD moments need to be calculated makes this approach a few times more expensive than the ground-state CCSD(T) approach.

**2.2. Real-Time TDDFT.** In real-time time-dependent density functional theory (RT-TDDFT), the time-dependent Kohn–Sham (KS) equations are explicitly propagated in time:

$$i \frac{\partial \psi_i(\mathbf{r}, t)}{\partial t} = \left[ -\frac{1}{2} \nabla^2 + v_{\text{KS}}[\rho](\mathbf{r}, t) \right] \psi_i(t) \quad (5)$$

$$= \left[ -\frac{1}{2} \nabla^2 + v_{\text{ext}}(\mathbf{r}, t) + v_{\text{H}}(\mathbf{r}, t) + v_{\text{XC}}[\rho](\mathbf{r}, t) \right] \psi_i(t) \quad (6)$$

where  $\rho(\mathbf{r}, t)$  is the charge density,  $v_{\text{ext}}(\mathbf{r}, t)$  is the external potential describing the nuclear–electron and applied field contributions,  $v_{\text{H}}(\mathbf{r}, t)$  is the electron–electron potential, and  $v_{\text{XC}}[\rho](\mathbf{r}, t)$  is the exchange–correlation potential, which is henceforth assumed to depend only on the instantaneous density (adiabatic approximation). In a Gaussian-orbital basis, it is simpler to work with density matrices rather than KS orbitals, in which case the evolution of the electronic density is governed by the von Neumann equation:

$$i \frac{\partial \mathbf{P}'}{\partial t} = [\mathbf{F}'(t), \mathbf{P}'(t)] \quad (7)$$

where the prime notation denotes matrices in the orthogonal molecular orbital (MO) basis and unprimed denotes matrices in the atomic orbital (AO) basis. Note that in eq 7, all matrices are complex quantities. The Fock matrix  $\mathbf{F}(t)$  is computed in the AO basis similar to ground state DFT, with the important distinction that in the absence of Hartree–Fock exchange (e.g., pure DFT);  $\mathbf{F}(t)$  is real symmetric and only depends on the real part of  $\mathbf{P}(t)$ . If HF exchange is included (e.g., hybrid functionals), it becomes complex Hermitian (see ref 47 for details of the NWChem RT-TDDFT implementation, derivations, and references).

There are numerous approaches taken to propagate eq 7. In this study, we use a second order Magnus scheme, which is equivalent to an exponential midpoint propagator

$$\mathbf{P}'(t + \Delta t) = e^{-i\mathbf{F}'(t + \Delta t/2)\Delta t} \mathbf{P}'(t) e^{i\mathbf{F}'(t + \Delta t/2)\Delta t} \quad (8)$$

where we compute the Fock matrix at the future time via linear extrapolation from the previous two values, followed by iterative interpolation until converged. This approach is extremely stable, as it maintains the idempotency of the density matrix and yields order  $(\Delta t)^2$  accuracy. In practice, this allows for time steps on the order of  $\Delta t = 0.1 \text{ au} = 2.42 \times 10^{-3} \text{ fs}$  with a minimal loss of accuracy. The exponentiation of eq 8 is done via contractive power series, where the operator is first divided by  $2^m$  such that the norm of the scaled operator is less than 1, performing the power series (which is guaranteed to converge well numerically since it is contractive), then squaring the result  $m$  times to recover the result. All real-time TDDFT simulations here used a time step of  $\Delta t = 0.2 \text{ au} = 0.0048 \text{ fs}$  and ran up to  $1500 \text{ au} = 36.3 \text{ fs}$ , which corresponds to 7500 time steps.

To obtain spectroscopic information, the system is excited via a linearly polarized ( $x, y, z$ ) narrow Gaussian electric field kick, which adds to the Fock matrix via dipole coupling:

$$\mathbf{E}(t) = \kappa \exp[-(t - t_0)^2/2w^2] \hat{\mathbf{d}} \quad (9)$$

where  $\hat{\mathbf{d}} = \hat{x}, \hat{y}, \hat{z}$  is the polarization,  $\kappa$  is the field maximum (dimensions of electric field),  $t_0$  is the center of the pulse, and  $w$  is the width, which is typically  $\sim \Delta t$ . This induces all electronic modes simultaneously, and the Fourier transform of the resulting

time-dependent dipole moment yields the absorption spectrum for that polarization. The sum of the three spectra gives the full absorption. In the limit of a small electric field perturbation, real-time TDDFT and linear-response yield essentially identical spectroscopic results. Unlike LR-TDDFT, RT-TDDFT is also valid in the strong perturbation regime, but the studies presented here are all the weak-field type and thus comparable to LR-TDDFT. All kick-type results here used a kick with  $\kappa = 0.002 \text{ au} = 1.0 \text{ V/nm}$ ,  $t_0 = 3.0 \text{ au} = 0.07 \text{ fs}$ , and  $w = 0.2 \text{ au} = 0.0048 \text{ fs}$ .

The true power of RT-TDDFT, however, lies in direct modeling of the electron dynamics in response to a realistic stimulus, such as a laser tuned to resonance with a particular electronic transition. For example, to excite the system into a particular state of interest, it is simplest to use a Gaussian enveloped monochromatic laser pulse of the form:

$$\mathbf{E}(t) = \kappa \exp[-(t - t_0)^2/2w^2] \cos(\omega_0 t) \hat{\mathbf{d}} \quad (10)$$

where  $\omega_0$  is the driving frequency and  $w$  is broad enough to encapsulate at least a few oscillations. In this case, the charge density can be visualized in 4D (three space + time), which yields detailed insight into the fundamental nature of the excitation. This is especially important as an intuitive metric for characterizing charge transfer excitations, and when elucidating the mechanism of excitations. In this paper, RT-TDDFT is used as a visual tool to assign longitudinal and transverse excitations into two distinct classes (ionic vs covalent, respectively) and to study the physical origin of the red-shift with acene length.

**2.3. Time-Dependent Semiempirical Methods.** A well-known alternative to first-principles approaches is semiempirical methods (e.g., PM3<sup>28</sup> and ZINDO<sup>48</sup>) which can be extended to a time-dependent formalism.<sup>30</sup> A minimal valence basis set is used, so that there are only four orbitals for each carbon atom. Typically, the Fock matrix has the generic Hartree–Fock-like form:

$$F_{ij} = h_{ij} + \sum_{kl} v_{ijkl} P_{ij} \quad (11)$$

where  $h_{ij}$  and  $v_{ijkl}$  are semiempirical one-body and interaction parameters, respectively. Unlike Hartree–Fock and DFT, however, the interaction parameters are restricted to be at most two-center. The calculations are done in an atomic basis (rather than molecular orbital basis, which earlier TD-semiempirical methods use) so that the calculation of the Fock matrix scales like  $N^2$ , where  $N$  is the number of orbitals.

After the initial SCF solution labeled as  $P_0$ , the same von Neumann equation as in TDDFT (eq 7) is propagated. While the same real-time approach as in eq 8 could have been used, here, however a different algorithm is found to be more efficient. The algorithm has been covered recently (see ref 30), so it will only be briefly reviewed. Basically, the linear-response von Neumann operator is constructed:

$$LZ \equiv \frac{dZ}{dt} = -i \frac{[F(P_0 + \eta Z), P_0 + \eta Z] - [F(P_0), P_0]}{\eta} \quad (12)$$

for the deviation from the initial density matrix:

$$Z \equiv \mathbf{P} - P_0 \quad (13)$$

and  $\eta$  is a small parameter ensuring linearity. Then, the time-dependent dynamics are represented by writing a Chebyshev

algorithm for the propagator:

$$\mathbf{Z}(t) = e^{Lt} \mathbf{Z}_0 = \sum_n (2 - \delta_{n0}) J_n(t\Delta H) T_n \left( \frac{L}{\Delta H} \right) \mathbf{Z}_0 \quad (14)$$

where we introduced the Bessel and modified Chebyshev operators, with the latter propagated as

$$T_n \left( \frac{L}{\Delta H} \right) \mathbf{Z}_0 = 2 \frac{L}{\Delta H} T_{n-1} \left( \frac{L}{\Delta H} \right) \mathbf{Z}_0 + T_{n-2} \left( \frac{L}{\Delta H} \right) \mathbf{Z}_0 \quad (15)$$

and

$$\mathbf{Z}_0 = -i[\mathbf{D}, \mathbf{P}_0] \quad (16)$$

where  $\mathbf{D}$  is the dipole moment matrix.  $\Delta H$  is half the spectrum width, so that  $(\Delta H)^{-1}$  is the effective time-step; it is quite large (almost 0.4 au), so that the overall number of iterations required is quite small (a few thousands even without any signal processing approaches). This approach minimizes the number of matrix multiplications, which in semiempirical calculations are the most time-consuming steps (scales as  $N^3$  unless sparse matrix algorithms are used). Further savings are obtained by Fourier transforming the time-dependent Bessel function coefficients in eq 14 analytically, thereby reducing the required number of iterations. As with RT-TDDFT, spectroscopic information is obtained via kick-type excitations.

### 3. RESULTS

In this section, we present acene vertical excitation energies (VEEs) for a wide range of theories: coupled cluster (EOMCCSD, CR-EOMCCSD(T)), linear response TDDFT with a global hybrid functional (B3LYP<sup>49</sup>) and a variety of range-separated functionals (CAM-B3LYP,<sup>50</sup> LC-BLYP, LC- $\omega$ PBE,<sup>51</sup> BNL<sup>52</sup>), real-time TDDFT with the BNL functional, and two semiempirical methods (TD-ZINDO, TD-PM3). Before discussing results, it is important to note that vertical excitation energies, which correspond to the energy difference between ground and excited states without a change in geometry, cannot be directly measured experimentally (see ref 21). As a good approximation, VEEs can be measured experimentally via the locations of experimental UV–vis absorption peaks, but the accuracy of this approximation varies depending on state and molecule, with deviations typically on the order of a few tenths of an electronvolt. To ensure meaningful comparisons between the computed VEEs and experimental results, we use the corrected acene experimental values from Grimme and Parac<sup>20</sup> (see ref 53 for the original experimental results). In a nutshell, these incorporate adjustments to the  $L_a$  and  $L_b$  states computed from TDDFT (B3LYP/TZVP) excitation energies with fully optimized excited state geometries (calculated for acenes  $N = 2, 3, 4$ ; extrapolated to  $N = 5, 6, 7$ ). This somewhat accounts for geometry relaxation effects, but significant theory–experiment discrepancies still arise from basis set quality and the level of theory, specifically the treatment of correlation effects.

The  $L_a$  and  $L_b$  vertical excitation energies for the set of acenes are summarized in Table 1, along with the corrected experimental values, and the mean average error (MAE) from the experiment, for the full set of acenes for each approach. These VEEs (for a few representative theories) are plotted against acene size in Figure 4. Qualitatively speaking, all methods capture most of the gross features, including the red-shift of the  $L_b$  (longitudinal) state with acene length and the steeper red-shift of the  $L_a$  (transverse) state with acene length. However, there is only mixed success in

describing the important experimentally observed crossover of the lowest energy state from  $L_a \rightarrow L_b$  around anthracene; this is discussed in more detail below.

**3.1. Equation-of-Motion Coupled Cluster.** Overall, CR-EOMCCSD(T) has the best agreement with experimental energies, with a MAE of 0.07 eV for the  $L_a$  state and 0.06 eV for the  $L_b$ . Most importantly, CR-EOMCCSD(T) simultaneously describes both states well and captures the crossover at the right energy (near anthracene). That is, it predicts that  $L_a$  is lower in energy than  $L_b$  for naphthalene. They are roughly equal for anthracene, and  $L_b$  is lower afterwards (see Figure 2). In contrast to the experimental vertical excitation energies, the EOMCCSD and CR-EOMCCSD(T) approaches predict for anthracene the reversed ordering of the  $L_a$  and  $L_b$  states. The CR-EOMCCSD(T) excitation energy for the  $L_b$  state is located 0.1 eV below the one corresponding to the  $L_a$  state. Similar reverse ordering has been reported in the context of multireference Møller–Plesset (MRPT) theory<sup>54,55</sup> calculations for low-lying excited states of anthracene.<sup>16</sup> In the case of the MRPT approach, the 0.17 eV separation between  $L_b$  and  $L_a$  states is slightly larger than 0.1 eV obtained with the CR-EOMCCSD(T) method for POL1 basis set. The CC2 model,<sup>56</sup> which is an approximation to the EOMCCSD formalism, predicts the  $L_a$  state to the lowest state, and the calculated separation between  $L_a$  and  $L_b$  states is around 0.2 eV.

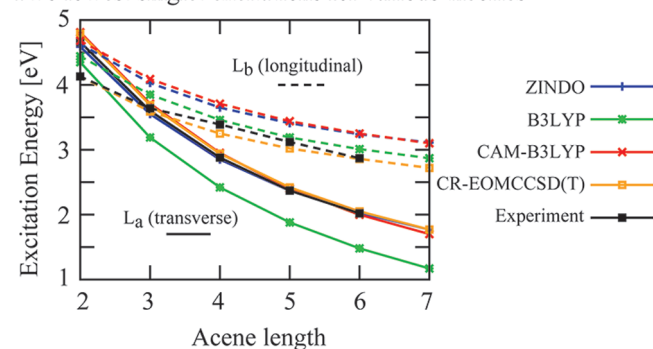
**3.2. Linear Response TDDFT.** The range-separation parameter for the CAM-B3LYP,<sup>50</sup> LC-BLYP, and BNL<sup>52</sup> functionals was taken to be 0.33 au<sup>-1</sup>; for LC- $\omega$ PBE,<sup>51</sup> it was 0.30 au<sup>-1</sup>. For the transverse charge-transfer-like  $L_a$  state (solid lines in Figure 2), all of the range-separated TDDFT results agree well with experimental results and EOMCC, with MAE typically around a few hundredths of an electronvolt. Real-time BNL results are essentially the same as the corresponding linear response ones, since the kick perturbation was small. Range-corrected TDDFT is less accurate for the  $L_b$  state, however, with MAEs of  $\sim 0.3$  eV, which is almost twice that of B3LYP. Thus, range-separated TDDFT excels at predicting the challenging charge-transfer-like  $L_a$  state, but using a range-separated functional significantly compromises the accuracy of the  $L_b$  state versus a global hybrid approach (e.g., B3LYP). To better understand the accuracy of RSH functionals, two versions of the CAM-B3LYP functional were studied: The first, denoted “CAM-B3LYP (I)”, has an asymptote of 0.65/ $r$  (i.e.,  $\alpha + \beta = 0.65$ ), while the second, denoted “CAM-B3LYP (II)”, has an asymptote of 1.0/ $r$ . The full Hartree–Fock asymptote in the exchange in CAM-B3LYP (II) improves the accuracy in the  $L_a$  state at a cost of slightly decreasing the accuracy of the  $L_b$  state. On another note, range-separated TDDFT correctly predicts the  $L_a \rightarrow L_b$  crossover (intersection of like-colored solid and dashed lines in Figure 2), albeit at a lower energy than the experiment. B3LYP, in contrast, fails to even qualitatively capture this crossover. In short, using range-separated functionals overcomes many of the failures of pure or hybrid DFT functionals in describing the transverse  $L_a$  state and the  $L_a \rightarrow L_b$  crossover, with overall accuracy rivaling that of CC2. The use of “tuned” RSHs, which has been pioneered by Baer and co-workers,<sup>57</sup> shows promise in further improving the accuracy of TDDFT for systems such as this.<sup>19</sup>

**3.3. Time-Dependent PM3 and ZINDO.** We performed time-dependent simulations with two typical semiempirical methods, PM3 and ZINDO. The latter is well-known to be better for spectra, as our results indicate. In order to parametrize the TD-ZINDO approach against the coupled CR-EOMCCSD(T) results for the charge-transfer-like  $L_a$ , we scaled down the strength of the  $\pi\pi'$  interaction potentials, as is commonly done in ZINDO. We

**Table 1.** The Two Lowest Singlet Excitation Energies in eV for the  $N = 2-7$  Series of Acenes for a Range of Theories and the Corresponding Mean Absolute Error (MAE) and Maximum Absolute Error (XAE) in eV from the Experimental Values<sup>a</sup>

N	ZINDO	ZINDO	CAM-		CAM-	LC-	LC-	BNL		CC2 <sup>20</sup>	EOM-	CR-EOM-	expt <sup>20</sup>	
	PM3	(I)	(II)	B3LYP	B3LYP (I)	B3LYP (II)	BLYP	$\omega$ PBE	BNL	BNL (real-time)	CCSD	CCSD(T)		
	L <sub>a</sub> state (transverse; bright)													
2	3.50	4.23	4.59	4.35	4.64	4.81	4.77	4.77	4.86	4.79	4.88	5.09	4.79	4.66
3	2.94	3.30	3.55	3.19	3.51	3.71	3.66	3.66	3.72	3.68	3.69	4.00	3.69	3.60
4	2.53	2.67	2.85	2.42	2.75	2.95	2.91	2.90	2.94	2.91	2.90	3.25	2.94	2.88
5	2.22	2.23	2.37	1.88	2.21	2.40	2.37	2.37	2.39	2.41	2.35	2.72	2.42	2.37
6	1.99	1.92	2.03	1.48	1.82	2.00	1.99	1.99	2.00	1.96	1.95	2.34	2.05	2.02
7	1.81	1.68	1.77	1.17	1.52	1.70	1.69	1.70	1.70	1.69	1.60	2.05	1.77	—
MAE	0.47	0.24	0.03	0.44	0.12	0.08	0.05	0.04	0.08	0.07	0.08	0.37	0.07	—
XAE	1.16	0.43	0.07	0.54	0.20	0.15	0.11	0.11	0.20	0.13	0.22	0.43	0.13	—
	L <sub>b</sub> state (longitudinal; dim)													
2	3.34	4.21	4.63	4.44	4.59	4.68	4.58	4.58	4.64	4.61	4.46	4.43	4.13	4.13
3	2.91	3.67	4.03	3.85	4.02	4.09	4.02	4.02	4.07	4.03	3.89	3.90	3.59	3.64
4	2.62	3.32	3.65	3.46	3.64	3.71	3.65	3.65	3.70	3.68	3.52	3.54	3.25	3.39
5	2.42	3.10	3.41	3.19	3.38	3.44	3.39	3.40	3.44	3.42	3.27	3.30	3.02	3.12
6	2.27	2.96	3.24	3.01	3.20	3.25	3.21	3.22	3.26	3.23	3.09	3.12	2.86	2.87
7	2.15	2.82	3.11	2.87	3.06	3.10	3.07	3.09	3.12	3.03	2.97	2.99	2.74	—
MAE	0.72	0.06	0.36	0.16	0.34	0.40	0.34	0.34	0.39	0.36	0.22	0.23	0.06	—
XAE	0.79	0.09	0.50	0.31	0.46	0.55	0.45	0.45	0.51	0.48	0.33	0.30	0.14	—

<sup>a</sup>The L<sub>a</sub> state corresponds to a transverse excitation with high oscillator strength (bright) and the L<sub>b</sub> state to a longitudinal excitation with low oscillator strength (dim). All TDDFT results are linear response unless noted otherwise.

**Figure 2.** Comparison between the two lowest singlet excitation energies of the set of acenes for a selection of theories, along with the experimental values. The solid lines correspond to the L<sub>a</sub> (transverse) excitation and the dashed lines to the L<sub>b</sub> (longitudinal) excitation.

**Figure 2.** Comparison between the two lowest singlet excitation energies of the set of acenes for a selection of theories, along with the experimental values. The solid lines correspond to the L<sub>a</sub> (transverse) excitation and the dashed lines to the L<sub>b</sub> (longitudinal) excitation.

found that a scaling factor of 0.64, which we denote “ZINDO (II)”, yielded the best fit, compared to the stock scaling factor of 0.70 (denoted “ZINDO (I)”). In the case of the general ZINDO (I), the L<sub>a</sub> is fairly poorly described (MAE of 0.24 eV), whereas the longitudinal L<sub>b</sub> is quite well described, akin to the B3LYP results. The L<sub>a</sub>-tuned ZINDO (II), however, is extremely accurate for the L<sub>a</sub> state, but as with range-separated TDDFT, the corresponding accuracy in the L<sub>b</sub> state suffers. One drawback of ZINDO, however, is that it fails to properly capture the crossover. ZINDO (I) predicts that L<sub>a</sub> and L<sub>b</sub> are roughly equal in energy at  $N = 2$ , whereas ZINDO (II) incorrectly predicts that L<sub>b</sub> is always higher in energy than L<sub>a</sub>. Of course, the excellent quality of ZINDO (II) results for the L<sub>a</sub> are a consequence of being fit to this particular

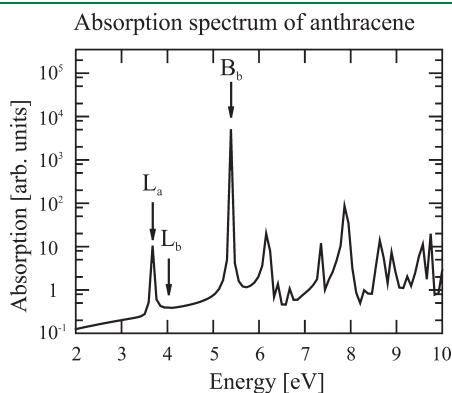
state, but it is still quite remarkable that with a single parameter it is possible to simultaneously fit six molecules so well. These results suggest that carefully parametrized semiempirical approaches are an excellent tool for modeling excitations in large polyaromatic hydrocarbons, where large system size makes coupled cluster, or even TDDFT, unfeasible.

### 3.4. Real-Time Visualization of the Excited Charge Density.

Next, to gain insight into the nature of the excitations, we present real-time real-space visualization of the excited state charge density for the (transverse) L<sub>a</sub> state. The (longitudinal) L<sub>b</sub> state has too small an oscillator strength to visualize clearly, so the major bright longitudinal UV B<sub>b</sub> absorption (see Figure 3) was chosen as an illustrative analogue (note this peak is not compared in Table 1). As before, the system was described using the BNL functional, and for speed the smaller 6-31G\*\* basis set was used instead of POL1. The spectra of the acenes with 6-31G\*\* basis sets were extremely similar to the POL1 spectra, save a slight blue-shift due to the smaller basis.

Figure 4 shows real-time TDDFT snapshots of the deviation of the charge density from the ground state for anthracene and heptacene after resonant excitation to the L<sub>a</sub> state. Unlike plots of molecular orbitals, which are strictly ground state quantities, Figure 4 corresponds to the actual charge density dynamics resulting from an excitation. For the longitudinal excitation (top), blue isosurfaces correspond to positive charge density deviation from the ground state,  $\rho(\mathbf{r},t) - \rho(\mathbf{r},t=0) = 10^{-6} \text{ \AA}^{-3}$ , and red isosurfaces to the corresponding negative deviation. In the transverse excitation (bottom), the isosurface values were  $10^{-7} \text{ \AA}^{-3}$ . The two excitations were induced via longitudinal or transverse polarized enveloped laser pulses (see eq 10), with  $w = 2\pi/\omega_0$  and  $t_0 = 5w$ ; the values of the driving frequencies  $\omega_0$  are shown in Figure 4, along with the time taken for half an oscillation to occur.

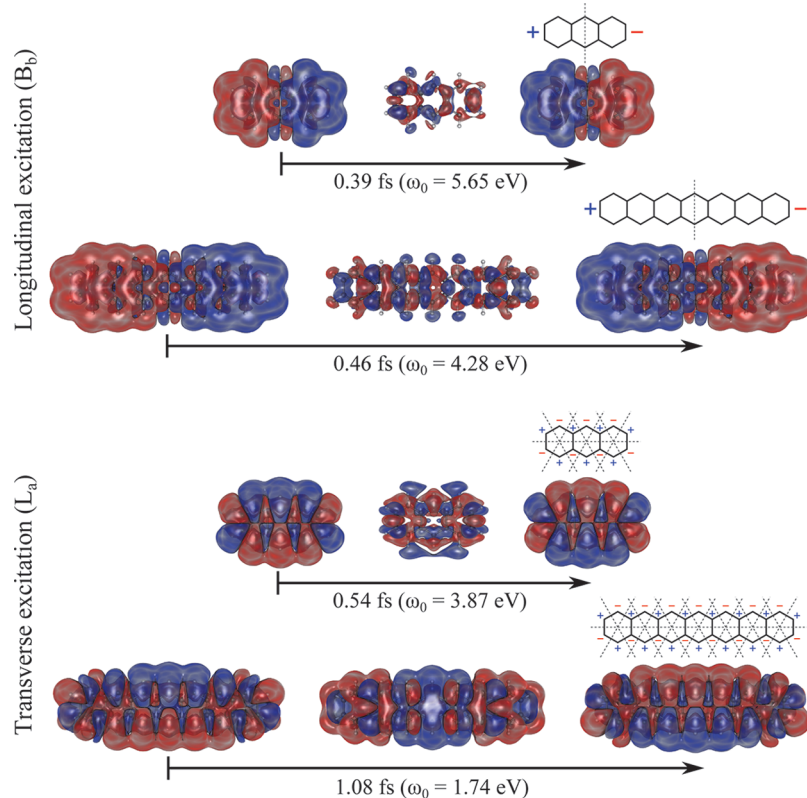
Two distinct mechanisms of excitation are clearly visible in Figure 4. The longitudinally excited charge density ( $B_b$  state; top) sloshes back and forth along the  $\pi^*$  orbitals along the acene backbone; at the extrema, the charge density has piled up at one end of the molecule, with corresponding depletion (hole) on the opposite end. After transverse excitation ( $L_a$  state; bottom), however, the density is driven from delocalized  $\pi$  orbitals across



**Figure 3.** Absorption spectrum of anthracene ( $N = 3$ ) obtained via RT-TDDFT (POLI/BNL). The bright  $L_a$  and dim  $L_b$  peaks correspond to transverse and longitudinal excitations, respectively. The intensely bright longitudinal UV  $B_b$  peak is visualized in Figure 4 but not compared in Table 1.

the acene and forced to populate the orbitals above and below the C–H bonds, which leads to alternating “fingers” of accumulated charge, and thus alternating  $\cdots C^{\delta+}C^{\delta-}C^{\delta+}C^{\delta-} \cdots$  atoms along the acene. In a valence bond picture, this is an ionic-like excitation, in agreement with previous analyses.<sup>18,20</sup> The intramolecular charge-transfer-like character (or charge transfer in disguise) is not due to a long-range pileup of charge but instead arises from this ionic-like character. Here, range-separated functionals perform well because they are able to capture interaction between these regions of alternating charge and hole. This is related to Kuritz et al.’s discussion, where a state is characterized as charge-transfer-like on the basis of minimal overlap of auxiliary orbitals.<sup>19</sup>

RT-TDDFT can also shed light on the origin of the red shifts. As the acenes increase in length, the time taken to oscillate increases (frequency decreases) for both the transverse and longitudinal excitations. Although not immediately obvious, the red shifts of both excitations can be rationalized in a similar way. The simplest physical description for this comes from the perimeter free electron orbital (PFEO) theory,<sup>1,58</sup> which models the  $\pi$  electrons as being confined in an oval-shaped infinite potential with no other electron–nuclear or electron–electron interactions. This leads to a particle-on-a-ring wave function for each  $\pi$  electron; a particular electronic state is then characterized by the total ring quantum number  $Q$ , which is the sum of the individual ring quantum numbers. The number of nodal planes



**Figure 4.** Real-time TDDFT (6-31G\*\*/BNL) isosurface snapshots of the deviation of the charge density from the ground state for anthracene ( $N = 3$ ) and heptacene ( $N = 7$ ), after resonant excitation (frequencies shown in eV). Positive deviation (more charge density than in the ground state) is shown in blue, while negative deviation (less charge density than ground state) is shown in red. The time for a charge oscillation (half period) is shown in femtoseconds. The longitudinal  $B_b$  state (note: not compared in Table 1) is covalent in nature. The ionic character of the  $L_a$  state is clearly visible from the alternation of charge buildup above/below the C–H bonds and charge depletion on the carbon atoms between. The corresponding perimeter free electron orbital (PFEO) theory structures are shown, confirming that the excited state densities at the oscillation (half period) maxima are extremely similar to those arising from  $\pi$  electrons confined to a ring. The densities were visualized using Blender.<sup>59</sup>

for a particular state is then  $Q$ , with alternating positive and negative charge buildup at each antinode. This is clearly visible in Figure 4, where the charge density deviations at the maxima of the oscillations (i.e., the excited electronic states) directly match up to the PFE0 predictions. In anthracene, for example, the excited state charge density of the  $B_b$  state corresponds to a  $Q = 1$  state (one node; high longitudinal dipole moment), whereas the  $L_a$  state corresponds to  $Q = 7$  (seven nodes; low but nonzero transverse dipole moment). The transition to  $Q = 7$  ( $L_a$ ) requires less energy than that to  $Q = 1$  ( $B_b$ ), which is a consequence of Hund's rule.<sup>1</sup> Larger acenes have larger circumferences, and thus their excitation energies are red-shifted.

#### 4. CONCLUSIONS

In summary, we have computed the  $L_a$  and  $L_b$  vertical excitation energies for the acenes ranging from anthracene to heptacene, using a broad spectrum of excited-state theoretical approaches. High accuracy coupled cluster calculations (CR-EOMCCST(T)) agree extremely well with experimental results for both states and thus serve as a baseline for validating the lower level theories. Global hybrid TDDFT (e.g. B3LYP) performs poorly for the  $L_a$  state, as expected, whereas range-separated hybrid (RSH) TDDFT (e.g. CAM-B3LYP, LC-BLYP, etc) better describes the ionic  $L_a$  state, at a cost of lost accuracy for the  $L_b$  state. Real-time RSH TDDFT visualization shows that the excited state charge densities are consistent with the predictions of perimeter free electron orbital (PFE0) theory, and the red shifts of the excitations are due to particle-on-a-ring-like confinements. For the semiempirical methods, with proper parametrization, ZINDO rivals range-separated hybrids in accuracy, at a fraction of the computational cost. This suggests a multitiered approach to modeling complicated acene derivatives, as well films and crystals of these molecules: high accuracy coupled cluster calculations validate RSH TDDFT calculations on small (perhaps pairs of) molecules, which in turn enables careful parametrization of semiempirical calculations capable of modeling large systems.

#### AUTHOR INFORMATION

##### Corresponding Author

\*E-mail: kenneth.lopata@pnnl.gov; dxn@chem.ucla.edu; niri.govind@pnnl.gov; karol.kowalski@pnnl.gov.

#### ACKNOWLEDGMENT

A portion of the research was performed using EMSL, a national scientific user facility sponsored by the U.S. Department of Energy's Office of Biological and Environmental Research and located at Pacific Northwest National Laboratory (PNNL). PNNL is operated for the Department of Energy by the Battelle Memorial Institute under Contract DE-AC06-76RLO-1830. K.L. acknowledges the William Wiley Postdoctoral Fellowship from EMSL. K.K. and N.G. acknowledge support from the Extreme Scale Computing Initiative, a Laboratory Directed Research and Development Program at Pacific Northwest National Laboratory. D.N. and R.R. gratefully acknowledge support by DOE-EFRC.

#### REFERENCES

(1) Pope, M.; Swenberg, C. E. *Electronic processes in organic crystals and polymers*, 2nd ed.; Oxford University Press: Oxford, U.K., 1999; Chapter 1, pp 7–12.

- (2) Picciolo, L. C.; Murata, H.; Kafafi, Z. H. *Appl. Phys. Lett.* **2001**, *78*, 2378.
- (3) Wolak, M. A.; Jang, B. B.; Palilis, L. C.; Kafafi, Z. H. *J. Phys. Chem. B* **2004**, *108*, 5492–5499.
- (4) Xu, Q.; Duong, H. M.; Wudl, F.; Yang, Y. *Appl. Phys. Lett.* **2004**, *85*, 3357.
- (5) Lin, C. Y.; Wang, Y. C.; Hsu, S. J.; Lo, C. F.; Diao, E. W. *J. Phys. Chem. C* **2009**, *114*, 687–693.
- (6) Lloyd, M. T.; Mayer, A. C.; Subramanian, S.; Mourey, D. A.; Herman, D. J.; Bapat, A. V.; Anthony, J. E.; Malliaras, G. G. *J. Am. Chem. Soc.* **2007**, *129*, 9144–9149.
- (7) Jiang, Y.; Okamoto, T.; Beceril, H. A.; Hong, S.; Tang, M. L.; Mayer, A. C.; Parmer, J. E.; McGehee, M. D.; Bao, Z. *Macromolecules* **2010**, *42*, 6361–6367.
- (8) Sheraw, C. D.; Zhou, L.; Huang, J. R.; Gundlach, D. J.; Jackson, T. N.; Kane, M. G.; Hill, I. G.; Hammond, M. S.; Campi, J.; Greening, B. K.; Francl, J.; West, J. *Appl. Phys. Lett.* **2002**, *80*, 1088–1090.
- (9) Merlo, J. A.; Newman, C. R.; Gerlach, C. P.; Kelley, T. W.; Muires, D. V.; Fritz, S. E.; Toney, M. F.; Frisbie, C. D. *J. Am. Chem. Soc.* **2005**, *127*, 3997–4009.
- (10) Tang, M. L.; Reichardt, A. D.; Miyaki, N.; Stoltenberg, R. M.; Bao, Z. *J. Am. Chem. Soc.* **2008**, *130*, 6064–6065.
- (11) Lin, Y. Y.; Gundlach, D. J.; Nelson, S. F.; Jackson, T. N. *Electron Device Lett., IEEE* **1997**, *18*, 606–608.
- (12) Klauk, H.; Halik, M.; Zschieschang, U.; Schmid, G.; Radlik, W.; Weber, W. *J. Appl. Phys.* **2002**, *92*, 5259.
- (13) Kim, G. H.; Shtein, M.; Pipe, K. P. *Appl. Phys. Lett.* **2011**, *98*, 093303.
- (14) Anthony, J. E. *Chem. Rev.* **2006**, *106*, 5028–5048.
- (15) Anthony, J. E. *Angew. Chem., Int. Ed.* **2008**, *47*, 452–483.
- (16) Kawashima, Y.; Hashimoto, T.; Nakano, H.; Hirao, K. *Theor. Chem. Acc.* **1999**, *102*, 49–64.
- (17) Wong, B. M.; Hsieh, T. H. *J. Chem. Theory Comput.* **2010**, *6*, 3704–3712.
- (18) Richard, R. M.; Herbert, J. M. *J. Chem. Theory Comput.* **2011**, *7*, 1296–1306.
- (19) Kuritz, N.; Stein, T.; Baer, R.; Kronik, L. *J. Chem. Theory Comput.* **2011**, *7*, 2408–2415.
- (20) Grimme, S.; Parac, M. *ChemPhysChem* **2003**, *4*, 292–295.
- (21) Bak, K. L.; Koch, H.; Oddershede, J.; Christiansen, O.; Sauer, S. P. A. *J. Chem. Phys.* **2000**, *112*, 4173–4185.
- (22) Peach, M. J. G.; Benfield, P.; Helgaker, T.; Tozer, D. J. *J. Chem. Phys.* **2008**, *128*, 044118.
- (23) Huang, L.; Rocca, D.; Baroni, S.; Gubbins, K. E.; Nardelli, M. B. *J. Chem. Phys.* **2009**, *130*, 194701.
- (24) Runge, E.; Gross, E. K. U. *Phys. Rev. Lett.* **1984**, *52*, 997.
- (25) Petersilka, M.; Gossmann, U. J.; Gross, E. K. U. *Phys. Rev. Lett.* **1996**, *76*, 1212–1215.
- (26) Casida, M. E. In *Recent Advances in Density Functional Methods*; Chong, D. P., Ed.; World Scientific Publishing: River Edge, NJ, 1995; Vol. 1, Chapter 5, pp 155–192.
- (27) Valiev, M.; Bylaska, E. J.; Govind, N.; Kowalski, K.; Straatsma, T. P.; Van Dam, H. J. J.; Wang, D.; Nieplocha, J.; Apra, E.; Windus, T. L.; de Jong, W. A. *Comput. Phys. Commun.* **2010**, *181*, 1477–1489.
- (28) Stewart, J. J. P. *J. Comput. Chem.* **1989**, *10*, 209–220.
- (29) MOPAC 6.0. <http://ccl.net/cca/software/MS-WIN95-NT/mopac6/index.shtml> (accessed September 2011).
- (30) Bartell, L. A.; Wall, M. R.; Neuhauser, D. *J. Chem. Phys.* **2010**, *132*, 234106.
- (31) Zerner, M. C. et al. *ZINDO-MN*, version 2011; Quantum Theory Project, University of Florida: Gainesville, FL; Department of Chemistry, University of Minnesota: Minneapolis, MN, 2011.
- (32) Hirata, S.; Head-Gordon, M. *Chem. Phys. Lett.* **1999**, *302*, 375–382.
- (33) Bartlett, R. J.; Musiał, M. *Rev. Mod. Phys.* **2007**, *79*, 291–352.
- (34) Kowalski, K.; Krishnamoorthy, S.; Villa, O.; Hammond, J. R.; Govind, N. *J. Chem. Phys.* **2010**, *132*, 154103.
- (35) Watts, J. D.; Bartlett, R. J. *J. Chem. Phys. Lett.* **1995**, *233*, 81–87.



- (36) Watts, J. D.; Bartlett, R. J. *Chem. Phys. Lett.* **1996**, *258*, 581–588.
- (37) Christiansen, O.; Koch, H.; Jørgensen, P. *J. Chem. Phys.* **1996**, *105*, 1451.
- (38) Hirata, S.; Nooijen, M.; Grabowski, I.; Bartlett, R. J. *J. Chem. Phys.* **2001**, *114*, 3919.
- (39) Shiozaki, T.; Hirao, K.; Hirata, S. *J. Chem. Phys.* **2007**, *126*, 244106.
- (40) Manohar, P. U.; Krylov, A. I. *J. Chem. Phys.* **2008**, *129*, 194105.
- (41) Kowalski, K.; Piecuch, P. *J. Chem. Phys.* **2004**, *120*, 1715.
- (42) Kowalski, K.; Piecuch, P. *J. Chem. Phys.* **2001**, *115*, 2966.
- (43) Włoch, M.; Lodriguito, M. D.; Piecuch, P.; Gour, J. R. *Mol. Phys.* **2006**, *104*, 2991.
- (44) Piecuch, P.; Gour, J. R.; Włoch, M. *Int. J. Quantum Chem.* **2009**, *109*, 3268–3304.
- (45) Kowalski, K. *J. Chem. Phys.* **2009**, *130*, 194110.
- (46) Raghavachari, K.; Trucks, G. W.; Pople, J. A.; Head-Gordon, M. *Chem. Phys. Lett.* **1989**, *157*, 479–483.
- (47) Lopata, K.; Govind, N. *J. Chem. Theory Comput.* **2011**, *7*, 1344–1355.
- (48) Anderson, W. P.; Cundari, T. R.; Zerner, M. C. *Int. J. Quantum Chem.* **1991**, *39*, 31–45.
- (49) Stephens, P. J.; Devlin, F. J.; Chabalowski, C. F.; Frisch, M. J. *J. Phys. Chem.* **1994**, *98*, 11623–11627.
- (50) Yanai, T.; Tew, D. P.; Handy, N. C. *Chem. Phys. Lett.* **2004**, *393*, 51–57.
- (51) Rohrdanz, M. A.; Martins, K. M.; Herbert, J. M. *J. Chem. Phys.* **2009**, *130*, 054112.
- (52) Baer, R.; Neuhauser, D. *Phys. Rev. Lett.* **2005**, *94*, 043002.
- (53) Lambert, W. R.; Felker, P. M.; Syage, J. A.; Zewail, A. H. *J. Chem. Phys.* **1984**, *81*, 2195.
- (54) Hirao, K. *Chem. Phys. Lett.* **1992**, *190*, 374–380.
- (55) Hirao, K. *Chem. Phys. Lett.* **1992**, *196*, 397–403.
- (56) Christiansen, O.; Koch, H.; Jørgensen, P. *Chem. Phys. Lett.* **1995**, *243*, 409–418.
- (57) Baer, R.; Livshits, E.; Salzner, U. *Annu. Rev. Phys. Chem.* **2010**, *61*, 85–109.
- (58) Platt, J. R. *J. Chem. Phys.* **1949**, *17*, 484–495.
- (59) *Blender*; The Blender Foundation: Amsterdam, The Netherlands, 2010.

# Dual Fluorescence of Fluorazene in Solution: A Computational Study

Ignacio Fdez. Galván,\* M. Elena Martín, Aurora Muñoz-Losa, and Manuel A. Aguilar

Química Física, Edif. José María Viguera Lobo, Universidad de Extremadura, Avda. de Elvas s/n, 06071 Badajoz, Spain

**ABSTRACT:** The fluorazene molecule presents dual fluorescence in polar solvents. Its absorption and emission properties in gas phase and in acetonitrile solution have been studied theoretically using the complete active space second-order perturbation//complete active space self-consistent field quantum methodology and average solvent electrostatic potential from molecular dynamics for the solvent effects. In gas phase, two optimized excited-state geometries were obtained, one of them corresponds to a local excitation (LE), and the other is an intramolecular charge transfer (ICT) and lies higher in energy. In acetonitrile solution, a second ICT structure where the molecule remains planar is found, and the energy differences are reduced. Fluorescence energies from LE and the planar ICT have a good agreement with the experimental bands, but emission from the bent ICT has too low an energy.

## 1. INTRODUCTION

A significant number of organic molecules with electron-donating and -withdrawing groups when immersed in polar solvents display what is known as dual fluorescence. In nonpolar solvents, as for most molecules, the fluorescence spectrum exhibits a single band, whose maximum is only slightly shifted as the solvent polarity increases, this is called the “normal” band. In polar solvents, a second fluorescence band appears in the spectrum, and the position of this second band varies more significantly with the solvent polarity, this is called the “anomalous” band. The relative intensity of the anomalous band increases with the polarity, so that in highly polar solvents, the normal band can disappear, and only the anomalous band is observed. This dual fluorescence phenomenon has been profusely studied in the literature since its discovery half a century ago.<sup>1,2</sup> Most of these studies are focused on the prototype molecule 4-(*N,N*-dimethylamino)benzonitrile (DMABN) or its derivatives, including experimental investigations<sup>3–8</sup> and theoretical works.<sup>9–16</sup> It was suggested early on that the origin of the anomalous fluorescence band is the existence of an intramolecular charge transfer (ICT) excited state, which is not normally accessible in nonpolar solvents but which is stabilized in polar solvents and can thus compete with the state responsible for the normal band, usually called a local excitation (LE) state.

The validity of this explanation for the dual fluorescence is still generally accepted. There is, however, a continuing controversy between the various groups that have investigated this subject, regarding the nature and the geometry of the ICT state, the mechanism through which the LE and ICT states are formed, the possible existence of further intermediate states, and practically every other detail of the dual fluorescence phenomenon.

Probably the most accepted models for the dual fluorescence in the DMABN molecule, and related compounds are the ones known as twisted ICT (TICT) and planar ICT (PICT). These models propose an ICT state where the donor and acceptor groups adopt, respectively, a perpendicular or coplanar conformation. Experimental evidence favoring one model or the other is usually derived from comparison of the properties of compounds with different geometric constraints and substituents.

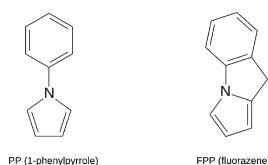
For example, compounds, like 3,5-dimethyl-4-(*N,N*-dimethylamino)benzonitrile, where the dimethylamino group is forced to be twisted, display only the ICT band in fluorescence, suggesting a TICT is responsible for the band. Other compounds where the twisting is hindered (like 6-cyano-1,2,3,4-tetrahydroquinoline, NTC6) can present dual fluorescence, which points to a PICT state. These apparently contradictory conclusions possibly indicate that the two models are not exclusive, and each particular system will favor one of them.

In recent years, a pair of closely related molecules has been studied for their dual fluorescence properties, see Figure 1. The two rings in 1-phenylpyrrole (PP) can freely rotate around the middle bond, while the methylene bridge in fluorazene (FPP) effectively locks the rings in an almost planar conformation. Both molecules display a very similar photophysical behavior, and in particular, both show dual fluorescence in polar solvents. One of the differences between the two molecules is that apparently FPP presents enhanced ICT emission compared to PP: the ICT band appears in less polar solvents, and its quantum yield is higher. This fact naturally leads to the conclusion that the PICT model applies better to these molecules.<sup>17,18</sup> However, most theoretical calculations predict a twisted structure for the ICT state of PP,<sup>19–23</sup> which seems unsatisfactory.

In a previous work,<sup>24</sup> we carried out a theoretical study on the absorption and fluorescence properties of the PP molecule, both in the gas phase and in acetonitrile solution. Our conclusion was that there are different molecular structures accessible for the PP and that the twisting of the rings is not necessary for reaching the emitting ICT state. In this work we present a similar study for the FPP molecule, where its electronic states are described with a multiconfigurational quantum method, and we used an explicit model of atomic detail for the solvent. By examining the relative energies, geometries, and emission energies of the different electronic states, we expect to obtain meaningful conclusions for the study of this system.

Received: July 28, 2011

Published: September 27, 2011



**Figure 1.** Two related compounds with dual fluorescence.

## 2. METHODS AND DETAILS

Solvent effects on the FPP UV–vis spectra were calculated with the average solvent electrostatic potential from molecular dynamics (ASEP/MD) method. This is a sequential quantum mechanics/molecular mechanics (QM/MM) method implementing the mean field approximation. It combines, alternately, a high-level QM description of the solute with a classical MM description of the solvent. One of its main features is the fact that the solvent effect is introduced into the solute’s wave function as an average perturbation. Details of the method have been described in previous papers,<sup>25–27</sup> so here we will only present a brief outline.

As mentioned above, ASEP/MD is a method combining QM and MM techniques, with the particularity that full QM and molecular dynamics (MD) calculations are alternated and not simultaneous. During the MD simulations, the intramolecular geometry and charge distribution of all molecules are considered fixed. From the resulting simulation data, the average electrostatic potential generated by the solvent on the solute (ASEP) is obtained. This potential is introduced as a perturbation into the solute’s quantum mechanical Hamiltonian, and by solving the associated Schrödinger equation, one gets a new charge distribution for the solute, which is used in the next MD simulation. This iterative process is repeated until the electron distribution of the solute and the solvent structure around it are mutually equilibrated.

The ASEP/MD framework can also be used to optimize the geometry of the solute molecule.<sup>28</sup> At each step of the ASEP/MD procedure, the gradient and Hessian on the system’s free energy surface (including the van der Waals contribution) can be obtained, and thus they can be used to search for stationary points on this surface by some optimization method. In the computation of the gradient and Hessian, the free energy gradient method<sup>29</sup> is used, with the incorporation of the mean field approximation to reduce the number of quantum calculations needed. In this way, after each MD simulation, the solute geometry is optimized within the fixed “average” solvent structure by using the free energy derivatives. In the next MD simulation, the new solute geometry and charge distribution are used. This approach allows the optimization of the solute geometry in parallel to the solvent structure.

For calculating transition energies, nonequilibrium solvation is assumed. The iterative process is only performed on the initial state of the transition (the ground state for absorption, the excited state for emission), i.e., the atomic charges for the MD and the energy derivatives for the geometry optimization of the solute are calculated with the initial state’s wave function. Then, with a frozen solvent model, the energies of the final states are obtained.

Once the different solute electronic states and the solvent structure around them have been optimized and equilibrated, the free energy differences between those states can be calculated, within the ASEP/MD framework, making use of the free energy

perturbation method.<sup>30,31</sup> The expression we use to calculate the free energy difference between two species in equilibrium in solution,  $\Delta G$ , is

$$\Delta G = \Delta E + \Delta G_{\text{int}} + \Delta V \quad (1)$$

where  $\Delta E$  is the difference in the internal quantum energy of the solute between the two species,  $\Delta G_{\text{int}}$  is the difference in the solute–solvent interaction energy, which is calculated classically with the free energy perturbation (FEP) method, and  $\Delta V$  is a term that includes the difference in the zero point energy and entropic contributions of the solute. The last term,  $\Delta V$ , is normally evaluated by applying the harmonic approximation to the vibrational modes of the solute in solution, and it needs the information provided by the Hessian matrix. In this work, obtaining an accurate enough Hessian matrix required computational resources that were too large, and we decided to approximate the results by neglecting this term. It must be noted that this  $\Delta V$  term refers only to the internal nuclear degrees of freedom of the solute; free energy contributions from the solvent around the solute are properly accounted for in the  $\Delta G_{\text{int}}$  term.

**2.1. Computational Details.** The quantum calculations on the solute molecule were done with the complete active space self-consistent field (CASSCF) method,<sup>32</sup> using the cc-pVDZ basis set and aug-cc-pVDZ in some selected cases. The active orbitals were the 6  $\pi$  and  $\pi^*$  valence orbitals of the phenyl ring, plus the 5  $\pi$  and  $\pi^*$  of the pyrrole ring, and 12 electrons were included in these orbitals, for a (12,11) total active space. All calculations were performed using a state average (SA) of the first five singlet states, with equal weights. It is known that, in order to obtain accurate transition energies, it is necessary to include the dynamic electron correlation in the quantum calculations, which we did with the complete active space second order perturbation (CASPT2) method,<sup>33,34</sup> using the SA(5)-CASSCF(12,11) wave functions as a reference. An ionization potential–electron affinity (IPEA) shifted zeroth-order Hamiltonian has been proposed for CASPT2 calculations,<sup>35</sup> which is supposed to reduce systematic overstabilization errors in open-shell systems (as is the case of the excited states studied here). We did all CASPT2 with the proposed IPEA shift of 0.25  $E_h$  (CASPT2(0.25)) as well as with no IPEA shift (CASPT2(0.00)). To minimize the appearance of intruder states, an additional imaginary shift of 0.1  $iE_h$  was used. No symmetry was assumed in any case.

The MD simulations were carried out with rigid molecules, with acetonitrile ( $\text{CH}_3\text{CN}$ ) as a solvent. Lennard-Jones parameters and solvent atomic charges were taken from the optimized potentials for liquid simulations, all atoms (OPLS-AA) force field,<sup>36</sup> solute atomic charges were calculated from the quantum calculations through a least-squares fit to the electrostatic potential obtained at the points where the solvent charges are located. The geometry of acetonitrile was optimized with Becke’s three-parameter Lee–Yang–Parr density functional (B3LYP) and the 6-311G\*\* basis set. A total of 375  $\text{CH}_3\text{CN}$  molecules and the solute were included at the experimental solvent density (779.3  $\text{kg}/\text{m}^3$ ). Periodic boundary conditions were applied, and spherical cutoffs were used to truncate the interatomic interactions at 12.75 Å. Long-range interactions were calculated using the Ewald sum technique. The temperature was fixed at 298.15 K by using the Nosé–Hoover thermostat. A time step of 0.5 fs was used during the simulations, and each one was run for 50 ps after 25 ps of equilibration.

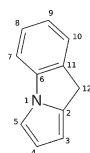


Figure 2. Atom numbering of the FPP molecule.

**Table 1.** Definition of geometric parameters for the FPP molecule.  $d$  is a bond length,  $a$  a bond angle, and  $D$  a dihedral angle. Point A is placed at  $C_6 + n(\text{Ph})$ , where  $n(\text{Ph})$  is the normal vector of the best-fit plane for the phenyl carbon atoms, except  $C_6$ . Point B is defined similarly for the N1 atom and the pyrrole ring (including the nitrogen)

$$\overline{\text{Ph}} = \frac{1}{6}(d(C_6C_7) + d(C_7C_8) + d(C_8C_9) + d(C_9C_{10}) + d(C_{10}C_{11}) + d(C_{11}C_6))$$

$$\overline{\text{Py}} = \frac{1}{5}(d(N_1C_2) + d(C_2C_3) + d(C_3C_4) + d(C_4C_5) + d(C_5N_1))$$

$$Q(\text{Ph}) = \frac{1}{4}(d(C_6C_7) + d(C_8C_9) + d(C_9C_{10}) + d(C_{11}C_6)) - \frac{1}{2}(d(C_7C_8) + d(C_{10}C_{11}))$$

$$Q'(\text{Ph}) = \frac{1}{4}(d(C_6C_7) + d(C_7C_8) + d(C_9C_{10}) + d(C_{10}C_{11})) - \frac{1}{2}(d(C_8C_9) + d(C_{11}C_6))$$

$$Q(\text{Py}) = \frac{1}{3}(d(N_1C_2) + d(C_3C_4) + d(C_5N_1)) - \frac{1}{2}(d(C_2C_3) + d(C_4C_5))$$

$$\text{Ph} - \text{Py} = d(N_1C_6)$$

$$\phi = a(\text{AC}_6\text{N}_1) - 90^\circ$$

$$\psi = a(\text{BN}_1\text{C}_6) - 90^\circ$$

$$\theta = D(\text{AC}_6\text{N}_1\text{B})$$

At each step of the ASEP/MD procedure, 500 configurations evenly distributed from the MD run were used to calculate the ASEP. The charges from each solvent molecule were kept explicitly whenever the molecule's center of mass was closer than  $9 a_0$  to any solute nucleus; the effect of the farther molecules was included in an additional shell of fitted charges. Each ASEP/MD run was continued until the energies and solute geometry and charges were stabilized for at least five iterations, results are reported as the average of these last five iterations.

For in solution calculations, a development version of the ASEP/MD software<sup>26</sup> was used. All quantum calculations were performed with Molcas-7.4.<sup>37</sup> All MD simulations were performed using Moldy.<sup>38</sup> The electrostatic potential generated by the solute was calculated with Molden.<sup>39</sup>

### 3. RESULTS AND DISCUSSION

**3.1. Gas Phase.** *3.1.1. Optimized Geometries.* The geometry of the FPP molecule was optimized in the gas phase at the SA(5)-CASSCF(12,11)/cc-pVDZ level for the electronic ground state and different singlet excited states. For describing and comparing the structures, we use some geometric parameters, such as the average bond length of the phenyl ring ( $\overline{\text{Ph}}$ ), the average bond length of the pyrrole ring ( $\overline{\text{Py}}$ ), the phenyl-pyrrole bond length (Ph–Py), or the phenyl pyrrole twist angle ( $\theta$ ). See Figure 2 and Table 1 for the atom numbering and parameter definitions.

**Table 2.** Geometrical parameters and dipole moments of the different optimized structures of FPP in the gas phase. Geometries optimized at the SA-CASSCF level, dipoles calculated at the CASPT2(0.00) level. The negative sign in the dipole indicates the negative charge is displaced toward the phenyl ring

	GS ( $S_0$ )	LE ( $S_1$ )	BQ ( $S_1$ )
$\overline{\text{Ph}}$ (Å)	1.400	1.433	1.421
$\overline{\text{Py}}$ (Å)	1.391	1.393	1.396
$Q(\text{Ph})$ (Å)	0.005	−0.003	0.069
$Q(\text{Py})$ (Å)	0.016	0.030	−0.105
Ph–Py (Å)	1.376	1.363	1.446
$\phi$ (°)	−0.1	0.0	30.1
$\psi$ (°)	−0.1	0.0	−4.9
$\theta$ (°)	0.0	0.0	2.1
$\mu$ (D)	1.18	0.10	−6.33

**Table 3.** Vertical Absorption Energies (in eV), Dipole Moments (in D), and Oscillator Strengths For the FPP Molecule in the Gas Phase at the GS Geometry<sup>a</sup>

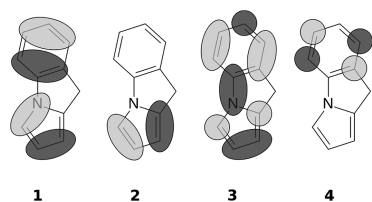
	vertical energies					
	CASSCF	CASPT2	expt <sup>b</sup>	$\mu$	$f$	
$S_0$				1.18		
$S_1$	4.63	4.64	4.26	4.22	0.013	
$S_2$	5.75	5.18	4.68	4.71	−3.65	0.276
$S_3$	5.91	5.56	5.18	(5.05)	−6.12	0.028
$S_4$	6.21	5.69	5.23		−7.99	0.174

<sup>a</sup>Dipole moments and oscillator strengths calculated at the CASPT2(0.00) level. <sup>b</sup>In *n*-hexane.<sup>18</sup>

The optimized ground state (GS) structure shows benzene and pyrrole rings with normal aromatic bond lengths ( $\overline{\text{Ph}} = 1.400$  Å,  $\overline{\text{Py}} = 1.391$  Å). The two rings are coplanar, as can be seen in the values of the angles  $\phi$ ,  $\psi$  and  $\theta$  in Table 2. The bond lengths are in general agreement with other published computed values,<sup>23,40</sup> although our Ph–Py length, 1.376 Å, is 0.02 Å shorter than the value reported in those works. This difference is probably due to the state averaging in our calculations.

At the ground state geometry, the first excited state corresponds mainly to a  $\pi \rightarrow \pi^*$  transition in the phenyl ring. Optimisation of this state leads to the LE (local excitation) geometry. In this structure the rings are also coplanar and Ph–Py is shorter than for the GS (1.363 Å). The local excitation character of this state is reflected in the significant increase of  $\overline{\text{Ph}}$  to 1.433 Å. These features agree with the results of Xu et al.,<sup>23</sup> with about the same difference in Ph–Py as with GS. He and Li,<sup>40</sup> however, report a LE geometry with more important differences, which they call “quinoid-like”, with a still shorter Ph–Py length (1.347 Å) and smaller  $\overline{\text{Ph}}$ ; this might be due to their use of a reduced active space (10 electrons in 9 orbitals). The dipole moment of this state is practically zero.

The higher excited states at the GS geometry have a marked charge transfer character. The electron density polarization is inverted with respect to the ground state and the negative charge is displaced toward the phenyl ring (this change of direction in the polarization is indicated with a negative sign in the dipole moment values in the tables). We optimized the geometry of a



**Figure 3.** Main active molecular  $\pi$  orbitals of FPP (simplified). In the dominant ground-state configuration, orbitals 1 and 2 are doubly occupied, while 3 and 4 are empty ( $1^2 2^2 3^0 4^0$ ).

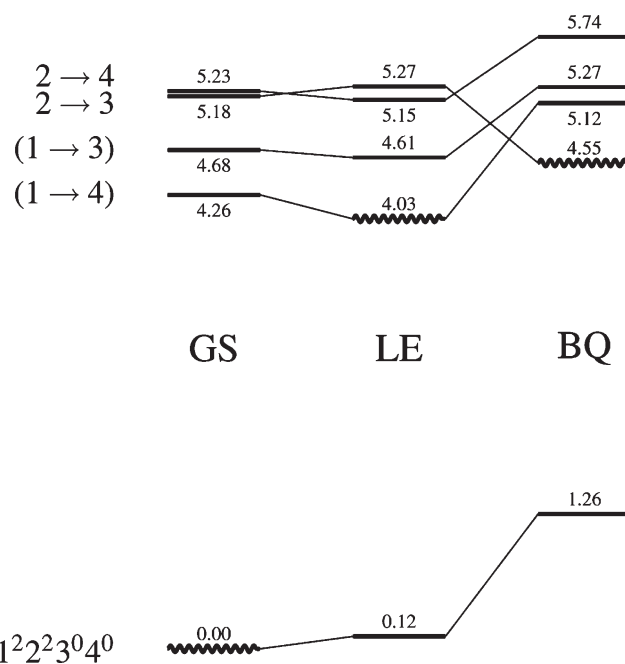
**Table 4.** Vertical Emission Energies (transitions to  $S_0$ , in eV), Dipole Moments (in D), and Oscillator Strengths for the FPP Molecule in the Gas Phase<sup>a</sup>

	vertical energies		$\mu$	$f$	$\Delta E$	
	CASPT2	expt <sup>b</sup>				
LE ( $S_1$ )	4.30	3.92	3.98	0.10	0.015	4.03
BQ ( $S_1$ )	3.55	3.29	−6.33	0.005	0.005	4.55

<sup>a</sup>  $\Delta E$  is the relative energy (in eV) with respect to the ground-state minimum, GS. Dipole moments, oscillator strengths, and  $\Delta E$  calculated at the CASPT2(0.00) level. <sup>b</sup> In *n*-hexane.<sup>18</sup>

charge transfer state in the gas phase, characterized by a quinoidal phenyl ring and a bend between the two aromatic rings, and therefore we will name it BQ (bent quinoidal). This structure has a pyramidalised  $C_6$  atom, which is also displaced out of the main phenyl plane. The two rings are distorted as described by the values of  $Q(\text{Ph})$  (positive) and  $Q(\text{Py})$  (negative), and the Ph–Py length is significantly larger than for the GS and LE structures. Xu et al.<sup>23</sup> report a similar structure for the ICT minimum, but He and Li<sup>40</sup> give a structure with “anti-quinoidal” phenyl ring (negative  $Q(\text{Ph})$ ). We could not obtain any different ICT minimum in our calculations in the gas phase, which may be due again to the different computational level employed.

**3.1.2. Absorption.** The vertical absorption properties of FPP calculated at the optimized ground-state geometry (GS) are summarized in Table 3. The CASSCF transition energies are included for comparison, but it is known that dynamic electron correlation must be included to obtain reliable results, and therefore, we will only discuss CASPT2 energies in the rest of the article. By comparing the two CASPT2 columns, it is clear that CASPT2(0.25) values are consistently 0.4 eV to 0.5 eV larger than CASPT2(0.00) values, a difference that has been found and discussed in other works.<sup>24,41–43</sup> Other properties like dipole moments or oscillator strengths do not show such variations, and only CASPT2(0.00) values are reported for them. From the values in Table 3, the  $S_0 \rightarrow S_2$  transition appears to be the most active in absorption, while the  $S_0 \rightarrow S_1$  transition should be much weaker, and the  $S_0 \rightarrow S_4$  transition could also be observed. The experimental absorption spectrum of FPP in *n*-hexane<sup>18</sup> has a strong band at 4.71 eV, a much weaker band at 4.22 eV, and a shoulder at around 5.05 eV. Although CASPT2-(0.25) results are generally less sensitive to basis set or active space changes and more similar to other methods of like quality, for the present calculations CASPT2(0.00) results are in better agreement with the experimental values. Therefore, to facilitate the discussion, in the rest of this work, we will refer in general to CASPT2(0.00) values. A single-point calculation with diffuse functions, using the aug-cc-pVDZ basis set, yielded very similar



**Figure 4.** Relative energies (CASPT2(0.00), in eV) of the calculated electronic states of FPP in the gas phase at the optimized geometries. The state for which each geometry is optimized is drawn as a wavy line. States of equivalent electron configuration are joined by lines. For the nature of the different states, labeled on the left, refer to Table 3, Figure 3, and the corresponding text.

results, with all absorption energies 0.1 eV to 0.2 eV lower, as observed in previous works when the basis set is enlarged.

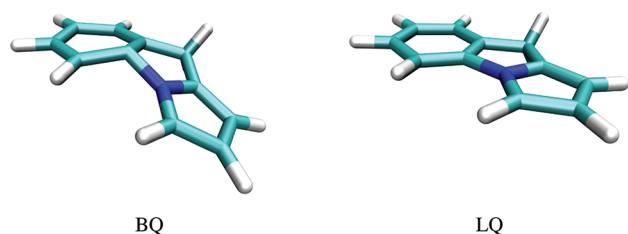
The electronic states  $S_3$  and  $S_4$  are dominated by single excitations from the pyrrole ring to the phenyl. In terms of the simplified molecular orbitals pictured in Figure 3,  $S_3$  is a  $2 \rightarrow 3$  transition, and  $S_4$  is a  $2 \rightarrow 4$  transition.  $S_1$  and  $S_2$  are not so clearly dominated by one configuration, but they have the larger contribution from  $1 \rightarrow 4$  and  $1 \rightarrow 3$  transitions, respectively. The electronic state optimized in the LE structure described above is equivalent to  $S_{11}$ , while the state optimized in the BQ ICT structure corresponds to  $S_3$  ( $2 \rightarrow 3$  transition), as suggested by the values of  $Q(\text{Ph})$  and  $Q(\text{Py})$ .

**3.1.3. Fluorescence.** The fluorescence energies from the two excited states optimized in the gas phase are shown in Table 4. The predicted emission from the LE state is 0.34 eV lower than the absorption with both CASPT2 variants, this agrees fairly well with the experimental Stokes shift of 0.24 eV in *n*-hexane. As occurred in the absorption, the best agreement with the experimental fluorescence is obtained with CASPT2(0.00). The  $\Delta E$  value of 4.03 eV can be compared with the experimental value obtained from the crossing point of the absorption and fluorescence spectra, which is 4.24 eV in *n*-hexane. Fluorescence at the BQ geometry is calculated to have a much lower energy (0.63 eV lower at the CASPT2(0.00) level), but the emitting state is 0.52 eV above the LE state and above the Franck–Condon  $S_1$  state at the GS geometry. This is probably a reason why there is no observed ICT fluorescence in nonpolar solvents. A scheme of the relative energies of the electronic states at the different geometries is presented in Figure 4. Again, aug-cc-pVDZ single-point calculations give very similar result, with fluorescence energies around 0.1 eV lower in all cases.

**Table 5.** Geometrical Parameters and Dipole Moments of the Different Optimized Structures of FPP in Acetonitrile Solution<sup>a</sup>

	GS (S <sub>0</sub> )	LE (S <sub>1</sub> )	BQ (S <sub>1</sub> )	LQ (S <sub>1</sub> )
$\overline{\text{Ph}}$ (Å)	1.400	1.433	1.419	1.413
$\overline{\text{Py}}$ (Å)	1.391	1.393	1.393	1.391
Q(Ph) (Å)	0.005	-0.003	0.064	0.052 <sup>b</sup>
Q(Py) (Å)	0.015	0.030	-0.101	-0.085
Ph-Py (Å)	1.380	1.364	1.462	1.451
$\phi$ (°)	-0.1	-0.1	26.0	0.6
$\psi$ (°)	-0.3	-0.1	-4.6	0.0
$\theta$ (°)	0.0	0.0	4.4	-0.3
$\mu$ (D)	1.75	0.27	-9.61	-12.06

<sup>a</sup> Geometries optimized at the SA-CASSCF level, dipoles calculated at the CASPT2(0.00) level. The negative sign in the dipole indicates the negative charge is displaced toward the phenyl ring. <sup>b</sup> Q'(Ph).

**Figure 5.** Perspective view of the two optimized ICT structures in acetonitrile.**Table 6.** Vertical Absorption Energies (in eV), Dipole Moments (in D), and Oscillator Strengths for the FPP Molecule in Acetonitrile at the GS Geometry

	vertical energies		$\mu$	$f$
	CASPT2	expt <sup>18</sup>		
S <sub>0</sub>			1.76	
S <sub>1</sub>	4.31	4.26	0.82	0.011
S <sub>2</sub>	4.75	4.73	-2.55	0.291
S <sub>3</sub>	5.28		-5.29	0.031
S <sub>4</sub>	5.41		-7.66	0.142

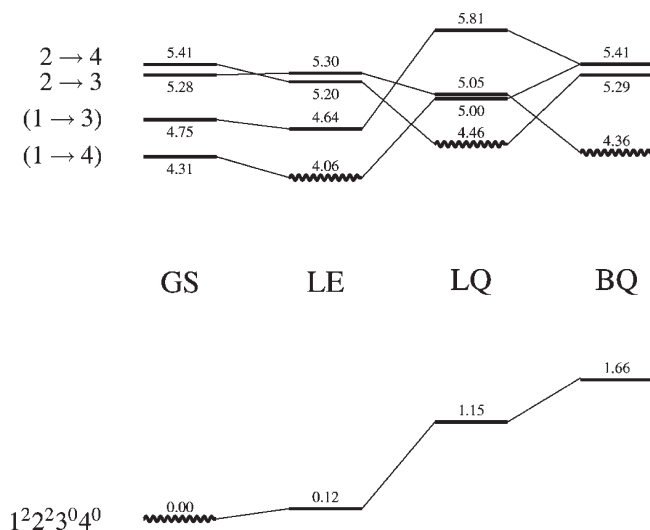
**3.2. Acetonitrile Solution.** **3.2.1. Optimized Geometries.** The different electronic states obtained in the gas phase for the FPP molecule were also optimized in acetonitrile solution, using the ASEP/MD method<sup>25–27</sup> to model the solvation process. The resulting geometries are given in Table 5. The changes in the geometry are small in all cases, and practically negligible for GS and LE. In the BQ structure the most significant change between the gas phase and acetonitrile is the lengthening of the Ph-Py bond and a slight planarization of the  $\phi$  angle. As expected in a polar solvent, dipole moments are enhanced, only slightly in GS and LE, and more significantly in BQ.

In addition to the minima already described, in solution it was possible to find another minimum in the S<sub>1</sub> surface with ICT character. This minimum is characterized by a planar structure ( $\phi$ ,  $\psi$ , and  $\theta$  angles close to zero) and a quinoidal phenyl ring (with two opposite bonds shorter than the other four), and therefore we name it linear quinoidal (LQ). It is interesting that

**Table 7.** Vertical Emission Energies (transitions to S<sub>0</sub>, in eV), Dipole Moments (in D), and oscillator Strengths for the FPP Molecule in Acetonitrile<sup>a</sup>

	vertical energies		$\mu$	$f$	$\Delta G$
	CASPT2	expt <sup>18</sup>			
LE (S <sub>1</sub> )	3.94	3.94	0.27	0.014	4.06
BQ (S <sub>1</sub> )	2.70		-9.60	0.003	4.36
LQ (S <sub>1</sub> )	3.31	3.29	-12.06	0.003	4.46

<sup>a</sup>  $\Delta G$  is the relative free energy (in eV) with respect to the ground-state minimum, GS.

**Figure 6.** Relative free energies (CASPT2(0.00), in eV) of the calculated electronic states of FPP in acetonitrile solution at the optimized geometries. The state for which each geometry is optimized is marked as a wavy line, this is also the state with which the solvent is in equilibrium. States of equivalent electron configuration are joined by lines. For the nature of the different states, labeled on the left, refer to Table 6 and Figure 3.

in this LQ structure the phenyl deformation does not happen along the C<sub>10</sub>-C<sub>11</sub> bond but along the C<sub>11</sub>-C<sub>6</sub> bond, so that it is best described with Q'(Ph) instead of Q(Ph) (see Table 1). The dipole moment of this structure is even larger than for BQ, and it can be noted that the dipole moments of the GS, LE, and LQ structures are in very good agreement with the experimental estimations<sup>17</sup> (1.7 D for the ground state, 1 D for the LE state, and -13 D for the ICT state). The two structures BQ and LQ are compared in Figure 5.

**3.2.2. Absorption.** The calculated absorption properties of FPP in acetonitrile are summarized in Table 6. All values are very close to the gas phase results, which is not surprising given the weak dipole moment of the ground state and the negligible change in the optimized GS geometry. In the two lowest transitions, a small blue shift is predicted, in accordance with the change of dipole moment between the states. This small blue shift is also observed experimentally when the absorptions in *n*-hexane and acetonitrile solutions are compared.<sup>18</sup>

**3.2.3. Fluorescence.** The results for the excited state emission properties from the different optimized structures of FPP in solution are shown in Table 7. Similarly to what was found for the absorptions, there is very little change in the LE emission from

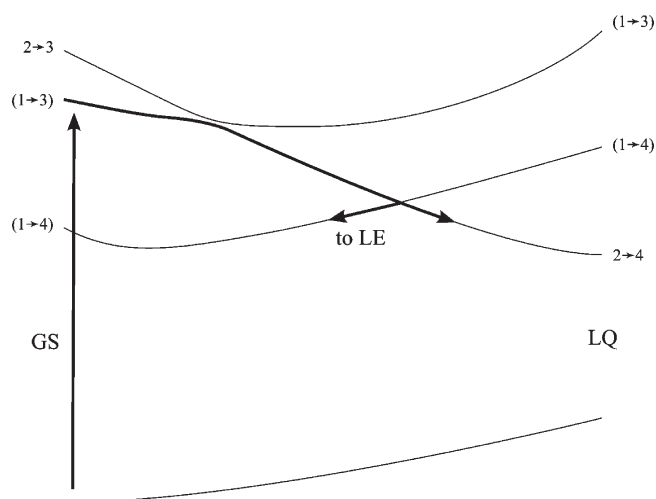
the gas phase to acetonitrile, this is consistent with the experiments, where the LE fluorescence band shows practically no solvatochromic shift from *n*-hexane to acetonitrile.<sup>18</sup>

The two optimized ICT structures have very different emission energies, with the value for LQ being 0.6 eV larger than for BQ. The observed red-shifted band of FPP in acetonitrile is centered at around 3.29 eV, which is in excellent agreement with the predicted LQ emission. On the basis of the fluorescence energies, the BQ structure can be ruled out as the main source of the ICT band. The relative free energy of the states is listed in the  $\Delta G$  column, where it can be seen that the difference in free energy between BQ and LQ is small (0.1 eV or 2.3 kcal/mol), and LE is around 0.3 eV (7 kcal/mol) below BQ. All three states are below the absorption Franck–Condon energy ( $S_2$  at GS). A scheme of the energies of the first five states at each structure is shown in Figure 6. Due to the different phenyl deformation in LQ, the equivalence between the states (the lines joining the horizontal lines) is only partial, and there is considerable mixture.

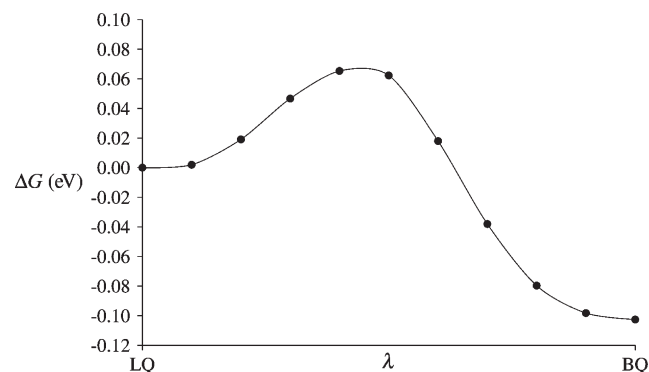
In a theoretical study of PP and FPP, using the polarizable continuum model (PCM) method to include solvent effects, Xu et al.<sup>23</sup> also found that the free energy difference between the ICT structures of FPP is relatively small. In agreement with our results, they obtained a significantly lower fluorescence energy for BQ than for a planar (symmetry-constrained) structure. However, the solvatochromic shift given by PCM is much weaker than ours (0.26 eV vs 0.59 eV for BQ), which leads them to conclude that the emission from BQ in acetonitrile is more similar to the observed ICT band, while the emission from the planar structure would be indistinguishable from the LE emission. In contrast, the fluorescence energies reported in this work indicate that BQ emission is too low to correspond to the experimental band, while the emission from the planar LQ structure is a much better candidate. The difference between our results and those of Xu et al. can be attributed to the absence of specific solute–solvent interactions in PCM and to the different active spaces used in both works.

Druzhinin et al. have estimated some thermodynamic quantities for the FPP system from the fluorescence properties;<sup>18</sup> in particular, from their data it can be concluded that the free energy difference between the emitting LE and ICT states is around 1 kcal/mol in acetonitrile at room temperature. Our results yield the two candidate ICT states about 8 kcal/mol higher in energy than the LE state. Here we must recall that we are using CASPT2-(0.00) values in this discussion because they make comparison of electron transition energies with experiments easier. As we indicated previously, transition energies with CASPT2(0.25) are 0.4 to 0.5 eV larger, but the relative stability of the different excited states does not change much. Nevertheless if we take CASPT2(0.25) values, then LE is only around 0.2 eV (4 kcal/mol) below BQ, the free energy difference between LQ and BQ is lower than 2 kcal/mol, and both ICT states lie very close to the  $S_1$  state at GS. Considering the errors, approximations and assumptions in the experiments, interpretations, and calculations, there is qualitative agreement with the recent experimental findings.

If, as we propose, the experimental ICT band corresponds to emission from the LQ structure, there must be a reason why the BQ structure is not formed or its fluorescence is not observed. Since, as seen in Figure 6 and Table 7, the free energies of LQ and BQ, and the oscillator strengths for their vertical emissions are very close, it can be interesting to analyze which ICT structure is reached first during the solute relaxation after the initial absorption. The structural similarity between GS and LQ suggests that



**Figure 7.** Qualitative scheme of the excited-state optimization of FPP in acetonitrile, starting from the Franck–Condon absorption at the GS structure. The electronic surfaces are labeled as in Figure 6. The bold lines and arrows indicate the path followed by the solute wave function.



**Figure 8.** Free energy profile for a geometry interpolation ( $\lambda$  parameter) between the LQ and BQ structures, in acetonitrile.

relaxation leads to LQ first, but a complete description of the process would require a study of the coupled dynamics of solute and solvent, which is beyond the scope of this work. However, within the mean field approximation of ASEP/MD, we can get a qualitative picture by following the gradient during the solute geometry optimization. We did this, optimizing the solute geometry starting with the  $S_2$  state at the GS geometry (Franck–Condon absorption) and observed that the solute structure tends to LQ. The process followed corresponds to that pictured in Figure 7, where there is initially a crossing between the surfaces of  $S_2$  and  $S_3$  and a change in the wave function nature occurs, such that after the crossing there is a clear ICT character in the  $S_2$  surface. Afterward, there is an intersection between the electronic surfaces corresponding to the LE and LQ states; after this intersection, the solute can proceed to either the LQ or LE structures, depending on which surface is followed, which would be determined by the system dynamics.

Once the LQ structure has been reached, there must be some barrier preventing the interconversion between LQ and BQ. We tried to estimate this barrier by performing a FEP calculation between the two ICT structures. The solute geometry was interpolated (in internal coordinates) in 10 steps between LQ

**Table 8. Main Experimental and Calculated Excited-State Absorption Bands of FPP (in eV)<sup>a</sup>**

	expt <sup>18</sup>		
<i>n</i> -hexane	1.50	2.92	
acetonitrile	1.55	2.05	3.40
MS-CASPT2, gas phase			
LE	1.58 (0.149)	3.39 (0.011)	
BQ	2.97 (0.008)	3.32 (0.094)	
MS-CASPT2, acetonitrile			
LE	1.56 (0.126)	2.24 (0.027)	
BQ	2.91 (0.008)	3.17 (0.044)	
LQ	2.31 (0.011)	3.75 (0.034)	

<sup>a</sup> Oscillator strength is in parentheses.

and BQ, and at each step, the solute wave function and the solvent were equilibrated. The resulting free energy profile (see Figure 8) shows a barrier of around 0.1 eV (2 kcal/mol) in the direction LQ → BQ. This value is only an upper bound for the barrier in equilibrium conditions, because the path followed is not optimized. The low value obtained for the barrier indicates that dynamical effects are probably important, and therefore further investigations are needed to elucidate the reasons why apparently no emission from BQ is observed experimentally.

It is interesting to compare the results obtained for FPP with those for the related PP.<sup>24</sup> In both systems we find an LQ ICT state, with a fluorescence energy matching the observed spectrum, and other bent structures (BQ) close in energy but with a predicted emission that is too low for the experimental band. These similarities between the two molecules can explain the parallels in their photophysical behavior. In the case of PP, a perpendicular state (PQ) was also found to be possible, which could be a route to nonradiative deactivation, thus decreasing the relative intensity of the ICT emission when compared to FPP, where this PQ structure is not available.

**3.3. Excited-State Absorption.** Druzhinin et al. have also measured the transient absorption spectra of FPP in *n*-hexane and acetonitrile<sup>18</sup> at different delay times, which has allowed them to assign certain absorption bands to the emitting states responsible for the two fluorescence bands. We have calculated the absorption energies from S<sub>1</sub> to higher excited states at the different optimized structures, with the goal of confirming the nature of the emitting states and their identity with the states probed in the transient absorption. For these calculations we had to include a larger number of states in the CASSCF state averaging (ten in total), and the multistate variant of CASPT2 was needed to separate the electronic states.<sup>44</sup>

The results are summarized in Table 8. Experimentally, the excited state absorption (ESA) spectrum of FPP in *n*-hexane is dominated by a band at 1.50 eV, with a minor band at 2.92 eV, which are attributed to the LE state, since this is the only state observed in the fluorescence spectrum. In acetonitrile, the band at 1.55 eV decreases over the first few ps, while the band at 3.40 eV increases and is therefore assigned to the ICT state.

The most intense absorption predicted by the present calculations occurs at around 1.58 eV, for the LE structure, and changes very little from the gas phase to acetonitrile. This value can be compared with the 1.50 and 1.55 eV bands observed in the experiments, confirming that the LE state can be the source of

these bands. The other absorption found in *n*-hexane should also correspond to the LE state, but to adequately reproduce it in the calculations, a higher number of states would probably be needed (the ninth root is still only 2.7 eV above S<sub>1</sub>).

For the ICT structures in acetonitrile solution, BQ, and LQ, we do not find any absorption of similar intensity, all oscillator strengths being significantly lower. This somewhat agrees with the ESA spectrum measured at longer delay times, which is relatively weak. The experimental band at 3.40 lies approximately between the two most intense absorptions predicted, at 3.17 and 3.40 eV. It can be tempting to assign the experimental band to either of the two predicted transitions (or a combination thereof), but both theoretical values correspond to high excited states and are therefore subject to significant errors. The fact is that the present results do not allow an unequivocal determination of the origin of the ESA spectrum of FPP in acetonitrile at long delays.

## 4. CONCLUSIONS

We have studied the ground and excited singlet states of fluorazene in the gas phase and in acetonitrile solution, using a high-level quantum method for the electronic structure and an explicit mean-field MM model for the solvent. The optimized structures for the GS and the LE state provide good agreement with the observed absorption bands and the higher-energy fluorescence band. These states are characterized by very low dipole moments and are only weakly affected by the solvent; in consequence, their photophysical properties show little change between the gas phase and solution. The agreement between the computed results for the LE state and the emission and excited-state absorption properties in *n*-hexane indicates that this state is adequately described by the present theoretical methods, and there is, in our opinion, little doubt on its nature and participation in the dual fluorescence of FPP.

The situation is less clear for the ICT state, responsible of the lower-energy fluorescence band. In the gas phase only a minimum is located, and this state is 0.5 eV higher in energy than the LE state. In acetonitrile solution we obtain two optimized structures for states of significant charge-transfer character, both structures being similar in energy. In one of these structures, LQ, the molecule skeleton is kept planar, and its emission energy and dipole moment are in good agreement with the experimental band, while in the other, BQ, the C<sub>6</sub> atom is pyramidalized, and its emission energy is around 0.6 eV lower. Our results therefore suggest that the experimental ICT fluorescence originates from LQ.

The excited-state absorption calculations for the different structures confirm the LE state as responsible for the 800 nm band, but they do not allow a conclusive assignment for the ICT signals.

Finally, why emission from BQ is not experimentally registered remains an open question, and to arrive to a definitive conclusion, possibly more sophisticated and accurate electronic structure methods are needed, along with the inclusion of further effects not considered in this work, such as the excited-state dynamics or vibronic coupling. With the current results, however, we can state that the twist between the electron-donor and -acceptor groups is not necessary for an ICT state to be stabilized.

## ■ AUTHOR INFORMATION

### Corresponding Author

\*E-mail: jellby@unex.es.



## ACKNOWLEDGMENT

This work was supported by the CTQ2008-06224/BQU Project from the Ministerio de Ciencia e Innovación of Spain, co-financed by the European Regional Development Fund (ERDF), and the PRI08A056 Project from the Consejería de Economía, Comercio e Innovación of the Junta de Extremadura. I.F.G. acknowledges the Junta de Extremadura and the European Social Fund for financial support, and A.M.L. acknowledges financial support from the Juan de la Cierva subprogramme of the Ministerio de Ciencia e Innovación of Spain. The authors also thank the Fundación Computación y Tecnologías Avanzadas de Extremadura (COMPUTAEX) for additional computational resources.

## REFERENCES

- (1) Lippert, E.; Lüder, W.; Moll, F.; Nägele, W.; Boos, H.; Prigge, H.; Seibold-Blankenstein, I. *Angew. Chem.* **1961**, *73*, 695–706.
- (2) Grabowski, Z. R.; Rotkiewicz, K.; Rettig, W. *Chem. Rev.* **2003**, *103*, 3899–4032.
- (3) Herbich, J.; Pérez Salgado, F.; Rettschnick, R. P. H.; Grabowski, Z. R.; Wojtowicz, H. *J. Phys. Chem.* **1991**, *95*, 3491–3497.
- (4) Peng, L. W.; Dantus, M.; Zewail, A. H.; Kemnitz, K.; Hicks, J. M.; Eisenthal, K. B. *J. Phys. Chem.* **1987**, *91*, 6162–6167.
- (5) Lommatzsch, U.; Gerlach, A.; Lahmann, C.; Brutschy, B. *J. Phys. Chem. A* **1998**, *102*, 6421–6435.
- (6) Chudoba, C.; Kummrow, A.; Dreyer, J.; Stenger, J.; Nibbering, E. T. J.; Elsaesser, T.; Zachariasse, K. A. *Chem. Phys. Lett.* **1999**, *309*, 357–363.
- (7) Druzhinin, S. I.; Demeter, A.; Galievsky, V. A.; Yoshihara, T.; Zachariasse, K. A. *J. Phys. Chem. A* **2003**, *107*, 8075–8085.
- (8) Druzhinin, S. I.; Mayer, P.; Stalke, D.; von Blow, R.; Noltemeyer, M.; Zachariasse, K. A. *J. Am. Chem. Soc.* **2010**, *132*, 7730.
- (9) Kato, S.; Amatatsu, Y. *J. Chem. Phys.* **1990**, *92*, 7241–7257.
- (10) Hayashi, S.; Ando, K.; Kato, S. *J. Phys. Chem.* **1995**, *99*, 955–964.
- (11) Serrano-Andrés, L.; Merchán, M.; Roos, B. O.; Lindh, R. *J. Am. Chem. Soc.* **1995**, *117*, 3189–3204.
- (12) Sudholt, W.; Arnulf Staib, A. L. S.; Domcke, W. *Phys. Chem. Chem. Phys.* **2000**, *2*, 4341–4353.
- (13) Mennucci, B.; Toniolo, A.; Tomasi, J. *J. Am. Chem. Soc.* **2000**, *122*, 10621–10630.
- (14) Rappoport, D.; Furche, F. *J. Am. Chem. Soc.* **2004**, *126*, 1277–1284.
- (15) Minezawa, N.; Kato, S. *J. Phys. Chem. A* **2005**, *109*, 5445–5453.
- (16) Gómez, I.; Mercier, Y.; Reguero, M. *J. Phys. Chem. A* **2006**, *110*, 11455–11461.
- (17) Yoshihara, T.; Druzhinin, S. I.; Zachariasse, K. A. *J. Am. Chem. Soc.* **2004**, *126*, 8535–8539.
- (18) Druzhinin, S. I.; Kovalenko, S. A.; Senyushkina, T. A.; Demeter, A.; Zachariasse, K. A. *J. Phys. Chem. A* **2010**, *114*, 1621–1632.
- (19) Parusel, A. B. *J. Phys. Chem. Chem. Phys.* **2000**, *2*, 5545–5552.
- (20) Proppe, B.; Merchán, M.; Serrano-Andrés, L. *J. Phys. Chem. A* **2000**, *104*, 1608–1616.
- (21) Zillberg, S.; Haas, Y. *J. Phys. Chem. A* **2002**, *106*, 1–11.
- (22) Schweke, D.; Baumgarten, H.; Haas, Y.; Rettig, W.; Dick, B. *J. Phys. Chem. A* **2005**, *109*, 576–585.
- (23) Xu, X.; Cao, Z.; Zhang, Q. *J. Phys. Chem. A* **2006**, *110*, 1740–1748.
- (24) Fdez. Galván, I.; Martín, M. E.; Muñoz-Losa, A.; Sánchez, M. L.; Aguilar, M. A. *J. Chem. Theory Comput.* **2011**, *7*, 1850–1857.
- (25) Sánchez, M. L.; Aguilar, M. A.; Olivares del Valle, F. J. *J. Comput. Chem.* **1997**, *18*, 313–322.
- (26) Fdez. Galván, I.; Sánchez, M. L.; Martín, M. E.; Olivares del Valle, F. J.; Aguilar, M. A. *Comput. Phys. Commun.* **2003**, *155*, 244–259.
- (27) Aguilar, M. A.; Sánchez, M. L.; Martín, M. E.; Fdez. Galván, I. An Effective Hamiltonian Method from Simulations: ASEP/MD. In *Continuum Solvation Models in Chemical Physics*, 1st ed.; Mennucci, B., Cammi, R., Eds. Wiley: Hoboken, NJ, 2007; Chapter 4.5, pp 580–592.
- (28) Fdez. Galván, I.; Sánchez, M. L.; Martín, M. E.; Olivares del Valle, F. J.; Aguilar, M. A. *J. Chem. Phys.* **2003**, *118*, 255–263.
- (29) Okuyama-Yoshida, N.; Nagaoka, M.; Yamabe, T. *Int. J. Quantum Chem.* **1998**, *70*, 95–103.
- (30) Zwanzig, R. W. *J. Chem. Phys.* **1954**, *22*, 1420–1426.
- (31) Fdez. Galván, I.; Aguilar, M. A.; Ruiz-López, M. F. *J. Phys. Chem. B* **2005**, *109*, 23024–23030.
- (32) Roos, B. O.; Taylor, P. R.; Siegbahn, P. E. M. *Chem. Phys.* **1980**, *48*, 157–173.
- (33) Andersson, K.; Malmqvist, P.-Å.; Roos, B. O.; Sadlej, A. J.; Wolinski, K. *J. Phys. Chem.* **1990**, *94*, 5483–5488.
- (34) Andersson, K.; Malmqvist, P.-Å.; Roos, B. O. *J. Chem. Phys.* **1992**, *96*, 1218–1226.
- (35) Ghigo, G.; Roos, B. O.; Malmqvist, P.-Å. *Chem. Phys. Lett.* **2004**, *396*, 142–149.
- (36) Jorgensen, W. L.; Maxwell, D. S.; Tirado-Rives, J. *J. Am. Chem. Soc.* **1996**, *118*, 11225–11236.
- (37) Aquilante, F.; De Vico, L.; Ferré, N.; Ghigo, G.; Malmqvist, P.-Å.; Neogrády, P.; Pedersen, T. B.; PitoYák, M.; Reiher, M.; Roos, B. O.; Serrano-Andrés, L.; Urban, M.; Veryazov, V.; Lindh, R. *Comput. Chem.* **2010**, *31*, 224–247.
- (38) Refson, K. *Comput. Phys. Commun.* **2000**, *126*, 310–329.
- (39) Schaftenaar, G.; Noordik, J. H. *J. Comput.-Aided Mol. Des.* **2000**, *14*, 123–134.
- (40) He, R.-X.; Li, X.-Y. *Chem. Phys.* **2007**, *332*, 325–335.
- (41) Fdez. Galván, I.; Martín, M. E.; Muñoz-Losa, A.; Aguilar, M. A. *J. Chem. Theory Comput.* **2009**, *5*, 341–349.
- (42) Valsson, O.; Filippi, C. *J. Chem. Theory Comput.* **2010**, *6*, 1275–1292.
- (43) Fdez. Galván, I.; Martín, M. E.; Aguilar, M. A. *J. Chem. Theory Comput.* **2010**, *6*, 2445–2454.
- (44) Finley, J.; Malmqvist, P.-Å.; Roos, B. O.; Serrano-Andrés, L. *Chem. Phys. Lett.* **1998**, *288*, 299–306.

# Accurate Anharmonic Vibrational Frequencies for Uracil: The Performance of Composite Schemes and Hybrid CC/DFT Model

Cristina Puzzarini,<sup>\*,†</sup> Malgorzata Biczysko,<sup>‡</sup> and Vincenzo Barone<sup>§</sup>

<sup>†</sup>Dipartimento di Chimica "G. Ciamician", Università di Bologna, Via F. Selmi 2, 40126 Bologna, Italy

<sup>‡</sup>Center for Nanotechnology Innovation @NEST, Istituto Italiano di Tecnologia, Piazza San Silvestro, 12 - 56127 Pisa, Italy

<sup>§</sup>Scuola Normale Superiore, Piazza dei Cavalieri 7, 56126 Pisa, Italy

**ABSTRACT:** The vibrational spectrum (frequencies as well as intensities) of uracil has been investigated at a high level of theory. The harmonic force field has been evaluated at the coupled-cluster (CC) level in conjunction with a triple- $\zeta$  basis set. Extrapolation to the basis set limit as well as inclusion of core-correlation and diffuse-function corrections have been considered by means of the second-order Møller–Plesset perturbation theory. To go beyond the harmonic approximation, a hybrid CC/DFT approach has been employed, which will be proved to provide state-of-the-art results. As the spectroscopic investigation of uracil is hampered by numerous Fermi resonances, models for explicitly taking them into account have been implemented and applied. On general grounds, the computational procedure presented is able to provide the proper accuracy to support experimental investigations of large molecules of biological interest.

## INTRODUCTION

Nowadays, spectroscopic techniques represent the most reliable and flexible approaches for the investigation of structural and dynamical properties of molecular and supra-molecular systems, either isolated or in condensed phases. However, interpretation of spectra is seldom straightforward, and integrated experimental/computational investigations are becoming more and more popular, thanks to the improved reliability and effectiveness of quantum mechanical (QM) computations.<sup>1,2</sup> For small molecules, the most refined QM methods provide such accurate results that any disagreement between computational data and experimental findings casts serious doubts on the reliability of the latter.<sup>3,4</sup> With regard to the topic of the present work, as a significant example, we mention the new IR spectrum assignment for the vinyl radical based on anharmonic force field computations<sup>5–7</sup> that was eventually confirmed by new purposely tailored experiments.<sup>8</sup>

For polyatomic molecules, effective computational solutions of the vibrational problem and simulation of IR and Raman spectra are among the most important tasks of contemporary computational chemistry.<sup>9</sup> While theoretical evaluations of vibrational frequencies and IR/Raman intensities within the harmonic approximation have become a routine tool for assisting the interpretation of spectroscopic experiments, in the past decade great effort has been made to go beyond the harmonic approximation and perform anharmonic computations by means of perturbative<sup>10–23</sup> or variational approaches.<sup>24–32</sup> An effective approach is obtained when the vibrational second-order perturbation theory (VPT2)<sup>10–14,23</sup> is applied to a fourth-order representation of the potential energy surface (PES). In particular, Density Functional Theory (DFT) using hybrid (especially B3LYP<sup>7,33–37</sup>) or double-hybrid (especially B2PLYP<sup>38–40</sup>) functionals in conjunction with medium-sized basis sets is known to provide rather accurate results and can be exploited for large systems. Further improvements in accuracy can be obtained by

computing the harmonic part of the force field at a more refined level, with the coupled cluster (CC) method providing the most effective route, at least in the absence of strong multireference character.<sup>33,37,41–44</sup> The inclusion of the CCSD(T) harmonic part in a DFT anharmonic force field leads to the definition of hybrid CC/DFT approaches,<sup>37,43,44</sup> which nowadays represent the method of choice for computing accurate vibrational spectra of semirigid systems.

The present work is part of a comprehensive research project aimed at extending composite schemes to the accurate prediction of molecular and spectroscopic properties for small- to medium-sized molecules. On this topic, we are particularly interested in building blocks of biomolecules in view of their relevance in several fields ranging from prebiotic systems to biosensors. Characterization of isolated molecules in the gas phase is a mandatory prerequisite for the subsequent analysis of the role of different effects (e.g., hydrogen bonding, environmental effects, etc.) in determining the overall behavior of these systems. However, from an experimental point of view, the structural characterization of the simplest building blocks of biomolecules (e.g., amino acids or nucleic acid bases) in the gas phase is not straightforward at all. In the field of vibrational spectroscopy, interpretation of spectra suffers from the overlapping of several bands and from the presence of strong resonances. In the case of uracil, the simplest nucleobase (see Figure 1), encouraging results have already been obtained for the molecular parameters and spectroscopic properties related to rotational spectroscopy by means of a composite QM scheme.<sup>45</sup> On the other hand, such an accurate approach can be used to benchmark less accurate but computationally cheap methods rooted in the density functional theory as well as to set up hybrid CC/DFT models.

**Received:** August 8, 2011

**Published:** September 28, 2011

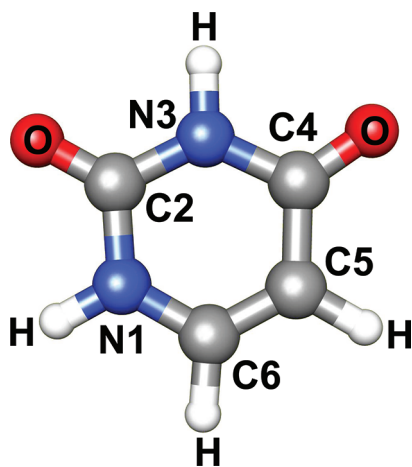


Figure 1. Molecular structure of uracil: atoms labeling.

Moving to the topic of the present work, vibrational spectra of uracil have been investigated in several theoretical and experimental works. The IR and Raman spectra of gaseous,<sup>46</sup> dissolved,<sup>47,48</sup> polycrystalline,<sup>49–51</sup> and matrix-isolated<sup>52–60</sup> uracil were recorded. A number of semiempirical and ab initio quantum mechanical approaches were also used to compute and analyze the vibrational spectra (see refs 61 and 62 and references therein). However, a number of unresolved questions are still open, and they mostly concern the intensities of the C–O stretching vibrations and the assignment of some other bands. As a consequence, there are still significant doubts on the interpretation of the vibrational spectrum of isolated uracil, the main problem being the large number of overtones and combination bands having intensities comparable to those of some fundamentals. This is mainly due to Darling–Dennison and Fermi resonances, which need to be properly taken into account in order to reliably reproduce the uracil IR spectrum. This implies that one should calculate the latter at the anharmonic level by means of an accurate quantum chemical approach, also including Fermi and Darling–Dennison resonances.

The paper is organized as follows. In the next section, the methodology used is explained together with the corresponding quantum-chemical details. Thereafter, the results are reported and discussed with particular emphasis on the proper account of interactions in the anharmonic frequency calculations. Harmonic frequencies and intensities are also reported and discussed.

## METHODOLOGY AND COMPUTATIONAL DETAILS

**Coupled-Cluster Computations.** The best-estimated harmonic force field for uracil in its electronic ground state has been evaluated by means of a composite scheme to account for electron-correlation and basis-set effects. This approach is based on the assumption of the additivity for various contributions. The second-order Møller–Plesset perturbation theory (MP2)<sup>63</sup> and CC singles and doubles approximation augmented by a perturbative treatment of triple excitations [CCSD(T)]<sup>64</sup> have been employed. Correlation-consistent basis sets, (aug)-cc-p(C)VnZ ( $n = T, Q$ ),<sup>65–67</sup> have been used in conjunction with the above-mentioned methods.

The harmonic force field has been computed at the MP2 and CCSD(T) levels employing different basis sets.<sup>68</sup> Following the procedure introduced in ref 17, the harmonic frequencies,  $\omega$ ,

have been extrapolated to the complete basis set (CBS) limit starting from the results obtained at the MP2/cc-pVTZ and MP2/cc-pVQZ levels. More precisely, the extrapolated correlation contribution has been added to the HF-SCF CBS limit, which is assumed to be reached at the HF/cc-pVQZ level. The latter assumption seems to be reasonable, as in most cases the differences in frequency between the HF/cc-pVTZ and HF/cc-pVQZ levels are smaller than  $1 \text{ cm}^{-1}$  and always much smaller than 0.5% in relative terms. Corrections due to core–valence (CV) correlation and effects due to diffuse functions (aug) in the basis set have then been evaluated at the MP2/cc-pCVTZ ( $\Delta\omega(\text{CV}) = \omega(\text{MP2}/\text{cc-pCVTZ}, \text{all}) - \omega(\text{MP2}/\text{cc-pCVTZ}, \text{fc})$ ) and MP2/aug-cc-pVTZ levels ( $\Delta\omega(\text{aug}) = \omega(\text{MP2}/\text{aug-cc-pVTZ}, \text{fc}) - \omega(\text{MP2}/\text{cc-pVTZ}, \text{fc})$ ), respectively. The latter correction has been introduced, as diffuse functions are required to properly describe electronegative atoms and thus to recover the limitations that are inherent in the extrapolation of the cc-pVTZ and cc-pVQZ basis sets to the CBS limit. In the expressions given in parentheses, “fc” and “all” stand for frozen-core approximation and all electrons (also 1s electrons of C, N, and O) correlated, respectively. Higher-order electron-correlation energy contributions ( $\Delta\omega((T))$ ) have also been considered. The corresponding corrections have been derived by comparing the harmonic frequencies at the MP2 and CCSD(T) levels, both in the cc-pVTZ basis set. Inclusion of all terms

$$\omega(\text{best}) = \omega(\text{CBS}(T, Q)) + \Delta\omega(\text{CV}) + \Delta\omega(\text{aug}) + \Delta\omega((T)) \quad (1)$$

finally provides the best estimated harmonic frequencies.

An analogous composite scheme has also been used to determine best estimates for the infrared intensities,  $I(\text{best})$ , within the harmonic approximation. As extrapolation schemes have not been formulated yet for such a property and diffuse functions are known to be important to correctly describe dipole moment derivatives, eq 1 has been rearranged as follows:

$$I(\text{best}) = I(\text{MP2}/\text{augVTZ}) + \Delta I(\text{CV}) + \Delta I(\text{QZ} - \text{TZ}) + \Delta I((T)) \quad (2)$$

where  $\Delta I(\text{QZ} - \text{TZ})$  is the correction due to the MP2/cc-pVQZ – MP2/cc-pVTZ difference.

The `CFOUR` program package<sup>69</sup> has been employed for all computations mentioned in this section.

**Density Functional Theory Computations.** Density Functional Theory has been employed to compute harmonic and anharmonic force fields. Within the DFT approach, the standard B3LYP functional<sup>70</sup> has been used in conjunction with the aug-N07D<sup>71</sup> and aug-cc-pVTZ<sup>65,66</sup> basis sets. Recently, the original polarized double- $\zeta$  basis set N07D<sup>71–74</sup> has been modified by consistent inclusion of diffuse s functions and then further augmented by one set of diffuse d functions on C and N atoms (already present on O). On general grounds, the DFT/N07D approach has been developed for spectroscopic studies of medium-to-large molecular systems and provides an excellent compromise between reliability and computational effort.<sup>7,37,39,75,76</sup> The two different basis sets employed allow us to monitor basis set effects as well as the performance of the newly developed aug-N07D set. Harmonic and anharmonic force fields have been computed starting from equilibrium structures optimized using tight convergence criteria.

As concerns anharmonic force fields, the third and semidiagonal fourth force constants have been obtained by numerical

Table 1. Harmonic Vibrational Frequencies ( $\text{cm}^{-1}$ ) of Uracil

assignment	B3LYP/aug-N07D	B3LYP/aug-cc-pVTZ	MP2/cc-pVTZ	MP2/cc-pVQZ	CBS	$\Delta\text{CV}$	$\Delta\text{aug}$	$\Delta(T)$	best <sup>d</sup> estimate
$\nu(\text{N1-H})$	3640.9	3634.4	3669.4	3666.4	3665.8	5.3	-14.9	-3.4	3652.7
$\nu(\text{N3-H})$	3596.4	3589.5	3616.4	3612.9	3612.3	5.4	-15.6	-0.0	3602.2
$\nu(\text{C5-H})$	3248.7	3243.9	3292.0	3289.2	3287.6	6.1	-10.5	-30.5	3252.8
$\nu(\text{C6-H})$	3207.7	3200.6	3246.5	3246.8	3245.9	5.9	-6.9	-27.5	3217.5
$\nu(\text{C2=O})$	1798.9	1792.7	1828.6	1818.3	1815.6	4.6	-24.3	-5.9	1790.0
$\nu(\text{C4=O})$	1764.3	1760.2	1790.4	1780.8	1779.1	5.0	-21.6	-1.0	1761.5
$\nu(\text{C5=C6})$	1674.0	1672.3	1686.2	1684.2	1684.8	6.0	-8.3	-5.0	1677.5
$\delta(\text{N1-H})$	1499.7	1497.8	1513.4	1513.5	1513.8	4.4	-4.9	-7.9	1505.4
$\delta(\text{C6-H})$	1403.5	1405.5	1427.2	1428.6	1428.7	4.3	-3.8	-2.1	1427.2
$\delta(\text{N3-H})$	1417.8	1422.3	1413.3	1413.8	1414.8	4.0	-3.9	-0.8	1414.0
$\delta(\text{C5-H})$	1383.2	1382.5	1389.2	1389.6	1390.1	3.2	-0.4	1.0	1394.0
$\nu(\text{ring})$	1229.0	1227.9	1246.8	1247.3	1248.6	5.0	-1.6	-3.9	1248.2
$\nu(\text{ring})$	1195.7	1192.5	1212.2	1212.4	1212.3	3.5	-2.5	-8.0	1205.3
$\nu(\text{ring})$	1086.8	1097.2	1095.3	1094.7	1094.3	3.5	-4.1	-9.5	1084.2
$\nu(\text{ring})$	992.8	993.7	992.7	994.7	995.1	3.3	-0.9	-0.9	995.4
$\nu(\text{ring})$	973.4	985.9	968.9	969.7	970.5	3.9	2.2	-9.0	967.7
$\nu(\text{ring})$	770.6	769.7	774.9	776.5	777.6	2.9	-1.5	-6.2	772.8
$\delta(\text{ring})$	557.5	574.0	562.2	564.5	565.6	3.9	3.2	-27.5	545.3
$\delta(\text{ring})$	542.4	543.6	538.1	539.4	540.3	2.5	-2.0	-0.0	540.8
$\delta(\text{ring})$	521.4	522.5	515.4	516.6	516.9	2.2	-1.0	-1.0	517.2
$\delta(\text{C=O})$	386.9	387.7	385.7	386.7	386.8	1.7	-1.5	0.4	387.4
$\gamma(\text{C6-H})$	965.3	964.0	977.0	978.3	978.6	3.3	-1.0	-7.6	973.3
$\gamma(\text{C5-H})$	823.5	828.2	817.8	816.6	815.7	3.6	3.1	-8.8	813.6
$\gamma(\text{C2=O})$	767.0	771.8	762.8	764.5	765.9	4.8	-1.9	-3.5	765.2
$\gamma(\text{C4=O})$	730.5	734.3	735.2	734.2	733.8	2.5	-0.4	-8.2	727.6
$\gamma(\text{N3-H})$	675.8	687.9	691.6	690.5	690.1	2.6	-3.5	-18.8	670.3
$\gamma(\text{N1-H})$	559.3	559.9	559.0	559.9	559.8	2.3	-2.1	-1.2	558.7
$\gamma(\text{ring})$	399.3	404.5	394.8	395.7	395.9	2.8	-11.5	-11.5	387.9
$\gamma(\text{ring})$	168.5	169.8	163.2	163.7	164.1	1.1	-0.3	-5.8	159.1
$\gamma(\text{ring})$	154.1	152.4	146.4	147.0	146.9	1.4	-1.5	-6.4	140.4

<sup>d</sup> From eq 1.

differentiation of the analytical second derivatives. The semidiagonal quartic force fields<sup>77</sup> have then been used to compute spectroscopic parameters and, in particular, anharmonic frequencies by means of the fully automated generalized second-order vibrational perturbation (GVPT2) approach,<sup>10,11</sup> as implemented in the Gaussian package.<sup>13,14,23</sup> As in perturbative treatments nearly resonant contributions should be removed, in the present work two possible approaches to defining Fermi and Darling–Dennison resonances have been followed. First, an automatic procedure (GVPT2, Fermi: auto) based on the criteria proposed by Martin et al.,<sup>15</sup> which are known to provide accurate results,<sup>35,61</sup> has been used to remove potentially divergent terms. In the current version of the code, an *ad hoc* procedure (GVPT2, Fermi: INP) to directly specify resonant terms has been also implemented. The latter allows us to directly compare with other theoretical approaches as well as to test the influence of any specific interaction on the overall results. In both cases, in a second step, all resonant terms are then treated variationally.<sup>11,13</sup> Finally, simple removal of resonant terms leads to the so-called deperturbed model, DVPT2.

All DFT computations have been performed employing the Gaussian suite of programs for quantum chemistry.<sup>78</sup>

**The Hybrid CC/DFT Approach.** A hybrid CCSD(T)/DFT approach<sup>33,37,41–44</sup> has also been used to evaluate anharmonic

frequencies. This model is based on the assumption that the differences between CCSD(T) and B3LYP anharmonic frequencies are solely due to the harmonic terms. In this way, prohibitively expensive computations of cubic and quartic force constants at the CCSD(T) level are avoided, and the hybrid CCSD(T)/DFT scheme therefore provides a viable route to extend accurate predictions of anharmonic frequencies to relatively large systems. In the present case, two possible approaches have been implemented. In the simplest one, the hybrid frequencies have been computed by means of *a posteriori* DFT corrections to the best-estimated harmonic frequencies:  $\nu_{\text{CC/DFT}} = \omega(\text{best}) + \Delta\nu_{\text{DFT}}$ . Such an approximation has already been validated for several closed- and open-shell systems (see, for instance, ref 37). On the other hand, in the second approach, the best-estimated harmonic frequencies are directly introduced into the GVPT2 computations along with the 3rd and 4th force constants obtained at the DFT level. It should be noted that the latter scheme can significantly improve the quality of the results when the discrepancy between harmonic frequencies computed at the DFT level and best estimates leads to an incorrect definition of Fermi resonances through automatic procedures.<sup>37,39</sup> However, a more general way to overcome such inconsistencies relies on generalized treatment that completely avoids divergent terms.<sup>79</sup>

Table 2. Anharmonic Vibrational Frequencies ( $\text{cm}^{-1}$ ) of Uracil from the CCSD(T)/B3LYP Hybrid Force Field<sup>a</sup>

symmetry/mode	assignment	harmonic best		GVPT2				B3LYP/6-31+G(d,p) <sup>c</sup> (Ten et al.)	experiment <sup>b</sup>		
		estimate	DVPT2	Fermi: DFT <sup>d</sup>	Fermi: CC <sup>e</sup>	Fermi: DFT+INP <sup>f</sup>	Fermi: CC+INP <sup>f</sup>		frequency	intensity	
A'	$\omega_1$	$\nu(\text{N1-H})$	3652.7	3484	3484	3485	3484	3485	3480	3485	166
A'	$\omega_2$	$\nu(\text{N3-H})$	3602.2	3436	3436	3436	3436*	3436*	3436	3435	100
A'	$\omega_3$	$\nu(\text{C5-H})$	3252.8	3117	3117	3117	3117	3117	3140		4
A'	$\omega_4$	$\nu(\text{C6-H})$	3217.5	3084	3072*	3072*	3083*	3072*	3082		
A'	$\omega_5$	$\nu(\text{C2=O})$	1790.0	1760	1762*	1762	1771*	1761*	1776	1764	680
A'	$\omega_6$	$\nu(\text{C4=O})$	1761.5	1735	1744*	1737*	1760*	1733*	1762	1706	291
A'	$\omega_7$	$\nu(\text{C5=C6})$	1677.5	1644	1644	1644	1644	1643	1643	1643	33
A'	$\omega_8$	$\delta(\text{N1-H})$	1505.4	1465	1461*	1465	1484*	1466*	1461	1472	83
A'	$\omega_9$	$\delta(\text{C6-H})$	1427.2	1382	1386*	1385	1392*	1388*	1394	1400	56
A'	$\omega_{10}$	$\delta(\text{N3-H})$	1414.0	1394	1391*	1386*	1384*	1384*	1374	1389	21
A'	$\omega_{11}$	$\delta(\text{C5-H})$	1394.0	1360	1360*	1353	1361*	1355*	1353	1359	13
A'	$\omega_{12}$	$\nu(\text{ring})$	1248.2	1220	1223*	1221	1226*	1221*	1210	1217	4
A'	$\omega_{13}$	$\nu(\text{ring})$	1205.3	1176	1176	1176	1176	1176	1171	1185	109
A'	$\omega_{14}$	$\nu(\text{ring})$	1084.2	1063	1061*	1061*	1064*	1061*	1073	1075	14
A'	$\omega_{15}$	$\nu(\text{ring})$	995.4	980	981*	978	995*	978*	961	980	
A'	$\omega_{16}$	$\nu(\text{ring})$	967.7	940	940	940	940	940	947	958	7
A'	$\omega_{17}$	$\nu(\text{ring})$	772.8	756	756	752*	756	751*	756	759	
A'	$\omega_{18}$	$\delta(\text{ring})$	545.3	552	549*	555*	545*	549*	550	562	17
A'	$\omega_{19}$	$\delta(\text{ring})$	540.8	533	533	530	529	530	534	537	7
A'	$\omega_{20}$	$\delta(\text{ring})$	517.2	511	511	510	511	510	515	516	23
A'	$\omega_{21}$	$\delta(\text{C=O})$	387.4	383	383	384	383*	386*	385	391	33
A''	$\omega_{22}$	$\gamma(\text{C6-H})$	973.3	955	955	954	955*	954*	950	987	2
A''	$\omega_{23}$	$\gamma(\text{C5-H})$	813.6	793	793	793	793	793	803	804	175
A''	$\omega_{24}$	$\gamma(\text{C2=O})$	765.2	746	746	746	746	746	749	757	125
A''	$\omega_{25}$	$\gamma(\text{C4=O})$	727.6	711	711	711	711*	711*	715	718	14
A''	$\omega_{26}$	$\gamma(\text{N3-H})$	670.3	654	654*	654	660*	654*	666	662	100
A''	$\omega_{27}$	$\gamma(\text{N1-H})$	558.7	546	546	549	546	555	567	551	25
A''	$\omega_{28}$	$\gamma(\text{ring})$	387.9	387	387	387	387	384	398	411	
A''	$\omega_{29}$	$\gamma(\text{ring})$	159.1	155	155	155	155	155	167	185	
A''	$\omega_{30}$	$\gamma(\text{ring})$	140.4	133	133	132	133	132	155		
MAE <sup>g</sup>			9	10	10	11	9	10	9		

<sup>a</sup>The asterisks denote the modes explicitly considered in the variational treatment of Fermi resonances. See text. <sup>b</sup>Refs 58–60. <sup>c</sup>Ref 62: Anharmonic frequencies and intensities at the B3LYP/6-31+G(d,p) level with Fermi resonances accounted for. <sup>d</sup>Fermi resonances identified through automatic procedure using harmonic frequencies and anharmonic force constants computed at the DFT level. <sup>e</sup>Fermi resonances identified through automatic procedure using harmonic frequencies computed at the CC level and anharmonic force constants computed at the DFT level. <sup>f</sup>All Fermi resonances reported in Ten et al.<sup>62</sup> have been considered. <sup>g</sup>MAE stands for Mean Absolute Error. See text.

## RESULTS AND DISCUSSION

**Harmonic Vibrational Frequencies.** Harmonic vibrational frequencies, as obtained from the composite scheme described in the Methodology and Computational Details section, are collected in Table 1. From this table, we first note that the MP2/cc-pVQZ level of theory already provides results close to the CBS limit, the differences being in most cases smaller than  $1 \text{ cm}^{-1}$ ; this is mostly related to the fact that the correlation contributions are already well converged at the MP2/cc-pVQZ level of theory. While core-correlation corrections are quite small, i.e., they range from 1 to  $6 \text{ cm}^{-1}$ , and tend to enlarge when the frequency value increases, the effects due to the inclusion of diffuse functions in the basis set are contradictory, as they range from being negligible ( $<1 \text{ cm}^{-1}$ ) to being large ( $>20 \text{ cm}^{-1}$  in absolute value terms). While CV corrections are always positive, those related to the

diffuse functions are in most cases negative. As concerns the effect of higher-order electron-correlation corrections, for which the inclusion of triples is expected to be the most relevant contribution, we note that they are rather large and mostly negative. The particularly large corrections observed for a few cases when accounting for the  $\Delta\omega(\text{aug})$  and  $\Delta\omega(\text{T})$  contributions deserve to be discussed a little bit more in detail. For the former, as already noticed for the molecular structure (see ref 45), diffuse functions in the basis set are important for correctly describing the oxygen atoms; in fact, large  $\Delta\omega(\text{aug})$  corrections are observed for the two C–O stretchings. As mentioned above, the inclusion of such corrections is expected to recover the limitations of our extrapolation based on the cc-pVTZ and cc-pVQZ basis sets. Less straightforward is how to understand the large effect due to higher excitations observed for the two C–H stretchings, two N–H out-of-plane vibrations, and one ring

Table 3. Combination Bands and Overtones ( $\text{cm}^{-1}$ ) of Uracil from the CCSD(T)/B3LYP Hybrid Force Field<sup>a</sup>

assignment	GVPT2					B3LYP/6-31+G(d,p) <sup>b</sup> (Ten et al.)	Experiment <sup>f</sup>
	DVPT2	Fermi: DFT <sup>c</sup>	Fermi: CC <sup>d</sup>	Fermi: DFT+INP <sup>e</sup>	Fermi: CC+INP <sup>e</sup>		
$\omega_{28} + \omega_{29}$	544	544	543	544	542	569	557
$\omega_{19} + \omega_{30}$	667	667	669	665*	669*	691	682
$\omega_{18} + \omega_{30}$	672	672	685	667*	685*	707	685
$\omega_{27} + \omega_{28}$	947	947	936	947*	936*	984	963
$2\omega_{19}$	1065	1067*	1069*	1053*	1069*	1069	1070
$2\omega_{27}$	1126	1126	1110*	1126	1110*	1141	1102
$\omega_{23} + \omega_{28}$	1180	1180	1179	1180	1179	1197	1192
$\omega_{24} + \omega_{27}$	1304	1304	1294	1304	1294	1317	1283
$2\omega_{26}$	1297	1297	1298	1311	1303	1325	1314
$\omega_{25} + \omega_{26}$	1369	1368*	1370*	1374*	1370*	1382	1366
$\omega_{14} + \omega_{21}$	1445	1445	1445	1445*	1445*	1448	1461
$2\omega_{17}$	1510	1510	1519	1510	1518	1519	1525
$\omega_{16} + \omega_{17}$	1695	1695	1700	1695*	1699*	1698	1699
$\omega_{13} + \omega_{19}$	1703	1703	1710*	1697*	1710*	1705	1707
$\omega_{12} + \omega_{20}$	1729	1729	1729*	1732*	1727*	1720	1720
$\omega_{12} + \omega_{19}$	1753	1749*	1757	1758*	1757*	1739	1733
$\omega_{22} + \omega_{23}$	1746	1746	1745	1746*	1745*	1750	1758
$\omega_{12} + \omega_{18}$	1760	1760	1774	1757*	1773*	1753	1762
$\omega_7 + \omega_{11}$	2998	2998	2998	2999	2998	2984	2970
$2\omega_6$	3456	3456	3431	3481	3442	3470	3477
MAE <sup>g</sup>	13	13	13	14	12	12	
MAE(all) <sup>h</sup>	11	11	11	12	11	12	

<sup>a</sup> The asterisks denote the modes explicitly considered in the variational treatment of Fermi resonances. See text. <sup>b</sup> Ref 62: Anharmonic frequencies and intensities at the B3LYP/6-31+G(d,p) level with Fermi resonances accounted for. <sup>c</sup> Fermi resonances identified through automatic procedure using harmonic frequencies and anharmonic force constants computed at the DFT level. <sup>d</sup> Fermi resonances identified through automatic procedure using harmonic frequencies computed at the CC level and anharmonic force constants computed at the DFT level. <sup>e</sup> All Fermi resonances reported in Ten et al.<sup>62</sup> have been considered. <sup>f</sup> Refs 58–60. <sup>g</sup> MAE stands for Mean Absolute Error. See text. <sup>h</sup> MAE computed considering fundamentals, overtones, and combination bands. See text.

deformation of  $A'$  symmetry. For instance, in the case of the C–H stretchings, such an effect might be related to the significant coupling between the two vibrational modes, which probably requires an improved correlation treatment. Finally, a brief discussion on the accuracy of the best estimated harmonic frequencies is warranted. First of all, we need to point out the role of higher-order effects in the correlation treatment beyond CCSD(T). From the literature available (see for example refs 80 and 81), the full CC singles, doubles, and triples method (CCSDT) is expected to provide no improvements with respect to CCSD(T). The corrections due to quadruple excitations seem to be non-negligible and opposite in sign with respect to core-correlation effects.<sup>81,82</sup> On the other hand, for these contributions, the literature available is very limited and mostly related to diatomics, while the importance of taking into account the effects of core correlation is well recognized.<sup>81–83</sup> On the basis of the approximations made, the estimates for neglected contributions (mainly due to higher excitations beyond CCSDT:  $-0.1\%$  to  $-0.3\%$ , in relative terms), the corrections included as well as the literature on this topic (see, for example, refs 17 and 81), we expect that the accuracy obtained is a few wavenumbers: from  $4\text{ cm}^{-1}$  to  $11\text{ cm}^{-1}$ , where the latter value essentially applies to the larger frequency values.

A comparison of best-estimated harmonic frequencies with DFT results (Table 1) confirms the overall good accuracy of the latter in the present case. In fact, such a comparison shows a mean

absolute error (MAE), with respect to best-estimated values, of about  $7\text{ cm}^{-1}$  and  $11\text{ cm}^{-1}$  for the B3LYP/aug-N07D and B3LYP/aug-cc-pVTZ levels of theory, respectively. The MAEs point out that the B3LYP/aug-N07D level of theory performs a little bit better than the B3LYP/aug-cc-pVTZ one, despite the significantly lower computational cost (232 vs 460 basis sets, respectively). Although this can be due to error compensation, the observed trend is quite general and suggests the B3LYP/augN07D level as the method of choice for spectroscopic investigations of large molecules, provided that DFT does not fail in describing the system under consideration.<sup>39,84</sup>

In the literature, some previous theoretical and experimental data are available for comparison. While the comparison to experiment is meaningful only once anharmonic corrections are accounted for (we therefore postpone such a comparison to the next section), a brief comment is deserved for what concerns theory. Among the literature papers available, we mention the works carried out by Barone et al.<sup>61</sup> and Ten et al.<sup>62</sup> Within the B3LYP approach, the former allows us to point out the improvement in the performance obtained by the aug-N07D set with respect to the 6-31G(d) basis set augmented by diffuse functions only on oxygen atoms. In fact, the discrepancies with respect to the best estimates decrease by a few (2–5) to somewhat greater (up to  $\sim 20$ ) wavenumbers. As concerns ref 62, rather good agreement between the B3LYP/6-31+G(d,p) harmonic frequencies

and our best estimated values is observed, further confirming the good performance of B3LYP for this specific molecule.

**Anharmonic Vibrational Frequencies.** Going beyond the harmonic approximation, as already mentioned in the Introduction, the GVPT2 model applied to anharmonic DFT force fields (in particular with Becke-family hybrid functionals) with the proper treatment of Fermi and Darling–Dennison resonances<sup>15</sup> is known to provide accurate results for semirigid systems.<sup>35,61</sup> As the use of a hybrid CC/DFT scheme improves the accuracy, we limit our discussion to the results obtained by means of the two hybrid models described in the Methodology and Computational Details section, with anharmonic corrections computed at the B3LYP/aug-N07D level. The anharmonic frequencies evaluated with the different models for taking into account resonances are compared in Tables 2 and 3 for fundamental transitions and for overtones and combinational bands, respectively; for comparison purposes, the deperturbed values are also given. In particular, as uracil is known to be a difficult case due to the presence of strong Fermi resonances, such a comparison provides us with additional insights into the performance of the DVPT2 scheme with respect to GVPT2. For GVPT2 calculations, two cases have been actually considered: that where all of the Fermi resonances included in the work by Ten et al.<sup>62</sup> have been considered (Fermi: INP) and that in which Fermi interactions are taken into account through an automatic procedure with the harmonic frequencies evaluated either at the DFT level or from the composite scheme (Fermi: DFT and Fermi: CC, respectively). We first note that all schemes provide very similar results. In Tables 2 and 3, our results are also compared to the available experimental<sup>58–60</sup> and theoretical<sup>62</sup> data. On average, the hybrid CC/DFT approach leads to discrepancies, with respect to experimental results, of about  $10\text{ cm}^{-1}$  and  $13\text{ cm}^{-1}$  (with the largest discrepancies being  $\sim 38\text{ cm}^{-1}$  and  $29\text{ cm}^{-1}$ ) for fundamentals and for overtones and combination bands, respectively. As concerns the results of ref 62, we note that the B3LYP/6-31+G(d,p) level of theory with the proper account of Fermi resonances shows good agreement with experimental results but a worse one than that noted for our hybrid CC/DFT approach. In fact, even though a MAE of  $\sim 12\text{ cm}^{-1}$  is observed, larger discrepancies (the maximum discrepancy is  $56\text{ cm}^{-1}$ ) are evident. As the B3LYP/6-31+G(d,p) anharmonic corrections can be considered as good as ours (B3LYP/aug-N07D) and Ten et al.<sup>62</sup> correctly accounted for all Fermi interactions, the slight improvement obtained should be mainly ascribed to the hybrid approach used, which includes accurate estimates for harmonic frequencies. Anyway, it is worth noting that the present system seems to be a fortunate case, as usually the performance of DFT at the harmonic level with respect to highly correlated methods is definitely worse.<sup>33,37,39,42,43</sup>

As we claim for our anharmonic frequencies an overall accuracy of about  $11\text{ cm}^{-1}$ , we can consider ours state-of-the-art results. For such a challenging case, systematic strategies to further reduce MAE below  $10\text{ cm}^{-1}$  are not yet available. Additionally, it should be noted that in many cases the accuracy of experimental data (in particular for the lower intensity transitions) is not sufficient to justify the effort to reach a  $1\text{ cm}^{-1}$  agreement between theory and experimental results. It is worth noting that the largest discrepancies are found for the  $1700\text{--}1800\text{ cm}^{-1}$  frequency range, where several intense transitions due to the Fermi interaction involving C–O stretching vibrations are exhibited. In fact, both ours and Ten et al.'s investigations predict a number of vibrational transitions in this zone, a

few of them with similar intensity.<sup>62</sup> Therefore, the proper assignment is cumbersome because of different possible interactions (see modes involved in Fermi resonances). The direct comparison between the computational simulated IR spectra and the experimental counterpart would be more meaningful. For uracil, such an attempt has been carried by Ten et al.,<sup>62</sup> but the analysis was hampered by the lack of experimental data in the proper numerical form as well as by the limited experimental resolution.

As the vibrational spectrum of uracil is dominated by a large number of resonances, these deserve to be discussed in some detail. As mentioned above, the automatic procedure used to define possible Fermi interactions (Fermi: DFT; Fermi: CC) has been compared to the results where all Fermi interactions (Fermi: INP) postulated to be important<sup>62</sup> have been considered. A total of 16 (DFT) or 13 (CC) Fermi-type interactions have been pointed out by the automatic procedure, and the corresponding vibrational energy levels have been included in the variational treatment. As the *ad hoc* definition of additional interactions (leading to a total of 33 and 41 resonances for DFT and CC, respectively) does not improve the agreement with experimental results, the present results confirm the reliability of our “black-box” procedure, which takes into account both the zero-order energy difference between two resonant states and the strength of the coupling. Furthermore, we note that in the present case, due to the overall good accuracy of B3LYP/aug-N07D harmonic frequencies, similar results are obtained either by the *a posteriori* approach or by including best estimates for harmonic frequencies in the perturbative treatment. However, this conclusion is not of general validity. As a matter of fact, previous studies<sup>37,39</sup> unambiguously showed that in difficult cases, i.e., when the DFT harmonic frequencies present significant discrepancies with respect to best estimates, proper inclusion of accurate harmonic frequencies into the perturbative expressions leads to significant improvements. In general, the hybrid scheme (in particular Fermi: CC) has to be recommended whenever feasible, since it guarantees reliable harmonic frequencies as well as a proper definition of Fermi resonances. For the latter problem, the development and validation of alternative, more general VPT2 approaches that completely avoid resonant terms is in progress.<sup>79</sup>

**Harmonic Vibrational Intensities.** As concerns IR intensities computed within the double harmonic approximation, it has been shown by Schaefer et al.<sup>85,86</sup> that improvements in the electron-correlation treatment lead to converged values and that quantitative IR intensity predictions can be obtained at the CCSD(T) level in conjunction with basis sets of at least aug-cc-pVTZ quality.<sup>86</sup> In the present work, the extent of various contributions to IR intensities (within the harmonic approximation) has been investigated by means of the composite scheme introduced in the Methodology and Computational Details section. The results, reported in Table 4, show a slow convergence to the CBS limit, unlike what was observed for harmonic frequencies. In fact, differences up to tens of kilometers per mole are obtained by comparing the MP2/cc-pVTZ and MP2/cc-pVQZ levels (see the sixth column). Even for this property, CV corrections appear to be small, being negligible in most cases. Even if not explicitly reported in Table 4, we note that the effects of diffuse functions (pointed out by comparing MP2-cc-pVTZ and MP2-aug-cc-pVTZ) are in most cases on the order of 10%, which means for the most intense transitions intensity enhancements of even  $60\text{--}90\text{ kcal/mol}$ . The same conclusion is drawn for higher-order correlation contributions, i.e.,  $\Delta I((T))$  corrections.

Table 4. Double Harmonic IR Intensities (km/mol) of Uracil

assignment	B3LYP/aug-N07D	B3LYP/aug-cc-pVTZ	MP2/aug-cc-pVTZ	$\Delta CV$	$\Delta(QZ-TZ)$	$\Delta(T)$	best <sup>d</sup> estimate
$\nu(N1-H)$	105.55	102.37	122.44	1.16	3.24	-14.40	112.44
$\nu(N3-H)$	68.51	65.58	77.11	1.01	3.89	-9.48	72.53
$\nu(C5-H)$	1.19	1.19	2.58	0.19	0.54	-1.22	2.09
$\nu(C6-H)$	2.50	2.31	2.08	-0.10	-0.46	0.28	1.80
$\nu(C2=O)$	626.64	599.84	753.68	-0.09	28.12	-55.76	725.95
$\nu(C4=O)$	770.74	774.96	568.37	6.84	56.71	88.39	720.31
$\nu(C5=C6)$	60.30	66.95	19.00	-0.42	-2.15	17.68	34.12
$\delta(N1-H)$	88.66	81.81	121.62	2.39	-5.10	-13.07	105.84
$\delta(C6-H)$	92.71	86.84	66.33	3.20	5.62	27.71	102.86
$\delta(N3-H)$	9.89	7.88	38.65	-4.72	-9.11	-20.23	4.59
$\delta(C5-H)$	35.17	56.65	9.86	-0.48	0.10	2.68	12.17
$\nu(\text{ring})$	2.76	10.86	15.54	1.14	2.71	0.17	19.57
$\nu(\text{ring})$	106.52	104.44	108.11	-2.04	2.21	2.36	110.64
$\nu(\text{ring})$	5.59	4.75	5.43	0.28	-0.59	1.30	6.41
$\nu(\text{ring})$	7.17	7.35	8.05	0.02	0.44	-0.24	8.28
$\nu(\text{ring})$	0.04	0.10	0.48	-0.04	-0.10	0.15	0.49
$\nu(\text{ring})$	3.27	3.65	3.67	-0.01	0.26	0.16	4.08
$\delta(\text{ring})$	36.66	41.70	40.88	0.80	1.03	-2.55	40.16
$\delta(\text{ring})$	6.52	7.39	6.57	-0.00	0.27	-0.27	6.57
$\delta(\text{ring})$	21.29	21.64	18.35	0.11	0.62	1.47	20.55
$\delta(C=O)$	20.97	20.74	20.68	0.13	0.40	1.32	22.53
$\gamma(C6-H)$	9.35	10.93	10.91	-0.14	0.41	-0.28	10.90
$\gamma(C5-H)$	61.73	60.69	49.90	0.74	-2.22	1.60	50.01
$\gamma(C2=O)$	30.98	34.35	30.31	0.64	0.98	-0.43	31.51
$\gamma(C4=O)$	17.79	15.48	9.61	0.55	1.47	0.59	12.22
$\gamma(N3-H)$	75.93	68.57	81.97	-1.20	-1.98	3.16	81.94
$\gamma(N1-H)$	4.64	4.53	3.00	0.03	0.15	0.38	3.56
$\gamma(\text{ring})$	24.03	20.98	23.41	-0.36	-1.67	2.58	23.97
$\gamma(\text{ring})$	0.24	0.08	0.35	-0.06	0.05	0.43	0.67
$\gamma(\text{ring})$	1.57	1.93	0.93	0.11	0.19	-0.33	0.90

<sup>d</sup> From eq 2.

Effects of diffuse functions and core correlation are not always in the same direction. For example, they are both positive for  $\nu(C4=O)$ , while they nearly cancel out for  $\nu(C2=O)$ . On general grounds, it might be of some interest to note that the vibrational modes that show the largest corrections are once again the C–O stretchings. Furthermore, while they exhibit comparable intensities when effects of triple excitations are included, very different intensities are observed at the MP2 level. On the whole, the present results confirm that the prediction of IR intensities is more demanding than band positions.<sup>9</sup> Since they depend on dipole moment derivatives, reliable estimates require an accurate description of the electronic charge density and its variation along normal modes. In turn, these conditions imply that IR intensities are sensitive to electron-correlation effects and basis-set extension. In view of the large extent of some contributions and the lack of literature on this topic, it is difficult to estimate the accuracy of our best estimated values. For these reasons, a benchmark investigation on IR intensities for small- to medium-sized molecules is in progress.

By comparing the results obtained at different levels of theory, it is worth noting that the MP2 level generally tends to overestimate IR intensities with respect to CCSD(T). DFT performs quite well with respect to our best estimated values, with a MAE of about 10 and 13 km/mol for the B3LYP/aug-N07D and

B3LYP/aug-cc-pVTZ levels, respectively. In Table 2, the available experimental data are reported. Even though our computed harmonic intensities are able to reproduce on average the experimental trend, a quantitative comparison is not possible at the present stage. Work along this direction is in progress, and it mainly concerns benchmark studies in view of verifying the accuracy obtainable by QM computations as well as the investigation of resonance effects on anharmonic contributions.

## CONCLUSIONS

Vibrational frequencies and intensities have been investigated within the harmonic approximation at a high level of theory. Following our previous work,<sup>45</sup> a composite scheme has been exploited in order to account for basis-set effects and to include core-correlation corrections. As the accuracy in vibrational frequencies evaluation is mostly related to the harmonic force field, the computational approach introduced paves the way toward spectroscopic accuracy for molecules with a larger and larger number of atoms. In particular, the approach presented is expected to be accurate and reliable for all types of systems, also when DFT fails in a proper description of the electronic structure. The comparison to available experimental data points out that the use of hybrid CC/DFT schemes with the proper



account of Fermi interactions allows one to obtain state-of-the-art anharmonic frequencies for a system as challenging as uracil. Within pure DFT approaches, for both frequencies and intensities, the present study confirms that the aug-N07D basis set is able to provide results close to or even better than those obtained with the larger aug-cc-pVTZ set. We also note that, in the present case, the accuracy of DFT results is at least comparable to that of MP2 computations employing extended basis sets, but with a considerable savings of computational time and thus a much better scaling with the dimensions of the system.

In conclusion, even if further developments are still required and under consideration, in our opinion, the results of the present investigation show the remarkable performance of composite schemes as well as of integrated CC and DFT approaches in the field of vibrational spectroscopy.

## AUTHOR INFORMATION

### Corresponding Author

\*E-mail: cristina.puzzarini@unibo.it

## ACKNOWLEDGMENT

This work was supported by Italian MIUR (PRIN 2008, PRIN 2009) and by the University of Bologna (RFO funds). The COST-CMTS Action CM1002 "COnvergent Distributed Environment for Computational Spectroscopy (CODECS)" is also acknowledged.

## REFERENCES

- Barone, V. *Computational Strategies for Spectroscopy, from Small Molecules to Nano Systems*; John Wiley & Sons, Inc.: Chichester, U. K., 2011.
- Jensen, P.; Bunker, P. R. *Computational Molecular Spectroscopy*; John Wiley & Sons: Chichester, U. K., 2000.
- Puzzarini, C.; Coriani, S.; Rizzo, A.; Gauss, J. *Chem. Phys. Lett.* **2005**, *409*, 118–123.
- Puzzarini, C.; Stanton, J. S.; Gauss, J. *Int. Rev. Phys. Chem.* **2010**, *29*, 273–367.
- Sattelmeyer, K. W.; Schaefer, H. F., III. *J. Chem. Phys.* **2002**, *117*, 7914–7416.
- Sharma, A. R.; Braams, B. J.; Carter, S.; Shepler, B. C.; Bowman, J. M. *J. Chem. Phys.* **2009**, *130*, 174301/1–9.
- Barone, V.; Bloino, J.; Biczysko, M. *Phys. Chem. Chem. Phys.* **2010**, *12*, 1092–1101.
- Nikow, M.; Wilhelm, M. J.; Dai, H.-L. *J. Phys. Chem. A* **2009**, *113*, 8857–8870.
- Cappelli, C.; Biczysko, M. In *Computational Strategies for Spectroscopy, from Small Molecules to Nano Systems*; Barone, V., Ed.; Wiley: Chichester, U. K., 2011; Chapter: Time Independent Approach to Vibrational Spectroscopies, pp 309–360.
- Mills, I. M. In *Molecular Spectroscopy: Modern Research*; Rao, K. N., Mathews, C. W., Eds.; Academic: New York, 1972.
- Amos, R. D.; Handy, N. C.; Green, W. H.; Jayatilaka, D.; Willets, A.; Palmieri, P. *J. Chem. Phys.* **1991**, *95*, 8323–8336.
- Gaw, F.; Willetts, A.; Handy, N.; Green, W. In *SPECTRO - A Program for Derivation of Spectroscopic Constants from Provided Quartic Force Fields and Cubic Dipole Fields*; Bowman, J. M., Ed.; JAI Press: Greenwich, CT, 1991; Vol. 1B, pp 169–185.
- Barone, V. *J. Chem. Phys.* **2005**, *122*, 014108/1–10.
- Barone, V. *J. Chem. Phys.* **2004**, *120*, 3059–3065.
- Martin, J. M. L.; Lee, T. J.; Taylor, P. R.; François, J.-P. *J. Chem. Phys.* **1995**, *103*, 2589–2591.
- Stanton, J. F.; Gauss, J. *J. Chem. Phys.* **1998**, *108*, 9218–9820.
- Tew, D. P.; Klopper, W.; Heckert, M.; Gauss, J. *J. Phys. Chem. A* **2007**, *111*, 11242–11248.
- Ruden, T. A.; Taylor, P. R.; Helgaker, T. *J. Chem. Phys.* **2003**, *119*, 1951–1960.
- Ruud, K.; Åstrand, P. O.; Taylor, P. R. *J. Chem. Phys.* **2000**, *112*, 2668–2683.
- Vázquez, J.; Stanton, J. F. *Mol. Phys.* **2006**, *104*, 377–388.
- Vázquez, J.; Stanton, J. F. *Mol. Phys.* **2007**, *105*, 101–109.
- Stanton, J. F.; Gauss, J. *Int. Rev. Phys. Chem.* **2000**, *19*, 61–95.
- Bloino, J.; Guido, C.; Lipparini, F.; Barone, V. *Chem. Phys. Lett.* **2010**, *496*, 157–161.
- Bowman, J. M. *Science* **2000**, *290*, 724–725.
- Bowman, J. M.; Carter, S.; Huang, X. *Int. Rev. Phys. Chem.* **2003**, *22*, 533–549.
- Carter, S.; Handy, N. *J. Chem. Phys.* **2000**, *113*, 987–993.
- Chaban, G.; Jung, J.; Gerber, R. *J. Chem. Phys.* **1999**, *111*, 1823–1829.
- Rauhut, G.; Hrenar, T. *Chem. Phys.* **2008**, *346*, 160–166.
- Christiansen, O. *Phys. Chem. Chem. Phys.* **2007**, *9*, 2942–2953.
- Norris, L. S.; Ratner, M. A.; Roitberg, A. E.; Gerber, R. B. *J. Chem. Phys.* **1996**, *105*, 11261–11268.
- Christiansen, O. *J. Chem. Phys.* **2003**, *119*, 5773–5781.
- Carter, S.; Sharma, A. R.; Bowman, J. M.; Rosmus, P.; Tarroni, R. *J. Chem. Phys.* **2009**, *131*, 224106–224121.
- Carbonniere, P.; Lucca, T.; Pouchan, C.; Rega, N.; Barone, V. *J. Comput. Chem.* **2005**, *26*, 384–388.
- Boese, A. D.; Martin, J. M. L. *J. Phys. Chem. A* **2004**, *108*, 3085–3096.
- Barone, V. *J. Phys. Chem. A* **2004**, *108*, 4146–4150.
- Barone, V. *Chem. Phys. Lett.* **2004**, *383*, 528–532.
- Puzzarini, C.; Biczysko, M.; Barone, V. *J. Chem. Theory Comput.* **2010**, *6*, 828–838.
- Neese, F.; Schwabe, T.; Grimme, S. *J. Chem. Phys.* **2007**, *126*, 124115/1–15.
- Biczysko, M.; Panek, P.; Scalmani, G.; Bloino, J.; Barone, V. *J. Chem. Theory Comput.* **2010**, *6*, 2115–2125.
- Kozuch, S.; Gruzman, D.; Martin, J. M. L. *J. Phys. Chem. C* **2010**, *114*, 20801–20808.
- Puzzarini, C.; Barone, V. *J. Chem. Phys.* **2008**, *129*, 084306/1–7.
- Puzzarini, C.; Barone, V. *Phys. Chem. Chem. Phys.* **2008**, *10*, 6991–6997.
- Begue, D.; Carbonniere, P.; Pouchan, C. *J. Phys. Chem. A* **2005**, *109*, 4611–4616.
- Begue, D.; Benidar, A.; Pouchan, C. *Chem. Phys. Lett.* **2006**, *430*, 215–220.
- Puzzarini, C.; Barone, V. *Phys. Chem. Chem. Phys.* **2011**, *13*, 7158–7166.
- Harsanyi, L.; Csaszar, P. *Acta Chim. Hung.* **1983**, *113*, 257–278.
- Lordand, R. C.; Thomas, G. *J. Spectrochim. Acta A* **1967**, *23*, 2551–2591.
- Aamouche, A.; Ghomi, C.; Coulombeau, M.; Jobic, H.; Grajcar, L.; Baron, M. H.; Baumruk, V.; Turpin, P. Y.; Henriët, C.; Berthier, G. *J. Phys. Chem.* **1996**, *100*, S224–S227.
- Susi, H.; Ard, J. S. *Spectrochim. Acta A* **1971**, *27*, 1549–1582.
- Florian, J.; Hroudá, V. *Spectrochim. Acta A* **1993**, *49*, 921–938.
- Wojcik, M. *J. Mol. Struct.* **1990**, *219*, 305–310.
- Barnes, A. J.; Stuckey, M. A.; Le Gall, L. *Spectrochim. Acta A* **1984**, *40*, 419–431.
- Wojcik, M. J.; Rostkowska, H.; Szczepaniak, K.; Person, W. B. *Spectrochim. Acta A* **1989**, *45*, 499–502.
- Nowak, M. *J. Mol. Struct.* **1989**, *193*, 35–49.
- Graindourze, M.; Grootaers, T.; Smets, J.; Zeegers-Huyskens, T.; Maes, G. *J. Mol. Struct.* **1991**, *243*, 37–60.
- Ivanov, A. Y.; Plokhotnichenko, A. M.; Radchenko, E. D.; Sheina, G. G.; Blagoi, Y. P. *J. Mol. Struct.* **1995**, *372*, 91–100.
- Maltese, M.; Passerini, S.; Nunziante-Cesaro, S.; Dobos, S.; Harsanyi, L. *J. Mol. Struct.* **1984**, *116*, 49–65.
- Graindourze, M.; Smets, J.; Zeegers-Huyskens, T.; Maes, G. *J. Mol. Struct.* **1990**, *222*, 345–364.
- Szczesniak, M.; Nowak, M. J.; Rostkowska, H.; Szczepaniak, K.; Person, W. B.; Shugar, D. *J. Am. Chem. Soc.* **1983**, *105*, 5969–5976.

- (60) Chin, S.; Scott, I.; Szczepaniak, K.; Person, W. B. *J. Am. Chem. Soc.* **1984**, *106*, 3415–3422.
- (61) Barone, V.; Festa, G.; Grandi, A.; Rega, N.; Sanna, N. *Chem. Phys. Lett.* **2004**, *388*, 279–283.
- (62) Ten, G. N.; Nechaev, V. V.; Krasnoshchekov, S. V. *Optics Spectrosc.* **2010**, *37*–44.
- (63) Møller, C.; Plesset, M. S. *Phys. Rev.* **1934**, *46*, 618–622.
- (64) Raghavachari, K.; Trucks, G. W.; Pople, J. A.; Head-Gordon, M. *Chem. Phys. Lett.* **1989**, *157*, 479–483.
- (65) Dunning, J., T. H. *J. Chem. Phys.* **1989**, *90*, 1007–1023.
- (66) Kendall, A.; Dunning, T. H., Jr.; Harrison, R. J. *J. Chem. Phys.* **1992**, *96*, 6796–6806.
- (67) Woon, D. E.; Dunning, J., T. H. *J. Chem. Phys.* **1995**, *103*, 4572–4585.
- (68) Stanton, J. F.; Gauss, J. *Chem. Phys. Lett.* **1997**, *276*, 70–77.
- (69) Stanton, J. F.; Gauss, J.; Harding, M. E.; Szalay, P. G. CFour, A quantum chemical program package, with contributions from Auer, A. A.; Bartlett, R. J.; Benedikt, U.; Berger, C.; Bernholdt, D. E.; Bomble, Y. J.; Christiansen, O.; Heckert, M.; Heun, O.; Huber, C.; Jagau, T.-C.; Jonsson, D.; Jusélius, J.; Klein, K.; Lauderdale, W. J.; Matthews, D.; Metzroth, T.; O'Neill, D. P.; Price, D. R.; Prochnow, E.; Ruud, K.; Schiffmann, F.; Stopkowitz, S.; Varner, M.; Vázquez, J.; Watts, J. D.; Wang, F. and the integral packages MOLECULE (Almloef, J.; Taylor, P. R.), PROPS (Taylor, P. R.), ABACUS (Helgaker, T.; Jensen, H. J. Aa.; Jørgensen, P.; Olsen, J.) and ECP routines by Mitin, A. V.; van Wuelen, C. For the current version, see <http://www.cfour.de> (accessed August 7, 2011).
- (70) Becke, D. *J. Chem. Phys.* **1993**, *98*, 5648–5652.
- (71) Double and triple- $\zeta$  basis sets of N07 family are available for download. See <http://idea.sns.it> (accessed June 30, 2011).
- (72) Barone, V.; Cimino, P.; Stendardo, E. *J. Chem. Theory Comput.* **2008**, *4*, 751–764.
- (73) Barone, V.; Cimino, P. *Chem. Phys. Lett.* **2008**, *454*, 139–143.
- (74) Barone, V.; Cimino, P. *J. Chem. Theory Comput.* **2009**, *5*, 192–199.
- (75) Barone, V.; Biczysko, M.; Bloino, J.; Borkowska-Panek, M.; Carnimeo, I.; Panek, P. *Int. J. Quantum Chem.* **2011**, [Online] DOI: 10.1002/qua.23224.
- (76) Carnimeo, I.; Biczysko, M.; Bloino, J.; Barone, V. *Phys. Chem. Chem. Phys.* **2011**, *13*, 16713–16727.
- (77) Note that only the  $K_{ijk}$ ,  $K_{iiii}$  and  $K_{ijij}$  force constants are used in the vibrational perturbative treatment (VPT2).
- (78) Frisch, M. J.; Trucks, G. W.; Schlegel, H. B.; Scuseria, G. E.; Robb, M. A.; Cheeseman, J. R.; Scalmani, G.; Barone, V.; Mennucci, B.; Petersson, G. A.; Nakatsuji, H.; Caricato, M.; Li, X.; Hratchian, H. R.; Izmaylov, A. F.; Bloino, J.; Zheng, G.; Sonnenberg, J. L.; Hada, M.; Ehara, M.; Toyota, K.; Fukuda, R.; Hasegawa, J.; Ishida, M.; Nakajima, T.; Honda, Y.; Kitao, O.; Nakai, H.; Vreven, T.; Montgomery, J. R., Jr.; Peralta, J. A.; Ogliaro, F.; Bearpark, M.; Heyd, J. J.; Brothers, E.; Kudin, K. N.; Staroverov, V. N.; Kobayashi, R.; Normand, J.; Raghavachari, K.; Rendell, A.; Burant, J. C.; Iyengar, S. S.; Tomasi, J.; Cossi, M.; Rega, N.; Millam, J. M.; Klene, M.; Knox, J. E.; Cross, J. B.; Bakken, V.; Adamo, C.; Jaramillo, J.; Gomperts, R.; Stratmann, R. E.; Yazyev, O.; Austin, A. J.; Cammi, R.; Pomelli, C.; Ochterski, J. W.; Martin, R. L.; Morokuma, K.; Zakrzewski, V. G.; Voth, G. A.; Salvador, P.; Dannenberg, J. J.; Dapprich, S.; Daniels, A. D.; Farkas, O.; Foresman, J. B.; Ortiz, J. V.; Cioslowski, J.; Fox, D. J. *Gaussian 09*, Revision B.01; Gaussian Inc.: Wallingford, CT, 2009.
- (79) Kuhler, K. M.; Truhlar, D. G.; Isaacson, A. D. *J. Chem. Phys.* **1996**, *104*, 4664–4671.
- (80) Feller, D. A.; Sordo, J. A. *J. Chem. Phys.* **2000**, *112*, 5604–5610.
- (81) Ruden, T. A.; Helgaker, T.; Jørgensen, P.; Olsen, J. *J. Chem. Phys.* **2004**, *121*, 5874–5884.
- (82) Martin, J. M. L. *Chem. Phys. Lett.* **1998**, *292*, 411–420.
- (83) Pawłowski, F.; Halkier, A.; Jørgensen, P.; Bak, K. L.; Helgaker, T.; Klopper, W. *J. Chem. Phys.* **2003**, *118*, 2539–2549.
- (84) Bloino, J.; Biczysko, M.; Santoro, F.; Barone, V. *J. Chem. Theory Comput.* **2010**, *6*, 1256–1274.
- (85) Thomas, J. R.; DeLeeuw, B. J.; Vacek, G.; Crawford, T. D.; Yamaguchi, Y.; Schaefer, H. F., III. *J. Chem. Phys.* **1993**, *99*, 403–412.
- (86) Galabov, B.; Yamaguchi, Y.; Remington, R. B.; Schaefer, H. F., III. *J. Phys. Chem. A* **2002**, *106*, 819–832.

# Polarizable Force Fields and Polarizable Continuum Model: A Fluctuating Charges/PCM Approach. 1. Theory and Implementation

Filippo Lipparini\* and Vincenzo Barone

Scuola Normale Superiore, Piazza dei Cavalieri 7, 56126 Pisa, Italy

**ABSTRACT:** We present a combined fluctuating charges–polarizable continuum model approach to describe molecules in solution. Both static and dynamic approaches are discussed: analytical first and second derivatives are shown as well as an extended lagrangian for molecular dynamics simulations. In particular, we use the polarizable continuum model to provide nonperiodic boundary conditions for molecular dynamics simulations of aqueous solutions. The extended lagrangian method is extensively discussed, with specific reference to the fluctuating charge model, from a numerical point of view by means of several examples, and a rationalization of the behavior found is presented. Several prototypical applications are shown, especially regarding solvation of ions and polar molecules in water.

## 1. INTRODUCTION

The astonishing development of computational resources during recent decades has made possible studies of larger and larger molecular systems together with the computation of accurate and complex physical–chemical properties. Both classical and quantum mechanical (QM) approaches have enormously increased either their range of application or their accuracy or both, allowing the study of several processes ranging from folding studies of huge biological systems to extremely accurate computations of spectroscopic parameters for medium-large molecules.

Complex systems, like nanostructured ones or solutions and, more in general, what is usually referred to as “Soft Matter” represent, nevertheless, an interesting challenge for theoretical and computational chemistry. The huge dimensionality of such systems, which can be considered microscopic but certainly not molecular, is still far beyond the possibilities of modern computational infrastructures: a computational study of a complex system by means of standard tools is nowadays still unfeasible when a QM treatment of the whole system is required. This can be both a curse and a blessing: the quantity of data arising from the direct study of such a system would be difficult to analyze and even more difficult to translate into chemically understandable information when a local property tuned by the chemical environment is the target of the study.

The unfeasibility of “brute force” approaches, on the other hand, is not to be considered as an insuperable obstacle. Chemical intuition is often the way to get a valuable answer at a reasonable price and is the driving force in the definition of *focused models*, where the target of a study is well-defined and distinguished from the environment, as complex as it might be, that surrounds it.

A prototypical focused model may use different levels of theory, from a very sophisticated QM approach to describe the core, to a cheaper one for its closest surroundings, to a classical but still atomistic one for the distant surroundings, to a continuum to describe the boundaries. In this paper, we will focus on the two latter shells and, in particular, on their interface. As the

continuum is concerned, the Polarizable Continuum Model (PCM)<sup>1,2</sup> is one of the most successful models, thanks to its generality and its versatility. The PCM represents a solvent, or other more complex matrices<sup>3</sup> such as an anisotropic medium or a weak ionic solution or even a metal nanoparticle, by means of a polarizable, infinite, dielectric medium which surrounds a molecular cavity that accommodates the “solute”. However, when dealing with solvents responsible for specific interactions like hydrogen bonds, a continuous approach may not be sufficient to achieve a correct description of the system: a mixed continuous–atomistic treatment of the solvent, using molecular mechanics (MM) to describe the atomistic portion, can be greatly beneficial.<sup>4–12</sup> On the other hand, the mixed strategy is advantageous with respect to a fully atomistic one as the PCM easily takes into account the long-range interactions that would require a huge number of solvent molecules, increasing significantly the computational cost of the simulation, and implicitly includes the statistical average of their configurations.

In this paper, we will present a combined PCM/MM description using a polarizable force field. The most popular approaches used to include polarization effects in MM include the induced point dipole method,<sup>13</sup> the classical Drude oscillator model,<sup>14</sup> and the fluctuating charges model.<sup>15–17</sup>

We find the FQ model particularly appealing in view of its strong connection both with quantum mechanics and classical electrostatics: the model is based on concepts, such as atomic hardness and electronegativity, which can be rigorously defined in the framework of density functional theory (DFT);<sup>18,19</sup> on the other, the electronic distribution is represented by effective atomic charges which interact classically. There is a strong connection with the formalism adopted by semiempirical methods, like the density functional tight binding<sup>20,21</sup> approach, and the FQ model, for they both treat the electronic polarization with some suitably defined—and QM derived—charges that are made self-consistent; on the other hand, the same strong formal

Received: June 6, 2011

Published: September 15, 2011

analogy holds between the FQ model and Apparent Surface Charge (ASC) methods like PCM, where the definition of the polarizable charges is classical. We find the smoothness in switching from a classical and continuous description to an atomistic and quantum mechanical one aesthetically fascinating and promising in the perspective of a complete, multiscale description of complex systems. Another advantage of the FQ model with respect to the point dipole method is that only the electrostatic potential needs to be calculated: as the electric field is discontinuous at the cavity surface, its use can be a source of numerical instabilities which are avoided using only the potential.

Two different possibilities are offered by the FQ model, eventually coupled with the PCM: It can be used to calculate molecular properties by means of response theory or analytical derivatives in a standard, static fashion or for molecular dynamics simulations. We will refer to the first approach as “time independent” and to the second as “time dependent”. In the first case, the FQ model can be used to describe the solute and (or) a few molecules of solvent in a QM/MM/PCM fashion: the standard machinery of computational chemistry can hence be employed to calculate structural and spectroscopic properties.<sup>5,8,12</sup> On the other hand, the PCM is an effective and physically suitable way to enforce nonperiodic boundary conditions (nPBC) in molecular dynamics (MD) simulations.<sup>22–27</sup> While the use of PBC is convenient to describe solids like crystals or metals or pure liquids, this is not always the case when dealing with intrinsically nonperiodic systems, like a molecule in solution: to avoid spurious interactions between the molecule and its copy in a neighboring cell, a large amount of molecules of solvent is to be used. This is especially true with charged solutes, as the Coulomb interaction decays very slowly with the distance.<sup>28–30</sup> The PCM can be successfully and efficiently used to impose boundary conditions by defining a suitable volume (cavity), enclosed by a regular surface like a sphere, a cylinder, or an ellipsoid, that accommodates the solute and the number of molecules of solvent needed to fill the volume. We would like to point out that the PCM treatment of the electrostatics is, in principle, exact: the electrostatic potential is hence well-defined in the whole space, and no discontinuity arises because of the definition of a cavity if this is regular enough.<sup>31</sup> On the other hand, confinement and nonelectrostatic interactions are a more delicate aspect. Nevertheless, it has been shown in the literature<sup>12,26,32–35</sup> that a tailored confinement potential or a proper buffer can provide accurate results also for the description of bulk liquids and avoid compenetrations between the continuous and atomistic portions of the solvent. This is of particular importance when dealing with a polarizable force field, as the penetration of a polarizable molecule in the polarizable dielectric could give rise to instabilities. The aspects of confinement and nonelectrostatics are of course closely related, as the repulsion interaction is responsible for the impossibility for different molecules to compenetrates. The PCM cavity can be kept fixed during the simulation so that the computational cost per step due to the PCM reduces to the calculation of the electrostatic potential at some representative surface point and a matrix/vector multiplication, as will be clarified in section 3. This constraint can be easily relaxed, but while this is necessary in order to do geometry optimizations, previous attempts<sup>22–24</sup> show that a fixed cavity is enough to have a satisfactory description of the solvent.

In this contribution, we will mainly focus on the time dependent approach and hence on molecular dynamics simulations; for

the sake of completeness, the FQ contribution to analytical first and second derivatives will be derived in the perspective of a future development in the time independent field of applications. To the best of our knowledge, this is the first derivation of analytical second derivatives for the FQ model.

This paper is organized as follows: In section 2, the FQ model is discussed, and the FQ/PCM equations are derived. Section 3 focuses on analytical derivatives of the FQ/PCM contribution to the energy and on the use of FQ/PCM in molecular dynamics simulations. Finally, in section 4, some preliminary numerical results are presented.

## 2. THEORY

The FQ model is based on the *electronegativity equalization principle*<sup>18,36</sup> (EEP), which states that, at equilibrium, the instantaneous electronegativity  $\tilde{\chi}$  of each atom has the same value. Considering an isolated atom, it is possible to expand in Taylor series its energy with respect to the net charge on the atom itself. To the second order:

$$E = E_0 + \frac{\partial E}{\partial Q} Q + \frac{1}{2} \frac{\partial^2 E}{\partial Q^2} Q^2 \quad (1)$$

The parameters—i.e., the energy derivatives—that appear in eq 1 hold a clear physical significance: the first derivative is in fact a Mulliken electronegativity, while the second is a chemical hardness.

$$\left. \frac{\partial E}{\partial Q} \right|_{Q=0} = \chi, \quad \left. \frac{\partial^2 E}{\partial Q^2} \right|_{Q=0} = 2\eta$$

It is possible to extend eq 1 to a molecular system by taking into account the interaction between charges located on different sites:

$$E = E_0 + \sum_i [\chi_i q_i + \eta_i q_i^2 + \sum_{j>i} J_{ij} q_i q_j] \quad (2)$$

where the hardness kernel  $J$  represents, as stated, the interaction, and the sum runs over the nuclei. The electronegativity of the  $i$ th atom is hence defined as the derivative of the energy with respect to the  $i$ th charge:

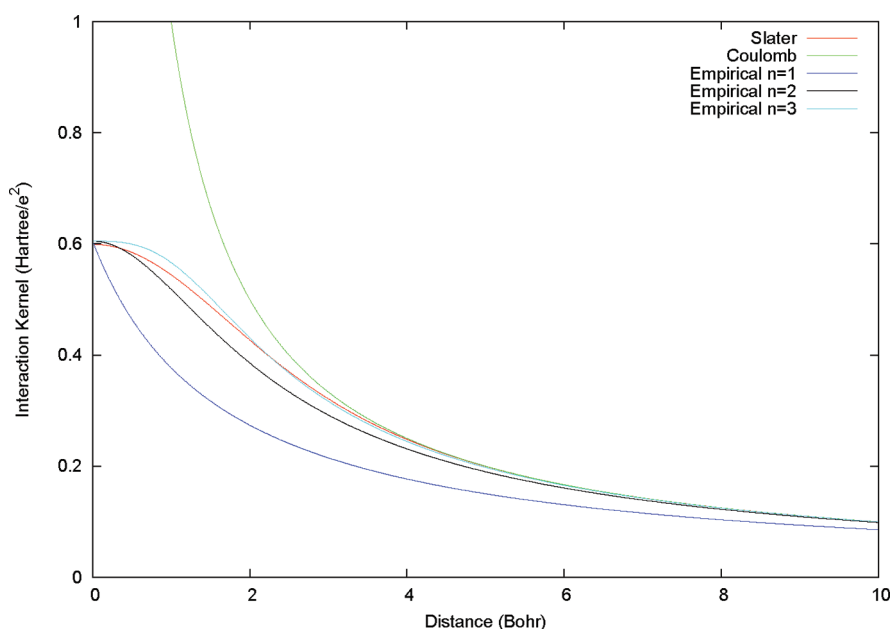
$$\tilde{\chi}_i = \frac{\partial E}{\partial q_i} = \chi_i + \sum_j J_{ij} q_j \quad (3)$$

where we put  $J_{ii} = 2\eta_i$ . The EEP states that the electronegativity of each atom in the molecule has the same value. This can be stated in an equivalent, but more advantageous, formulation defining the atomic partial charges as the constrained minimum of functional

$$F(\mathbf{q}, \lambda) = E_0 + \sum_i [\chi_i q_i + \eta_i q_i^2 + \sum_{j>i} J_{ij} q_i q_j] + \lambda (\sum_i q_i - Q_{\text{tot}})$$

where the constraint, imposed by means of a Lagrange multiplier, is meant to preserve the total charge. When more than a molecule is present, there are two possible different strategies to impose the charge constraint:

1. The entire system is constrained to have charge  $Q_{\text{tot}}$ , and no constraint is imposed on each molecule. This allows



**Figure 1.** Comparison between different expressions for the interaction kernel.

intermolecular charge transfer and makes, at the equilibrium, the electronegativity of each atom the same.

- Each molecule is constrained to assume a fixed, total charge  $Q_\alpha$  (requiring of course these charges to sum to  $Q_{\text{tot}}$ ). Following this second possibility, the electronegativity of each atom in the same molecule will be the same but will have in general different values among different molecules.

It has been pointed out<sup>15</sup> that the first choice allows for charge transfer even when it is unphysical—i.e., when two molecules are separated by a large distance. In the literature, some models have been proposed<sup>37–42</sup> to take into account in a correct way charge transfer; however, we will not deal with such phenomena here and adopt the second choice for constraints. Dropping the constant term, the functional to be minimized reads thus

$$\begin{aligned}
 F(\mathbf{q}, \boldsymbol{\lambda}) &= \sum_{\alpha, i} q_{\alpha i} \chi_{\alpha i} + \frac{1}{2} \sum_{\alpha, i} \sum_{\beta, j} q_{\alpha i} J_{\alpha i, \beta j} q_{\beta j} + \sum_{\alpha} \lambda_{\alpha} \sum_i (q_{\alpha i} - Q_{\alpha}) \\
 &= \mathbf{q}^\dagger \boldsymbol{\chi} + \frac{1}{2} \mathbf{q}^\dagger \mathbf{J} \mathbf{q} + \boldsymbol{\lambda}^\dagger \mathbf{q}
 \end{aligned} \quad (4)$$

where the Greek indexes run on molecules and the Latin ones on atoms of each molecule.

The interaction kernel  $\mathbf{J}$  describes the Coulomb repulsion between two atoms and, from a quantum mechanical point of view, can be conveniently described in terms of the Coulomb interaction between two charge distributions represented by spherical (s) Slater orbitals:

$$J_{ij}(r_{ij}) = \int_{\mathbb{R}^3} d\mathbf{r} \int_{\mathbb{R}^3} d\mathbf{r}' \frac{|\phi_i(\mathbf{r} - \mathbf{r}_i)|^2 |\phi_j(\mathbf{r}' - \mathbf{r}_j)|^2}{|\mathbf{r} - \mathbf{r}'|} \quad (5)$$

where

$$\phi_i(r) = \mathcal{N}_i r^{n_i - 1} e^{-\xi_i r}$$

$\mathcal{N}_i$  is a normalization constant,  $\mathbf{r}_i$  is the position of the  $i$ th nucleus,  $n_i$  is the principal quantum number, and  $\xi_i$  is the Slater exponent. This choice was the one proposed by Rappé and

co-workers in their pioneering work<sup>36</sup> and is widely pursued in the literature.<sup>15,40,43–46</sup> As the integral in eq 5 goes to the bare coulomb interaction when two sites are sufficiently distant, this choice is usually reserved to the intramolecular terms, while the intermolecular ones are described as classical Coulomb interactions between point charges. This choice is particularly pursued when rigid molecules are considered: the integrals are to be computed only at the beginning of the simulation, and only the intermolecular contributions need to be updated on the fly, which can be done very efficiently using linear scaling techniques. A different, but in principle similar, choice describes the interaction by means of Gaussian orbitals,<sup>19,38,47,48</sup> which allow an easy and efficient computation of the interaction kernel, exploiting the machinery of common quantum chemistry codes. The integrals that need to be calculated are

$$\begin{aligned}
 J_{ij}(r_{ij}) &= \int_{\mathbb{R}^3} d\mathbf{r} \int_{\mathbb{R}^3} d\mathbf{r}' \frac{|\phi_i(\mathbf{r} - \mathbf{r}_i)|^2 |\phi_j(\mathbf{r}' - \mathbf{r}_j)|^2}{|\mathbf{r} - \mathbf{r}'|} \\
 &= \frac{1}{r_{ij}} \operatorname{erf} \left( \frac{r_{ij}}{\sqrt{R_i^2 + R_j^2}} \right)
 \end{aligned} \quad (6)$$

where

$$\varphi(\mathbf{r} - \mathbf{r}_i) = \frac{1}{(R_i^2 \pi)^{3/2}} e^{-r - r_i/R_i^2}$$

and  $R_i$  is the width of the distribution. Some authors<sup>19,47</sup> generalized the ansatz of s-type functions to s- and p-type Gaussian type orbitals: this allows the model to accurately reproduce out of plane contributions to polarizabilities for planar molecules.

While the use of basis functions is general and elegant, especially when the relation between FQ and the semiempirical model is considered, the computation of the integrals at each step of a molecular dynamics simulation with flexible molecules would be overwhelmingly costly. A different strategy has been

proposed by several authors<sup>49–51</sup> who approximate the Slater integral by means of an empirical Coulombic interaction formula. Let  $\eta_{ij}$  be a parameter related to the interaction of the atoms  $i$  and  $j$ —for instance, this might be the value of the Slater integral at zero distance, but more in general it can be considered a parameter. The approximate formula is

$$J_{ij}(r_{ij}) = \frac{\eta_{ij}}{[1 + \eta_{ij}^n r_{ij}^{1/n}]^{1/n}} \quad (7)$$

In the literature, three different exponents ( $n = 1$ ,<sup>49</sup>  $n = 2$ ,<sup>50</sup> and  $n = 3$ <sup>51</sup>) have been proposed. We report in Figure 1 the interaction kernel element between an oxygen and a hydrogen atom  $J_{\text{OH}}$ , as a function of the distance between the two centers, calculated with both eq 5 and eq 7 with the three exponents proposed, in comparison to the bare Coulomb interaction. All of these expressions show the correct asymptotic behavior, and it is possible to see that the approximate formula, eq 7, closely resembles the Slater integral, especially with exponents  $n = 2$  and  $n = 3$ . Throughout this work, we will always adopt the  $n = 2$  choice.

The hardness parameters  $\eta_{ij}$  lack a precise physical significance when used to approximate eq 5 but make perfect mathematical sense as the limit for small interatomic distances of the aforementioned integral:

$$\eta_{ij} = \lim_{r_{ij} \rightarrow 0} J_{ij}(r_{ij}) \quad (8)$$

As a consequence of their definition, in principle, one should define as many hardness parameters as the number of couples of different atoms. It has been proposed by some authors<sup>50,51</sup> to define the off-diagonal elements as the arithmetical or geometrical averages of the diagonal ones: this is, of course, an approximation, but it reduces greatly the number of parameters to be considered. We will examine the effects of this approximation in section 4 with a numerical example. The obvious advantage of the use of an empirical formula is in terms of computational effort: no integral needs to be evaluated, making the FQ method suitable for flexible molecule MD simulations. This definition of the interaction kernel is the one used in the FQ CharMM force field<sup>52,53</sup> and in the reactive ReaxFF force field.<sup>54</sup>

By taking the derivative of eq 4 with respect to the charges and to the Lagrange multipliers, one obtains the constrained minimum condition:

$$\begin{cases} \sum_{\beta,j} J_{\alpha i, \beta j} q_{\beta j} + \lambda_{\alpha} = -\chi_{\alpha i} \\ \sum_i q_{\alpha i} = Q_{\alpha} \end{cases} \quad (9)$$

This equation can be recast in a more compact formalism introducing the extended  $\mathbf{D}$  matrix:

$$\begin{pmatrix} \mathbf{J} & \mathbf{1}_{\lambda} \\ \mathbf{1}_{\lambda}^{\dagger} & \mathbf{0} \end{pmatrix}$$

where  $\mathbf{1}_{\lambda}$  is a rectangular matrix which accounts for the Lagrangians. The linear system of equation reads now:

$$\mathbf{D}\mathbf{q}_{\lambda} = -\mathbf{C} \quad (10)$$

where  $\mathbf{C}$  is a vector containing atomic electronegativities and total charge constraints, whereas  $\mathbf{q}_{\lambda}$  is a vector containing charges and Lagrange multipliers. It is important to notice that, although

$\mathbf{J}$  is positive definite,  $\mathbf{D}$  is not, so that it is not possible to solve eq 10 by means of standard minimization procedures.

**2.1. Coupling the FQ Force Field with the Polarizable Continuum Model.** The PCM solves Poisson's equation with suitable boundary conditions in the presence of a dielectric medium with a cavity, where the solute is accommodated. As we are mainly interested in polar solvents involved in hydrogen bond formation, we will restrict our discussion to the conductor-like model (C-PCM);<sup>55–58</sup> a generalization to the IEF-PCM<sup>59–61</sup> model, which we are not concerned about at the moment, is however straightforward. The coupling of a polarizable force field with the PCM requires one to take into account the mutual polarization of the atomistic and continuous part. This problem has recently been solved by Steindal and co-workers<sup>9</sup> in the framework of QM/MM calculation, with the polarization of the force field described by means of Thole's point dipole method.<sup>13</sup> The coupling between the polarizable force field and the PCM is handled by defining the proper extended matrix. We will follow here a slightly different strategy. The recently introduced variational formalism for the PCM<sup>62,63</sup> (V-PCM) recasts the PCM problem in a calculus of variations fashion: the energy of the solvated system is defined as the (unconstrained) minimum of a suitable, strictly convex functional. We report here its expression for the C-PCM (see Appendix ; all of the details can be found in ref 62):

$$\mathcal{G}(\boldsymbol{\sigma}) = \frac{1}{2f(\epsilon)} \boldsymbol{\sigma}^{\dagger} \mathbf{S} \boldsymbol{\sigma} + \boldsymbol{\sigma}^{\dagger} \mathbf{V}[\rho] \quad (11)$$

where  $\boldsymbol{\sigma}$  is the vector containing the apparent surface charges that represent the polarization of the dielectric,  $\mathbf{S}$  is the coulomb interaction matrix,  $\mathbf{V}$  is the electrostatic potential produced by the solute *in vacuo*, and  $\epsilon$  is the solvent's dielectric permittivity. The scaling factor  $f(\epsilon) = (\epsilon - 1)/\epsilon$  is used to adjust the results of the conductor-like model to dielectrics.<sup>56</sup> We point out that an equivalent functional was originally proposed by Klamt and Schuurmann,<sup>55</sup> but it was not fully exploited in its variational aspect.

As discussed in the Introduction, for molecular dynamics related calculations, we will restrict ourselves to fixed cavities surrounded by a regular surface, in particular a sphere, filled with the solute and a suitable number of solvent molecules. This means that the PCM response matrix, which depends only on the geometry of the cavity, can be computed and inverted once for ever at the beginning of the calculation without the need of performing costly linear algebra at each step. It is assumed that the surface of the cavity is partitioned in a suitable mesh and the surface elements are provided with some basis function. The traditional choice<sup>2</sup> is to use piecewise constant functions on the surface elements, which corresponds to reproducing the polarization surface density of charge by means of point charges  $\sigma_i$ . Recently, a new discretization scheme was presented,<sup>64</sup> where the PCM apparent surface charge is expanded in terms of Gaussian functions. This is very convenient when the molecular cavity follows the motion of the solute, as it provides continuous energy and gradients. As in MD simulations, we are only dealing with fixed, regular cavities; we will not adopt the aforementioned scheme both because we would not exploit its advantages and because we are treating the FQ as point charges. A generalization to the continuous surface charge (CSC) formalism of Scalmani and Frisch<sup>64</sup> is straightforward and consistent with the definition of the interaction kernel in terms of the coulomb overlap of

Gaussian orbitals. From the perspective of using the standard PCM cavity to compute molecular properties in a time independent fashion, that is by means of analytical derivatives, the CSC will be necessary to avoid numerical instabilities and discontinuities.

Exploiting the new variational formalism, it is possible to define a (free) energy functional of both the FQ and the PCM polarization charges. We point out that although the PCM free energy functional is strictly convex, this is not the case for the FQ one, as the electroneutrality constraint needs to be imposed; nevertheless, the energy of the coupled system can be written as the constrained minimum of the functional

$$\Lambda(\mathbf{q}, \boldsymbol{\sigma}, \boldsymbol{\lambda}) = \mathbf{q}^\dagger \boldsymbol{\chi} + \frac{1}{2} \mathbf{q}^\dagger \mathbf{J} \mathbf{q} + \sum_{\alpha=1}^{NM} \lambda_\alpha \sum_{j=1}^{NA} (q_j^\alpha - Q_{\text{tot}}^\alpha) + \frac{1}{2f(\varepsilon)} \boldsymbol{\sigma}^\dagger \mathbf{S} \boldsymbol{\sigma} + \boldsymbol{\sigma}^\dagger \Phi(\mathbf{q}) \quad (12)$$

where  $\Phi(\mathbf{q})$  is the electrostatic potential due to the FQ at the PCM cavity's surface elements. Let  $V$  be the potential due to the PCM charges at the atoms. We can express those potentials in a symmetric fashion as

$$\Phi_i = \sum_{j=1}^N \frac{q_j}{|\mathbf{r}_j - \mathbf{s}_i|} \stackrel{\text{def}}{=} \sum_{j=1}^N \Omega_{ij} q_j$$

$$V_i = \sum_{j=1}^{\text{NTs}} \frac{\sigma_j}{\mathbf{s}_j - \mathbf{r}_i} = \sum_{j=1}^{\text{NTs}} \Omega_{ij}^\dagger \sigma_j$$

where  $N$  is the total number of atoms and NTs is the number of PCM surface elements.

These equations define an interaction kernel between the PCM charges and the FQ  $\Omega$ . Imposing the minimum condition, one has to solve a coupled linear system

$$\begin{pmatrix} \mathbf{D} & \Omega^\dagger \\ \Omega & \mathbf{S}/f(\varepsilon) \end{pmatrix} \begin{pmatrix} \mathbf{q}_\lambda \\ \boldsymbol{\sigma} \end{pmatrix} = \begin{pmatrix} -\mathbf{C} \\ \mathbf{0} \end{pmatrix} \quad (13)$$

which is consistent with what found by Steindal and co-workers<sup>9</sup> for point dipoles.

It is interesting to spend a few lines on eq 13. As was already pointed out in the Introduction, the EEP leads to a mixed classical/quantum model where the FQs arise from quantum atomic theory but interact in a classical way. On the other hand, PCM is a fully classical model, derived from electrostatics: it describes the polarization response of a classical dielectric medium due to a classical source. From a more formal point of view, the main difference between these two models is that while the fundamental equation of the PCM is Poisson's equation, which relates an electrostatic potential with a classical source, the same does not hold for the FQ model: it is in fact not possible to identify a "source", in the Maxwellian meaning of the term. The source of the FQ is nested in the electronic structure of the system, represented in an approximate fashion by the EEP. These considerations can be extracted by the formal structure of eq 13: the right-hand side of the equation shows a source term for the FQ, the electronegativities, and no external source for the PCM part, which arises from the *classical* density of charge produced by the FQs themselves.

### 3. ANALYTICAL DERIVATIVES

The gradients of the energy functional are easily derived from eq 4. Differentiating once and using the chain rule:

$$F^{(x_i)}(\mathbf{q}, \boldsymbol{\lambda}) = \frac{dF}{dx_i} = \frac{\partial F}{\partial x_i} + \frac{\partial F}{\partial \mathbf{q}} \frac{\partial \mathbf{q}}{\partial x_i} + \frac{\partial F}{\partial \boldsymbol{\lambda}} \frac{\partial \boldsymbol{\lambda}}{\partial x_i}$$

Assuming that eq 10 has been solved, the last two terms vanish; that is, the total derivative and the partial derivative of the functional coincide. This can be seen as a classical equivalent to the Hellman–Feynman theorem for variational methods. Hence:

$$F^{(x_i)}(\mathbf{q}, \boldsymbol{\lambda}) = \frac{\partial F(\mathbf{q}, \boldsymbol{\lambda})}{\partial r_{ij}} \frac{\partial r_{ij}}{\partial x_i} = \frac{1}{2} \mathbf{q}^\dagger \frac{\partial \mathbf{J}}{\partial r_{ij}} \mathbf{q} \frac{x_i}{r_{ij}} \quad (14)$$

In a MMPol-PCM calculation, there is also a contribution arising from the interaction between the FQ and the PCM charges to be added. Following the same arguments, only the partial derivative of eq 12 needs to be calculated:

$$\Lambda^{(x_i)}(\mathbf{q}, \boldsymbol{\lambda}, \boldsymbol{\sigma}) = \left[ \frac{1}{2} \mathbf{q}^\dagger \frac{\partial \mathbf{J}}{\partial r_{ij}} \mathbf{q} + \boldsymbol{\sigma}^\dagger \frac{\partial \Omega}{\partial r_{ij}} \mathbf{q} + \frac{1}{2f(\varepsilon)} \boldsymbol{\sigma}^\dagger \frac{\partial \mathbf{S}}{\partial r_{ij}} \boldsymbol{\sigma} \right] \frac{x_i}{r_{ij}} \quad (15)$$

where, when we work with a fixed cavity, the contribution involving the derivatives of the  $\mathbf{S}$  matrix vanishes.

The second derivatives of eq 4 can be obtained differentiating once again eq 14. Using the chain rule and adopting a more compact notation:

$$F^{(xy)} = \frac{d}{dy} \frac{\partial F}{\partial x} = \frac{\partial^2 F}{\partial x \partial y} + \frac{\partial^2 F}{\partial x \partial \mathbf{q}} \frac{\partial \mathbf{q}}{\partial y} + \frac{\partial^2 F}{\partial x \partial \boldsymbol{\lambda}} \frac{\partial \boldsymbol{\lambda}}{\partial y} \quad (16)$$

the last term of eq 16 vanishes, but it is necessary to compute the derivative of the charges:

$$F^{(xy)} = \frac{1}{2} \mathbf{q}^\dagger \mathbf{J}^{(xy)} \mathbf{q} + \mathbf{q}^\dagger \mathbf{J}^{(x)} \mathbf{q}^{(y)} \quad (17)$$

Differentiating once eq 10:

$$(\mathbf{D} \mathbf{q}_\lambda)^{(y)} = \mathbf{D}^{(y)} \mathbf{q}_\lambda + \mathbf{D} \mathbf{q}_\lambda^{(y)} = 0$$

The derivatives of the charges can thus be obtained solving a *coupled-perturbed* system of equations, which we will call CPFQ in analogy to the coupled-perturbed Hartree–Fock equations:

$$\mathbf{D} \mathbf{q}_\lambda^{(y)} = -\mathbf{D}^{(y)} \mathbf{q}_\lambda \quad (18)$$

Once again, in a FQ/PCM calculation, the PCM contributions need to be taken into account. A similar set of equations can be obtained differentiating eq 15:

$$\Lambda^{(xy)} = \frac{1}{2} \mathbf{q}^\dagger \mathbf{J}^{(xy)} \mathbf{q} + \frac{1}{2f(\varepsilon)} \boldsymbol{\sigma}^\dagger \mathbf{S}^{(xy)} \boldsymbol{\sigma} + \boldsymbol{\sigma}^\dagger \Omega^{(xy)} \mathbf{q} + \mathbf{q}^\dagger \mathbf{J}^{(x)} \mathbf{q}^{(y)} + \boldsymbol{\sigma}^\dagger \Omega^{(x)} \mathbf{q}^{(y)} + \mathbf{q}^\dagger \Omega^{(x)} \boldsymbol{\sigma}^{(y)} + \frac{1}{f(\varepsilon)} \boldsymbol{\sigma}^\dagger \mathbf{S}^{(x)} \boldsymbol{\sigma}^{(y)} \quad (19)$$

$$\quad (20)$$

where the derivatives of both the FQ and the PCM charges are obtained solving a set of coupled perturbed equations obtained

differentiating eq 13:

$$\begin{pmatrix} \mathbf{D} & \mathbf{\Omega}^+ \\ \mathbf{\Omega} & \mathbf{S}/f(\varepsilon) \end{pmatrix} \begin{pmatrix} \mathbf{q}_\lambda \\ \boldsymbol{\sigma} \end{pmatrix}^{(y)} = - \begin{pmatrix} \mathbf{D} & \mathbf{\Omega}^+ \\ \mathbf{\Omega} & \mathbf{S}/f(\varepsilon) \end{pmatrix}^{(y)} \begin{pmatrix} \mathbf{q}_\lambda \\ \boldsymbol{\sigma} \end{pmatrix} \quad (21)$$

We point out that solving eq 18 or 21 and therefore calculating the derivatives of the FQs gives access to several second order properties, among which are the IR absorption intensities, computed as the derivatives of the dipole moment with respect to the normal modes. Explicit expressions of the first and second derivatives of the interaction kernel are reported in Appendix.

To perform a molecular dynamics simulation, only the energy and its first derivatives are required. Nevertheless, solving eq 10 or 13 at each propagation step would be computationally very demanding; it is however possible to exploit the extended Lagrangian method,<sup>65–68</sup> considering the FQ as independent degrees of freedom endowed with a suitable fictitious mass, like is done in Car–Parrinello<sup>67</sup> MD:

$$\mathcal{L}(\mathbf{R}, \dot{\mathbf{R}}, \mathbf{q}, \dot{\mathbf{q}}) = \frac{1}{2} \mathbf{M} \dot{\mathbf{R}}^2 + \frac{1}{2} \boldsymbol{\mu} \dot{\mathbf{q}}^2 - U(\mathbf{R}) - F(\mathbf{q}, \lambda) \quad (22)$$

The force acting on the *i*th charge in the  $\alpha$ th molecule is obtained differentiating once eq 4 with respect to the charge:

$$\mu \ddot{q}_i^\alpha = \frac{\partial F}{\partial q_{\alpha i}} = -\tilde{\chi}_j^\alpha - \lambda_\alpha \quad (23)$$

The Lagrangian multipliers can be determined imposing the conservation of the total charge for each molecule, that is:

$$\sum_{i=1}^{NA} \ddot{q}_i^\alpha = 0$$

where NA is the number of atoms in the molecule, and hence

$$\lambda_\alpha = -\frac{1}{NA} \sum_{i=1}^{NA} \tilde{\chi}_i^\alpha \quad (24)$$

Substituting eq 24 into eq 23, one gets

$$\mu \ddot{q}_i^\alpha = -\frac{1}{NA} \sum_{j=1}^{NA} (\tilde{\chi}_i^\alpha - \tilde{\chi}_j^\alpha) \quad (25)$$

Equation 25 shows that the dynamics of the charges are governed by the electronegativity equalization principle: the forces that act on the charges arise from differences in the local electronegativities and vanish when the EEP is satisfied.

If a FQ/PCM simulation is done, there is an additional term to be added to the forces on the charge, that is

$$\tilde{\chi} = \mathbf{Jq} + \boldsymbol{\chi} + \mathbf{V} = \mathbf{Jq} + \boldsymbol{\chi} + \mathbf{\Omega}\boldsymbol{\sigma} \quad (26)$$

This means that the PCM equations need to be solved to calculate the forces on the charges (and on the nuclei). As we work with a fixed cavity, this can be done easily inverting the PCM matrix separately at the first iteration: the PCM charges can be calculated when necessary, calculating the interaction potential  $\Phi = \mathbf{\Omega}^+ \mathbf{q}$  with the FQ and then multiplying it with minus the inverse of the scaled PCM matrix:

$$\boldsymbol{\sigma} = -f(\varepsilon) \mathbf{S}^{-1} \mathbf{\Omega}^+ \mathbf{q}$$

An extended Lagrangian approach for the PCM charges too has been proposed<sup>62,69</sup> and is currently under investigation.

**Table 1. Comparison between Analytical and Numerical IR Frequencies ( $\text{cm}^{-1}$ ) and Intensities ( $\text{km/mol}$ ) of NMA**

mode	frequencies		intensities	
	analytical	numerical	analytical	numerical
1	88.2116	88.2158	0.0048	0.0049
2	163.5498	163.5521	0.9680	0.9680
3	202.5621	202.5628	0.3784	0.3784
4	303.7496	303.7498	1.4539	1.4539
5	443.8597	443.8598	0.2519	0.2519
6	591.8727	591.8729	2.2844	2.2844
7	635.4398	635.4399	7.8228	7.8230
8	805.8804	805.8805	0.1881	0.1881
9	824.2376	824.2395	86.4033	86.4029
10	965.8140	965.8141	1.5627	1.5627
11	1046.8707	1046.8711	2.0386	2.0386
12	1056.1803	1056.1807	0.3581	0.3581
13	1067.9730	1067.9732	1.9047	1.9047
14	1103.3990	1103.3992	27.8512	27.8511
15	1265.7315	1265.7320	2.5165	2.5165
16	1394.6985	1394.6990	0.0000	0.0000
17	1400.8901	1400.8899	0.2803	0.2803
18	1411.7677	1411.7682	6.6289	6.6289
19	1414.1528	1414.1527	1.7540	1.7540
20	1451.4045	1451.4048	5.7040	5.7040
21	1528.0293	1528.0294	7.6152	7.6152
22	1623.0639	1623.0640	4.3400	4.3400
23	1681.4400	1681.4401	0.0750	0.0750
24	2837.1613	2837.1610	21.5101	21.5101
25	2864.5112	2864.5108	15.3084	15.3084
26	2968.5008	2968.5004	0.1301	0.1301
27	2971.2358	2971.2355	2.5114	2.5114
28	2979.4347	2979.4343	2.6738	2.6738
29	2981.1227	2981.1224	1.3473	1.3473
30	3180.8164	3180.8157	33.2933	33.2934

## 4. NUMERICAL RESULTS

All of the calculations have been performed with a locally modified development version of the Gaussian<sup>70</sup> suite of programs. As a first numerical test, we report in Table 1 the vibrational frequencies and intensities calculated with eq 17 and by numerical differentiation of the energy gradient for a N-methyl acetamide (NMA) molecule. We employed the AMBER<sup>71</sup> force field endowed with FQs; the electrostatic parameters used for NMA were taken from ref 52 and slightly adjusted to better reproduce the molecule's gas phase dipole moment (see Table 2). We point out that we are not expecting to obtain accurate spectroscopic data, but we are only testing our second derivatives implementation. The discrepancies between the results obtained analytically and numerically are always negligible and reasonably due to the truncation error in numerical differentiation, which is the result we expected. We will not go any further on the time independent approach and focus on the dynamics.

It has been pointed out in the literature<sup>15,42,52</sup> that the time evolution of the extended system defined by the Lagrangian in eq 22 behaves differently whether the molecules are kept rigid or



Table 2. Parameters Used for the NMA/Water Simulations

atoms	$\chi$	$\eta$
C	379.83	240.34
H	367.20	501.42
H	367.20	501.42
H	370.10	501.42
C	379.55	214.44
O	430.09	230.06
N	390.88	260.00
H	313.34	517.26
C	380.09	240.34
H	373.02	501.42
H	363.53	501.42
H	373.02	501.42

not. While in the first case little or no thermal coupling is observed between the FQ's dynamics and the atomic one, a strong coupling is observed when intramolecular motions are allowed to take place. To better understand the behavior of the system, we have performed several short simulations with different fictitious masses for the FQs. To avoid as much as possible numerical errors due to the time propagation, we have implemented two very accurate symplectic integrators that, while being very expensive, are known for their long-term stability and accuracy. The time evolution of a classical system can be conveniently described using the Liouville formalism: if  $A(\mathbf{q}, \mathbf{p})$  is any property of the system that only depends implicitly on time, one has

$$\frac{dA(\mathbf{q}, \mathbf{p})}{dt} = \sum_i \left( \dot{q}_i \frac{\partial A}{\partial q_i} + \dot{p}_i \frac{\partial A}{\partial p_i} \right) = \{H, A\}$$

where  $q_i$  and  $p_i$  are the coordinates and momenta of the particles,  $H$  is its Hamiltonian, and  $\{\cdot, \cdot\}$  is a Poisson bracket. If one introduces the Liouvillian operator as

$$iL A = \{H, A\}$$

it is possible to write the formal solution of the equations of motion as

$$A(t) = e^{iL t} A(0)$$

Unfortunately, the exponential map cannot be computed explicitly; on the other hand, if we can write

$$iL = iL_T + iL_V$$

that is, if the Hamiltonian is separable into a kinetic and a potential contribution, we have two maps that we know how to explicitly compute. Symplectic integrators approximate the exact time evolution with

$$\exp(i(L_T + L_V)\delta t) \approx \prod_{j=1}^k \exp(i c_j L_T \delta t) \exp(i d_j L_V \delta t) + O(\delta t^{k+1}) \quad (27)$$

The order  $k$  of the expansion determines the order of the integrator. The first order integrator is also known as the symplectic Euler method, while the second order corresponds to the Verlet method. Following the work of Yoshida,<sup>72</sup> we also implemented a fourth order and a sixth order symplectic

Table 3. Coefficients for the Symplectic Integrator

order	$c_i$	$d_i$
2	1/2	1
	1/2	0
4	$1/(2(2 - 2^{1/3}))$	$1/(2 - 2^{1/3})$
	$(1 - 2^{1/3})/(2(2 - 2^{1/3}))$	$2^{1/3}/(2 - 2^{1/3})$
	$c_2$	$d_1$
	$c_1$	0
6	0.78451361047756	0.39225680523978
	0.23557321335936	0.51004341191846
	-1.1776799841789	-0.47105338540976
	1.31518632038390	0.068753168252520
	$c_3$	$d_4$
	$c_2$	$d_3$
	$c_1$	$d_2$
	0	$d_1$

integrator, which require respectively three and eight force evaluations per step and are respectively  $O(\delta t^5)$  and  $O(\delta t^7)$  accurate. We report the coefficients  $c_i$  and  $d_i$  obtained by Yoshida in his paper in Table 3.

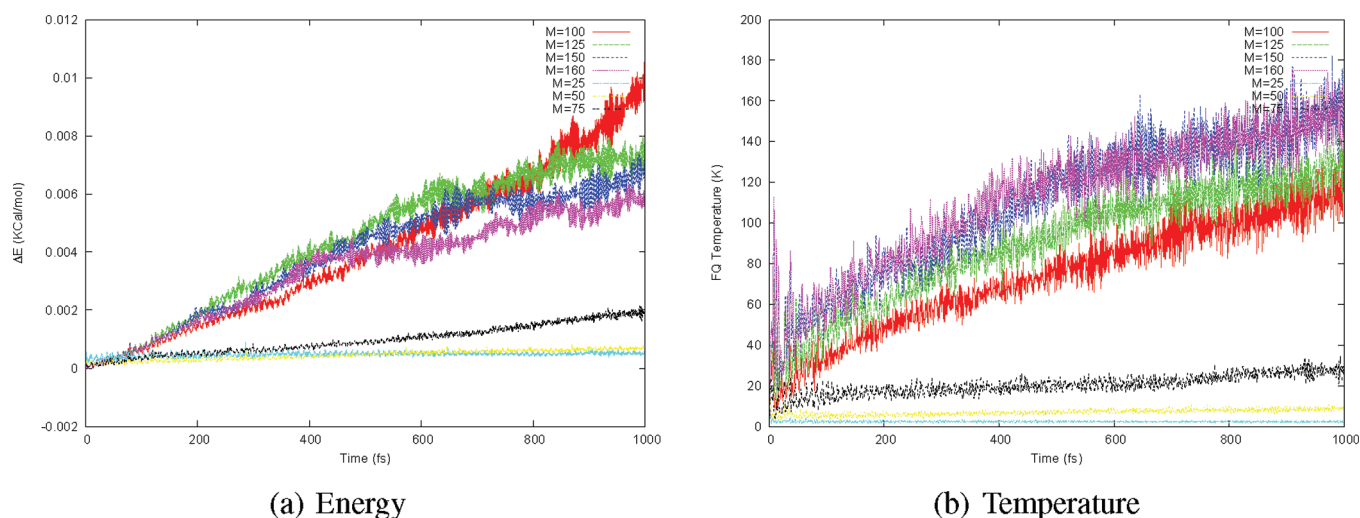
In Figure 2, we report the total energy and the temperature of the charges during the propagation for the first picosecond of a simulation. The propagation was obtained using the fourth order symplectic integrator with a 0.1 fs time step. The test system used is a cluster of 171 water molecules, described with the AMBER-(TIP3P)<sup>71,73</sup> (flexible) force field endowed with fluctuating charges. The parameters used to define the FQ part of the force field are  $\eta_O = 367.0$  kcal/mol<sup>2</sup>,  $\eta_H = 392.2$  kcal/mol<sup>2</sup>, and  $\chi_{OH} = 130.0$  kcal/mol. The first two parameters are those given in ref 15, while the electronegativity difference between oxygen and hydrogen has been adjusted to reproduce the dipole moment of liquid water. An extensive and thorough work of parametrization will be the object of a future communication.

We notice that, as the fictitious mass of the FQ is increased, a more intense thermal coupling between the nuclear and FQ dynamics occurs: while with a very small mass (below 50 atomic units) there is almost no coupling and the FQ temperature remains stable around a few degrees Kelvin, it increases rapidly, reaching values near 300–400 K when a bigger mass is used. The same simulation was also carried out without the PCM embedding; we do not report the results, as no significant difference is seen. The behavior of the energy conservation and of temperature fluctuations stabilize and remain mostly unchanged after the first picosecond of simulation.

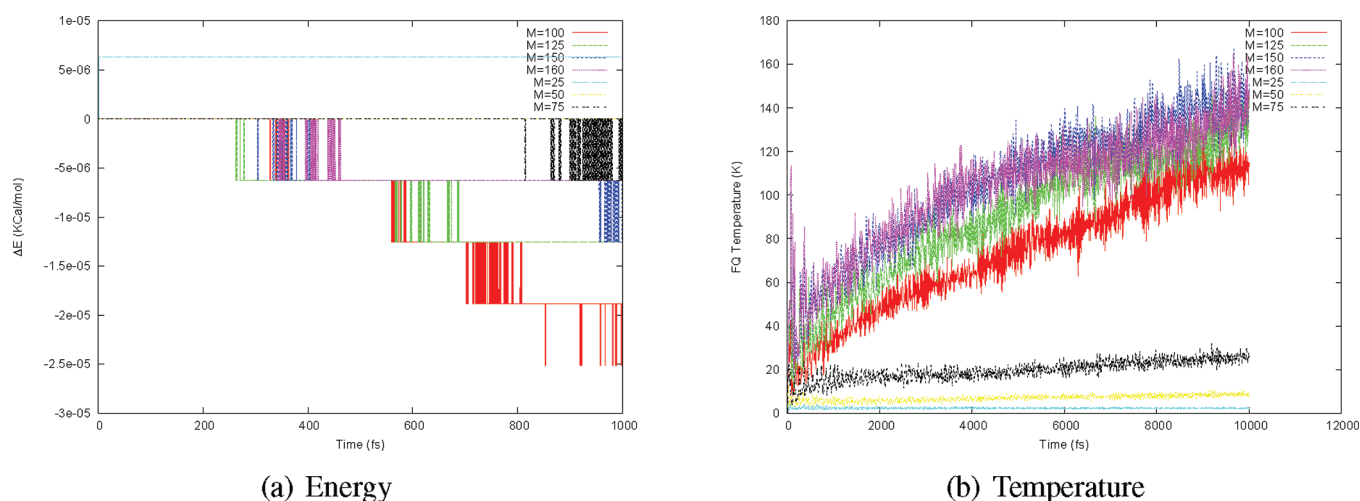
To be sure that the thermal coupling is not (at least, not only) due to numerical errors during the integration of the equations of motion, we repeated these tests also with the sixth order symplectic integrator: the results are reported in Figure 3

We notice that, in any simulation, the energy drifts are always very small (below 0.1 kcal/mol) and almost not noticeable (below  $3 \times 10^{-5}$  kcal/mol) when using the sixth order integrator; nevertheless, the thermal coupling is always strong when the FQ fictitious mass is above 50 atomic units.

These results seem to suggest that a small mass has to be used in order to avoid an excessive energy transfer from the nuclei to the charges; on the other hand, a small mass limits the propagation to a very small time step. We report in Figure 4 the energy conservation for a fictitious mass of 25 au and 150 au for different



**Figure 2.** Energy conservation and temperature of the FQ with different fictitious masses for the FQ—fourth order integrator.



**Figure 3.** Energy conservation and temperature of the FQ with different fictitious masses for the FQ—sixth order integrator.

time steps. With the smaller mass, the propagation could not be carried out with a time step longer than 0.25 fs, while it remained reasonably stable even with a 0.4–0.5 fs time step when using the heavier mass.

To overcome this difficulty,<sup>15,52</sup> the thermal coupling can be removed by thermostating separately the nuclei and the charges at two different temperatures. This allows one to propagate the system using a large fictitious mass for the charges and hence a reasonable time step. Nevertheless, the time step chosen to propagate the charges cannot be too small if molecules are not kept rigid, as the polarization energy should be considered a “fast” motion: the order of magnitude of its variation with respect to the molecular geometry is comparable to that of the bonding energy terms in a force field. This is shown in Figure 5, where the O–H stretching energy is compared with the electrostatic energy as a function of the O–H distance.

This behavior can be rationalized by considering the motion of the FQs at a fixed geometry. Following the elegant analysis of Olano and Rick,<sup>46</sup> it is possible to calculate a frequency of oscillation for some suitably defined linear combination of the charges, considered dynamical variables. Following Olano and

Rick, let us consider an isolated molecule of water at its equilibrium geometry ( $C_{2v}$  symmetry). Imposing explicitly the constraint on the total charge, i.e.

$$q_O = -2q_H$$

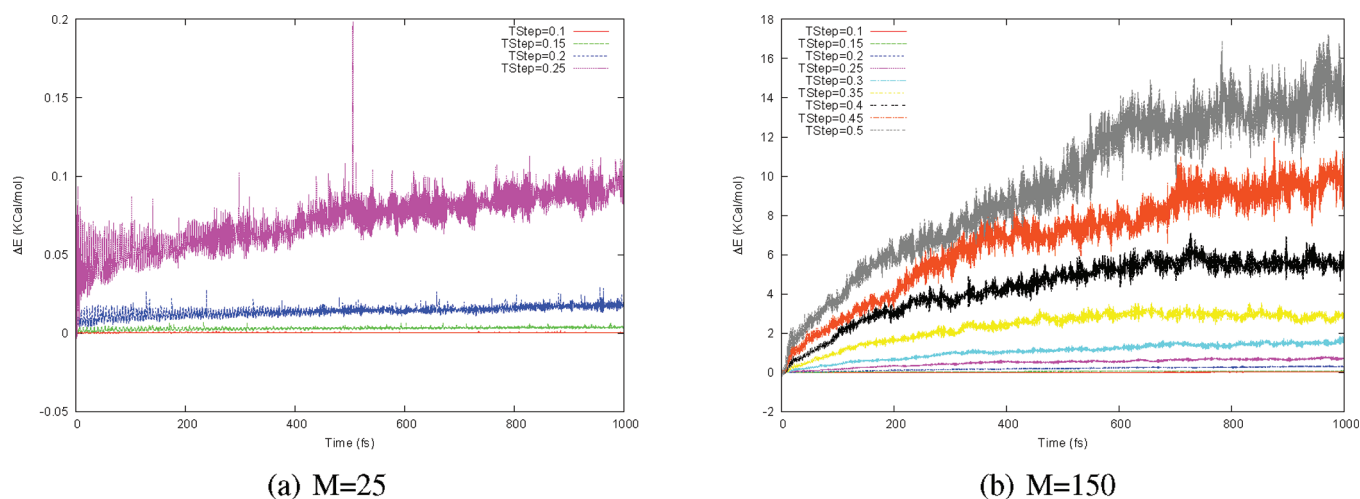
where the hydrogens carry the same charge for symmetry reasons, the electrostatic energy reads:

$$U(q_H) = \frac{1}{2} \begin{pmatrix} q_{H1} \\ q_{H2} \end{pmatrix}^\dagger \begin{pmatrix} \alpha & \beta \\ \beta & \alpha \end{pmatrix} \begin{pmatrix} q_{H1} \\ q_{H2} \end{pmatrix} + \Delta\chi + q \quad (28)$$

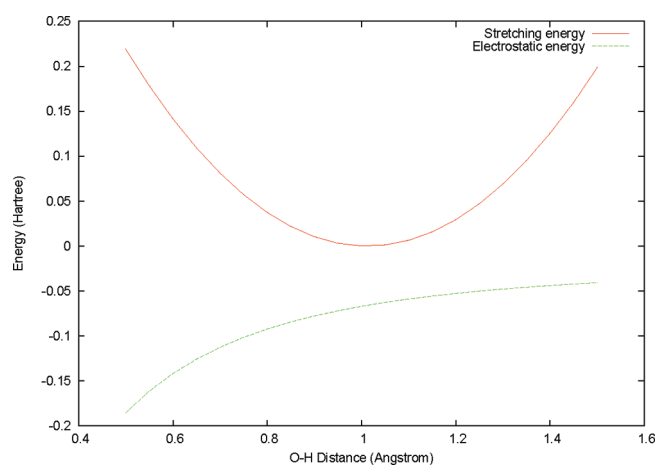
where

$$\Delta\chi_i = \chi_i - \chi_O, \alpha = J_{HH}(0) + J_{OO}(0) - 2J_{OH}(r_{OH}), \\ \beta = J_{HH}(r_{HH}) + J_{OO}(0) - 2J_{OH}(r_{OH})$$

The eigenvalues of this matrix correspond, when divided by the charge fictitious mass, to the oscillation frequency of the charge “normal modes”. With our values, the eigenvalues are  $\lambda_1 = 0.4298$  au, corresponding to a “symmetric” oscillation, and 0.31378 au, corresponding to an asymmetric one. In Figure 6, we report the

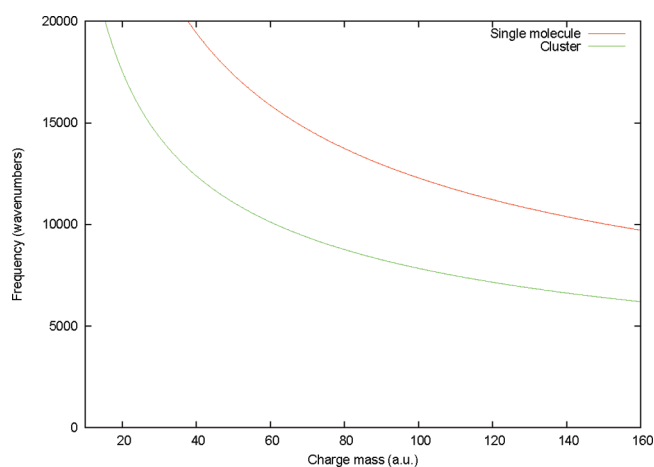


**Figure 4.** Energy conservation with different FQ fictitious masses and time steps.



**Figure 5.** Stretching and electrostatic energy with respect to O–H distance.

frequency as a function of the fictitious mass for the lowest eigenvalue. We see that, while no value of the fictitious mass makes the frequency close to a typical vibrational frequency (which could cause a resonance and, hence, a numerical catastrophe), with large values of the fictitious mass, the separation becomes smaller. On the other hand, small masses correspond to very high characteristic frequencies, hence the necessity of using a very small time step. If we consider a system of interacting molecules, the analysis becomes trickier. Without performing a detailed and accurate treatment of the problem in terms of normal modes, it is possible to estimate the frequency of the lowest energy vibration by the lowest eigenvalue of the  $J$  matrix. We report in the same Figure 6 the results for the cluster containing 171 molecules that we used for the tests reported at the beginning of this section. Although the numbers are not quantitatively accurate, they certainly provide a qualitatively correct picture of the dynamics of the charges. When the biggest mass is used, the oscillation frequency of the lowest energy normal mode is relatively close to an O–H stretching: this explains the strong thermal coupling between the polarization and atomic degrees of freedom.



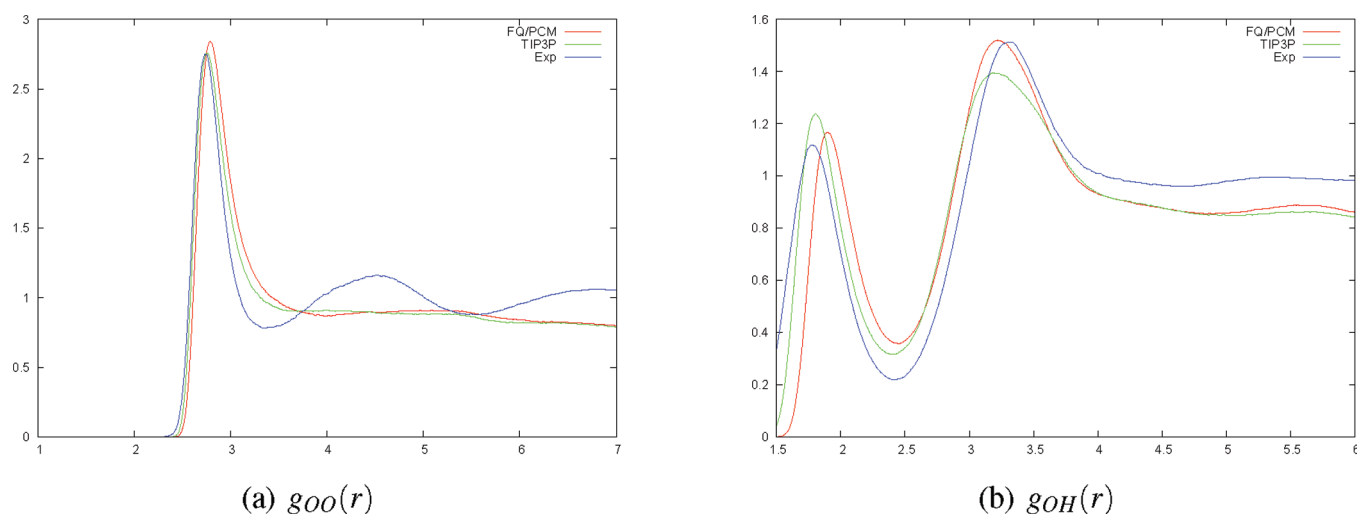
**Figure 6.** Oscillation frequency for the lowest energy charge normal mode as a function of the fictitious mass for a single molecule of water and for a cluster of water molecules.

To test our implementation, we tried to reproduce the first peak of the pair correlation function for water. A 100 ps simulation was performed at 298 K and unitary density, using the velocity Verlet integrator, with a 0.25 fs time step. The system was composed of 457 molecules of water in a spherical box of radius  $r_c = 14.5 \text{ \AA}$ . The starting configuration was obtained with a standard  $(N,V,T)$  simulation at 298 K.

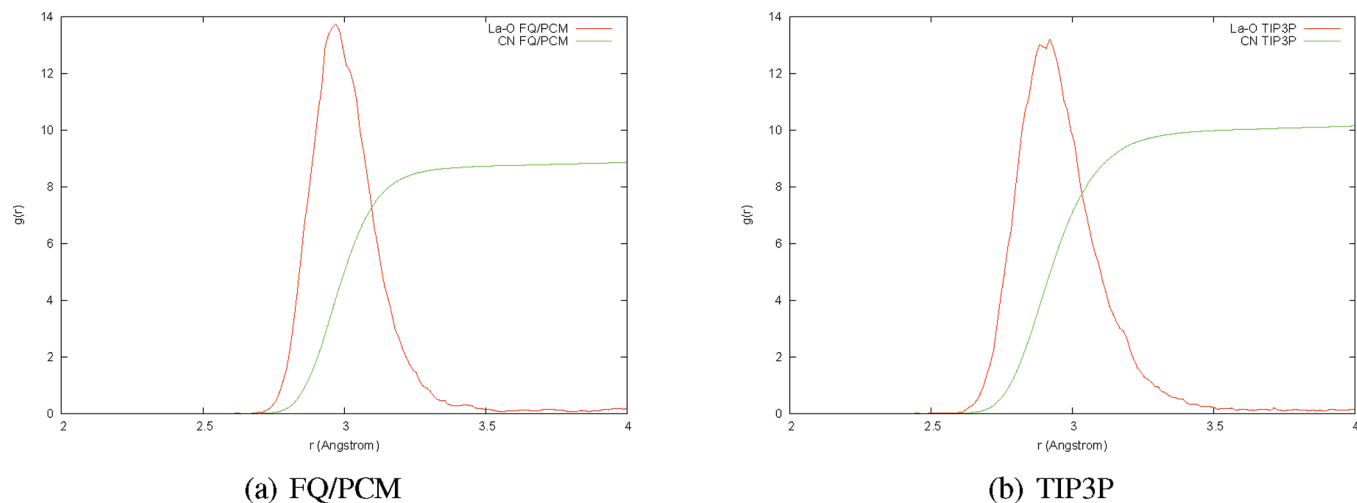
When enforcing nPBCs, spurious boundary effects can occur, which should be cured employing either effective potentials or proper buffer regions.<sup>12,26,32–35</sup> In the present study, we have employed a simple radial potential  $V_c$ , which is sufficient for the description of local effects far from the boundaries. More reliable effective potentials will be implemented after the optimization of polarizable flexible solvent force fields. In particular, we adopt for the  $V_c$  potential a sixth-degree polynomial expression:

$$V_c(r) = \begin{cases} 0 & r \leq r_c \\ k(r - r_c)^6 & r > r_c \end{cases}$$

The temperature was kept constant using a stochastic Andersen



**Figure 7.** O–O and O–H (intermolecular) radial distribution function for water.



**Figure 8.** First peak of the La–O radial distribution function and its integral (coordination number).

thermostat<sup>65</sup> for the nuclei and rescaling the velocities for the charges so that their temperature remained fixed to 4 K. We used the AMBER/TIP3P<sup>71,73</sup> force field for the nonelectrostatic part of the potential energy. The system was surrounded by a PCM spherical cavity of radius 15.5 Å. We used a larger radius for the PCM cavity to avoid contact between the FQ and the PCM apparent surface charge and, hence, the potential polarization catastrophe. The simple potential adopted was able to enforce the confinement, allowing a very small penetration of the  $r > r_c$  region. The size of the PCM cavity with respect to the confinement radius will be the object of further study.

In Figure 7, we report the radial distribution function we obtained with a FQ/PCM simulation and a nonpolarizable TIP3P simulation in comparison to the experimental one of Soper and Phillips.<sup>74</sup> A simulation carried out with 1116 molecules of water confined in a 20.5 Å sphere surrounded by a PCM spherical cavity of radius 21.5 Å produced comparable results. The agreement is only qualitatively acceptable; on the other hand, since a thorough parametrization was not carried out, as regards either the force field or the confinement potential, we did

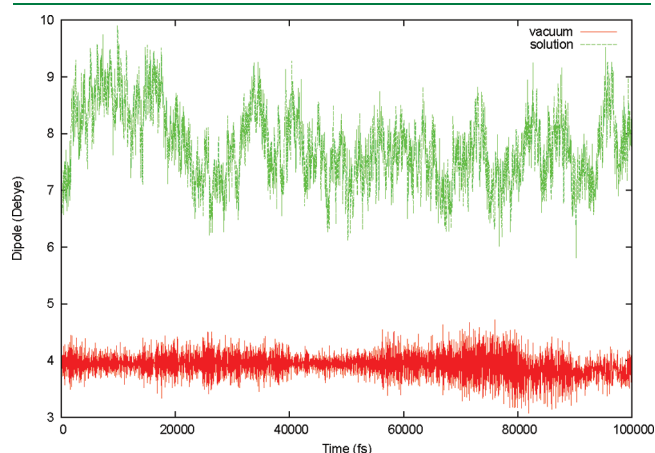
not expect to perfectly reproduce the properties of the bulk liquid.

On the other hand, we felt that the agreement was satisfactory enough to try to reproduce the solvation properties of a cation like lanthanum. It was shown by Duvail and co-workers<sup>75,76</sup> that a polarizable force field is necessary to correctly describe the coordination of lanthanides in water. We hence tried to reproduce their results, adopting the same Buckingham-6 potential they propose to describe the nonelectrostatic interaction of the cation with the solvent and our FQ force field. A 100 ps simulation was carried out on a 12.5 Å spherical box, obtained cutting the cluster used for the water simulation, and creating a hollow, 3 Å in radius cavity in the middle that accommodates a lanthanum ion, surrounded by 241 water molecules. A 0.25 fs time step and the velocity Verlet integrator were employed. The system was thermostatted at 298 K (4 K for the FQ as before), and data were collected after a 2 ps of equilibration, as suggested by Duvail and co-workers.<sup>75</sup> Confinement was kept as before with a sixth degree polynomial potential, and the system was surrounded with a 13.5 Å spherical PCM cavity. The La–O radial

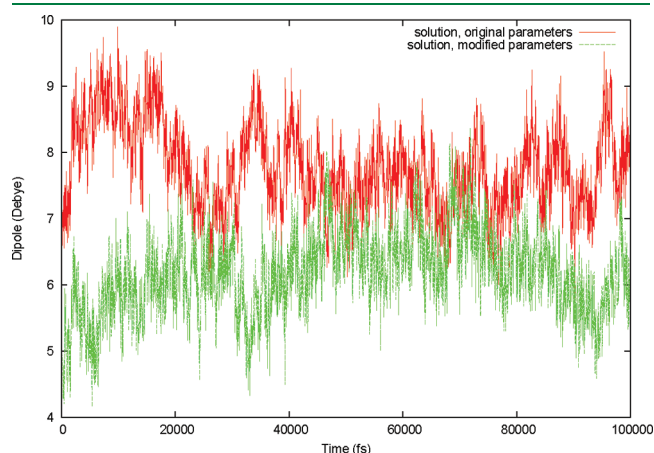
distribution function was calculated and the coordination number determined as the area of the first peak. We report in Figure 8 the first peak of the  $g_{\text{La-O}}(r)$  function and its integral obtained with the FQ polarizable force field and PCM embedding in comparison to the ones obtained with the TIP3P force field. The FQ/PCM calculation gives a coordination number for the  $\text{La}^{3+}$  ion of slightly less than 9 (8.78), which is in good agreement with both the result obtained by Duvaal et al. and the experimental data they report. On the other hand, the unpolarizable force field gives a coordination number slightly bigger than 10 (10.1), again in agreement with what Duvaal and co-worker observed.

As a last pilot application of the method implemented, we studied the solvation of a N-methyl acetamide (NMA) molecule in water. The starting configuration was obtained as for the lanthanum ion, cutting a 12.5 Å sphere of water molecules from the cluster used for the pure water simulation and creating a 3 Å cavity to accommodate the NMA molecule. The system was equilibrated for 2 ps, and a 100 ps simulation was run with a 0.25 fs time step, using the velocity Verlet integrator. We employed the AMBER<sup>71,73</sup> force field (TIP3P for water) endowed with fluctuating charges. The FQ parameters were taken from ref 52; we slightly adjusted the electronegativities to better reproduce the NMA's gas phase dipole moment (see Table 2). Confinement

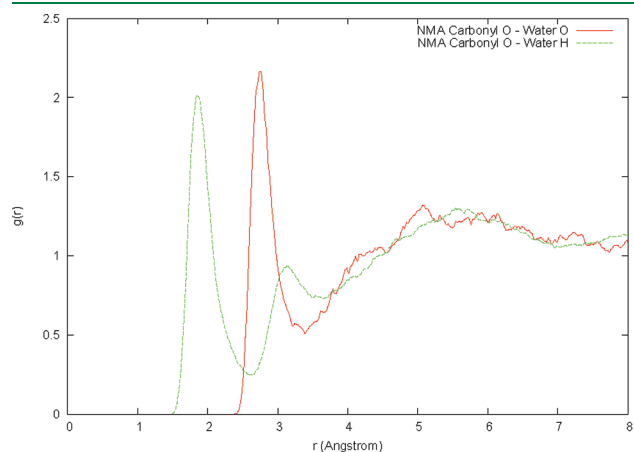
was maintained with the same sextic potential, and a 13.5 Å PCM cavity surrounded the system. We report in Figure 9 the dipole moment of the NMA molecule along the simulation and in Figure 10 the radial distribution function between the NMA carbonyl oxygen and water and between the NMA amide hydrogen and water, respectively. The results obtained are consistent with those of Patel and co-workers,<sup>52</sup> which was to be expected, as we used a only slightly modified version of their parameters. The radial distribution functions show the presence of a strong hydrogen bond between NMA and water. As Patel and co-workers have already pointed out, the dipole moment of the NMA molecule in water is higher (the average is slightly lower than 8 D) than the one obtained by *ab initio* computation, with implicit or explicit solvation (5.18 to 5.33 D). In our case, this overpolarization effect can be attributed to two different facts. As we have already pointed out, this work was not concerned with the parametrization of a FQ polarizable force field: we used parameters taken from the literature, eventually making some adjustment, but without a thorough effort to produce a general and reliable set of parameters. On the other hand, we think that one of the assumptions made is probably a considerable source of error. The definition of the off-diagonal elements of the hardness matrix in eq 8 is an approximation, and



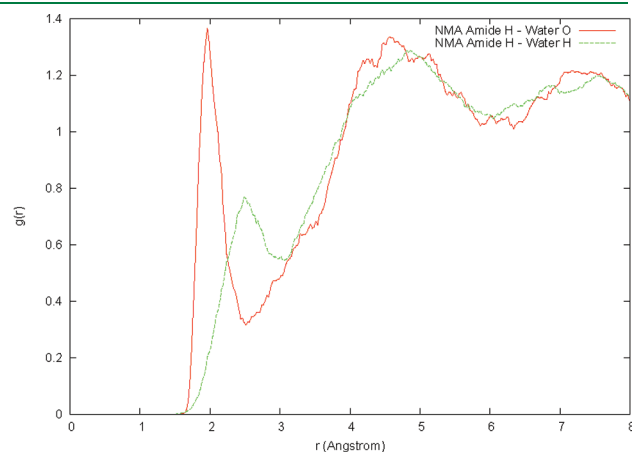
**Figure 9.** Instantaneous dipole moment (Debye) of the NMA molecule in solution during the simulation.



**Figure 11.** Instantaneous dipole moment (Debye) of the NMA molecule in solution during the simulation—original vs modified parameters.

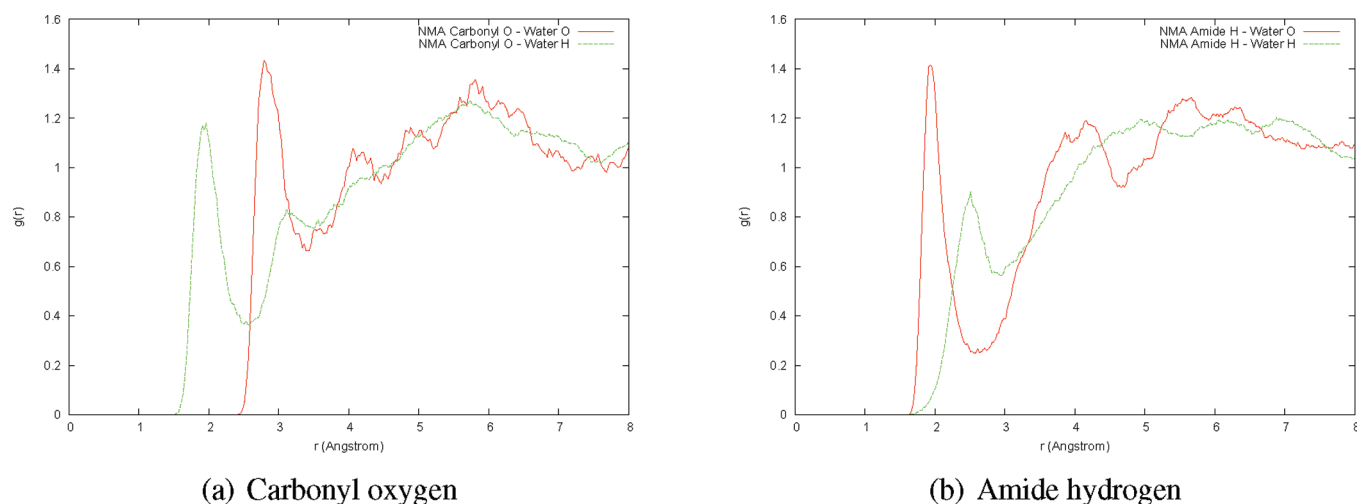


(a) Carbonyl oxygen



(b) Amide hydrogen

**Figure 10.** Radial distribution functions for NMA's hydrogen bond donors/acceptors and water's atoms.



**Figure 12.** Radial distribution functions for NMA's hydrogen bond donors/acceptors and water's atoms with modified parameters.

we feel it to be a little too simplistic, especially when strong, specific interactions are involved. As a matter of fact, the whole model does not consider covalent or “quantum” interactions at all: when dealing with hydrogen bonds, it is probably necessary to add some degrees of freedom in the parametrization. To show this, we run a simulation changing the off diagonal hardness elements for the atoms involved in hydrogen bond formation, i.e., the amide hydrogen with the water oxygen and the carbonyl oxygen with the water hydrogen. We lowered the values obtained as the arithmetical mean of the atomic ones from  $\eta_{\text{H}_{\text{NMA}}\text{O}_{\text{wat}}} = 311.14$  kcal/mol to  $\eta_{\text{H}_{\text{NMA}}\text{O}_{\text{wat}}} = 300.00$  kcal/mol and from  $\eta_{\text{H}_{\text{NMA}}\text{O}_{\text{wat}}} = 442.20$  kcal/mol to  $\eta_{\text{H}_{\text{NMA}}\text{O}_{\text{wat}}} = 420.00$  kcal/mol. We reiterate that this is only a conceptual experiment and that an accurate determination of the optimal parameters has not yet been done. The results are reported in Figures 11 and 12. We see that with the modified parameters, the instantaneous dipole moment oscillates around a value (slightly larger than 6.0 D) which is much closer to the one obtained by *ab initio* calculations; on the other hand, the simulation provides a much weaker formation of hydrogen bonds, as shown by the heights of the peaks in Figure 12. A small extension of the manifold of parameters hence seems definitely worth the effort.

## 5. CONCLUSIONS AND PERSPECTIVES

We have implemented a polarizable force field using the fluctuating charges model, and we have combined it with the polarizable continuum model as a tool to perform molecular dynamics simulations with nonperiodic boundary conditions in NVE/NVT ensembles. Extensive numerical testing has been performed to better understand the behavior of the MD simulations when the charges are propagated in a Car–Parrinello-like fashion, and a rationalization has been proposed on the basis of the analysis of the frequencies of oscillation that characterize the charge dynamics. Several prototypical applications have been shown: calculation of IR frequencies and intensities and simulation of pure water and of charged and uncharged atomic and molecular solutes have been performed and discussed. We notice that the lack of an accurate parametrization of the whole force field is probably the first problem that needs to be addressed. From the NMA simulations, we found that the manifold of parameters that need to be taken into

account for the electrostatics can have a crucial effect on the quality of the results: its extension with parameters chosen to describe specific intermolecular interactions might introduce major improvements. As pointed out by Verstraelen and co-workers<sup>41</sup> in their excellent work, parametrization is a complex matter, and the literature presents sets of parameters widely varying in one respect to another, all because the cost functions used for calibration present flat and elongated minima. Another important extension of the model that we have not considered yet is charge transfer.<sup>15,40,43–46</sup> It has been shown that this phenomenon can be of great importance in determining the properties of liquid water, in which, as we aim to model aqueous solutions, we are very interested. Concerning more specifically molecular dynamics, other ensembles, especially the NPT one,<sup>77</sup> need to be considered, and a more exhaustive study on the effect of different thermostats is under investigation. In particular, the use of the Andersen barostat combined with Nosé–Hoover–Poincaré thermostats may be convenient, as it allows the propagation of the trajectory with an arbitrary order symplectic integrator. The perspectives of a FQ/PCM implementation are many. The FQ model provides the electrostatic interactions to the ReaxFF reactive force field:<sup>54</sup> reactivity in solution by means of classical MD simulations is an interesting development. On the other hand, an interface with the QM world seems the most promising target to pursue, especially considering that our implementation is rooted in a very QM-oriented computational package. From this point of view, the FQ polarizable force field is the missing term in the GLOB<sup>22–26</sup> model, where the core and the continuum are polarizable but the atomistic layer of the solvent is not. A fully polarizable QM/MM/PCM layered method, suitable both for accurately reproducing the solvent effect on molecular properties and for mixed classical/*ab initio* molecular dynamics, is, after the aforementioned developments are complete, an important target we wish to pursue.

## APPENDIX A. THE C-PCM AS A VARIATIONAL PROBLEM

The conductor-like PCM describes the solvent as a conducting medium which occupies all the space but a hollow cavity that accommodates the solute. Let  $C$  be such cavity, which we will assume to be a bounded, simply connected open subset of  $\mathbb{R}^3$

with regular enough boundary  $\Gamma = \partial C$ , and let  $\rho$  be the molecular density of charge, which we will assume to be supported inside  $C$ . As we are assuming that the cavity is surrounded by a conductor, the potential at the boundary must vanish, and hence we need to solve the Poisson equation

$$\nabla^2 \phi = -4\pi\rho, \phi = 0 \quad \forall \mathbf{r} \notin C \quad (29)$$

Linearity allows us to write the potential as a sum of a molecular term  $\Phi$  and a reaction term  $V_r$ , due to the polarization of the surroundings:

$$\phi = \Phi + V_r$$

The molecular term is the potential produced by the density  $\rho$  in *vacuo*, while the reaction term can be modeled as the potential produced by an apparent surface charge  $\sigma$ , which represents the polarization of the medium:

$$\Phi(\mathbf{r}) = \int_C \frac{\rho(\mathbf{r}')}{|\mathbf{r} - \mathbf{r}'|} d\mathbf{r}', V_r(\mathbf{r}) = \int_\Gamma \frac{\sigma(\mathbf{s})}{|\mathbf{r} - \mathbf{s}|} ds$$

As the total potential must vanish at the boundary, the C-PCM integral equation will read:

$$V_r(\mathbf{s}) = \int_\Gamma \frac{\sigma(\mathbf{s}')}{|\mathbf{s} - \mathbf{s}'|} ds' \stackrel{\text{def}}{=} (\mathcal{J}\sigma)(\mathbf{s}) = -\Phi(\mathbf{s}) \quad (30)$$

where we have introduced the integral operator  $\mathcal{J}$ . It is known<sup>31</sup> that if  $\Gamma$  is regular enough,  $\mathcal{J}$  is a self-adjoint, positive definite, coercive operator in a suitable Hilbert space  $V$ : this means that the (unique) minimum of the functional

$$J(\sigma) = \frac{1}{2} \langle \sigma, \mathcal{J}\sigma \rangle_V + \langle \sigma, \Phi|_\Gamma \rangle_V \quad (31)$$

(we denote with  $\langle \cdot, \cdot \rangle_V$  the scalar product in  $V$ ) is also the solution of the integral equation.<sup>78</sup> Adding the C-PCM scaling factor  $1/(f(\epsilon))$  and discretizing, eq 11 is recovered.

## APPENDIX B. DERIVATIVES OF THE FQ CONTRIBUTION TO THE ENERGY: EXPLICIT CONTRIBUTIONS

During the calculation of the gradients and of the Hessian of the FQ contribution to the energy, both first and second partial derivatives of the interaction kernel eq 7 are to be computed. We report here their expressions.

$$\frac{\partial J_{ij}}{\partial r_k^\mu} = (1 - \delta_{ij}) \frac{\partial J_{ij}}{\partial r_{ij}} \frac{r_i^\mu - r_j^\mu}{r_{ij}} (\delta_{ik} - \delta_{jk}) \quad (32)$$

where  $r_k^\mu$  is the  $\mu$ th Cartesian coordinate of the position vector of the  $k$ th particle and

$$r_{ij}^2 = \sum_{\mu=1}^3 (r_i^\mu - r_j^\mu)^2 = \sum_{\mu=1}^3 (r_{ij}^\mu)^2$$

We will need to calculate

$$\frac{\partial J_{ij}}{\partial r_{ij}} = \frac{\partial}{\partial r_{ij}} \frac{\eta_{ij}}{(1 + r_{ij}^n \eta_{ij}^n)^{1/n}} = - \frac{\eta_{ij}^{n+1} r_{ij}^{n-1}}{(1 + \eta_{ij}^n r_{ij}^n)^{1/n+1}} \quad (33)$$

In particular, if we choose  $n = 2$ , eq 32 becomes

$$\frac{\partial J_{ij}}{\partial r_i^\mu} = - \frac{\eta_{ij}^3 r_{ij}^\mu}{(1 + \eta_{ij}^2 r_{ij}^2)^{3/2}} \quad (34)$$

The partial second derivative has a rather cumbersome expression:

$$\frac{\partial^2 J_{ij}}{\partial r_k^\mu \partial r_l^\nu} = \frac{\partial^2 J_{ij}}{\partial r_{ij}^2} \frac{\partial r_{ij}}{\partial r_k^\mu} \frac{\partial r_{ij}}{\partial r_l^\nu} + \frac{\partial J_{ij}}{\partial r_{ij}} \frac{\partial^2 r_{ij}}{\partial r_k^\mu \partial r_l^\nu} \quad (35)$$

The second derivative of the interaction kernel with respect to the interatomic distance is

$$\frac{\partial^2 J_{ij}}{\partial r_i^\mu \partial r_j^\nu} = \frac{\eta_{ij}^{n+1} r_{ij}^{n-2} (2r_{ij}^n \eta_{ij}^n - n + 1)}{(1 + \eta_{ij}^n r_{ij}^n)^{1/n+2}}$$

With our  $n = 2$  choice, eq 35 becomes

$$\frac{\partial^2 J_{ij}}{\partial r_k^\mu \partial r_l^\nu} = (1 - \delta_{ij}) \left\{ \frac{3\eta_{ij}^5 r_{ij}^\mu r_{ij}^\nu}{(1 + \eta_{ij}^2 r_{ij}^2)^{5/2}} - \delta_{\mu\nu} \frac{\eta_{ij}^3}{(1 + \eta_{ij}^2 r_{ij}^2)^{3/2}} \right\} [\delta_{kl}(\delta_{ik} + \delta_{jk}) - (1 - \delta_{kl})(\delta_{ik}\delta_{jl} + \delta_{il}\delta_{jk})] \quad (36)$$

## AUTHOR INFORMATION

### Corresponding Author

\*E-mail: flipparini@sns.it.

## ACKNOWLEDGMENT

The authors gratefully acknowledge financial support from Gaussian, Inc.

## REFERENCES

- (1) Tomasi, J.; Persico, M. *Chem. Rev.* **1994**, *94*, 2027–2094.
- (2) Tomasi, J.; Mennucci, B.; Cammi, R. *Chem. Rev.* **2005**, *105*, 2999–3093.
- (3) Mennucci, B. *J. Phys. Chem. Lett.* **2010**, *1*, 1666–1674 and references therein.
- (4) Vreven, T.; Mennucci, B.; da Silva, C.; Morokuma, K.; Tomasi, J. *J. Chem. Phys.* **2001**, *115*, 62–72.
- (5) Rega, N.; Cossi, M.; Barone, V. *J. Am. Chem. Soc.* **1998**, *120*, 5723–5732.
- (6) Cui, Q. *J. Chem. Phys.* **2002**, *117*, 4720–4728.
- (7) Pedone, A.; Biczysko, M.; Barone, V. *ChemPhysChem* **2010**, *11*, 1812–1832.
- (8) Barone, V.; Bloino, J.; Monti, S.; Pedone, A.; Prampolini, G. *Phys. Chem. Chem. Phys.* **2010**, *12*, 10550–10561.
- (9) Steindal, A. H.; Ruud, K.; Frediani, L.; Aidas, K.; Kongsted, J. *J. Phys. Chem. B* **2011**, *115*, 3027–3037.
- (10) Rega, N.; Cossi, M.; Barone, V. *J. Am. Chem. Soc.* **1997**, *119*, 12962–12967.
- (11) Benzi, C.; Improta, R.; Scalmani, G.; Barone, V. *J. Comput. Chem.* **2002**, *23*, 341–350.
- (12) Barone, V.; Biczysko, M.; Brancato, G. Extending the Range of Computational Spectroscopy by QM/MM Approaches: Time-Dependent and Time-Independent Routes. In *Combining Quantum Mechanics and Molecular Mechanics. Some Recent Progresses in QM/MM Methods*; Sabin, J. R., Canuto, S., Eds.; Academic Press: Waltham, MA, 2010; Vol. 59, pp 17–57.
- (13) Thole, B. *Chem. Phys.* **1981**, *59*, 341–350.
- (14) Lamoureux, G.; Roux, B. *J. Chem. Phys.* **2003**, *119*, 3025–3039.
- (15) Rick, S. W.; Stuart, S. J.; Berne, B. J. *J. Chem. Phys.* **1994**, *101*, 6141–6156.
- (16) Rick, S. W.; Berne, B. J. *J. Am. Chem. Soc.* **1996**, *118*, 672–679.
- (17) Rick, S. W.; Stuart, S. J.; Bader, J. S.; Berne, B. J. *J. Mol. Liq.* **1995**, *65–66*, 31–40.

- (18) Mortier, W. J.; Van Genechten, K.; Gasteiger, J. *J. Am. Chem. Soc.* **1985**, *107*, 829–835.
- (19) Chelli, R.; Procacci, P. *J. Chem. Phys.* **2002**, *117*, 9175–9189.
- (20) Elstner, M.; Porezag, D.; Jungnickel, G.; Elsner, J.; Haugk, M.; Frauenheim, T.; Suhai, S.; Seifert, G. *Phys. Rev. B* **1998**, *58*, 7260–7268.
- (21) Trani, F.; Barone, V. *J. Chem. Theory Comput.* **2011**, *7*, 713–719.
- (22) Brancato, G.; Rega, N.; Barone, V. *J. Chem. Phys.* **2006**, *124*, 214505.
- (23) Brancato, G.; Nola, A. D.; Barone, V.; Amadei, A. *J. Chem. Phys.* **2005**, *122*, 154109.
- (24) Rega, N.; Brancato, G.; Barone, V. *Chem. Phys. Lett.* **2006**, *422*, 367–371.
- (25) Brancato, G.; Barone, V.; Rega, N. *Theor. Chem. Acc.* **2007**, *117*, 1001–1015.
- (26) Brancato, G.; Rega, N.; Barone, V. *J. Chem. Phys.* **2008**, *128*, 144501.
- (27) Rega, N.; Brancato, G.; Petrone, A.; Caruso, P.; Barone, V. *J. Chem. Phys.* **2011**, *134*, 074504.
- (28) Smith, P. E.; Pettitt, B. M. *J. Chem. Phys.* **1996**, *105*, 4289–4293.
- (29) Hünenberger, P. H.; McCammon, J. A. *J. Chem. Phys.* **1999**, *110*, 1856–1872.
- (30) Bergdorf, M.; Peter, C.; Hünenberger, P. H. *J. Chem. Phys.* **2003**, *119*, 9129–9144.
- (31) Cancès, E. In *Continuum Solvation Models in Chemical Physics*; Mennucci, B., Cammi, R., Eds.; Wiley: New York, 2007.
- (32) Warshel, A.; King, G. *Chem. Phys. Lett.* **1985**, *121*, 124–129.
- (33) King, G.; Warshel, A. *J. Chem. Phys.* **1989**, *91*, 3647–3661.
- (34) Luzhkov, V.; Warshel, A. *J. Comput. Chem.* **1992**, *13*, 199–213.
- (35) Sham, Y. Y.; Warshel, A. *J. Chem. Phys.* **1998**, *109*, 7940–7944.
- (36) Rappe, A.; Goddard, W. J. *Phys. Chem.* **1991**, *95*, 3358–3363.
- (37) Chen, J.; Martinez, T. *J. Chem. Phys. Lett.* **2008**, *463*, 288.
- (38) Chelli, R.; Pagliai, M.; Procacci, P.; Cardini, G.; Schettino, V. *J. Chem. Phys.* **2005**, *122*, 074504.
- (39) Nistor, R. A.; Polihronov, J. G.; Müser, M. H.; Mosey, N. J. *J. Chem. Phys.* **2006**, *125*, 094108.
- (40) Lee, A. J.; Rick, S. W. *J. Chem. Phys.* **2011**, *134*, 184507.
- (41) Verstraelen, T.; Bultinck, P.; Van Speybroeck, V.; Ayers, P. W.; Van Neck, D.; Waroquier, M. *J. Chem. Theory Comput.* **2011**, *7*, 1750–1764.
- (42) Ando, K. *J. Chem. Phys.* **2001**, *115*, 5228–5237.
- (43) Chelli, R.; Ciabatti, S.; Cardini, G.; Righini, R.; Procacci, P. *J. Chem. Phys.* **1999**, *111*, 4218–4229.
- (44) Llanta, E.; Ando, K.; Rey, R. *J. Phys. Chem. B* **2001**, *105*, 7783–7791.
- (45) Llanta, E.; Rey, R. *Chem. Phys. Lett.* **2001**, *340*, 173–178.
- (46) Olano, L. R.; Rick, S. W. *J. Comput. Chem.* **2005**, *26*, 699–707.
- (47) York, D. M.; Yang, W. *J. Chem. Phys.* **1996**, *104*, 159–172.
- (48) Verstraelen, T.; Speybroeck, V. V.; Waroquier, M. *J. Chem. Phys.* **2009**, *131*, 044127.
- (49) Nishimoto, K.; Mataga, N. *Z. Phys. Chem. (Frankfurt)* **1957**, *12*, 335.
- (50) Ohno, K. *Theor. Chem. Acc.* **1964**, *2*, 219–227.
- (51) Louwen, J. N.; Vogt, E. T. C. *J. Mol. Catal. A* **1998**, *134*, 63–77.
- (52) Patel, S.; Brooks, C. *J. Comput. Chem.* **2004**, *25*, 1–15.
- (53) Patel, S.; Mackerell, A.; Brooks, C. *J. Comput. Chem.* **2004**, *25*, 1504–1514.
- (54) van Duin, A. C. T.; Dasgupta, S.; Lorant, F.; Goddard, W. A. *J. Phys. Chem. A* **2001**, *105*, 9396–9409.
- (55) Klamt, A.; Schuurmann, G. *J. Chem. Soc., Perkin Trans. 2* **1993**, 799–805.
- (56) Barone, V.; Cossi, M. *J. Phys. Chem. A* **1998**, *102*, 1995–2001.
- (57) Cossi, M.; Rega, N.; Scalmani, G.; Barone, V. *J. Comput. Chem.* **2003**, *24*, 669–681.
- (58) Scalmani, G.; Barone, V.; Kudin, K.; Pomelli, C.; Scuseria, G.; Frisch, M. *Theor. Chem. Acc.* **2004**, *111*, 90–100.
- (59) Cancès, E.; Mennucci, B.; Tomasi, J. *J. Chem. Phys.* **1997**, *107*, 3032–3041.
- (60) Cancès, E.; Mennucci, B. *J. Math. Chem.* **1998**, *23*, 309–326.
- (61) Mennucci, B.; Cancès, E.; Tomasi, J. *J. Phys. Chem. B* **1997**, *101*, 10506–10517.
- (62) Lipparini, F.; Scalmani, G.; Mennucci, B.; Cancès, E.; Caricato, M.; Frisch, M. *J. Chem. Phys.* **2010**, *133*, 014106.
- (63) Lipparini, F.; Scalmani, G.; Mennucci, B.; Frisch, M. *J. Chem. Theory Comput.* **2011**, *7*, 610–617.
- (64) Scalmani, G.; Frisch, M. *J. Chem. Phys.* **2010**, *132*, 114110.
- (65) Andersen, H. C. *J. Chem. Phys.* **1980**, *72*, 2384–2393.
- (66) Parrinello, M.; Rahman, A. *Phys. Rev. Lett.* **1980**, *45*, 1196–1199.
- (67) Car, R.; Parrinello, M. *Phys. Rev. Lett.* **1985**, *55*, 2471–2474.
- (68) Crescenzi, O.; Pavone, M.; De Angelis, F.; Barone, V. *J. Phys. Chem. B* **2005**, *109*, 445–453, PMID: 16851035.
- (69) Caricato, M.; Scalmani, G.; Frisch, M. J. In *Continuum Solvation Models in Chemical Physics*; Mennucci, B., Cammi, R., Eds.; Wiley: New York, 2007.
- (70) Frisch, M. J.; Trucks, G. W.; Schlegel, H. B.; Scuseria, G. E.; Robb, M. A.; Cheeseman, J. R.; Scalmani, G.; Barone, V.; Mennucci, B.; Petersson, G. A.; Nakatsuji, H.; Caricato, M.; Li, X.; Hratchian, H. P.; Izmaylov, A. F.; Bloino, J.; Zheng, G.; Sonnenberg, J. L.; Hada, M.; Ehara, M.; Toyota, K.; Fukuda, R.; Hasegawa, J.; Ishida, M.; Nakajima, T.; Honda, Y.; Kitao, O.; Nakai, H.; Vreven, T.; Montgomery, J. A., Jr.; Peralta, J. E.; Ogliaro, F.; Bearpark, M.; Heyd, J. J.; Brothers, E.; Kudin, K. N.; Staroverov, V. N.; Kobayashi, R.; Normand, J.; Raghavachari, K.; Rendell, A.; Burant, J. C.; Iyengar, S. S.; Tomasi, J.; Cossi, M.; Rega, N.; Millam, J. M.; Klene, M.; Knox, J. E.; Cross, J. B.; Bakken, V.; Adamo, C.; Jaramillo, J.; Gomperts, R.; Stratmann, R. E.; Yazyev, O.; Austin, A. J.; Cammi, R.; Pomelli, C.; Ochterski, J. W.; Martin, R. L.; Morokuma, K.; Zakrzewski, V. G.; Voth, G. A.; Salvador, P.; Dannenberg, J. J.; Dapprich, S.; Daniels, A. D.; Farkas, Ö.; Foresman, J. B.; Ortiz, J. V.; Cioslowski, J.; Fox, D. J. *Gaussian Development Version, Revision H.10*; Gaussian Inc.: Wallingford, CT, 2010.
- (71) Cornell, W.; Cieplak, P.; Bayly, C.; Gould, I.; Merz, K.; Ferguson, D.; Spellmeyer, D.; Fox, T.; Caldwell, J.; Kollman, P. *J. Am. Chem. Soc.* **1995**, *117*, 5179–5197.
- (72) Yoshida, H. *Phys. Lett. A* **1990**, *150*, 262–268.
- (73) Jorgensen, W. L.; Chandrasekhar, J.; Madura, J. D.; Impey, R. W.; Klein, M. L. *J. Chem. Phys.* **1983**, *79*, 926–935.
- (74) Soper, A. K.; Phillips, M. G. *Chem. Phys.* **1986**, *107*, 47–60.
- (75) Duval, M.; Souaille, M.; Spezia, R.; Cartailier, T.; Vitorge, P. *J. Chem. Phys.* **2007**, *127*, 034503.
- (76) Duval, M.; Vitorge, P.; Spezia, R. *J. Chem. Phys.* **2009**, *130*, 104501.
- (77) Brancato, G.; Rega, N.; Barone, V. *Chem. Phys. Lett.* **2009**, *483*, 177–181.
- (78) Ern, A.; Guermond, J.-L. *Theory and Practice of Finite Elements*; Springer: Berlin, 2004; Vol. 159, pp 81–84.



# Insights into the Solvation and Mobility of the Hydroxyl Radical in Aqueous Solution

Edelsys Codorniu-Hernández and Peter G. Kusalik\*

Department of Chemistry, University of Calgary, 2500 University Drive NW, Calgary, T2N1N4, Alberta, Canada

**ABSTRACT:** A detailed description of the local solvation structure and mobility of hydroxyl radicals ( $\text{OH}^*$ ) in aqueous solution near ambient conditions is provided by Car–Parrinello molecular dynamics simulations. Here, we demonstrate that for HCTH/120 and BLYP functionals, smaller systems (i.e.,  $31 \cdot \text{H}_2\text{O}-\text{OH}^*$ ) are contaminated by system size effects, being biased for the presence of a three-electron two-centered hemibond structure between the oxygen atoms of a water molecule and the radical. Radial and spatial distribution functions of relatively large  $63 \cdot \text{H}_2\text{O}-\text{OH}^*$  systems reveal the existence of a 4-fold coordinated “inactive”  $\text{OH}^*$  structure with three H-bond donating neighbors and a strongly coordinated H-bond accepting neighbor. The local hydration structure around the radical exhibits more H-bond ordering than has been predicted by recent simulations employing classical force fields. Local structural fluctuations can end with spontaneous H-transfer reactions from the nearest H-bond donor water molecule, facilitated by the formation of an “active”  $\text{OH}^*$  state, resembling the proton transfer mechanism of hydrated  $\text{OH}^-$  (i.e., slight polarization of the  $(\text{H}_3\text{O}_2)^*$  complex). A comparison of the free energy barriers for the H-transfer reaction obtained by both DFT functionals and for both system sizes is also provided, demonstrating that this can be a very rapid process in water.

## 1. INTRODUCTION

The hydroxyl radical ( $\text{OH}^*$ ) has posed significant challenges to theoretical and experimental studies due to its high reactivity and very short lifetime.<sup>1</sup> However, it is still the target molecule of numerous investigations owing, in part, to its biological and atmospheric significance as well as its crucial role in industrial applications.<sup>2,3</sup> Described as the atmospheric “vacuum cleaner”, this radical is responsible for many of the reactions that remove volatile organic compounds from the air.<sup>4,5</sup> It oxidizes approximately 83% of annual methane emissions, making  $\text{OH}^*$  the most important processor of greenhouse gases.<sup>4,5</sup> Serious ailments such as cancer and Parkinson’s disease have also been related to  $\text{OH}^*$ ,<sup>6</sup> and it is recognized as the most reactive of the so-called reactive oxygen species (ROS). The only means to protect important cellular structures from its action is the use of antioxidants because, contrary to other oxidants,  $\text{OH}^*$  cannot be eliminated by an enzymatic reaction. This makes it a very dangerous compound to an organism, but interestingly,  $\text{OH}^*$  is also essential to the body’s natural defense mechanisms.<sup>6</sup>

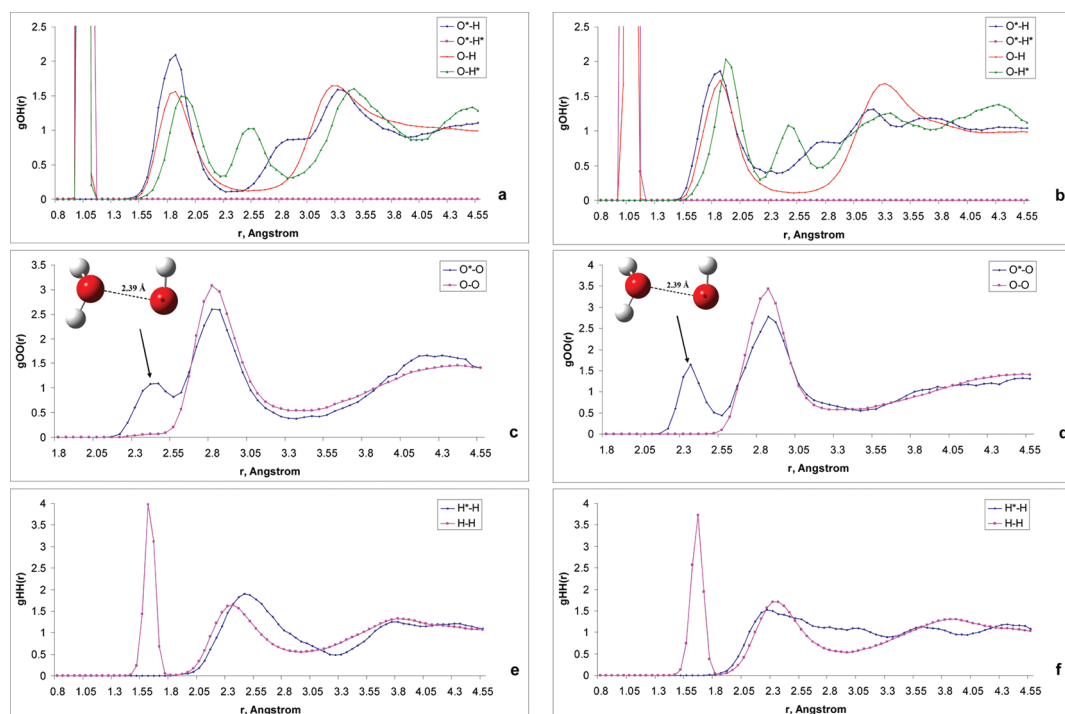
Water apparently has a crucial role in  $\text{OH}^*$  chemistry. Recent studies have pointed out its effect on modulating  $\text{OH}^*$  reactivity, providing a clear stabilization of transition states and higher reactivity via hydrogen bonding.<sup>7,8</sup> High level *ab initio* calculations of gas phase  $\text{OH}^*(\text{H}_2\text{O})_n$  clusters<sup>9–28</sup> have also been conducted. Their focus has been on the possible existence of a  $\text{H}_2\text{O}-\text{OH}^*$  complex which is speculated to influence strongly the diffusion and oxidative capacity of the radical. Only a few studies have considered the solvation of  $\text{OH}^*$  in liquid water,<sup>29–34</sup> aimed at demonstrating an expected  $\text{OH}^*$  ability to diffuse in water via hydrogen exchange analogous to the proton-exchange reaction in the case of  $\text{OH}^-$ .<sup>30,35</sup> Prior to our work,<sup>36</sup> this key reaction had not been directly detected by either experimental or theoretical studies due to the large challenges posed by  $\text{OH}^*$  to both fields. In addition, a recent spectroscopic observation<sup>37</sup>

made during the irradiation of  $\text{OH}^-$  in aqueous solutions has been attributed to ultrafast H-transfer reactions from neighboring water molecules to  $\text{OH}^*$ . These authors<sup>37</sup> had attempted to compare their proposed mechanism with previous Car–Parrinello MD simulations<sup>31</sup> in which a three-electron two-center hemibond structure between the oxygen atom of the radical and the oxygen atom of one water molecule was found to be a particularly stable structure.<sup>30,31</sup> In the presence of the hemibond, the supposed diffusion mechanism of  $\text{OH}^*$  in liquid water via a hydrogen exchange reaction is effectively impeded.<sup>30</sup> Afterward, Vandevonle et al.<sup>33</sup> claimed that BLYP and all GGA functionals overestimate the hemibond structure and, using self-interaction corrected methods, reported that  $\text{OH}^*$  acts as a good hydrogen bond donor but accepts fewer than two hydrogen bonds on average. However, the previous apparent inability of Car–Parrinello MD to determine a  $\text{OH}^*$  H-transfer reaction<sup>29–34</sup> seems inconsistent with the low reaction barrier (around 4.2 kcal/mol)<sup>27,38</sup> predicted for the hydrogen transfer reaction in the gas phase, which is in good agreement with the available experimentally derived data.<sup>11</sup>

As we have shown in a very recent paper,<sup>36</sup> Car–Parrinello molecular dynamics simulations of a larger ( $63 \cdot \text{H}_2\text{O}-\text{OH}^*$ ) system provide a different picture with respect to previous simulations using the smaller system ( $31 \cdot \text{H}_2\text{O}-\text{OH}^*$ ).<sup>29–34</sup> Here, we present a comparison between these two systems, providing a detailed description of the solvation and electronic behavior of  $\text{OH}^*$  in aqueous solution and demonstrating that smaller systems are contaminated by system size effects. Both the HCTH/120 and BLYP density functionals are employed. Important insights into the features of the H-transfer reaction ( $\text{OH}^* + \text{H}_2\text{O} \leftrightarrow \text{H}_2\text{O} + \text{OH}^*$ ) are also provided.

Received: June 18, 2011

Published: October 03, 2011

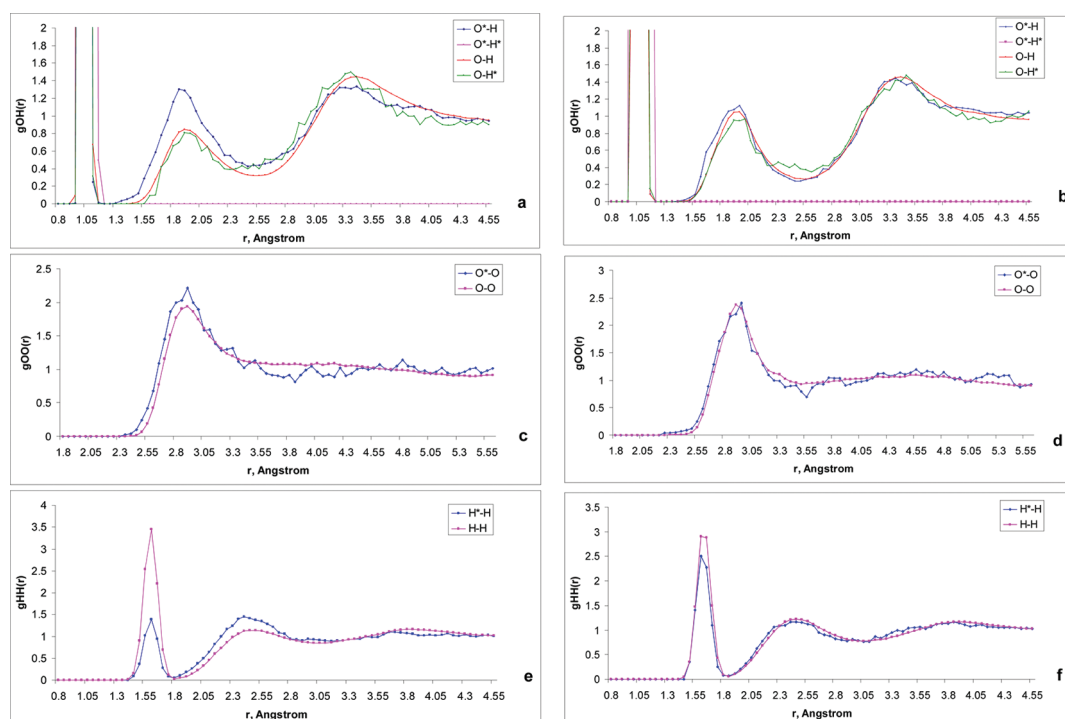


**Figure 1.** Radial distribution functions obtained with the HCTH/120 (left column) and BLYP (right column) functionals for total simulation times of 160 ps of  $31 \cdot \text{H}_2\text{O}-\text{OH}^*$  systems at a temperature of 310 K. (a, b) Oxygen–hydrogen RDFs in which  $\text{O}^*\text{H}$  is represented by a blue solid line,  $\text{O}^*\text{H}^*$  by a magenta solid line,  $\text{OH}$  by a red solid line, and  $\text{OH}^*$  by a green solid line. (c, d) Oxygen–oxygen RDFs in which  $\text{O}^*\text{O}$  is represented by a blue solid line and  $\text{OO}$  by a magenta solid line. (e, f) Hydrogen–hydrogen RDFs in which  $\text{H}^*\text{H}$  is represented by blue solid line and  $\text{H}-\text{H}$  by a magenta solid line. A representative structure for the  $\text{O}^*\text{O}$  hemibond interaction is shown in c and d.

## 2. SIMULATION DETAILS

**2.1. Car–Parrinello Molecular Dynamics Simulations.** The standard Car–Parrinello<sup>39</sup> DFT-based *ab initio* MD method was used with the CPMD<sup>40</sup> code to study  $31 \cdot \text{H}_2\text{O}-\text{OH}^*$  and  $63 \cdot \text{H}_2\text{O}-\text{OH}^*$  systems within periodic boundary conditions. The cubic simulation cells have lengths of 9.85 Å and 12.56 Å, respectively. The resulting density of 1 g/cm<sup>3</sup> corresponds to the density of water under ambient conditions. The local spin density approximation (LSDA) was employed to account for the unpaired electron at the hydroxyl radical. Two different density functionals were utilized and compared, the gradient-corrected exchange–correlation energy functionals of Becke, Lee, Yang, and Parr<sup>41,42</sup> (BLYP) and the HCTH/120.<sup>43</sup> The HCTH/120<sup>43</sup> functional was employed because it has been reported to describe accurately the properties of liquid water.<sup>43–45</sup> This particular functional is a highly parametrized GGA functional which was fit to a large set of empirical molecular properties. For our simulations, we applied the HCTH/120 functional and the Troullier–Martins norm-conserving pseudopotential,<sup>46</sup> where the valence electronic wave function is described with a plane wave basis with an energy cutoff of 90 Ry, which provides a reasonable basis set convergence for this particular system. In addition, we compare the results obtained from the HCTH/120 functional with those obtained with the BLYP functional with the Goedecker<sup>47</sup> norm-conserving pseudopotential and a valence electronic wave function described with a plane wave basis with a 75 Ry energy cutoff since this scheme was previously applied in the study of systems with 31  $\text{H}_2\text{O}$  molecules and a hydroxyl radical.<sup>29–31</sup> We use this scheme with the BLYP functional for both small ( $31 \cdot \text{H}_2\text{O}-\text{OH}^*$ ) and large ( $63 \cdot \text{H}_2\text{O}-\text{OH}^*$ ) systems for

consistency in our exploration of finite size effects. Tests with the BLYP functional in combination with either pseudopotential gave the same results. The simulations were carried out with a fictitious mass of 600 au and a simulation temperature of 310 K, whereas previous calculations employed 600 au and 800 au.<sup>31</sup> The supporting information of ref 36 explores the possible impact of the selected fictitious electronic mass for this system. We have performed a benchmarking exercise for the current methodology against non-DFT and all-electron DFT methods demonstrating that the electronic approach used in the present simulations is reasonable for our purposes. The fictitious electron kinetic energy and the dynamics of the atoms were controlled by a chain of three Nose–Hoover thermostats<sup>48</sup> operating at characteristic frequencies of 6000 cm<sup>-1</sup> and 2000 cm<sup>-1</sup>, respectively. The average fictitious kinetic energy was maintained at levels of 0.035 Ha for the  $31 \cdot \text{H}_2\text{O}-\text{OH}^*$  system and 0.06 Ha for the  $63 \cdot \text{H}_2\text{O}-\text{OH}^*$  system, both remaining stable during the whole simulation. The spin distribution function of two other systems,  $23 \cdot \text{H}_2\text{O}-\text{OH}^*$  and  $95 \cdot \text{H}_2\text{O}-\text{OH}^*$ , with dimensions of 8.4 and 14.4 Å, respectively, was calculated after geometry optimization. The input structures for the CPMD simulations were taken from a large liquid water system previously equilibrated by classical MD simulations. For these systems, a hydrogen atom was removed from the water molecule that was closest to the center of mass of the system. The first 2 ps of CPMD dynamics following the proper equilibration of the energies were also considered as equilibration and discarded. The total simulation times were 160 ps for the  $31 \cdot \text{H}_2\text{O}-\text{OH}^*$  systems and 50 ps for the  $63 \cdot \text{H}_2\text{O}-\text{OH}^*$  systems. The time step was set to 0.1 fs. Bin sizes of 0.05 Å and 5° were used for the radial and angular distribution functions, respectively. Taking into account the



**Figure 2.** Radial distribution functions obtained with the HCTH/120 (left column) and BLYP (right column) functionals for total simulation times of 50 ps of  $63 \cdot \text{H}_2\text{O}-\text{OH}^*$  systems at a temperature of 310 K. (a, b) Oxygen–hydrogen RDFs in which  $\text{O}^*\text{H}$  is represented by a blue solid line,  $\text{O}^*\text{H}^*$  by a magenta solid line,  $\text{OH}$  by a red solid line, and  $\text{OH}^*$  by a green solid line. (c, d) Oxygen–oxygen RDFs in which  $\text{O}^*\text{O}$  is represented by a blue solid line and  $\text{OO}$  by a magenta solid line. (e, f) Hydrogen–hydrogen RDFs in which  $\text{H}^*\text{H}$  is represented by a blue solid line and  $\text{HH}$  by a magenta solid line.

importance of obtaining measures of average spatial structures<sup>49</sup> for this kind of study, a custom code was used for calculating the probability distributions of oxygen and hydrogen atoms around the  $\text{O}^*$  (or  $\text{H}^*$ ) over all system trajectories. Spherical polar coordinates were used for averaging. VMD<sup>50</sup> was utilized for visualization of configurations as well as for displaying isosurfaces. Henceforth, hydrogen and oxygen atoms from water molecules are denoted as H and O, while those from the radical are denoted as  $\text{H}^*$  and  $\text{O}^*$ , respectively.

**2.2. Constrained MD.** The free energy of the hydrogen transfer reaction between the hydroxyl radical and the water molecules was performed using constrained molecular dynamics simulations.<sup>51</sup> The difference between the  $\text{O}^*\text{H}$  and  $\text{OH}$  distances was chosen as the constraint variable  $R$ , where H is the hydrogen atom being transferred from a neighboring water molecule to the radical. For each increment ( $\sim 0.15 \text{ \AA}$  for the  $63 \cdot \text{H}_2\text{O}-\text{OH}^*$  systems and  $\sim 0.05 \text{ \AA}$  for the  $31 \cdot \text{H}_2\text{O}-\text{OH}^*$  systems), the average constraint force was measured over a 3 ps trajectory. From such simulations, the free energy profile may be obtained from a straightforward thermodynamic integration over the coordinate  $R$ , and the symmetry in the results was numerically imposed for the  $63 \cdot \text{H}_2\text{O}-\text{OH}^*$  system.

### 3. RESULTS AND DISCUSSION

#### 3.1. Solvation Structure of the $\text{OH}^*$ in Aqueous Solution.

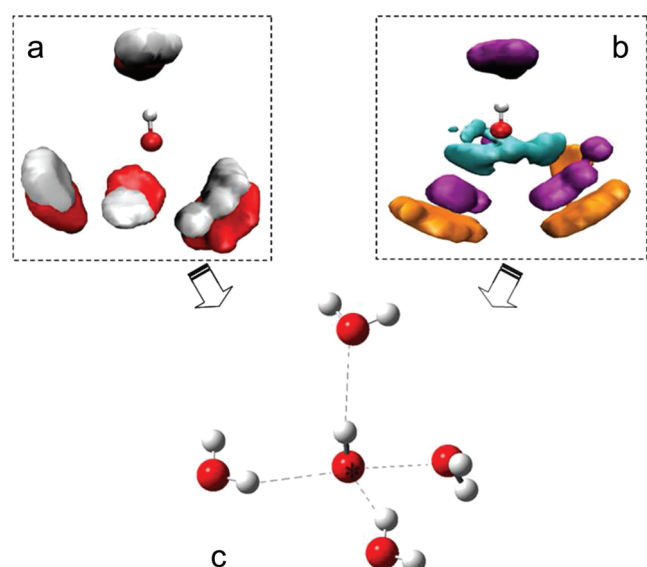
Significant differences can be observed in the radial, angular, and spatial distribution functions for the two studied system sizes (small,  $31 \cdot \text{H}_2\text{O}-\text{OH}^*$  and large,  $63 \cdot \text{H}_2\text{O}-\text{OH}^*$ ). The absence of a three-electron two-center hemibonded structure (between the oxygen atom of the radical and the oxygen atom

**Table 1.** Coordination Numbers,  $n(r)$ , for the Hydroxyl Radical in Which the Coordination of Hydrogen and Oxygen Atoms around the Radical Oxygen Has Been Measured ( $\text{O}^*\text{H}$  and  $\text{O}^*\text{O}$ , respectively) As Well As the Coordination of Oxygen Atoms around the Radical Hydrogen ( $\text{H}^*\text{O}$ )<sup>a</sup>

$r$ (Å)	$31 \text{ H}_2\text{O}-\text{OH}^*$		$63 \text{ H}_2\text{O}-\text{OH}^*$	
	BLYP	HCTH/120	BLYP	HCTH/120
$\text{O}^*\text{H}$	1.2		0.8	1
peak 1 ( $r$ )	2.4 (2.45)	2.1 (2.55)	3.1 (2.75)	3.2 (2.71)
peak 2 ( $r$ )	4.4 (2.90)	4.2 (3.05)		
	4.5	25.2	24.4	24.5
$\text{O}-\text{H}$	1.2	1.9	1.9	1.9
	2.5	3.9	3.9	3.9
	4.5	24.9	25.2	24.8
$\text{O}^*\text{O}$	2.6	0.99	1.07	0.1
	3.4	4.92	4.24	4.5
	4.5	12.2	12.9	11.3
$\text{O}-\text{O}$	3.4	4.2	4.1	4.4
	4.5	12.21	12.2	11.3
$\text{H}^*\text{O}$	2.2	0.9	0.7	1.6
	2.7	1.7	1.1	2.1
	3.6	5.8	6.1	6.5
	4.5	12.9	11	12.3

<sup>a</sup>The coordination numbers of hydrogen and oxygen atoms around the water oxygen are also provided as  $\text{OH}$  and  $\text{OO}$ , respectively.

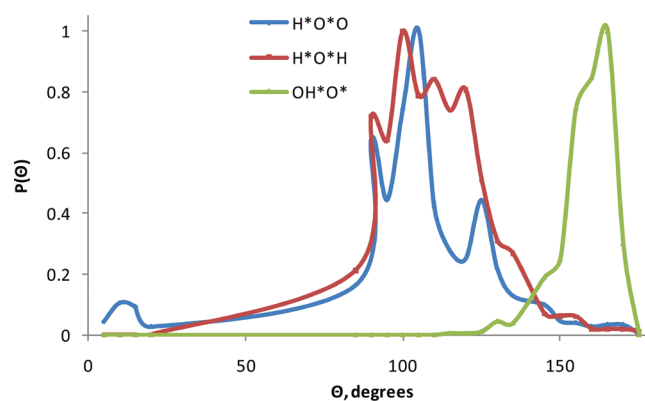
of one water molecule) in the  $63 \cdot \text{H}_2\text{O}-\text{OH}^*$  system is immediately apparent from an analysis of the radial distribution



**Figure 3.** (a) Spatial distributions of hydrogen (white isosurfaces, threshold 2.1) and oxygen (red isosurfaces, threshold 3.4) atoms around the OH\* for the  $31 \cdot \text{H}_2\text{O}-\text{OH}^*$  system obtained with the HCTH/120 functional from 160 ps trajectories. We remark that the results obtained with the BLYP functional appear very similar. A 4-fold coordination of OH\* is evident in which two water neighbors donate H bonds to the radical, one water interacts through the O–O\* hemibond, and another water accepts a H bond from OH\*. (b) Spatial distributions of oxygen atoms around the OH\* for the  $31 \cdot \text{H}_2\text{O}-\text{OH}^*$  system with respect to different O\*–O separations. Color legend: cyan ( $2.3 \pm 0.25 \text{ \AA}$ ), purple ( $2.8 \pm 0.35 \text{ \AA}$ ), orange ( $4.05 \pm 0.8 \text{ \AA}$ ). The cyan isosurfaces (threshold 2.8) demonstrate that the hemibond structure is not always formed with the same water molecule, the purple isosurfaces (threshold 3.2) show the 4-fold coordination of H-bonding oxygen atoms, and the orange isosurfaces (threshold 2.4) show the spatial structure of oxygen atoms in the second solvation shell. (c) A representative snapshot configuration exhibits a typical spatial arrangement of water molecules around the radical for this smaller system.

functions (RDF) shown in Figures 1 and 2. An unusual peak in the  $g_{\text{O}^*\text{O}}(r)$  function at around  $2.4 \text{ \AA}$  confirms the existence of this kind of interaction from both the HCTH/120 (Figure 1c) and BLYP (Figure 1d) functionals for the small  $31 \cdot \text{H}_2\text{O}-\text{OH}^*$  system; the peak does not appear with the increased size of the periodic simulation box for the  $63 \cdot \text{H}_2\text{O}-\text{OH}^*$  system (Figure 2c and d). Comparing Figures 1 and 2, we can see that the presence of the hemibond in the same system significantly alters the local structure around the radical.

The similarity of the BLYP and HCTH RDFs for both systems (Figures 1 and 2) is important to note. In the case of the  $31 \cdot \text{H}_2\text{O}-\text{OH}^*$  system, our results for the first 80 ps utilizing the BLYP functional reproduced the RDFs reported by Khalack and Lyubartsev;<sup>31</sup> however, the RDFs reported in Figure 1 averaged over a 160 ps trajectory show a more defined  $g_{\text{O}^*\text{O}}(r)$  hemibond peak (Figure 1d) relative to the shoulder apparent in ref 31, with a consequent distortion of the  $g_{\text{O}^*\text{H}}(r)$  in Figure 1. This confirms the persistence (i.e., stability) of the hemibond structure in  $31 \cdot \text{H}_2\text{O}-\text{OH}^*$  systems. The integration of the RDF peaks (Table 1) for this smaller system further supports the existence of a local structure very similar to those reported by Vassilev et al.<sup>30,33</sup> and Khalack and Lyubartsev,<sup>31</sup> with two H-bond donating neighbors, a H-bond accepting neighbor, and a fourth water molecule forming a hemibond. Such a structure is

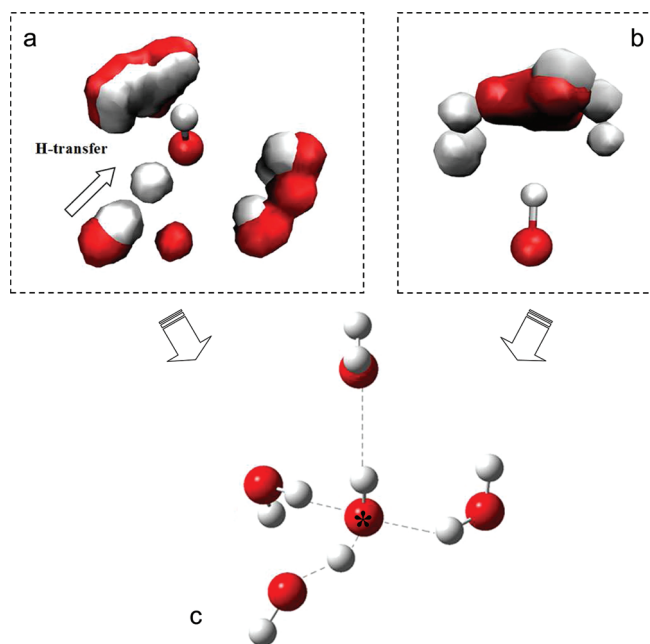


**Figure 4.** Distribution of the angles H\*O\*O (blue solid line), H\*O\*H (red solid line), and OH\*O\* (green solid line) for the first solvation shell molecules of the  $63 \cdot \text{H}_2\text{O}-\text{OH}^*$  and the HCTH/120 DFT functional averaged for 50 ps of simulation. The BLYP functional yields similar results. The blue and red lines represent the angular distribution of H-bond donating neighbors of OH\*, and the green line represents the angular distribution of the H-bond accepting neighbor of OH\*.

visualized in Figure 3, in which spatial distribution functions for the most probable locations of the nearest neighbors of the OH\* (averaged over the full 160 ps trajectories) are shown along with a representative configuration. The persistence of the hemibond appears to be rather unfavorable for the H-transfer process since no spontaneous H-transfers were observed in either of the 160 ps small system simulations.

A different local hydration structure of the radical is evident for the  $63 \cdot \text{H}_2\text{O}-\text{OH}^*$  system. It is noteworthy that in Figure 2 the OH\*–water distribution functions are generally similar to the corresponding water–water RDFs, and there is no evidence of the hemibond in the  $g_{\text{O}^*\text{O}}(r)$  peaks (Figure 2c and d). Again, the BLYP and HCTH/120 DFT functionals yield very similar RDFs. From the integration of the RDF peaks (Table 1), it is apparent that the OH\* accepts three H bonds and donates one with water molecules in its first hydration shell. The formation of an almost tetrahedral configuration around the radical is suggested by the angular distribution functions presented in Figure 4. This structure, identified as an “inactive” OH\* state in our recent report,<sup>36</sup> can be visualized with spatial distribution functions of the most probable locations of the nearest neighbors of the OH\* (Figure 5a and b). As we shall detail below, local structural fluctuations of this “inactive” OH\*(H<sub>2</sub>O)<sub>4</sub> can end with spontaneous H-transfer reactions in these larger systems. We note that in Figure 5a, the hydrogen atom transferred to the OH\* appears as one of the most probable locations. A representative configuration for the transfer process is shown in Figure 5c. These larger system results clearly impact our understanding of the mobility and solvation of OH\* in aqueous solution.

In addition to previous Car–Parrinello MD results, three classical MD studies have been published<sup>52–54</sup> looking to provide a description of the behavior of OH\* in aqueous solution. Classical potential models were employed and so were unable to capture the formation of the (H<sub>3</sub>O<sub>2</sub>)\* complex and the H-transfer reaction. Unfortunately, the potential models used apparently were not able to provide an adequate description of the specific features observed in the interaction of OH\* with water molecules. Consequently, Campo and Grigera<sup>52</sup> obtained a rather different local hydration structure for OH\* with only one H-bond water donor (53% of the time) and one H-bond water acceptor

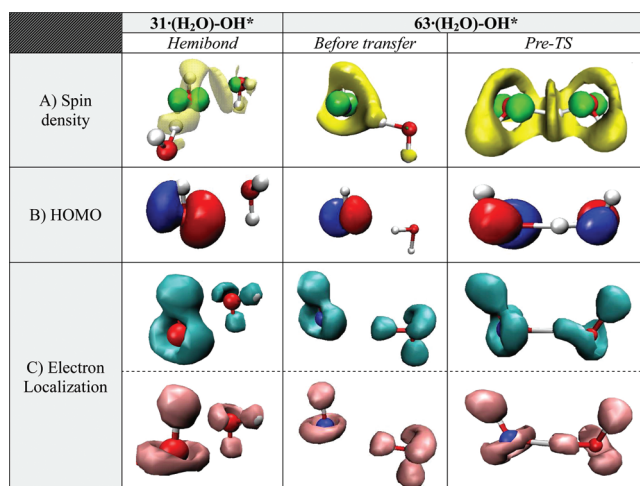


**Figure 5.** Spatial distributions of hydrogen (white isosurfaces, threshold 2.1) and oxygen (red isosurfaces, threshold 3.4) atoms around the  $\text{OH}^*$  for the  $63 \cdot \text{H}_2\text{O}-\text{OH}^*$  system obtained with the HCTH/120 functional from 50 ps trajectories. We remark that the results obtained with the BLYP functional appear very similar. (a) The isosurfaces corresponding to H-bond donating neighbors of  $\text{OH}^*$  are shown. We note that at shorter  $\text{O}^*\text{H}$  distances, the hydrogen atom position for H transfer is evident. (b) The isosurfaces for the H-bond accepting neighbor of  $\text{OH}^*$  are shown. (c) A representative snapshot configuration shows a typical spatial arrangement of water molecules around the radical during a H transfer.

(88% of the time). Svishchev and Plugatyr<sup>53</sup> reported that  $\text{OH}^*$  occupies “holes” in the tetrahedral arrangement of water molecules, with no indication of H-bond-like arrangements. Pabis et al.<sup>54</sup> developed a flexible potential for  $\text{OH}^*$  (derived from *ab initio* gas phase  $\text{OH}^*$ –water dimer energies) and observed that the radical tends to occupy cavities in the hydrogen-bonded network of the water molecules with hydration shells of 13–14 water molecules. These striking differences from the present results again point to the importance of an adequate description of the interactions of  $\text{OH}^*$  in a water environment.

A particular RDF feature worthy of comment is the flatness of the first peak of the  $g_{\text{OO}}(r)$ , particularly for the HCTH/120 functional (Figure 2c). A previous comparison of the oxygen–oxygen RDF obtained from *ab initio* MD (CPMD) for the HCTH/120, BLYP, and BPZ functionals with experimental results has been reported by Boese et al.<sup>43</sup> These authors found the position and depth of the first minimum more shallow and shifted to a slightly smaller distance with respect to experimental data. Comparing their RDF with Figure 2c, it can be seen that the distribution function obtained here is somewhat flatter and less pronounced than that for pure liquid water. Apparently, the presence of the  $\text{OH}^*$  has perturbed the water structure somewhat.

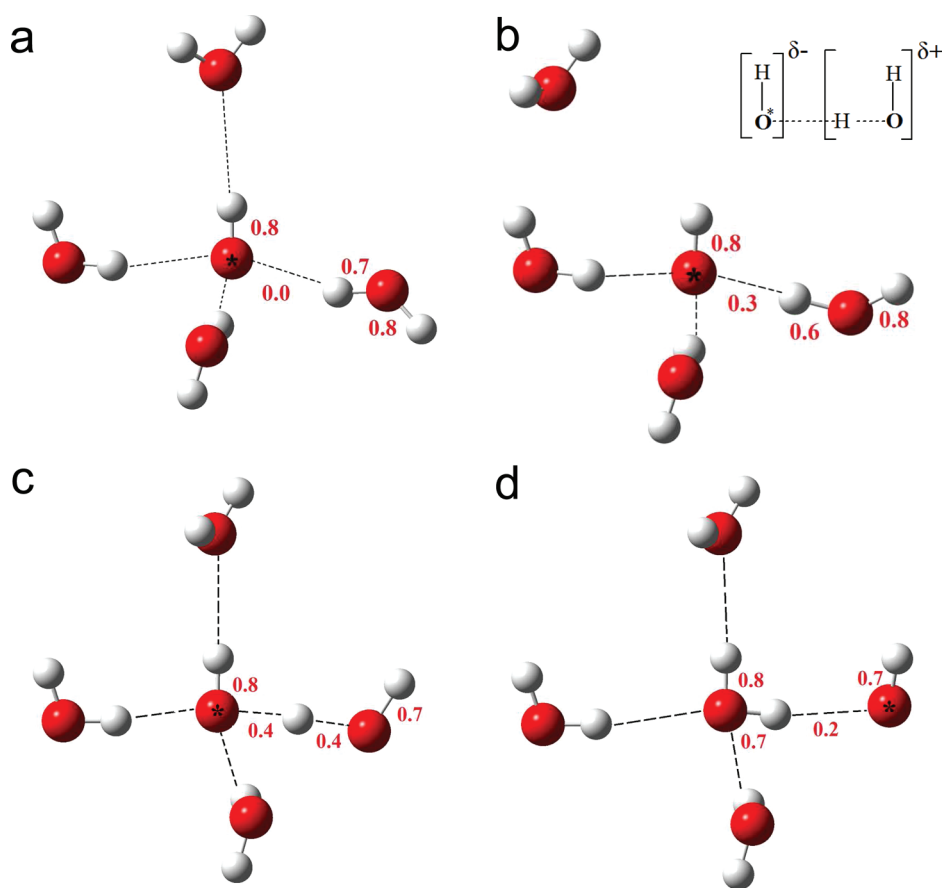
**3.2. Size Effects and Electronic Features of  $\text{OH}^*$  in Aqueous Solution.** The marked differences observed between the RDF and SDF results for the small and large systems indicate that those from the former are contaminated by system size effects. An examination of the electronic features for both systems yields further interesting results. Different from calculations in the gas



**Figure 6.** Electronic features of the three-electron two-center hemibond structure obtained with the  $31 \cdot \text{H}_2\text{O}-\text{OH}^*$  system and the spontaneous H-transfer reaction obtained with the  $63 \cdot \text{H}_2\text{O}-\text{OH}^*$  system<sup>36</sup> from the HCTH/120 functional. The BLYP functional yields similar results. Row A represents the spin density distribution with yellow (value of +0.0004) and green (value of  $-0.03$ ) isosurfaces. Row B shows the evolution of the HOMO where red and blue isosurfaces have values of  $-0.03$  and  $+0.04$ , respectively. For the smaller system, the negative HOMO isosurface points toward the hemibond oxygen. For the larger system, the HOMO is localized on the  $\text{OH}^*$  before the transfer and is perpendicular to the H bond with the nearest neighboring water molecule. In the pre-transition state, both the HOMO and HOMO–1 orbitals are shown due to the existence of  $\alpha$  and  $\beta$  degenerated states centered on the water and  $\text{OH}^*$  oxygens. Row C presents the ELF- $\alpha$  as cyan isosurfaces (threshold 0.85) and the ELF- $\beta$  as pink isosurfaces (threshold 0.85).

phase, periodic boundary conditions in MD simulations can affect certain systems, which is a possible explanation for this effect. In a  $31 \cdot \text{H}_2\text{O}-\text{OH}^*$  system, the  $\text{OH}^*$  must share its second hydration shell with its image in the periodic simulation box (being 9.84 Å in width). The highly reactive  $\text{OH}^*$  can potentially “sense” its image, allowing electronic artifacts and “undesirable” structures for the hydrated  $\text{OH}^*$  to arise during the simulation. We note that at the end of the geometry optimization step there is a spin delocalization in smaller systems ( $23 \cdot$  and  $31 \cdot \text{H}_2\text{O}-\text{OH}^*$ ), while for the  $63 \cdot$  and  $95 \cdot \text{H}_2\text{O}-\text{OH}^*$  systems, the spin density is primarily concentrated on the  $\text{OH}^*$ . Changes to the plane wave cutoff (90 Ry, 100 Ry, 120 Ry) do not alter this behavior, supporting the conjecture of a system size effect as the origins for this artifact. In our case, the selection of the  $63 \cdot \text{H}_2\text{O}-\text{OH}^*$  system implies a necessary compromise between accuracy and computational expense.

The finite size effects in MD simulations of the smaller  $31 \cdot \text{H}_2\text{O}-\text{OH}^*$  system are manifested as a persistent hemibonded structure, which is apparently a rather unfavorable structure for H abstraction. Further to the description of the electronic features during the H-transfer reaction we provided very recently,<sup>36</sup> here we focus on a comparison of the electronic properties of the  $63 \cdot \text{H}_2\text{O}-\text{OH}^*$  system with those of the  $31 \cdot \text{H}_2\text{O}-\text{OH}^*$  system (Figure 6). The spin density isosurface obtained for the smaller system evidences that the singly occupied  $p\pi$  MO of the  $\text{OH}^*$  is tied down in this hemibond; both the water and the  $\text{OH}^*$  share negative and positive spin density, consistent with a stabilizing resonance in which the electron pair is on the water oxygen atom



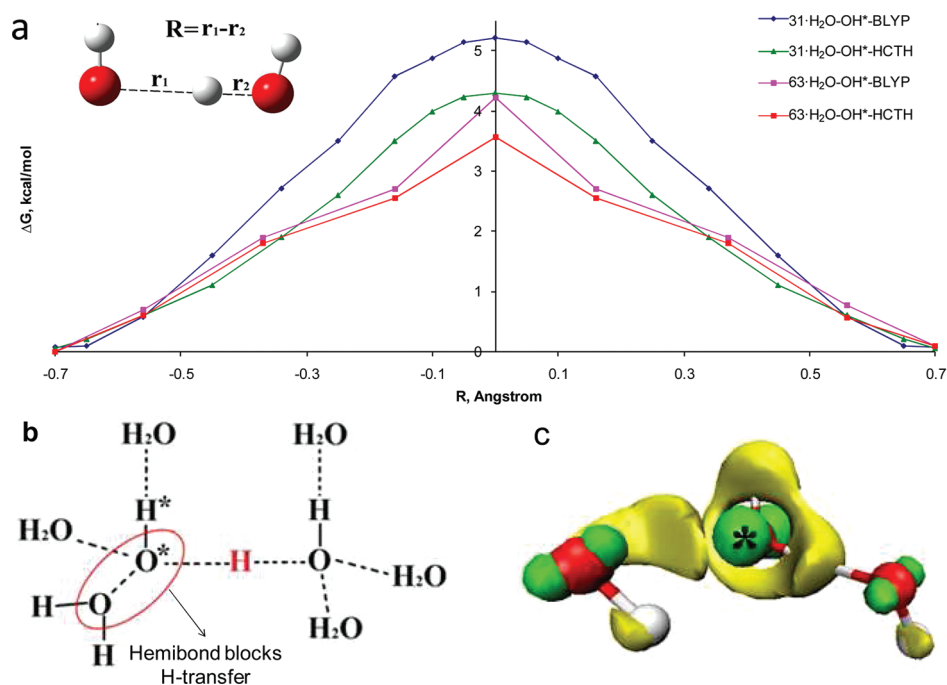
**Figure 7.** Schematic molecular configurations for different states during the spontaneous H-transfer reaction for the  $63 \cdot \text{H}_2\text{O} - \text{OH}^*$  system. (a) The initial state; (b) the pre-transition state, where the polarization of the structure is represented in the inset; (c) the transition state, a  $(\text{H}_3\text{O}_2)^*$  complex; and (d) the post-transfer state. The whole reaction (conversion from a–d) occurs in approximately 0.7 ps. The red numbers are the bond orders for the radical and the water molecule involved in the reaction.

while the unpaired electron is on the  $\text{OH}^*$  oxygen, or vice versa. In the case of the larger system, the spin density is primarily located on the  $\text{OH}^*$  prior to the H-transfer reaction. During the initial part of the reaction, a small portion of positive spin density is shared with the water molecule involved in the reaction (Figure 6). As we described in ref 36, a  $(\text{H}_3\text{O}_2)^*$  complex appears prior to the transition state (Pre-TS, Figure 6) in which positive and negative spin density is shared across both oxygen atoms.

The highest occupied molecular orbital (HOMO) provides an alternate description of the H-transfer reaction mechanism. The early movement of the electron was suggested<sup>36</sup> by the existence of essentially degenerated HOMO and HOMO–1, located on the  $\text{OH}^*$  oxygen and the water oxygen, respectively, in the pre-transition state. In the case of the smaller system, the existence of the hemibond seems not to relate with the location of the HOMO: Although the spin density shows that the unpaired electron is shared between the  $\text{OH}^*$  and the water molecule, the HOMO and HOMO–1 are both located only on the  $\text{OH}^*$  oxygen (see Figure 6). The electronic localization functions (ELF) presented in Figure 6 for the smaller system appear very similar to those of the larger system. These functions appear to remain localized on the  $\text{OH}^*$  and water molecule, further supporting the claim that resonance structures stabilize the hemibond rather than a chemical bonding. The  $\beta$  function (ELF\_BETA) for the  $31 \cdot \text{H}_2\text{O} - \text{OH}^*$  systems (Figure 6) exhibits a continuous ring around  $\text{OH}^*$  very similar to that for the  $63 \cdot \text{H}_2\text{O} - \text{OH}^*$

system before the transfer, and very similar to the picture obtained for the hydroxyl anion.<sup>35</sup> A p-like function for the unpaired electron appears in the ELF\_ALPHA for the  $\text{OH}^*$  for both small and larger systems, which becomes somewhat modified in the pre-transition state for the  $63 \cdot \text{H}_2\text{O} - \text{OH}^*$  system.

We have suggested<sup>36</sup> that the H-transfer reaction has characteristics of a hybrid mechanism apparently involving aspects of a hydrogen-atom transfer (HAT) and an electron–proton transfer (EPT). As the proton and the electron come from the same bond in this reaction, a HAT mechanism is expected. However, the evolution of the HOMO, ELF, and spin density suggests that an early electron movement occurs in the pre-transition state when the hydrogen atom (or proton) is still closer to the water oxygen. A schematic representation of the local structure of the  $\text{OH}^*$  during different states of this reaction is presented in Figure 7, along with the bond orders for the  $\text{OH}^*$  and the water molecule involved in the transfer. The three spontaneous events observed during the two simulations of the larger system all exhibited the same basic structural pattern. Figure 7a shows the “inactive”  $\text{OH}^*(\text{H}_2\text{O})_4$  state already introduced in Figure 5a. A key point in this reaction appears to be related with a change of the  $\text{OH}^*$  hydration structure from the “inactive” to the “active”  $\text{OH}^*(\text{H}_2\text{O})_3$  form (see Figure 7b). In the “active” state, the H-bond of the accepting neighbor to the  $\text{OH}^*$  becomes significantly weakened.<sup>36</sup> Interestingly, a slight polarization of the  $(\text{H}_3\text{O}_2)^*$  complex is evident at this point, while the hydrogen



**Figure 8.** Results of constrained molecular dynamic simulation for the H-transfer reaction between OH\* and a neighboring water molecule. (a) Free energy profiles using the BLYP density functional for the 31·H<sub>2</sub>O–OH\* system (blue line) and the 63·H<sub>2</sub>O–OH\* (magenta line)<sup>36</sup> and the HCTH/120 density functional for the 31·H<sub>2</sub>O–OH\* (green line) and 63·H<sub>2</sub>O–OH\* (red line)<sup>36</sup> utilizing  $R$  as the displacement coordinate. The estimated error is in the range of 0.1–0.2 kcal/mol, obtaining by comparing forward and reverse reaction pathways as well as considering the symmetry of these functions. (b) The transition state structure obtained with BLYP (31·H<sub>2</sub>O–OH\*) shows that the hemibond structure persists even with the use of structural constraints and apparently contributes to the higher value of the energy barrier for this small system. (c) The spin density corresponding to the transition state in b is shown as yellow (value of +0.0004) and green (value of -0.03) isosurfaces and attests to its delocalization among the radical, the hemibonded water, and the water involved in the H transfer.

atom (or proton) is still not fully shared between the two oxygens (see schematic representation in Figure 7c). Consequently, the “active” states of both OH\* and OH<sup>-</sup> are similar, where the H-transfer mechanism for OH\* has key elements that resemble the proton transfer mechanism of hydrated OH<sup>-</sup>. Yet, from an analysis of constrained MD trajectories with the smaller systems for either of the functionals used (see details below), it is evident that this polarization of the (H<sub>3</sub>O<sub>2</sub>)\* complex is not present; neither is the weakening of the H bond to the accepting neighbor. Again, the presence of the hemibond is significantly altering the observed behavior during the (imposed) H transfer.

**3.3. System Size Effects in the Free Energy Barrier for H Transfer.** As already stated, the H-transfer reaction is crucial to understanding the mobility and reactivity of the OH\* in aqueous solution. Figure 8a shows the free energy profiles obtained after averaging the values of the mean forces for the forward and reverse processes in both 31·H<sub>2</sub>O–OH\* and 63·H<sub>2</sub>O–OH\* systems for both DFT functionals. A small free energy barrier is predicted in all cases, where a value of about 4 kcal/mol was obtained for the 63·H<sub>2</sub>O–OH\* systems in good agreement with experimentally derived values<sup>11</sup> and high-level *ab initio* calculations in the gas phase (4.2 kcal/mol).<sup>27,38</sup> We note that tests using a larger step in the thermodynamic integration (i.e., ~0.15 Å for the 31·H<sub>2</sub>O–OH\* system) do not yield substantial differences in the resulting barrier height. For the small system modeled with the BLYP functional, the free energy barrier for the reaction is somewhat higher than those obtained for the larger systems. After examining the structural features of this smaller system during the imposed H transfer, it is possible to see (Figure 8b)

that the hemibond structure between the oxygen atoms of the radical and water molecules persists to the transition state. As shown in Figure 8c, the spin density is shared across the radical, the hemibonded water, and the water involved in the imposed H transfer. The hemibond predicted by HCTH/120 for the small systems seems to be less stable (c.f. the “shoulder” in the  $g_{OO}(r)$  of Figure 1c with the more defined peak in Figure 1d). With the HCTH/120 functional, the hemibond is sufficiently weak that it does not appear to survive during the imposed H transfer, thereby allowing the smaller system to exhibit a free energy barrier similar to those obtained for the larger system. Metadynamics results previously reported<sup>36</sup> confirm that the barrier for this reaction is indeed small, having an upper bound of 6 kcal/mol.

## CONCLUSIONS

Car–Parrinello molecular dynamics simulations of OH\* in liquid water, utilizing different system sizes with 31 and 63 water molecules, reveal significant insight into the hydration and mobility of OH\* in solution. Analysis of radial and spatial distribution functions demonstrates the existence of system size effects with consequent electronic implications within MD results when using smaller systems (i.e., 31·H<sub>2</sub>O–OH\*). Smaller systems (i.e., 31·H<sub>2</sub>O–OH\*) show the presence of a three-electron two-centered hemibond structure between the oxygen atoms of a water molecule and the radical. Simulations with 63·H<sub>2</sub>O–OH\* systems show two main states in the OH\* solvation, a 4-fold coordination OH\*(H<sub>2</sub>O)<sub>4</sub> as an “inactive” state in which OH\* is donating one H bond and accepting other three H bonds from

water molecules and an “active” state with three H-bond donating neighbors and a weakly coordinated H-bond accepting neighbor. Previously studied classical models seem to underestimate the interaction of water molecules with the radical. The H-transfer reaction is apparently a very rapid process in water with a relatively small free energy barrier which can contribute significantly to OH\* mobility in aqueous solution. Further spectroscopic characterization of this reaction, critical in various scientific fields, by modern ultrafast experimental techniques is clearly warranted.

## AUTHOR INFORMATION

### Corresponding Author

\*Telephone: (403)-220-6244. E-mail: pkusalik@ucalgary.ca.

## ACKNOWLEDGMENT

We are grateful for the financial support of the Natural Sciences and Engineering Research Council of Canada and the Canadian Foundation for Innovation. We also acknowledge computational resources made available via WestGrid (www.westgrid.ca) and the University of Calgary. E.C.-H. wants to acknowledge Prof. Dr. Alberto Rolo-Naranjo for his help with the SDF code and Dr. Daniel Boese for useful discussions.

## REFERENCES

- (1) Sies, H. *Eur. J. Biochem.* **1993**, *215*, 213–219.
- (2) Isaksen, I. S. A.; Dalsøren, S. B. *Science* **2011**, *331*, 38–39.
- (3) Ikai, H.; Nakamura, K.; Shirato, M.; Kanno, T.; Iwasawa, A.; Sasaki, K.; Niwano, Y.; Kohno, M. *Antimicrob. Agents Chemother.* **2010**, *54*, 5086–5091.
- (4) Seinfeld, J. H.; Pandis, S. N. *Atmospheric Chemistry and Physics: From air pollution to climate change*, 1st ed.; John Wiley and Sons, Inc.: New York, 1998; pp 204–206.
- (5) Allodi, M. A.; Dunn, M. E.; Livada, J.; Kirschner, K. N.; Shields, G. C. *J. Phys. Chem. A* **2006**, *110*, 13283–13289.
- (6) Manda, G.; Nechifor, M.-T.; Neagu, T.-M. *Curr. Chem. Biol.* **2009**, *3*, 342–366.
- (7) Mitroka, S.; Zimmeck, S.; Troya, D.; Tanko, J.-M. *J. Am. Chem. Soc.* **2010**, *132*, 2907–2913.
- (8) Vöhringer-Martinez, E.; Hansmann, B.; Hernandez-Soto, H.; Francisco, J. S.; Troe, J.; Abel, B. *Science* **2007**, *315*, 497–501.
- (9) Sennikov, P. G.; Ignatov, S. K.; Schrems, O. *ChemPhysChem* **2005**, *6*, 392–412.
- (10) Chipman, D. M. *J. Phys. Chem. A* **2008**, *112*, 13372–13381.
- (11) Dubey, M. K.; Mohrschladt, R.; Donahue, N. M.; Anderson, J. G. *J. Phys. Chem. A* **1997**, *101*, 1494–1500.
- (12) Tsuji, K.; Shibuya, K. *J. Phys. Chem. A* **2009**, *113*, 9945–9951.
- (13) Soloveichik, P.; O'Donnell, B. A.; Lester, M. I.; Francisco, J. S.; McCoy, A. B. *J. Phys. Chem. A* **2010**, *114*, 1529–1538.
- (14) Cooper, P. D.; Kjaergaard, H. G.; Langford, V. S.; McKinley, A. J.; Quickenden, T. I.; Schofield, D. P. *J. Am. Chem. Soc.* **2003**, *125*, 6048–6049.
- (15) Crespo-Otero, R.; Sánchez-García, E.; Suardíaz, R.; Montero, L. A.; Sander, W. *Chem. Phys.* **2008**, *353*, 193–201.
- (16) Du, S.; Francisco, J. S.; Schenter, G. K.; Iordanov, T. D.; Garrett, B. C.; Dupuis, M.; Li, J. *J. Chem. Phys.* **2006**, *124*, 224318–15.
- (17) Galano, A.; Narciso-López, M.; Francisco-Marquez, M. *J. Phys. Chem. A* **2010**, *114*, 5796–5809.
- (18) Autrey, T.; Brown, A. K.; Camaioni, D. M.; Dupuis, M.; Foster, N. S.; Getty, A. *J. Am. Chem. Soc.* **2004**, *126*, 3680–3681.
- (19) Hamad, S.; Lago, S.; Mejías, J. A. *J. Phys. Chem. A* **2002**, *106*, 9104–9113.
- (20) Du, S.; Francisco, J. S.; Schenter, G. K.; Garret, B. C. *J. Chem. Phys.* **2008**, *128*, 084307–8.
- (21) Engdahl, A.; Karlström, G.; Nelander, B. *J. Chem. Phys.* **2003**, *118*, 7797–7802.
- (22) Kim, K. S.; Kim, H. S.; Jang, J. H.; Kim, H. S.; Mhin, B. J.; Xie, Y.; Schaefer, H. F. *J. Chem. Phys.* **1991**, *94*, 2057–2062.
- (23) Nanayakkara, A. A.; Balint-Kurti, G. G.; Williams, I. H. *J. Phys. Chem.* **1992**, *96*, 3662–3669.
- (24) Xie, Y.; Schaefer, H. F. *J. Chem. Phys.* **1993**, *98*, 8829–8834.
- (25) Wang, B.; Hou, H.; Gu, Y. *Chem. Phys. Lett.* **1999**, *303*, 96–100.
- (26) Zhou, Z.; Qu, Y.; Fu, A.; Du, B.; He, F.; Gao, H. *Int. J. Quantum Chem.* **2002**, *89*, 550–558.
- (27) Uchimaru, T.; Chandra, A. K.; Tsuzuki, S.; Sugie, M.; Sekiya, A. *J. Comput. Chem.* **2003**, *24*, 1538–1548.
- (28) Gonzalez, J.; Caballero, M.; Aguilar-Mogas, A.; Torrent-Sucarrat, M.; Crehuet, R.; Sole, A.; Gimenez, X.; Ollivella, S.; Bofill, J. M.; Anglada, J. M. *Theor. Chem. Acc.* **2011**, *128*, 579–592.
- (29) Vassilev, P.; Louwse, M. J.; Baerends, E. *J. Chem. Phys. Lett.* **2004**, *398*, 212–216.
- (30) Vassilev, P.; Louwse, M. J.; Baerends, E. *J. Phys. Chem. B* **2005**, *109*, 23605–23610.
- (31) Khalack, J. M.; Lyubartsev, A. P. *J. Phys. Chem. A* **2005**, *109*, 378–386.
- (32) Cabral do Couto, P.; Guedes, R. C.; Martinho-Simões, J. A.; Costa Carbra, B. *J. Chem. Phys.* **2003**, *119*, 7344.
- (33) VandeVondele, J.; Sprik, M. *Phys. Chem. Chem. Phys.* **2005**, *7*, 1363–1367.
- (34) Chalmet, S.; Ruiz-López, M. F. *J. Chem. Phys.* **2006**, *124*, 194502–6.
- (35) Tuckerman, M. E.; Marx, D.; Parrinello, M. *Nature* **2002**, *417*, 925–929.
- (36) Codorniu-Hernández, E.; Kusalik, P. G. Mobility mechanism of hydroxyl radicals in aqueous solution via hydrogen transfer. Manuscript submitted.
- (37) Iglev, H.; Fisher, M. K.; Gliserin, A.; Laubereau, A. *J. Am. Chem. Soc.* **2011**, *133*, 790–795.
- (38) Masgrau, L.; Gonzalez-Lafont, A.; Lluch, J. M. *J. Phys. Chem. A* **1999**, *103*, 1044–1053.
- (39) Car, R.; Parrinello, M. *Phys. Rev. Lett.* **1985**, *55*, 2471–2474.
- (40) CPMD; IBM Corp.: Armonk, NY, 2006; MPI für Festkörperforschung: Stuttgart, Germany, 2001.
- (41) Becke, A. *Phys. Rev. A* **1998**, *38*, 3098–3100.
- (42) Lee, C.; Yang, W.; Parr, R. *Phys. Rev. B* **1998**, *37*, 785–789.
- (43) Boese, A. D.; Doltsinis, N. L.; Handy, N. C.; Sprik, M. *J. Chem. Phys.* **2000**, *112*, 1670–1678.
- (44) Boese, A. D.; Martin, J. M. L. *J. Chem. Phys.* **2004**, *121*, 3405–3416.
- (45) Izvekov, S.; Voth, G. A. *J. Chem. Phys.* **2005**, *123*, 134105–13.
- (46) Troullier, N.; Martins, J. L. *Phys. Rev. B* **1991**, *43*, 1993–2006.
- (47) Goedecker, S.; Teter, M.; Hutter, J. *Phys. Rev. B* **1996**, *54*, 1703–1710.
- (48) Martyna, G. J.; Klein, M. L.; Tuckerman, M. *J. Chem. Phys.* **1992**, *97*, 2635.
- (49) Kusalik, P. G.; Svishchev, I. M. *Science* **1994**, *265*, 1219–1221.
- (50) Humphrey, W.; Dalke, A.; Schulten, K. *J. Mol. Graphics* **1996**, *14*, 33–38.
- (51) Sprik, M.; Ciccotti, G. *J. Chem. Phys.* **1998**, *109*, 7737–7745.
- (52) Campo, M. G.; Grigera, J. R. *J. Chem. Phys.* **2005**, *123*, 084507–084507–6.
- (53) Svishchev, I. M.; Plugatyr, A. Y. *J. Phys. Chem. B* **2005**, *109*, 4123–4128.
- (54) Pabis, A.; Szala-Bilnik, J.; Swiatla-Wojcik, D. *Phys. Chem. Chem. Phys.* **2011**, *13*, 9458–9468.



# Accurate Molecular Crystal Lattice Energies from a Fragment QM/MM Approach with On-the-Fly Ab Initio Force Field Parametrization

Shuhao Wen and Gregory J. O. Beran\*

Department of Chemistry, University of California, Riverside, California 92521, United States

Supporting Information

**ABSTRACT:** We combine quantum and classical mechanics in a fragment-based many-body interaction model to predict organic molecular crystal lattice energies. Individual molecules in the central unit cell and their short-range pairwise interactions are modeled quantum mechanically, while long-range pairwise and many-body interactions are approximated classically. The classical contributions are evaluated using an accurate ab initio force field that is constructed on-the-fly from quantum mechanical calculations on the individual molecules in the unit cell. The force field parameters include ab initio distributed multipole moments, distributed polarizabilities, and isotropic two- and three-body atomic dispersion coefficients. This QM/MM fragment model reproduces full periodic MP2 lattice energies to within a couple kJ/mol at substantially reduced cost. When high-level electronic structure methods are coupled with the ab initio force field, molecular crystal lattice energies are predicted to within 2 kJ/mol of their experimental values for six of the seven crystals examined here. Finally, Axilrod–Teller–Muto three-body dispersion energy plays a nontrivial role in several of the molecular crystals studied here.

## 1. INTRODUCTION

Organic molecular crystals play a fundamental role in pharmaceuticals, agrochemicals, pigments, dyestuffs, foods, explosives, and organic electronic materials. Molecular crystal properties are strongly affected by the crystal packing. Multiple crystal packing arrangements, or polymorphs, are often thermodynamically accessible in real crystals and can have major real-world consequences.<sup>1</sup>

For example, a change in the crystal packing of rubrene, a promising organic semiconductor material, utterly destroys its high charge-carrier mobility.<sup>2</sup> Or consider that the appearance of a low-solubility polymorph of ritonavir, an anti-HIV drug, forced its temporary removal from the market. This prevented patients from receiving treatment and cost its maker an estimated \$250 million in lost sales.<sup>3</sup> Clearly, substantial scientific and financial interest lies in knowing the structures and the properties of stable crystal polymorphs.

Over the past decade, substantial progress has been made toward the dream of predicting molecular crystal structures starting from only a single molecule, as evidenced by recent major improvements in the results of the blind crystal structure prediction tests.<sup>4,5</sup> Two advances instrumental to this progress are the development of robust, anisotropic force fields that include distributed multipolar expansions of the molecular charge distribution, induction, and dispersion<sup>6,7</sup> and the application of quantum mechanical models to crystal structure prediction, either to help determine intramolecular conformations or for fully periodic density functional theory (DFT) calculations.<sup>6,8–16</sup>

Unfortunately, traditional density functionals suffer from well-known difficulties in describing van der Waals dispersion interactions,<sup>17,18</sup> which are critical to modeling intermolecular interactions in molecular crystals. A number of strategies to correct this deficiency have been adopted, ranging from empirical corrections to the development of nonlocal density functions.<sup>19–23</sup>

These techniques often work extremely well for molecular crystals, but counter-examples also exist.<sup>24,25</sup>

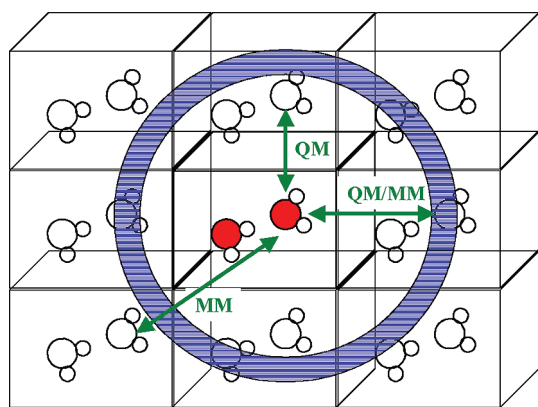
Moreover, the energy spacing between crystal polymorphs can be as little as 1 kJ/mol or less.<sup>24,26</sup> Differences among the lattice energy predictions from various density functionals often exceed this threshold. The absence of a clear strategy for systematically improving DFT calculations makes robust predictions difficult. Recent developments in periodic Møller–Plesset perturbation theory (MP2) are also very promising for crystal structure modeling,<sup>27–32</sup> but those calculations remain relatively computationally expensive.

The past few years have seen considerable interest in strategies that model molecular crystals through a hierarchical<sup>33,34</sup> or fragment-based scheme. The advantage of fragment-based models is that one can systematically improve the quality of the electronic structure method used to describe the fragments and their interactions. Many of these models, such as the fragment molecular orbital method,<sup>35</sup> are based on the many-body interaction expansion. The key distinguishing features between such methods lie in how they handle the long-range two-body (interactions between a pair of molecules) and the many-body (involving three molecules or more) contributions.

The most straightforward approach would be to simply neglect these terms, but they contribute too much to ignore. Long-range and many-body contributions typically account for ~5–10% of the lattice energy, but they can contribute as much as ~25%! One can do better if three-body terms are included explicitly, as has been demonstrated with highly accurate symmetry-adapted perturbation theory (SAPT) calculations, for example.<sup>36,37</sup> Unfortunately, the steep computational scaling of most electronic structure methods makes the explicit inclusion of three-body terms costly.

Received: August 2, 2011

Published: October 20, 2011



**Figure 1.** Pictorial representation for the treatment of two-body terms in periodic HMBI. Short-range dimers interactions are modeled quantum mechanically, while long-range ones are treated classically (MM). A linear combination of QM and MM is used in the blue region to transition smoothly between the two regimes.

The fragment molecular orbital method and related approaches incorporate many-body electrostatic induction effects via the use of an embedding potential in the one- and two-body terms.<sup>38–42</sup> However, this embedding potential complicates the evaluation of the nuclear derivatives required for structure optimization.<sup>39,43</sup> Alternatively, one can approximate the long-range/many-body terms in some fashion. A number of groups have used Hartree–Fock (HF) or DFT to capture these effects.<sup>44–48</sup> Both methods capture the many-body induction effectively, but they become computationally expensive for large unit cells. Both of these approaches have traditionally omitted many-body dispersion effects, which are sometimes important.

In our approach, quantum mechanics (QM) is used to treat the short-range interactions, while a polarizable force field (MM) is used to approximate the long-range two- and many-body interactions.<sup>25,49,50</sup> This hybrid QM/MM many-body interaction (HMBI) model differs from conventional QM/MM models in that it partitions different classes of interactions as either QM or MM based on their importance in the many-body interaction expansion, rather than by defining specific QM and MM regions of space.

We have recently demonstrated that this hybrid model enables the prediction of several small-molecule crystal lattice energies to within 4–5 kJ/mol, so-called chemical accuracy.<sup>25</sup> However, chemical accuracy is probably insufficient to discriminate among crystal polymorphs separated by only a kJ/mol or less. That earlier work used the Amoeba polarizable force field for the MM portion of the model. In this paper, we demonstrate that even better results are obtained when we replace the Amoeba force-field with a high-quality ab initio force field whose parameters are calculated “on-the-fly” via separate electronic structure calculations performed for each molecule in the crystal unit cell.

Like our earlier work on molecular clusters,<sup>50</sup> this force field includes electrostatic and induction effects based on distributed multipole moments and polarizabilities. Here, we augment those terms with atomic dispersion coefficients to describe long-range two-body dispersion and Axilrod–Teller–Muto three-body dispersion. We also implement an Ewald summation-based treatment of multipolar electrostatics and induction for the periodic crystals (see also Supporting Information).

This force field model is analogous to those used in high-quality MM crystal modeling (e.g., refs 5, 7, and 51), with the

force field parameters recalculated for each molecule/geometry to capture the variations in the properties (particularly multipole moments<sup>5,50,51</sup>) with geometry. No rigid monomer approximation is needed. Similar long-range terms are included in the “systematic fragmentation” model,<sup>34</sup> though the model described here differs in many details, including its use of a multipolar Ewald sum for long-range electrostatics, distributed polarizabilities, and the inclusion of three-body dispersion.

We demonstrate that this hybrid QM/MM approach reproduces periodic, fully quantum mechanical calculations to within a couple kJ/mol. More importantly, combining this ab initio force field with high-level electronic structure calculations reproduces experimental crystal lattice energies to within 2 kJ/mol for most of the crystals examined here. In other words, these predictions lie within the typical experimental error bars for molecular crystal lattice energies. Finally, we observe that the three-body Axilrod–Teller–Muto dispersion contribution is surprisingly important, even in some crystals where many-body induction would normally be expected to dominate.

## 2. THEORY

**2.1. Periodic Hybrid Many-Body Interaction Model.** The details of our hybrid QM/MM fragment approach for both clusters<sup>49,50</sup> and periodic systems<sup>25</sup> have been given previously, so we provide only a brief summary here. This fragment method decomposes a system into interacting molecules using a many-body interaction expansion. The intramolecular interactions and the most important intermolecular interactions are modeled quantum mechanically, while weaker intermolecular interactions are approximated classically. Specifically, for an infinite periodic molecular crystal, individual molecules in the central unit cell and shorter-range pairwise interactions are modeled quantum mechanically, while longer-range two-body interactions and all many-body interactions are treated using a polarizable force field (see Figure 1):

$$E_{\text{PBC}}^{\text{HMBI}} = E_{\text{PBC}}^{\text{MM}} + \sum_i (E_i^{\text{QM}} - E_i^{\text{MM}}) + \sum_{ij} d_{ij}^{\text{smooth}} (\Delta^2 E_{ij}^{\text{QM}} - \Delta^2 E_{ij}^{\text{MM}}) + \frac{1}{2} \sum_i \sum_{\vec{k}}^{\text{images}} d_{i\vec{k}}^{\text{smooth}} (\Delta^2 E_{i\vec{k}}^{\text{QM}} - \Delta^2 E_{i\vec{k}}^{\text{MM}}) \quad (1)$$

Here,  $i$  and  $j$  run over molecules in the central unit cell, while  $\vec{k}$  runs over all periodic image molecules within some cutoff distance of molecule  $i$ ,  $E_i$  corresponds to the energy of monomer  $i$ , and  $\Delta^2 E_{ij}$  is the interaction energy between monomers  $i$  and  $j$ . Both can be calculated either quantum mechanically (QM) or with a force field (MM).  $E_{\text{PBC}}^{\text{MM}}$  refers to the force field energy of the entire periodic crystal. To ensure smooth and continuous potential energy surfaces, the transition from short-range quantum to long-range classical treatments is spread over a finite region using a smoothing function  $d_{ij}^{\text{smooth}}$  that decays from 1 at radius  $r_1$  to 0 at radius  $r_0$ :<sup>52</sup>

$$d_{ij}^{\text{smooth}}(R) = \begin{cases} 1 & \text{if } x \leq r_1 \\ \frac{1}{1 + e^{2|r_1 - r_0|/(r_1 - R)} - |r_1 - r_0|/(R - r_0)} & \text{if } r_1 < x < r_0 \\ 0 & \text{if } x \geq r_0 \end{cases} \quad (2)$$

where  $R$  is the shortest intermolecular distance between any two atoms in the pair of molecules  $i$  and  $j$ . For any dimer where the shortest intermolecular separation is less than or equal to  $r_1$ , the two-body interaction is treated quantum mechanically ( $d^{\text{smooth}} = 1$ ). If the shortest intermolecular distance is greater than or equal to  $r_0$ , it is approximated classically ( $d^{\text{smooth}} = 0$ ). For dimers whose shortest intermolecular separation lies within the damping region between  $r_1$  and  $r_0$ , the dimer interaction energy is a linear combination of quantum and classical interactions ( $0 < d^{\text{smooth}} < 1$ ).

**2.2. Nature of the Ab Initio Force Field (AIFF) in Periodic Systems.** The success of this fragment QM/MM approach depends critically on the quality of the force field used to approximate the long-range two- and the many-body intermolecular interactions. The force field used here includes long-range two-body electrostatics, induction (both many-body and long-range two-body), long-range two-body dispersion, and three-body dispersion:

$$E^{\text{MM}} = E_{\text{es}} + E_{\text{ind}} + E_{2\text{-body disp}} + E_{3\text{-body disp}} \quad (3)$$

Note that the addition of dispersion terms and the incorporation of periodic boundary conditions distinguish this force field from an earlier version.<sup>50</sup>

The force field is parametrized with atom-centered distributed multipole moments, distributed static polarizabilities, and isotropic atomic dispersion coefficients. The isotropic dispersion coefficients are computed from the isotropic frequency-dependent polarizabilities. The multipole moments and polarizabilities are represented in a spherical tensor formalism and can be computed for each monomer in the unit cell.<sup>53</sup> The computational time required to determine these parameters is typically small compared to the time required to evaluate the QM interactions in the system.

The force field also requires short-range induction and dispersion damping function parameters which are unfortunately more difficult to obtain from first principles. As described below, the damping parameters are obtained empirically. The following sections describe each of the force-field terms in greater detail.

The following notation is used below: The letters A, B, and C refer to molecules, while  $a$ ,  $b$ , and  $c$  refer to atoms in those molecules. The letters  $t$  and  $u$  refer to spherical tensor components of the multipole moments/polarizabilities.

**2.2.1. Long-Range Two-Body Electrostatics.** The force fields adopts a distributed multipole representation of the interacting molecular charge densities.<sup>54–56</sup> Heavy atom densities are represented with a rank 4 expansion (up to hexadecapole moments), while hydrogen atoms include up to rank 2 (quadrupole moments). As described in ref 53, the interaction between two molecules A and B is given by

$$E_{\text{es}}^{\text{AB}} = \sum_{a \in \text{A}} \sum_{b \in \text{B}} \sum_{tu} Q_t^a T_{tu}^{ab} Q_u^b \quad (4)$$

where  $Q_t^a$  represents the  $t$ -th multipole moment component on atom  $a$ , and  $T_{tu}^{ab}$  contains the distance and orientation dependence of the interaction (the multipole moments are generally anisotropic and are typically represented in a local molecular coordinate system). Using real spherical tensors,  $t$  and  $u$  run over the 25 rank 4 components: charge ( $t = 00$ ), dipole ( $t = 10, 11c, 11s$ ), quadrupole ( $t = 20, 21c, 21s, 22c, 22s$ ), octopole ( $t = 30, 31c, 31s, 32c, 32s, 33c, 33s$ ), and hexadecapole ( $t = 40, 41c, 41s, 42c, 42s, 43c, 43s, 44c, 44s$ ) moments. Electrostatic interactions up to  $R^{-5}$  are included in eq 4.

Calculating the electrostatics for an infinite crystal formally requires a lattice summation between each of the central unit cell molecules A and all other molecules B, including an infinite number of periodic images:

$$E_{\text{es}}^{\text{lattice}} = \sum_{\text{A}} \sum_{\text{B} \neq \text{A}} E_{\text{es}}^{\text{AB}} \quad (5)$$

We evaluate this expression via multipolar Ewald summation, drawing heavily from ref 57. The resulting total two-body lattice interaction energy is given by

$$\begin{aligned} E_{\text{es}}^{\text{lattice}} = & \sum_{\text{A}} \sum_{\text{B} \neq \text{A}} \sum_{\text{N}} \sum_{ab} \sum_{tu} Q_t^a (\mathcal{F}_{tu}^{ab} + \tilde{\mathcal{F}}_{tu}^{ab}) Q_u^b \\ & - \frac{\gamma}{\sqrt{\pi}} \sum_{\text{A}} \sum_a (Q_{00}^a)^2 - \sum_{\text{A}} \sum_a \sum_{a' \neq a} \sum_{tu} Q_t^a T_{tu}^{aa'} Q_u^{a'} \\ & + \sum_{\text{A}} \sum_{\text{B} \neq \text{A}} \sum_{ab} \sum_{t+u=2} Q_t^a \mathcal{F}_{tu}^{ab} Q_u^b \end{aligned} \quad (6)$$

where  $N$  refers to the image cell index in the Ewald summation. In the first term of this equation, the  $\mathcal{F}_{tu}^{ab}$  and  $\tilde{\mathcal{F}}_{tu}^{ab}$  are the interaction functions in direct and reciprocal space, respectively.  $\mathcal{F}_{tu}^{ab}$  and  $\tilde{\mathcal{F}}_{tu}^{ab}$  are analogous to the  $T_{tu}^{ab}$  terms in eq 4 but with extra components arising from the Ewald summation. They include the orientation dependence between site–site vectors and lattice vectors in direct and reciprocal space, the site–site distance dependence, and the coefficient that controls the length scales in the direct and reciprocal space portions Ewald summation. Explicit expressions for  $\mathcal{F}_{tu}^{ab}$  and  $\tilde{\mathcal{F}}_{tu}^{ab}$  are given in the Supporting Information.

The first term in eq 6 gives the basic Ewald summation in direct and reciprocal space. However, the Ewald method introduces a self-interaction energy (i.e., the interaction of an atomic site with itself,  $\text{A} = \text{B}$  and  $a = b$ ), which is explicitly subtracted out by the second term in eq 6. Only the charge–charge self-interaction term needs to be corrected for in a spherical harmonic formulation.<sup>57</sup> The Ewald sum here also includes terms corresponding to interactions between pairs of atoms within a single molecule. These unwanted intramolecular electrostatic terms are eliminated by the third term. Finally, the fourth term corresponds to a boundary condition term for interactions with total multipole moment of two (dipole–dipole and charge–quadrupole interactions). Boundary condition terms with total multipole moment less than two, which correspond to a crystal shape-dependent surface contribution in polar/ionic crystals, are omitted in this formulation (i.e., tinfoil boundary conditions are adopted). See ref 57 and references cited therein for details. In practice, we perform the Ewald sum on the total multipole moments (permanent plus induced), as described below.

**2.2.2. Induction.** Many-body induction can also be important in determining molecular crystal structures and energetics, especially when the crystals contain polar molecules and/or hydrogen bonds.<sup>58</sup> For this reason, the force field includes self-consistent induction for both long-range two-body interactions (relatively unimportant) and many-body interactions (important).

Multipole moments on nearby molecules B induce multipole moments  $\Delta Q_t^a$  on atom  $a$  in molecule A:

$$\Delta Q_t^a = - \sum_{\text{B} \neq \text{A}} \sum_b \sum_{a'} \sum_{t'u} \alpha_{tt'}^{aa'} f_n(R, \beta) T_{t'u}^{a'b} (Q_u^b + \Delta Q_u^b) \quad (7)$$

where  $\alpha_{tt'}^{aa'}$  is the polarizability of atom  $a$  and  $f_n(R, \beta)$  is a short-range electrostatic damping function. We typically compute

atom-centered distributed polarizabilities up to rank 2 (quadrupole–quadrupole) on heavy atoms and up to rank 1 (dipole–dipole) on hydrogen atoms.<sup>59,60</sup> The distributed polarizabilities are computed according to the Williams–Stone–Misquitta procedure.<sup>61,62</sup>

Of course, the multipole moments on molecule A also induce multipole moments  $\Delta Q_u^b$  on the atoms in molecule B. Thus, eq 7 must be iterated to self-consistency on all atoms. We iterate the induced multipole moments until the energy converges to  $10^{-5}$  kJ/mol. The final induction energy, which includes many-body induction, is given by

$$E_{\text{ind}} = \frac{1}{2} \sum_A \sum_{B \neq A} \Delta Q_t^a f_n(R, \beta) T_{t'u}^{a'b} Q_u^b \quad (8)$$

We apply the Tang–Toennies damping function  $f_n(R, \beta)$ :

$$f_n(R, \beta) = 1 - \sum_{k=0}^n \left( \frac{(\beta R)^k}{k!} \right) e^{-\beta R} \quad (9)$$

The subscript  $n$  in  $f_n(R, \beta)$  refers to the order of the electrostatic interaction  $R^{-n}$  in  $T_{t'u}^{a'b}$ . The constant  $\beta$  is determined empirically for each type of molecule, as described previously.<sup>50</sup>

Generalizing this treatment of many-body induction to infinite periodic systems with high-order multipoles is complicated by two key issues. First, we need to determine the self-consistent induced moments in the context of an infinite lattice. In principle, the iteration of the induced moments to self-consistency could be coupled with the Ewald summation for the permanent electrostatics. Although that is conceptually straightforward, the complexity of the multipolar Ewald summation makes it messy in practice. Furthermore, it is unnecessary: distant molecules do not contribute significantly to the induced multipole moments in the central unit cell. They will, however, induce multipole moments on other molecules which are closer to the central unit cell and therefore interact in a many-body fashion.

Therefore, we evaluate the induced multipole moments in a finite cluster that is large enough to capture these effects. In practice, we find that including all molecules within 25 Å of the central unit cell molecules converges the induced multipole moments to within 0.001 au. At each iteration, we evaluate only the induced moments on the molecules in the central unit cell. Induction effects between groups of molecules outside the central unit cell are not included. Rather, the induced moments on the periodic image molecules are then set equal to those of the reference central cell molecules, mimicking the infinite crystal. This process is repeated until the induced moments reach self-consistency.

The use of a finite cluster introduces slight asymmetries in the induced moments on symmetry-equivalent atoms, with the induced multipole moments typically varying by a few percent or less. The larger the finite cluster, the smaller the errors. In principle, these errors could be eliminated entirely through a proper treatment of space group symmetry, though we do not do so here.

Second, the induction interactions must be damped at short-range to avoid the “polarization catastrophe,” particularly when evaluating the many-body induction terms. Damping is trivial to apply when determining the multipole moments in a finite cluster, but again it complicates the Ewald summation equations. Here, we include it by recognizing that damping is important only within short ranges (<10 Å). We perform the Ewald sum using undamped interaction energies and then correct the resulting induction energy with the difference between the damped and undamped interactions in the finite cluster used above to

determine the induced multipole moments. As long as the finite cluster is larger than the length scale on which the damping function operates, this approach introduces no additional errors.

To summarize, this approach for computing the induction energy in infinite periodic systems has four steps that can be implemented easily:

1. Determine the self-consistent induced multipole moments on all atoms in the central unit cell by interacting them with a finite number of periodic image molecules (a “cluster”). Short-range damping is applied while iterating to self-consistency. Compute the damped induction energy for the central unit-cell molecules interacting with this finite cluster,  $E_{\text{ind}}^{\text{cluster}}(\text{damped})$ .
2. Use the converged induced multipole moments from step 1 to compute the induction energy without short-range damping,  $E_{\text{ind}}^{\text{cluster}}(\text{undamped})$  (i.e., eq 8 with  $f_n = 1$ ). Compute the correction due to short-range damping  $\delta E_{\text{damp}}^{\text{cluster}} = E_{\text{ind}}^{\text{cluster}}(\text{damped}) - E_{\text{ind}}^{\text{cluster}}(\text{undamped})$ .
3. Replace the permanent multipole moments in eq 6 with the total multipole moments (permanent plus induced) and evaluate the total lattice energy,  $E_{\text{es+ind}}^{\text{lattice}}(\text{undamped})$ .
4. Correct the total lattice energy for short-range induction damping:

$$E_{\text{es+ind}}^{\text{lattice}}(\text{damped}) = E_{\text{es+ind}}^{\text{lattice}}(\text{undamped}) + \delta E_{\text{damp}}^{\text{cluster}} \quad (10)$$

Our strategy differs moderately from the one in ref 7, but the two approaches probably give similar results. The computational cost of this approach is small compared to the cost of the quantum mechanical calculations in the hybrid fragment model.

**2.2.3. Long-Range Two-Body Dispersion.** Two-body van der Waals dispersion makes an important contribution to molecular crystals, but the bulk of this interaction is typically captured in the QM portion of the fragment model. Nevertheless, we include two-body dispersion in the force-field to capture the long-range contributions that are missed in the short-range QM treatment.

Two-body dispersion is included in the force-field using isotropic atomic  $C_6$  and  $C_8$  dispersion coefficients:

$$E_{2\text{-body disp}} = - \sum_A \sum_{B \neq A} \sum_a \sum_b \left( f_6 \frac{C_6^{ab}}{R_{ab}^6} + f_8 \frac{C_8^{ab}}{R_{ab}^8} + \dots \right) \quad (11)$$

where the  $f_n$  are again Tang–Toennies damping functions (eq 9). The damping parameter  $\beta$  in  $f_n(R, \beta)$  is determined empirically from atomic van der Waals radii.<sup>63</sup> In this case, the damping function is fairly unimportant, since the force field is used only to describe long-range dispersion. The lattice sum is evaluated explicitly with a large cutoff (e.g., 15 or 20 Å).

The dispersion coefficients  $C_n^{ab}$  for atoms  $a$  and  $b$  (in atomic units) are determined via Casimir–Polder integration over the isotropic frequency-dependent polarizabilities:<sup>53</sup>

$$C_6^{ab} = \frac{3}{\pi} \int_0^\infty \alpha_{11}^a(iv) \alpha_{11}^b(iv) dv \quad (12)$$

$$C_8^{ab} = \frac{15}{2\pi} \left( \int_0^\infty \alpha_{11}^a(iv) \alpha_{22}^b(iv) dv + \int_0^\infty \alpha_{22}^a(iv) \alpha_{11}^b(iv) dv \right) \quad (13)$$

where  $\alpha_{11}$  and  $\alpha_{22}$  are the isotropic dipole–dipole and quadrupole–quadrupole frequency-dependent polarizabilities, respectively. Note that for a given geometry, the frequency-dependent polarizabilities need only be determined once for each atom in the unit cell. The dispersion coefficients and interaction energy between any pair of atoms can then be computed quickly by performing the one-dimensional integral over imaginary frequency.

The Casimir–Polder integral is evaluated via 10-point Gauss–Legendre quadrature after performing a change of variables that maps  $v$  to  $t$  according to  $v = v_0(1+t)/(1-t)$ , with  $v_0 = 0.5$ .<sup>64</sup> This transformation converts the semi-infinite integral to one between  $-1$  and  $1$ .

**2.2.4. Three-Body Dispersion.** Many-body dispersion is usually expected to be small compared to other intermolecular interactions, and it is often neglected in molecular crystal calculations. However, the leading many-body contribution, three-body Axilrod–Teller–Muto dispersion, makes a significant contribution in crystals containing nonpolar molecules, such as in benzene<sup>65</sup> or rare gases.<sup>66</sup> As we demonstrate below, it can also contribute nontrivially to the lattice energy of even some polar, hydrogen-bonded molecular crystals. The magnitude of the three-body dispersion contribution depends strongly on the orientation of the interacting bodies,<sup>65</sup> and it can be important for discriminating between putative crystal polymorphs.<sup>63</sup>

The AIFF incorporates the Axilrod–Teller–Muto triple-dipole three-body intermolecular dispersion term.<sup>67,68</sup> For a given set of three molecules ABC, this is given by

$$E_{3\text{-body disp}}^{\text{ABC}} = \sum_{a \in A} \sum_{b \in B} \sum_{c \in C} f_9 C_9^{abc} \frac{(1 + 3 \cos \hat{a} \cos \hat{b} \cos \hat{c})}{R_{ab}^3 R_{bc}^3 R_{ac}^3} \quad (14)$$

where  $C_9^{abc}$  is the dispersion coefficient for atom triplet  $abc$ ,  $R_{ij}$  is the distance between atoms  $i$  and  $j$ , and  $\hat{a}$ ,  $\hat{b}$ , and  $\hat{c}$  are the angles of the triangle formed by the three atoms. The damping function  $f_9$  is written as a product of three two-body Tang–Toennies damping functions.<sup>63,69,70</sup>

A similar atom–atom triple-dipole dispersion formulation for molecules has been used, for example, by von Lilienfeld and Tkatchenko.<sup>63</sup> They demonstrated that it reproduces SAPT three-body dispersion energies fairly well. Our implementation differs from theirs primarily in how the  $C_9$  coefficients are obtained. Three-body dispersion corrections based on coupled Kohn–Sham theory, such as the approach used here, are known to predict the asymptotic dispersion contributions accurately.<sup>37,69,71</sup>

The total three-body dispersion contribution of the lattice is given by summing eq 14 over all possible triplets of molecules. Only one of these molecules needs lie in the central unit cell. The other two may either reside in the unit cell or be periodic image molecules. We perform the lattice sum explicitly up to a user-defined cutoff (e.g., 10 Å). In practice, most of the contribution comes from cases with one molecule in the unit cell and the other two outside it. Of course, the details vary with the number and the chemical nature of the molecules in the unit cell.

In the formulation used here, the three atoms  $a$ ,  $b$ , and  $c$  lie on different molecules. This means that only the intermolecular three-body dispersion contribution is included. Important interatomic three-body dispersion contributions involving only one or two molecules need to be captured by the QM portion of the model.

The  $C_9$  coefficient can be calculated using Casimir–Polder integration:

$$C_9^{abc} = \frac{3}{\pi} \int_0^\infty \alpha_{11}^a(iv) \alpha_{11}^b(iv) \alpha_{11}^c(iv) dv \quad (15)$$

or it can be estimated from two-body  $C_6$ :

$$C_9^{abc} \approx \frac{2S^a S^b S^c (S^a + S^b + S^c)}{(S^a + S^b)(S^b + S^c)(S^c + S^a)} \quad (16)$$

where  $S^a = C_6^{aa}((\alpha_{11}^b(0)\alpha_{11}^c(0))/(\alpha_{11}^a(0)))$  and  $\alpha_{11}(0)$  is static dipole–dipole polarizabilities. We tested both approaches on a handful of systems and found that the  $C_9$  coefficients estimated via eq 16 are typically  $\sim 5\%$  larger than those computed from eq 15. Since we already have the frequency-dependent polarizabilities, we adopt eq 15.

This formalism assumes that any important interatomic three-body dispersion terms involving only one or two monomers (between two atoms on one monomer and one on another, for example) are handled in the quantum mechanical two-body interaction terms. Of course, three-body dispersion terms are *not* captured at the MP2 level,<sup>72</sup> in which case one might consider including additional terms.

**2.2.5. Terms Not Included in the Force Field.** Before discussing the results, it is worth considering some of the terms that are not included in the force field. First, in the two-body case, the force field does not include exchange/repulsion, penetration, and charge-transfer effects. These effects are primarily short-range in nature, so they are handled in the quantum mechanical portion of the model. The absence of these terms from the force-field is therefore expected to be unimportant.

At the many-body level, we only explicitly include some of the most important terms. Many-body induction is performed self-consistently, but many-body dispersion is only approximated via the leading Axilrod–Teller–Muto term. The use of distributed multipolar expansions and asymptotic dispersion interaction formulations is potentially problematic at short ranges, which is where these terms are most important. As described above, empirical damping factors are needed to compensate at short-range. Furthermore, various other many-body terms (exchange, exchange–induction, induction–dispersion, etc.) are not included explicitly.

Despite these issues, the model performs very well, as we demonstrate below. We attribute this success to a combination of fortuitous error cancellation and to the use of empirical damping factors. In particular, the induction damping parameter for each type of intermolecular interaction is determined by fitting the AIFF induction contribution to the MP2 many-body contribution in a small test set of trimers (none of which are taken from the crystal). So it really includes some effect of the other three-body contributions present in MP2.

Finally, we note that one might circumvent some of these issues by using a hybrid of explicit three-body QM calculations (using supermolecular or SAPT approaches) theory at short ranges and the asymptotic expressions at longer ranges.<sup>37</sup> Of course, including any quantum mechanical treatment of trimers would substantially increase the computational cost of the model, especially for larger molecules.

### 3. COMPUTATIONAL DETAILS

We perform lattice energy calculations for seven different molecular crystals (ice, formamide, acetamide, imidazole, benzene,

ammonia, and carbon dioxide) using a procedure that is analogous to the one used in ref 25. For the first five crystals, the geometries are identical to those used in ref 25. They were optimized under the constraint of frozen experimental lattice parameters with the HMBI method using RI-MP2/aug-cc-pVDZ and the Amoeba force field. The geometries for ammonia and carbon dioxide come from ref 30.

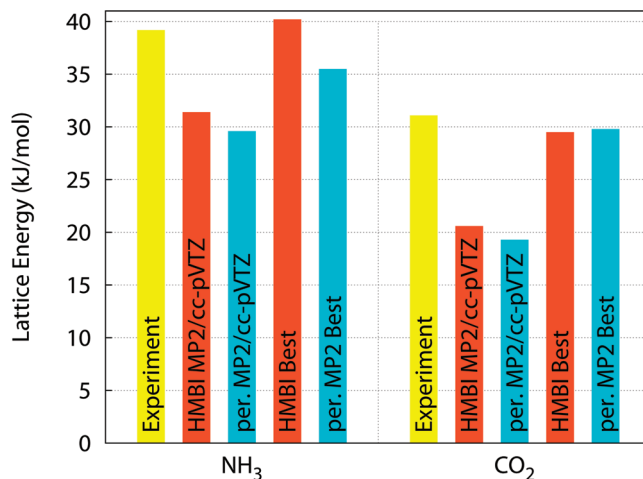
The quantum mechanical calculations were performed using counterpoise-corrected dual basis RI-MP2<sup>73</sup> and the Dunning aug-cc-pVDZ, aug-cc-pVTZ, and aug-cc-pVQZ basis sets<sup>74</sup> for the QM part. Energies at the complete basis set limit were estimated by separately extrapolating the HF and MP2 correlation triple and quadruple- $\zeta$  basis results.<sup>75,76</sup> Then, a post-MP2 correction was computed using CCSD(T) and a smaller, practical basis set (typically aug-cc-pVDZ, except aug-cc-pVTZ for ice and 6-31+G\*<sup>77,78</sup> for acetamide and benzene). All quantum calculations were performed using a development version of Q-Chem 3.1,<sup>79</sup> except for the CCSD(T) ones, which were performed with PSI3.<sup>80</sup>

For the AIFF, the distributed multipole moments and polarizabilities (both isotropic and anisotropic) were calculated with CamCASP<sup>81</sup> using asymptotically corrected PBE and the Sadlej basis set.<sup>82,83</sup> A single induction damping factor was determined empirically for each unique molecule type by optimizing the many-body induction against MP2 many-body induction in a set of 10 trimers at various configurations, none of which were taken from the crystal structure. The induction damping factors used here are: 1.45 (ice and carbon dioxide), 1.40 (formamide and benzene), and 1.35 (acetamide, ammonia, and imidazole) bohr<sup>-1</sup>. Further details of this procedure can be found in ref 50.

As in our previous work,<sup>25</sup> the smoothing region that transitions from QM to MM is conservatively set at 9–10 Å, except for water, for which 6–7 Å can be used. Shorter cutoffs may be feasible with the improved force field, but we have not investigated that.

Because each fragment is defined as a single molecule, specification of the crystal for the calculation is straightforward with our software. The geometries are specified in Cartesian coordinates, with atoms grouped by molecule. Separate sections of the input file define the lattice parameters, the QM job parameters, the AIFF property calculation parameters, and the various AIFF force field cutoffs, etc. Our software then automatically creates input files for each job (e.g., Q-Chem and CamCASP), distributes and runs the jobs across a user-defined number of parallel processors, collects the results, evaluates the AIFF contributions, and finally computes the HMBI energy according to eq 1.

Finally, for four of the crystals, we estimated the effect of relaxing the experimental lattice parameters, using the same technique we adopted previously.<sup>25</sup> In particular, we generated a one-dimensional potential energy scan by isotropically scaling the lattice parameters  $a$ ,  $b$ , and  $c$  in increments of 1%. For each set of lattice parameters, the atoms in the unit cell were optimized with planewave DFT, as described previously. HMBI single-point energies were computed with dual-basis RI-MP2/aug-cc-pVTZ and the AIFF at each point, and a cubic spline was used to estimate the optimal lattice parameters and the change in the lattice energy. This  $\Delta E_{\text{lattice}}^{\text{relax}}$  contribution is added to the calculated lattice energy to obtain our best estimate. As noted previously, DFT cost and convergence issues prevented us from applying this procedure to acetamide. For ammonia and carbon dioxide, a similar procedure was already used to determine the geometries, so we did not repeat it here.<sup>30</sup> Implementation of analytic HMBI



**Figure 2.** Comparison of HMBI with fully QM periodic MP2 on ammonia and carbon dioxide crystals. Results are presented for both models in the cc-pVTZ basis. We also present our best prediction (complete-basis MP2 +  $\Delta^{\text{CCSD(T)}}$ ) and the largest basis MP2 results from ref 30.

lattice gradients is in progress, so we hope to fully optimize the structures in the future.

## 4. RESULTS AND DISCUSSION

The performance of the HMBI model for molecular crystals will be evaluated in two ways: First, to gauge the quality of the AIFF approximation, we examine how faithfully the QM/MM approach used here reproduces fully QM results. Specifically, we compare with benchmark periodic local MP2 lattice energy predictions for the ammonia and carbon dioxide crystals. Second, we determine how accurately molecular crystal lattice energies can be predicted compared to experiment for seven different molecular crystals. Finally, we decompose the different force field contributions to identify the important interactions. We find that three-body dispersion interactions are surprisingly important in a number of cases, even in some hydrogen-bonded molecular crystals where induction would be expected to dominate.

### 4.1. Comparison with Local Periodic MP2 Lattice Energies.

As mentioned in the Introduction Section, periodic local MP2 calculations on small-molecule organic crystals are now feasible. Those calculations enable the benchmarking of the HMBI approach with a fully QM treatment. We examine two molecular crystals for which large-basis periodic MP2 results exist: ammonia and carbon dioxide.<sup>30</sup>

To match the results of ref 30 as closely as possible, we performed HMBI calculations using local TRIM-MP2<sup>84</sup> with the identical cc-pVTZ basis set and crystal structure. Their periodic MP2 calculations use a different (Saebø–Pulay style)<sup>85</sup> local MP2 approximation, but both models should perform similarly for individual intermolecular interactions. For the purposes of this comparison, we also omit the AIFF three-body dispersion terms, which are not present in MP2 (they first appear in MP3).<sup>72</sup>

As shown in Figure 2, the HMBI model reproduces the full MP2 results well: The predicted lattice energies for NH<sub>3</sub> and CO<sub>2</sub> differ by only 1.8 and 1.3 kJ/mol, respectively from the periodic MP2 results. The error introduced into the predicted lattice energies by these two crystals by the HMBI fragment approach is similar to or smaller than the difference between canonical and local MP2 in the QM portion of our model!

Table 1. HMBI-Predicted Crystal Lattice Energies (kJ/mol)

QM level	ice	formamide	acetamide	imidazole	benzene	NH <sub>3</sub>	CO <sub>2</sub>
DB-RI-MP2/aug-cc-pVDZ	52.8	70.1	72.2	96.4	60.7	33.4	22.1
DB-RI-MP2/aug-cc-pVTZ	56.7	74.9	76.6	100.2	60.6	37.2	26.1
DB-RI-MP2/aug-cc-pVQZ	58.3	76.7	78.4	100.8	62.8	38.4	27.9
DB-RI-MP2/CBS	59.9	78.6	79.8	102.8	61.6	39.3	29.1
$\Delta^{\text{CCSD(T)a}}$	0.4	1.8	-0.1	-14.2	-10.4	0.9	0.3
DB-RI-MP2/CBS + $\Delta^{\text{CCSD(T)}}$	60.2	80.4	79.7	88.6	51.2	40.2	29.5
est. lattice param. relax., $\Delta E_{\text{lattice}}^{\text{relax}}$	0.2	0.0		2.2	2.8		
est. change in lattice parameters	-0.9%	-0.3%		-2.6%	-3.4%		
best estimate <sup>b</sup>	60.4	80.4	79.7	90.8	54.0	40.2	29.5
experiment	59	82 ± 0.3 <sup>c</sup>	86 ± 2 <sup>c</sup>	91 ± 4 <sup>c</sup>	52 ± 3 <sup>c</sup>	39 <sup>d</sup>	31 <sup>e</sup>

<sup>a</sup> Post-MP2 correction,  $\Delta^{\text{CCSD(T)}} = E_{\text{lattice}}^{\text{CCSD(T)}} - E_{\text{lattice}}^{\text{MP2}}$ , using the basis sets described in the text. <sup>b</sup> Best estimate =  $E^{\text{DB-RI-MP2/CBS}} + \Delta^{\text{CCSD(T)}} + \Delta E_{\text{lattice}}^{\text{relax}}$ . <sup>c</sup> Reported errors are the standard deviation among the set of extrapolated 0 K lattice energies. Actual experimental errors may be larger. See ref 25. <sup>d</sup> See Supporting Information. <sup>e</sup> From ref 30.

We can also compare our predicted HF lattice energy in the same basis with the fully periodic HF value. In this case, we omit all dispersion terms from the AIFF and reoptimized the AIFF induction damping factor using a set of trimer many-body energies at the HF level instead of the MP2 level. The latter step is unimportant for carbon dioxide, which exhibits minimal many-body induction, but the smaller parameter of  $\beta = 1.20 \text{ bohr}^{-1}$  improves the HF prediction by almost 1 kJ/mol.

Compared with periodic HF, HMBI overestimates the ammonia lattice energy by 1.7 kJ/mol, while the carbon dioxide lattice energies are identical.<sup>30</sup> The fact that the AIFF induction damping factor differs between HF and MP2 supports the idea mentioned in Section 2.2.5 that the empirical damping factor incorporates some of the missing many-body effects into the AIFF induction term.

The good agreement between HMBI and the periodic models can also be partially attributed to the fact that the force field contributions to these crystal energies are rather small: -3.6 and -0.5 kJ/mol for NH<sub>3</sub> and CO<sub>2</sub>, respectively, when three-body dispersion is neglected. In both cases, repulsive three-body dispersion would add roughly +1 kJ/mol to the total energy. In any case, the force field captures the long-range and many-body interactions fairly accurately. Furthermore, the HMBI fragment approach is much less computationally expensive than periodic MP2. On the other hand, the fragment approach used here assumes that the crystal can be partitioned into separate molecular fragments, which is not always true (e.g., polynitrogen crystals).<sup>86</sup>

**4.2. Comparison with Experimental Lattice Energies.** Having demonstrated that the HMBI fragment model nearly reproduces fully quantum mechanical results, the next step is to determine how accurately such predictions reproduce experimental lattice energies. We have already demonstrated that HMBI with the Amoeba force field reproduces lattice energies to within 4–5 kJ/mol on 5 molecular crystals: ice, formamide, acetamide, imidazole, and benzene.<sup>25</sup> Here, we revisit those crystals with the improved force field, and we also add ammonia and carbon dioxide to the test set. These seven crystals include a representative range of intermolecular interactions, ranging from hydrogen bonding (ice, formamide, acetamide, and ammonia) to dispersion (benzene, carbon dioxide) or both (imidazole). For carbon dioxide, we use the experimental lattice energy quoted in ref 30. For ammonia, we use a revised version of the lattice energy cited in ref 30 in which we have made an improved estimate of the

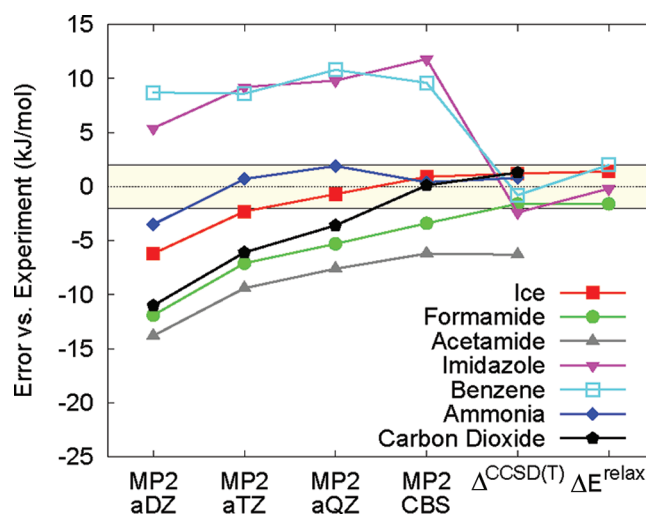
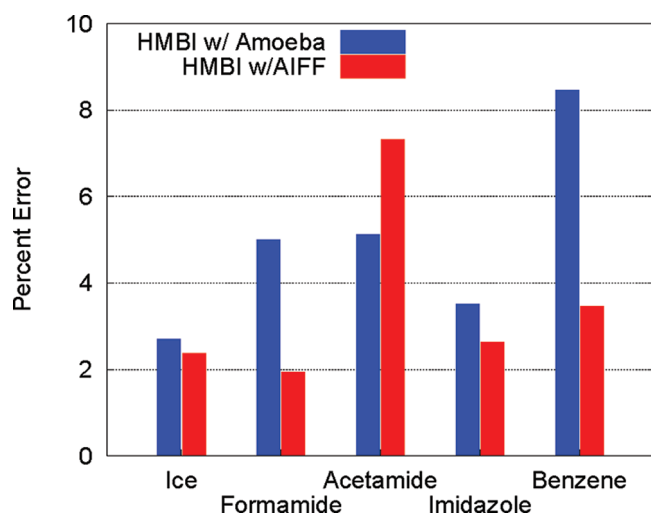


Figure 3. Convergence of the predicted lattice energies toward the experimental values. The yellow band highlights an error of  $\pm 2$  kJ/mol.

zero-point energy contribution. See the Supporting Information for details. For the other five, we use our earlier estimates for the experimental lattice energy.<sup>25</sup>

The calculated lattice energies are listed in Table 1 and plotted as errors relative to experiment in Figure 3. Bear in mind that the experimental lattice energies themselves are probably in error by a couple kJ/mol or more.<sup>25,87,88</sup> Increasing the quality of the wave function used for the QM calculations systematically converges the predicted HMBI lattice energies toward the experimental values. Post-MP2 correlation,  $\Delta^{\text{CCSD(T)}}$ , is particularly important for benzene and imidazole.<sup>25</sup> The estimated lattice parameter relaxation effects,  $\Delta E_{\text{lattice}}^{\text{relax}}$ , are small, with the largest shift of 2.8 kJ/mol coming from benzene. The estimated change in the lattice parameters is also mostly small. The largest change occurs for benzene, where the parameters shrink by an estimated 3.4%. Overall, in six of the seven cases, our best HMBI lattice energy prediction lies within 2 kJ/mol of the experimental value, which is on par with typical experimental errors!

These lattice energy predictions compare favorably with the other calculations found in the literature that have been summarized in ref 25. For example, periodic density functional theory predictions with empirical dispersion corrections often predict



**Figure 4.** Comparison of the percent errors in the best estimate HMBI lattice energy predictions for the Amoeba and AIFF force fields.

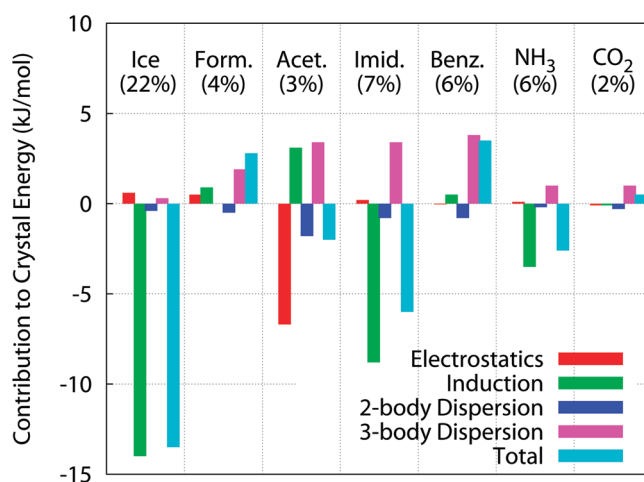
lattice energies with errors ranging from a couple kJ/mol to several times larger than that. For  $\text{NH}_3$  and  $\text{CO}_2$ , the 1–2 kJ/mol errors in our best predictions are comparable to or better than the best periodic MP2 predictions.

For additional perspective, consider that relative energy differences between closely spaced molecular crystal polymorphs are often on the order of  $\sim 1$  kJ/mol. One anticipates some degree of error cancellation between the “absolute” lattice energies when comparing relative polymorph energies. Thus, the ability to predict lattice energies to  $\sim 2$  kJ/mol bodes well for the possibility of reliably distinguishing between crystal polymorphs.

For all of the cases for which Amoeba is parametrized, except acetamide, the AIFF reduces the error in the predicted lattice energy, as shown in Figure 4. On the other hand, the 6 kJ/mol error for the best acetamide prediction is much worse than those for the other crystals. In fact, the use of the AIFF increases the error in the predicted lattice energy slightly compared to the Amoeba force field. The reasons for this behavior are unclear, but the predictions underestimate the experimental lattice energy for acetamide. Relaxing the experimental lattice parameters would increase the calculated lattice energy and potentially improve the prediction. Implementation of analytical derivatives of the AIFF and crystal lattice derivatives is in progress, so we hope to investigate this possibility in the near future. Alternatively, perhaps larger-basis  $\Delta^{\text{CCSD(T)}}$  correlation corrections are needed. The large size of the acetamide unit cell (18 monomers, 162 atoms, and 1008 significant dimers without symmetry) limited the CCSD(T) calculations to the small 6-31+G\* basis. Finally, another source of error might be the moderately small aug-cc-pVDZ basis used in optimizing the crystal structure. The degree of pyramidalization in  $\text{NH}_2$  groups can be quite sensitive to the electronic structure treatment, for example.

#### 4.3. Analysis of the Ab Initio Force Field Contributions.

Finally, we examine the AIFF contributions in more detail. As shown in Figure 4 the AIFF significantly improves upon Amoeba for capturing the long-range and many-body contributions. To provide further insight into this behavior, Figure 5 decomposes the AIFF energy into its individual contributions: long-range two-body electrostatic, induction (both long-range two- and three-body), long-range two-body dispersion, and three-body dispersion.



**Figure 5.** Ab initio force field contributions per molecule to the total crystal energy of the seven crystals considered. The numbers in parentheses indicate the fraction of the net force field contribution relative to the total lattice energy. Note that attractive (negative) contributions increase the lattice energy, while repulsive (positive) ones decrease it.

The electrostatics and induction terms provide a major contribution to the force field energy, particularly for hydrogen-bonded crystals like ice or ammonia. The importance of induction to describing hydrogen-bond cooperativity and organic crystals is well-known.<sup>58,89</sup>

On the other hand, the long-range two-body dispersion terms in the force field contribute very little in all cases. This does not mean that total two-body dispersion is unimportant. Rather, the important two-body dispersion contributions occur at shorter ranges and are captured in the QM part of the model. The fact that the force field dispersion terms only describe long-range contributions means that the  $C_6$  term ( $R^{-6}$  decay) is much more important than the  $C_8$  term ( $R^{-8}$  decay). In the cases examined here,  $C_6$  provides 97–99% of the total long-range dispersion, while the  $C_8$  coefficient contributes the remaining 1–3%. The  $C_{10}$  contributions are yet another order of magnitude smaller, so they are not included.

The contribution of three-body Axilrod–Teller–Muto dispersion is particularly interesting. Conventional wisdom would argue that many-body dispersion is only significant in nonpolar/aromatic species, where induction is unimportant. For this reason, its contribution is often ignored in systems containing polar or hydrogen-bonded molecules.

As expected, three-body dispersion contributes significantly for benzene, imidazole, and carbon dioxide, while its contribution is negligible for ice. For the benzene crystal structure used here, for example, the repulsive three-body dispersion term contributes 4.6 kJ/mol. This is fairly similar to the SAPT(DFT) result of 6.5 kJ/mol for the benzene crystal at the experimental geometry.<sup>37</sup>

Contrary to conventional wisdom, however, we observe that even for the hydrogen-bonded formamide and acetamide crystals, three-body dispersion contributes several kJ/mol to the overall lattice energy. More importantly, it is similar to or larger in magnitude than the induction contribution! Combining these results with the evidence that three-body dispersion can be important in ranking different crystal polymorphs<sup>63</sup> suggests that the many-body dispersion may be more important than has often been thought. Further study is clearly needed.



Together, the inclusion of three-body dispersion and the improved treatment of induction account for most of the differences between the results with the Amoeba and ab initio force fields. The long-range two-body electrostatics and dispersion contributions in the two force fields typically differ by much less than 1 kJ/mol.

The need to determine the AIFF parameters on-the-fly obviously makes the AIFF much more computationally expensive than Amoeba or other conventional force fields. On the other hand, the force field parameters are evaluated separately for each monomer in the central unit cell, so the number of parameters that need to be computed grows only linearly with the number of molecules in the unit cell. Trivial parallelization of the computational effort can be achieved easily by calculating the AIFF parameters for each monomer on a separate processor. The computational cost of evaluating these monomer properties is small compared to the cost of calculating many high-level QM dimer interaction energies. Overall, these AIFF-based crystal calculations are not significantly more expensive than the Amoeba-based ones for which timings have been reported previously.<sup>25</sup>

## 5. CONCLUSIONS

In summary, the HMBI fragment QM/MM model used here provides an accurate and computationally affordable means of predicting molecular crystal lattice energies. In particular, we have demonstrated that a polarizable force field based on distributed multipoles, distributed polarizabilities, and atomic dispersion coefficients which are calculated on-the-fly from DFT provides an accurate treatment of long-range two-body and many-body interactions.

For two different crystals, the model reproduces periodic MP2 results to within a couple kJ/mol. It also predicts six of the seven crystal lattice energies examined to within experimental error. The ability to systematically improve the predictions along with the standard hierarchy of conventional electronic structure methods and basis sets is critical to achieving these high accuracies.

The favorable computational scaling inherent to fragment methods makes it feasible to apply such high-level electronic structure methods to molecular crystals. In particular, the approximation of long-range and many-body intermolecular interactions using a polarizable force field makes the model described here linear scaling for all unit-cell sizes. Furthermore, fragment methods are naturally suited for massively parallel computing, and very high efficiencies can be obtained with hundreds or more processors.

All of these features make this approach very promising for molecular crystal structure prediction. Future work will focus on implementing nuclear gradients of this improved force field that will enable full crystal structure optimization. Other computational speed-ups could be obtained by exploiting crystal space group symmetry and by new developments in accurate, low-cost electronic structure methods.

## ■ ASSOCIATED CONTENT

**S** **Supporting Information.** More detailed expressions for the Ewald summation intermediates and discussion of the ammonia lattice energy are provided. This material is available free of charge via the Internet at <http://pubs.acs.org>.

## ■ AUTHOR INFORMATION

### Corresponding Author

\*E-mail: [gregory.beran@ucr.edu](mailto:gregory.beran@ucr.edu).

## ■ ACKNOWLEDGMENT

Financial support from the National Science Foundation (CHE-1112568) and supercomputer time from the Teragrid (TG-CHE090099 and TG-CHE110064) are gratefully acknowledged. The authors also thank Alston Misquitta and Anthony Stone for providing their CamCASP software and for helpful discussions on its use, and Lorenzo Maschio and co-workers for providing their optimized crystal structures for ammonia and carbon dioxide.

## ■ REFERENCES

- (1) Price, S. L. *Acc. Chem. Res.* **2009**, *42*, 117–26.
- (2) Haas, S.; et al. *Phys. Rev. B* **2007**, *76*, 115203.
- (3) Bauer, J.; Spanton, S.; Quick, R.; Quick, J.; Dziki, W.; Porter, W.; Morris, J. *Pharm. Res.* **2001**, *18*, 859–866.
- (4) Day, G. M.; et al. *Acta Cryst. B* **2009**, *65*, 107–125.
- (5) Kazantsev, A. V.; Karamertzanis, P. G.; Adjiman, C. S.; Pantelides, C. C.; Price, S. L.; Galek, P. T. a.; Day, G. M.; Cruz-Cabeza, A. J. *Int. J. Pharm.* **2011**, *418*, 168–178.
- (6) Price, S. L. *Int. Rev. Phys. Chem.* **2008**, *27*, 541–568.
- (7) Price, S. L.; Leslie, M.; Welch, G. W. A.; M. Habgood, L. S. P.; Karamertzanis, P. G.; Day, G. M. *Phys. Chem. Chem. Phys.* **2010**, *12*, 8478–8490.
- (8) Li, T.; Feng, S. *Pharm. Res.* **2006**, *23*, 2326–2332.
- (9) Kleis, J.; Lundqvist, B. I.; Langreth, D. C.; Schröder, E. *Phys. Rev. B* **2007**, *76*, 1002001.
- (10) Neumann, M. A.; Perrin, M. A. *J. Phys. Chem. B* **2005**, *109*, 15531–15541.
- (11) Neumann, M. A.; Leusen, F. J. J.; Kendrick, J. *Angew. Chem., Int. Ed.* **2008**, *47*, 2427–2430.
- (12) Civalleri, B.; Zicovich-Wilson, C. M.; Valenzano, L.; Ugliengo, P. *CrystEngComm* **2008**, *10*, 405–410.
- (13) Karamertzanis, P. G.; Day, G. M.; Welch, G. W. A.; Kendrick, J.; Leusen, F. J. J.; Neumann, M. A.; Price, S. L. *J. Chem. Phys.* **2008**, *128*, 244708.
- (14) Sorescu, D. C.; Rice, B. M. *J. Phys. Chem. C* **2010**, *114*, 6734–6748.
- (15) Balu, R.; Byrd, E. F. C.; Rice, B. M. *J. Phys. Chem. B* **2011**, *115*, 803–10.
- (16) Shimojo, F.; Wu, Z.; Nakano, A.; Kalia, R. K.; Vashishta, P. *J. Chem. Phys.* **2010**, *132*, 094106.
- (17) Kristyan, S.; Pulay, P. *Chem. Phys. Lett.* **1994**, *229*, 175–180.
- (18) Riley, K. E.; Pitonák, M.; Jurecka, P.; Hobza, P. *Chem. Rev.* **2010**, *110*, 5023–63.
- (19) Grimme, S. *J. Comput. Chem.* **2004**, *25*, 1463–1473.
- (20) Lu, D.; Li, Y.; Rocca, D.; Galli, G. *Phys. Rev. Lett.* **2009**, *102*, 206411.
- (21) Li, Y.; Lu, D.; Nguyen, H.-V.; Galli, G. *J. Phys. Chem. A* **2010**, *114*, 1944–1952.
- (22) Dion, M.; Rydberg, H.; Schröder, E.; Langreth, D. C.; Lundqvist, B. I. *Phys. Rev. Lett.* **2004**, *92*, 246401.
- (23) Thonhauser, T.; Cooper, V. R.; Li, S.; Puzder, A.; Hyldgaard, P.; Langreth, D. C. *Phys. Rev. B* **2007**, *76*, 125112.
- (24) Hongo, K.; Watson, M. A.; Sanchez-Carrera, R. S.; Iitaka, T.; Aspuru-Guzik, A. *J. Phys. Chem. Lett.* **2010**, *1*, 1789–1794.
- (25) Beran, G. J. O.; Nanda, K. *J. Phys. Chem. Lett.* **2010**, *1*, 3480–3487.
- (26) Rivera, S. A.; Allis, D. G.; Hudson, B. S. *Cryst. Growth Des.* **2008**, *8*, 3905–3907.
- (27) Usvyat, D.; Maschio, L.; Manby, F. R.; Casassa, S.; Pisani, C.; Schütz, M. *Phys. Rev. B* **2007**, *76*, 075102.
- (28) Pisani, C.; Maschio, L.; Casassa, S.; Halo, M.; Schütz, M.; Usvyat, D. *J. Comput. Chem.* **2008**, *29*, 2113–2124.
- (29) Erba, A.; Pisani, C.; Casassa, S.; Maschio, L.; Schütz, M.; Usvyat, D. *Phys. Rev. B* **2010**, *81*, 165108.

- (30) Maschio, L.; Usvyat, D.; Schütz, M.; Civalleri, B. *J. Chem. Phys.* **2010**, *132*, 134706.
- (31) Usvyat, D.; Civalleri, B.; Maschio, L.; Dovesi, R.; Pisani, C.; Schutz, M. *J. Chem. Phys.* **2011**, *134*, 214105.
- (32) Marsman, M.; Grueneis, A.; Paier, J.; Kresse, G. *J. Chem. Phys.* **2009**, *130*, 184103.
- (33) Manby, F. R.; Alfe, D.; Gillan, M. J. *Phys. Chem. Chem. Phys.* **2006**, *8*, 5178–5180.
- (34) Addicoat, M. a.; Collins, M. a. *J. Chem. Phys.* **2009**, *131*, 104103.
- (35) Fedorov, D. G.; Kitaura, K. *J. Phys. Chem. A* **2007**, *111*, 6904–6914.
- (36) Podeszwa, R.; Bukowski, R.; Rice, B. M.; Szalewicz, K. *Phys. Chem. Chem. Phys.* **2007**, *9*, 5561–5569.
- (37) Podeszwa, R.; Rice, B. M.; Szalewicz, K. *Phys. Rev. Lett.* **2008**, *101*, 115503.
- (38) Nagayoshi, K.; Ikeda, T.; Kitaura, K.; Nagase, S. *J. Theory Comput. Chem.* **2003**, *2*, 233–244.
- (39) Hirata, S. *J. Chem. Phys.* **2008**, *129*, 204104.
- (40) Sode, O.; Keceli, M.; Hirata, S.; Yagi, K. *Int. J. Quantum Chem.* **2009**, *109*, 1928–1939.
- (41) Dahlke, E. E.; Truhlar, D. G. *J. Phys. Chem. B* **2006**, *3*, 10595–10601.
- (42) Dahlke, E. E.; Truhlar, D. G. *J. Chem. Theory Comput.* **2007**, *3*, 46–53.
- (43) Nagata, T.; Brorsen, K.; Fedorov, D. G.; Kitaura, K.; Gordon, M. S. *J. Chem. Phys.* **2011**, *134*, 124115.
- (44) Tschumper, G. S. *Chem. Phys. Lett.* **2006**, *427*, 185–191.
- (45) Dahlke, E. E.; Truhlar, D. G. *J. Chem. Theory Comput.* **2007**, *3*, 1342–1348.
- (46) Stoll, H.; Paulus, B.; Fulde, P. *J. Chem. Phys.* **2005**, *123*, 144108.
- (47) Hermann, A.; Schwerdtfeger, P. *Phys. Rev. Lett.* **2008**, *101*, 183005.
- (48) Bludsky, O.; Rubes, M.; Soldan, P. *Phys. Rev. B* **2008**, *77*, 092103.
- (49) Beran, G. J. O. *J. Chem. Phys.* **2009**, *130*, 164115.
- (50) Sebetci, A.; Beran, G. J. O. *J. Chem. Theory Comput.* **2010**, *6*, 155–167.
- (51) Kazantsev, A. V.; Karamertzanis, P. G.; Adjiman, C. S.; Pantelides, C. C. *J. Chem. Theory Comput.* **2011**, *7*, 1998–2016.
- (52) Subotnik, J. E.; Sodt, A.; Head-Gordon, M. *J. Chem. Phys.* **2008**, *128*, 034103.
- (53) Stone, A. J. *The Theory of Intermolecular Forces*; Clarendon Press: Oxford, U.K., 2002; chpt. 3–4, pp 7–9.
- (54) Stone, A. J. *Chem. Phys. Lett.* **1981**, *83*, 233–239.
- (55) Stone, A. J.; Alderton, M. *Mol. Phys.* **1985**, *56*, 1047–1064.
- (56) Stone, A. J. *J. Chem. Theory Comput.* **2005**, *1*, 1128–1132.
- (57) Leslie, M. *Mol. Phys.* **2008**, *106*, 1567–1578.
- (58) Welch, G. W. A.; Karamertzanis, P. G.; Misquitta, A. J.; Stone, A. J.; Price, S. L. *J. Chem. Theory Comput.* **2008**, *4*, 522–532.
- (59) Stone, A. J.; Misquitta, A. J. *Int. Rev. Phys. Chem.* **2007**, *26*, 193–222.
- (60) Welch, G. W. A.; Karamertzanis, P. G.; Misquitta, A. J.; Stone, A. J.; Price, S. L. *J. Chem. Theory Comput.* **2008**, *4*, 522–532.
- (61) Misquitta, A. J.; Stone, A. J. *J. Chem. Theory Comput.* **2008**, *4*, 7–18.
- (62) Misquitta, A. J.; Stone, A. J.; Price, S. L. *J. Chem. Theory Comput.* **2008**, *4*, 19–32.
- (63) von Lilienfeld, O. A.; Tkatchenko, A. *J. Chem. Phys.* **2010**, *132*, 234109.
- (64) Misquitta, A.; Stone, A. *Mol. Phys.* **2008**, *106*, 1631–1643.
- (65) Podeszwa, R. *J. Phys. Chem. A* **2008**, *112*, 8884–8885.
- (66) Schwerdtfeger, P.; Assadollahzadeh, B.; Hermann, A. *Phys. Rev. B* **2010**, *82*, 205111.
- (67) Axilrod, P. M.; Teller, E. *J. Chem. Phys.* **1943**, *11*, 299–300.
- (68) Muto, Y. *Proc. Phys.-Math. Soc. Jpn.* **1943**, *17*, 629–631.
- (69) Lotrich, V. F.; Szalewicz, K. *J. Chem. Phys.* **1997**, *106*, 9688–9702.
- (70) Cencek, W.; Jeziorska, M.; Akin-Ojo, O.; Szalewicz, K. *J. Phys. Chem. A* **2007**, *111*, 11311–9.
- (71) Podeszwa, R.; Szalewicz, K. *J. Chem. Phys.* **2007**, *126*, 194101.
- (72) Chalasinski, G.; Szczesniak, M. M.; Kendall, R. A. *J. Chem. Phys.* **1994**, *101*, 8860–8869.
- (73) Steele, R. P.; Distasio, R. A.; Shao, Y.; Kong, J.; Head-Gordon, M. *J. Chem. Phys.* **2006**, *125*, 074108.
- (74) Dunning, T. H. *J. Chem. Phys.* **1989**, *90*, 1007–1023.
- (75) Karton, A.; Martin, J. M. L. *Theor. Chem. Acc.* **2006**, *115*, 330–333.
- (76) Helgaker, T.; Klopper, W.; Koch, H.; Noga, J. *J. Chem. Phys.* **1997**, *106*, 9639–9646.
- (77) Hehre, W. J.; Ditchfield, R.; Pople, J. A. *J. Chem. Phys.* **1972**, *56*, 2257–2261.
- (78) Hariharan, P. C.; Pople, J. A. *Theor. Chim. Acta* **1973**, *28*, 213–222.
- (79) Shao, Y.; et al. *Phys. Chem. Chem. Phys.* **2006**, *8*, 3172–3191.
- (80) Crawford, T. D.; Sherrill, C. D.; Valeev, E. F.; Fermann, J. T.; King, R. A.; Leininger, M. L.; Brown, S. T.; Janssen, C. L.; Seidl, E. T.; Kenny, J. P.; Allen, W. D. *J. Comput. Chem.* **2007**, *28*, 1610–1616.
- (81) Misquitta, A. J.; Stone, A. J. *CamCASP*, v5.6; University of Cambridge: Cambridge, U.K., 2011; <http://www-stone.ch.cam.ac.uk/programs.html>. Accessed February 23, 2011.
- (82) Sadlej, A. J. *Collect. Czech. Chem. Commun.* **1988**, *53*, 1995–2016.
- (83) Sadlej, A. J. *Theor. Chim. Acta* **1991**, *79*, 123–140.
- (84) Lee, M. S.; Maslen, P. E.; Head-Gordon, M. *J. Chem. Phys.* **2000**, *112*, 3592–3601.
- (85) Saebo, S.; Pulay, P. *Annu. Rev. Phys. Chem.* **1993**, *44*, 213–236.
- (86) Erba, A.; Maschio, L.; Salustro, S.; Casassa, S. *J. Chem. Phys.* **2011**, *134*, 074502.
- (87) Shipman, L. L.; Burgess, A. W.; Scheraga, H. A. *J. Phys. Chem.* **1976**, *80*, 52–54.
- (88) Chickos, J. S.; Acree, W. E. *J. Phys. Chem. Ref. Data* **2002**, *31*, 537–698.
- (89) Steiner, T. *Angew. Chem., Int. Ed.* **2002**, *41*, 48–76.

# The Nature of the Binding of Au, Ag, and Pd to Benzene, Coronene, and Graphene: From Benchmark CCSD(T) Calculations to Plane-Wave DFT Calculations

Jaroslav Granatier,<sup>†,||</sup> Petr Lazar,<sup>†,||</sup> Michal Otyepka,<sup>\*,‡</sup> and Pavel Hobza<sup>\*,†,‡,§</sup>

<sup>†</sup>Institute of Organic Chemistry and Biochemistry (IOCB), Academy of Sciences of the Czech Republic and Center for Biomolecules and Complex Molecular Systems, Flemingovo nam. 2, 166 10 Prague, Czech Republic

<sup>‡</sup>Department of Physical Chemistry, Faculty of Science, Regional Centre of Advanced Technologies and Materials (RCPTM), Palacky University Olomouc, tr. 17. listopadu 12, 771 46 Olomouc, Czech Republic

<sup>§</sup>Department of Chemistry, Pohang University of Science and Technology (POSTECH), San 31, Hyojadong, Namgu, Pohang 790-784, Korea

**ABSTRACT:** The adsorption of Ag, Au, and Pd atoms on benzene, coronene, and graphene has been studied using post Hartree–Fock wave function theory (CCSD(T), MP2) and density functional theory (M06-2X, DFT-D3, PBE, vdW-DF) methods. The CCSD(T) benchmark binding energies for benzene–M (M = Pd, Au, Ag) complexes are 19.7, 4.2, and 2.3 kcal/mol, respectively. We found that the nature of binding of the three metals is different: While silver binds predominantly through dispersion interactions, the binding of palladium has a covalent character, and the binding of gold involves a subtle combination of charge transfer and dispersion interactions as well as relativistic effects. We demonstrate that the CCSD(T) benchmark binding energies for benzene–M complexes can be reproduced in plane-wave density functional theory calculations by including a fraction of the exact exchange and a nonempirical van der Waals correction (EE+vdW). Applying the EE+vdW method, we obtained binding energies for the graphene–M (M = Pd, Au, Ag) complexes of 17.4, 5.6, and 4.3 kcal/mol, respectively. The trends in binding energies found for the benzene–M complexes correspond to those in coronene and graphene complexes. DFT methods that use empirical corrections to account for the effects of vdW interactions significantly overestimate binding energies in some of the studied systems.

## 1. INTRODUCTION

Metals are used as interfaces between graphene and conventional electronics; consequently, it is important to understand the nature of the interactions between metals and graphene if nanoelectronics and nanodevices are to reach their full potential.<sup>1</sup> In addition, nanoparticles of gold and palladium on graphene have found an increasing number of applications as biosensors, highly active catalysts, and energy storage devices.<sup>2–8</sup> Unfortunately, the theoretical description of the interactions between a graphene surface and transition metals is complicated by the large (infinite) number of carbon atoms in the graphene sheet and by the complex electronic structure of the transition metals, which is influenced by relativistic effects and both static and dynamic electron correlation. The size of the systems necessitates the use of periodic boundary conditions (i.e., the description of the electronic structure with a plane-wave basis set). Consequently, studies on the interactions between graphene and transition metals have relied heavily on various plane-wave density functional theory (DFT) methods. Surprisingly, the simple local density approximation (LDA) method still finds widespread use,<sup>9,10</sup> reflecting the fact that this method frequently provides better results (due to cancelation of errors) than fundamentally more accurate generalized gradient approximation (GGA) methods.<sup>11–14</sup> For example, the LDA reproduces the available experimental results for the adsorption of Au on graphite surface better than the other GGA, which underpredicts

the strength of binding to Au.<sup>15</sup> However, even the early experiments conducted in the 1970s<sup>16</sup> indicated that the binding of gold on carbon surfaces is heavily dependent on van der Waals (vdW) interactions. This is problematic because neither the LDA nor the various common DFT approaches can describe nonlocal correlation effects, such as vdW interactions. It is worth noting that the physical and chemical nomenclature is not unified; in the physical literature and in this paper, the term “vdW interaction” refers specifically to the London dispersion interaction, which is a weak noncovalent force arising from nonlocal electron correlation. Thus, while the LDA provides a fairly good estimate of the binding energy of gold,<sup>15</sup> it does so for the wrong reason. It is likely that this will have influenced the results obtained in other studies on the adsorption of gold on carbon surfaces,<sup>10,17–19</sup> the adsorption of various metal atoms (including Au, Pd, Fe, and Ti) on graphene,<sup>13</sup> and the adsorption of hydrogen on Pd-decorated graphene.<sup>20</sup> All of these studies were performed using DFT methods that do not account for the contributions of dispersion. It is possible that the adsorption of metals other than gold on carbon surfaces is not governed by the dispersion energy. However, it is impossible to calculate the energy changes involved in the binding of metal atoms to carbon surfaces with thermochemical accuracy (i.e., with errors below 1 kcal/mol) using methods that do not account for dispersion

Received: September 7, 2011

Published: October 05, 2011

energy. Moreover, it is well-known that dispersion energy is an important component of the overall stabilization energy in various types of noncovalent complexes such as those held together by hydrogen bonding,  $\pi$ -stacking, halogen bonding, and other noncovalent interactions.<sup>21</sup> In general, the use of DFT techniques that do not account for dispersion energy causes binding energies to be underestimated.<sup>21</sup>

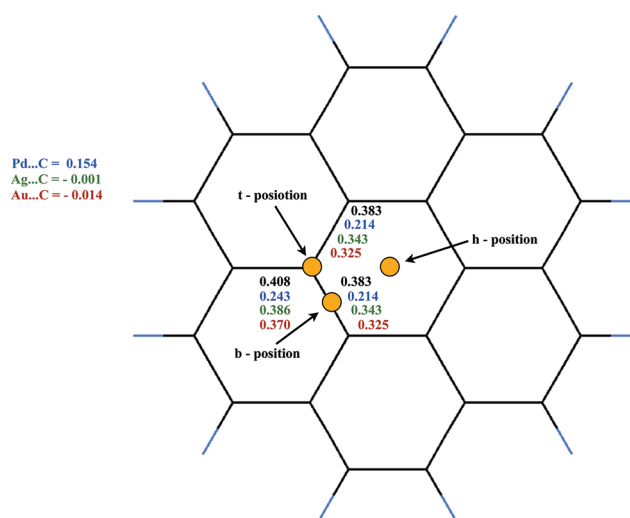
The aim of the study reported herein was to investigate the interaction of graphene with three different transition metals: gold, silver, and palladium. Since the number of quantum chemical methods that can be used to study infinite graphene sheets is rather limited, we initially studied two smaller systems as models of the graphene surface: benzene and coronene. Because the benzene–M (M = Pd, Au, Ag) complexes are comparatively small, they can be studied using even very accurate and computationally expensive wave function theory (WFT) methods based on the coupled cluster technique with iterative evaluation of the contributions of single and double electron excitations and perturbative evaluation of the contributions of triple excitations (CCSD(T)).<sup>22–24</sup> When used in conjunction with an extended basis set, this method provides stabilization energies for various types of noncovalent complexes with chemical or even higher accuracy ( $\pm 1$  or  $\pm 0.1$  kcal/mol)<sup>21</sup> and is therefore used to ‘benchmark’ the performance of less computationally expensive WFT and DFT techniques that account for dispersion interactions in some way. Our first aim was to identify a computational method that is less computationally demanding than CCSD(T) and uses a local basis set but yields good agreement with the CCSD(T) benchmark data. We then planned to use this method to study binding in coronene–M complexes; accurate calculations on these two groups of complexes would provide insights into the nature of the binding of the three different adatoms to carbon surfaces. Specifically, we sought to investigate the performance of the second-order Møller–Plesset (MP2),<sup>25</sup> DFT-D3,<sup>26</sup> and M06-2X<sup>27–29</sup> methods. The DFT-D3 method models the effects of dispersion forces using an additional empirical term that is proportional to  $R^{-6}$ , while the M06-2X functional achieves the same objective by incorporating modified parameters into its exchange–correlation functional. Our second aim was to compare the performance of DFT methods utilizing a plane-wave basis set to that of CCSD(T) in the benzene–M model systems. This comparison was performed to identify a DFT method that can be used to accurately model the interactions of transition-metal atoms with graphene. It was anticipated that the results obtained would make it possible to develop general guidelines for the efficient and accurate modeling of extended systems involving vdW interaction.

## 2. SYSTEMS INVESTIGATED

Benzene–M, coronene–M, and graphene–M (M = Pd, Au, Ag) complexes were investigated. The metal atoms were modeled as being adsorbed at one of three different positions: (t) a ‘top’ site directly above a C atom, (b) a ‘bridge’ site above the midpoint of a C–C bond, and (h) the ‘hollow’ site above the center of the aromatic ring. In the case of coronene, the analogous positions above the central benzenoid ring were considered (Figure 1).

## 3. CALCULATIONS

Benchmarking calculations on the benzene–M complexes were carried out at the spin-adapted CCSD(T) level with a



**Figure 1.** Coronene molecule, showing the three potential sites for the adsorption of metal atoms. Figure also shows the charge distribution in the bonds of the free coronene molecule (black), the coronene–Pd complex (blue), the coronene–Ag complex (green), and the coronene–Au complex (red), as calculated using the M06-2X method. Geometries of the complexes were optimized at the M06-2X level, starting from geometries in which the metal was adsorbed at the (t) position. All final optimized geometries have been bonded on the coronene in (t) position.

restricted closed-/open-shell Hartree–Fock (HF) reference function.<sup>22–24,30</sup> Because of the high computational demands of CCSD(T), the MP2 method<sup>25</sup> was also used. The  $(n-1)p^6$  ( $n-1$ ) $d^{10}$  shells of palladium and  $(n-1)p^6$  ( $n-1$ ) $d^{10}$   $ns^1$  shells of silver and gold were correlated. With the exception of the  $1s^2$  electrons of the carbon atoms, all of the electrons in benzene and coronene were correlated.

Relativistic effects, which are important in heavy transition metals (especially gold) and their complexes,<sup>31</sup> were modeled using the scalar one-component Douglas–Kroll–Hess approximation<sup>32,33</sup> in all wave function methods. All relativistic MP2 and CCSD(T) calculations were performed with ANO-RCC basis sets.<sup>34,35</sup> These basis sets contain diffuse and polarization functions, which are important when studying noncovalent interactions. Another advantage of these basis sets is that they are available with various degrees of contraction. All benchmark CCSD(T) calculations on the benzene–M complexes were performed with the VTZP contraction. MP2 calculations were performed using the VDZP and VTZP contractions as well as with a combination denoted VDZP/VTZP (VDZP for benzene and VTZP for the metal). To compare the relativistic and nonrelativistic CCSD(T) binding energies, calculations were also performed using the relativistic Pol-DK<sup>36</sup> basis sets and the otherwise-equivalent nonrelativistic Pol basis sets,<sup>37</sup> both of which are suitable for calculating molecular electronic properties and the interaction energies of noncovalent complexes.<sup>38</sup> This was done because comparisons of the relativistic and nonrelativistic stabilization energies can provide helpful insights into the nature of the bonding between an aromatic system and a metal atom.<sup>39–41</sup> Throughout this paper, the interaction energy is defined as the difference between the energy of a complex and the sum of the energies of its components; it is negative when the components are attracted to one another. The binding energy is defined as the absolute value of the interaction energy and is

therefore always positive. All calculated WFT interaction energies were corrected for the basis set superposition error (BSSE) using the counterpoise correction.<sup>42</sup> RHF/ROHF, MP2, and CCSD(T) energies were calculated using the MOLCAS 7.2 program package.<sup>43</sup>

The DFT-D3/TPSS/def2-QZVP<sup>26</sup> and M06-2X/lanl2dz<sup>27–29</sup> methods were also used to evaluate the interaction energies of the studied complexes. The DFT-D3 method uses an empirical correction term to describe the dispersion energy, while the M06-2X method accounts for dispersion using a reparameterized exchange–correlation functional. Both of the DFT techniques are substantially less computationally demanding than CCSD(T), making them applicable to large molecular systems.

The structures of benzene and coronene were optimized at the MP2/cc-pVTZ level, and their geometries were assumed to be frozen in all subsequent WFT and DFT calculations, with the exception of the M06-2X calculations on the coronene–M complexes, for which the change in the geometry of the coronene induced by adatom adsorption was studied by full reoptimization of the complex. The DFT-D3 calculations were performed using Turbomole 6.0,<sup>44</sup> and the M06-2X calculations were performed using Gaussian 09.<sup>45</sup>

Plane-wave DFT calculations for an infinite graphene surface were performed using the Vienna Ab initio Simulation Package (VASP) which makes use of the projector augmented wave (PAW) construction for the pseudopotential.<sup>46,47</sup> The GGA of Perdew–Burke–Ernzerhof (PBE)<sup>48</sup> was used to parametrize the exchange–correlation functional. All calculations were carried out using scalar relativistic approximation, i.e., without spin–orbit coupling (except one test calculation for benzene–Au complex, which is discussed later in the text). The structural parameters of benzene and graphene were relaxed by minimizing the forces acting on the atoms using a conjugate gradient algorithm. The energy cutoff for the plane-wave expansion of the eigenfunctions was set to 500 eV. The periodically repeating benzene molecules were separated by at least 8 Å of vacuum in the plane containing the benzene ring and 18 Å of vacuum in the perpendicular direction. The graphene sheet was modeled using a 4 × 4 supercell, i.e. each supercell contained 32 carbon atoms, using the calculated C–C bond length of 1.44 Å. The repeated sheets were separated from each other by 18 Å of vacuum, and the shortest distance between metal atoms was 10 Å. This construction minimizes electrostatic interactions between repeated images. A  $\Gamma$ -centered 5 × 5 × 1 *k*-point mesh was found to provide converged total energies and was consequently used for Brillouin zone integration. Spin polarization was taken into account in all calculations. Long range vdW (dispersion) interactions, which are absent in standard DFT, were included by means of the vdW density functional (vdW-DF)<sup>49</sup> for PBE-optimized geometries. The core of the vdW-DF method is a fully nonlocal expression for the correlation energy  $E_c^{\text{nl}}$ , which takes the following form:

$$E_c^{\text{nl}} = \int \text{d}r^3 \text{d}r'^3 n(r) \Phi(r, r') n(r') \quad (1)$$

Here,  $n(r)$  is the electron density obtained from a standard DFT calculation and the kernel  $\Phi(r, r')$  is a function that depends on  $r - r'$  and the magnitudes and gradients of the electron densities at the points  $r$  and  $r'$ . We used the JuNoLo program to evaluate the vdW term, with PBE electron densities serving as inputs.<sup>50</sup> The vdW-DF method uses standard semilocal GGA functionals

**Table 1.** DK Relativistic and Nonrelativistic Values of the IP, EA, and Dipole Polarizability ( $\alpha$ ) for Metal Atoms<sup>a</sup>

	IP (eV)		EA (eV)		$\alpha$
	DK rel.	nonrel.	DK rel.	nonrel.	DK rel.
Pd (MP2)	8.781		0.248		24.581
Pd (CCSD(T))	8.372		0.521		
Pd (expt)	8.3369 <sup>56</sup>		0.5621 <sup>57</sup>		
Ag (MP2)	7.615	7.013	1.109	0.880	
Ag (CCSD(T))	7.553	6.990	1.279	1.064	52.46 <sup>58</sup>
Ag (expt)	7.57623 <sup>59</sup>		1.304481 <sup>60</sup>		
Au (MP2)	9.342	7.108	2.248	1.043	
Au (CCSD(T))	9.137	7.072	2.250	1.191	36.06 <sup>58</sup>
Au (expt)	9.22553 <sup>61</sup>		2.308664 <sup>62</sup>		

<sup>a</sup> These calculations were been performed using the aug-cc-pVTZ and aug-cc-pVTZ-DK basis sets.<sup>54,55</sup>

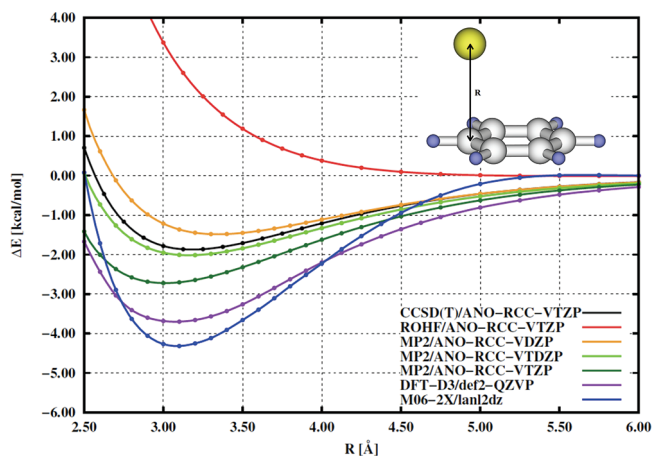
to describe the exchange energy. We chose to use the PBE exchange functional, since it was the functional used to calculate the input electron densities. The total energy was then calculated using the expression:

$$E_{\text{tot}}^{\text{nl}} = E_{\text{tot}}^{\text{DFT}} - E_c^{\text{PBE}} - E_x^{\text{PBE}} + (E_x^{\text{PBE}} + E_c^{\text{LDA}} - E_c^{\text{nl}}) \quad (2)$$

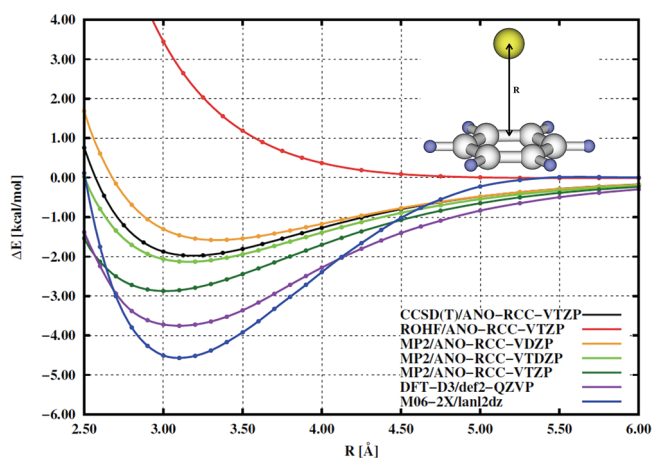
We refer to this method as PBE+vdW. The  $E_x^{\text{PBE}}$  terms are written out explicitly to emphasize the point that the PBE exchange energy inside the parentheses could in principle be replaced by that calculated using some other semilocal formulation; the revised Perdew–Burke–Ernzerhof (revPBE) was suggested in the original formulation of the vdW-DF method by Dion et al.,<sup>49</sup> and other exchange functionals have also been considered.<sup>51,63</sup> In this paper, we propose a different approach; in the spirit of the hybrid screened exchange functionals, we replaced one-quarter of  $E_x^{\text{PBE}}$  with the exact Hartree–Fock exchange,  $E_x^{\text{HF}}$ , which was evaluated in VASP using one-electron Kohn–Sham orbitals. The resulting total energy is denoted as EE+vdW. Notice that  $E_x^{\text{HF}}$  does not match the local density exchange in the constant density limit and so one should not simply exchange  $E_x^{\text{PBE}}$  for  $E_x^{\text{HF}}$ . A rationale for mixing one-quarter of  $E_x^{\text{HF}}$  with the approximate local density exchange was provided by Perdew et al.,<sup>48</sup> who showed that this hybrid matches the LDA in value, slope, and second derivative and is therefore readily embedded into the DFT scheme.

## 4. RESULTS AND DISCUSSION

**4.1. WFT and DFT Calculations on Benzene(Coronene)–M Complexes.** **4.1.1. Isolated Systems.** The DK relativistic and nonrelativistic CCSD(T) and MP2 one-electron properties of all three metal atoms are presented in Table 1. Benzene and coronene are electron donors, while the metal atoms are electron acceptors. Because of its electron affinity, Au is a much stronger electron acceptor than Ag and Pd. Relativistic effects significantly increase the electron affinity and the ionization potential of the gold atom and decrease its dipole polarizability; these effects are much smaller in the other metals considered. Consequently, it was expected that charge-transfer stabilization would be most important in the gold complexes. Ag has the greatest polarizability, followed by Au and Pd. Consequently, it was expected that the dispersion interaction would be strongest in the



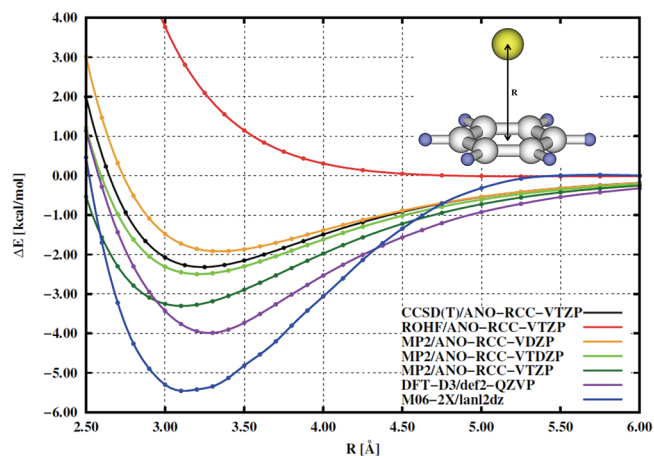
**Figure 2.** Relativistic WFT (BSSE corrected RHF/ROHF, MP2, and CCSD(T)) and DFT (DFT-D3 and M06-2X) potential energy curves for the benzene–Ag complex with the metal adsorbed at the (t) position.



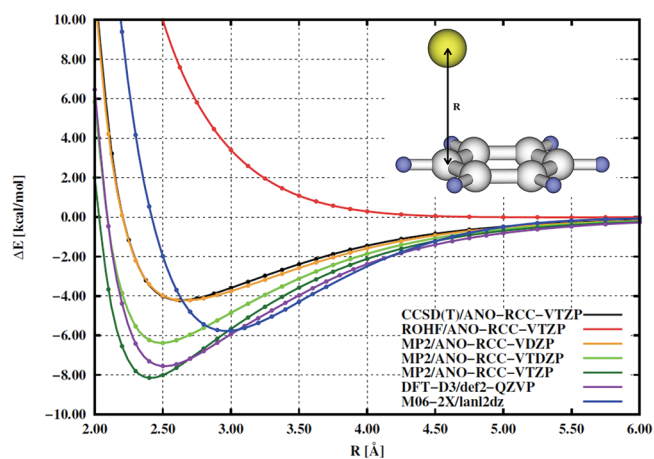
**Figure 3.** Relativistic WFT (BSSE corrected RHF/ROHF, MP2, and CCSD(T)) and DFT (DFT-D3 and M06-2X) potential energy curves for the benzene–Ag complex with the metal adsorbed at the (b) position.

benzene(coronene)–Ag complexes and would become progressively smaller in the corresponding Au and Pd species.

**4.1.2. Benzene–M Complexes.** Figures 2–10 and Table 2 show the characteristics of all complexes investigated in this work. The benzene–Au complex with the Au atom positioned over a carbon atom (t) was energetically similar but slightly more stable than that in which the metal atom was positioned over a C–C bond (b); both were more stable than that in which the gold atom occupied the ‘hollow’ site (h) above the center of the ring. The same relative order was given by all methods investigated. The benchmark (DK rel. CCSD(T)/ANO-RCC-VTZP) binding energies for the (t), (b), and (h) positions were 4.2, 4.1, and 3.2 kcal/mol, respectively. The DK-MP2/ANO-RCC-VDZP method yielded similar binding energies to CCSD(T) for all positions, but MP2 calculations using the larger VTDZP and VTZP basis sets (cf. Figures 5–7) overestimated the binding energies. M06-2X and DFT-D3 systematically overestimated the binding energies by 40–100%. For the (t) and (b) positions, the DFT-D3 energies were in worse agreement with the benchmark data than those obtained with M06-2X, but M06-2X strongly



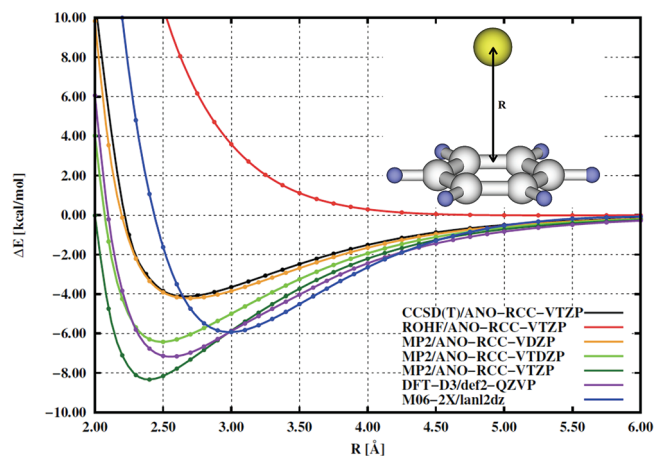
**Figure 4.** Relativistic WFT (BSSE corrected RHF/ROHF, MP2, and CCSD(T)) and DFT (DFT-D3 and M06-2X) potential energy curves for the benzene–Ag complex with the metal adsorbed at the (h) position.



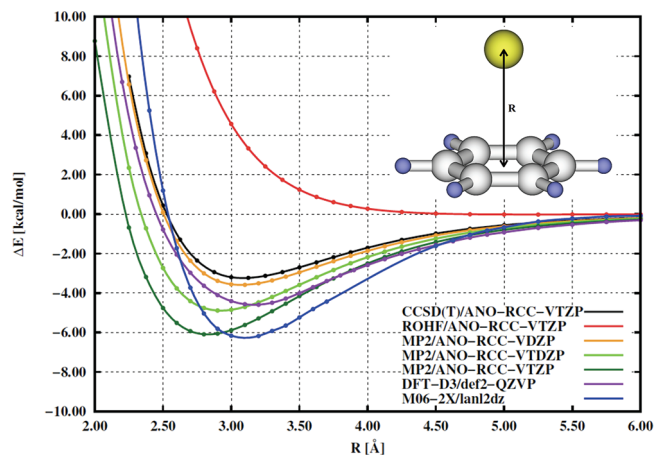
**Figure 5.** Relativistic WFT (BSSE corrected RHF/ROHF, MP2, and CCSD(T)) and DFT (DFT-D3 and M06-2X) potential energy curves for the benzene–Au complex with the metal adsorbed at the (t) position.

overestimates the stabilization for the (h) position. The M06-2X results were also qualitatively inconsistent with the CCSD(T) benchmarks in that they predict the complex with the gold atom in the (h) site to be the most stable.

The situation changes somewhat on switching from Au to Ag. Specifically, the calculated CCSD(T) energies for all three Ag adsorption positions were similar; the species generated by adsorption above the ‘hollow’ (h) was the most favorable but was only 25% more stable than the least favorable, which was generated by adsorption over a carbon atom (t). Similar trends were observed with all of the computational methods examined. The benchmark binding energies for the (h), (b) and (t) positions (2.3, 2.0, and 1.9 kcal/mol, respectively) are smaller than the corresponding values for the benzene–Au complexes by about 30% for (h) and 50% for the (t) and (b) positions. Additionally, the equilibrium distances between the metal atom and the ring were more than 0.5 Å larger in the Ag species than in their Au counterparts for the (b) and (t) positions. As was the case with the Au species, DK-MP2/ANO-RCC-VDZP provided binding energies that mirrored the benchmark results fairly



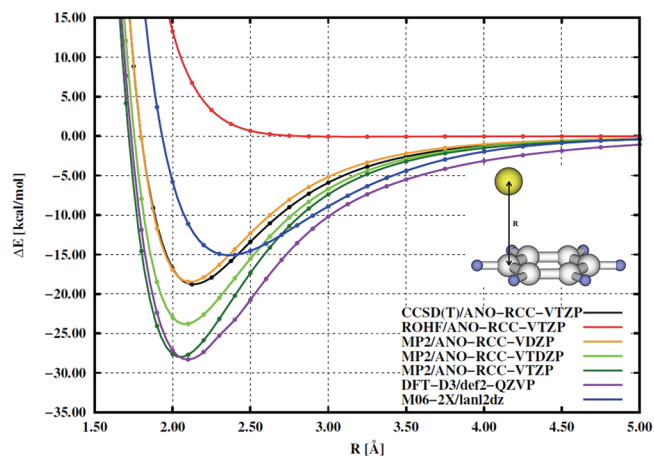
**Figure 6.** Relativistic WFT (BSSE corrected RHF/ROHF, MP2, and CCSD(T)) and DFT (DFT-D3 and M06-2X) potential energy curves for the benzene–Au complex with the metal adsorbed at the (b) position.



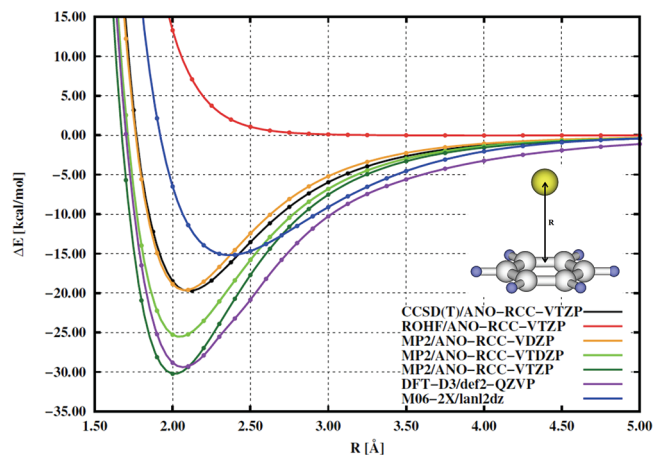
**Figure 7.** Relativistic WFT (BSSE corrected RHF/ROHF, MP2, and CCSD(T)) and DFT (DFT-D3 and M06-2X) potential energy curves for the benzene–Au complex with the metal adsorbed at the (h) position.

closely, while MP2 with triple- $\zeta$  basis set overestimated the binding energies (cf. Figures 2–4). Neither of the DFT methods examined provided reliable binding energies; both DFT-D3 and M06-2X strongly overestimated the stabilization for all three positions.

The low binding energies for Au and Ag are indicative of noncovalent binding. The binding energies for Pd were an order of magnitude higher, suggesting that in this case, the interaction between the metal and the arene is partially covalent. The (b) and (t) positions, which are similar in energy, are preferred to (h), and all methods examined yielded the same order of energies. The benchmark binding energies for the (b), (t), and (h) positions were 19.7, 18.8, and 12.8 kcal/mol, respectively. These higher binding energies were associated with considerably shorter internuclear distances between the Pd and C atoms than was the case in the Au and Ag complexes; adsorption of Pd in the (t) position resulted in an internuclear distance of only 2.1 Å, which is similar to the length of covalent C–Pd bonds. As in both of the preceding cases, DK-MP2/ANO-RCC-VDZP was the method



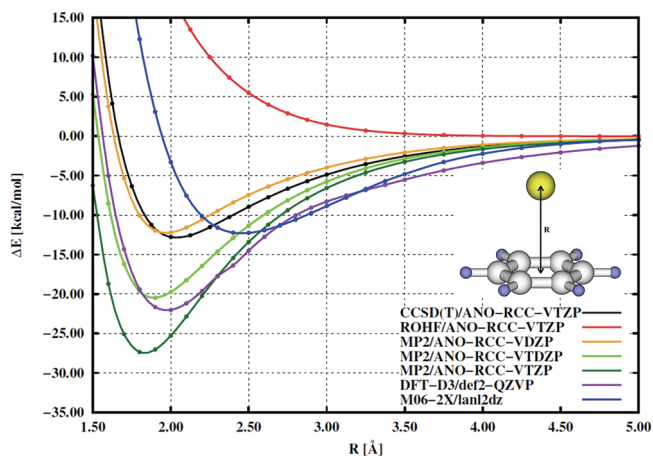
**Figure 8.** Relativistic WFT (BSSE-corrected RHF/ROHF, MP2, and CCSD(T)) and DFT (DFT-D3 and M06-2X) potential energy curves for the benzene–Pd complex with the metal adsorbed at the (t) position.



**Figure 9.** Relativistic WFT (BSSE-corrected RHF/ROHF, MP2, and CCSD(T)) and DFT (DFT-D3 and M06-2X) potential energy curves for the benzene–Pd complex with the metal adsorbed at the (b) position.

whose energies were in best agreement with the benchmark values, with the other MP2 methods once again significantly overestimating the binding energies for all three positions (cf. Figures 8–10). DFT-D3 also significantly overestimates the binding energies (by 35% or more), but M06-2X provides binding energies that agree quite well with the benchmark values, although the (t) and (b) sites are slightly underbound.

These results clearly demonstrate that the interactions of Pd atoms with benzene differ significantly from those of Au and Ag atoms. The binding energies of Pd are much higher than those of Au and Ag, and the corresponding internuclear distances are much shorter. The Au and Ag binding energies are in the range typically associated with noncovalent interactions, whereas the Pd binding energies encroach on ranges more commonly associated with covalent bonds. All three binding sites yield broadly similar binding energies for the adsorption of Au and Ag, but the (b) and (t) positions are clearly favored over the (h) site in the case of Pd adsorption. Of the computational methods tested, DK-MP2/ANO-RCC-VDZP provided the best



**Figure 10.** Relativistic WFT (BSSE-corrected RHF/ROHF, MP2, and CCSD(T)) and DFT (DFT-D3 and M06-2X) potential energy curves for the benzene–Pd complex with the metal adsorbed at the (h) position.

agreement with the benchmark CCSD(T) energies and can thus reasonably be expected to provide accurate results when applied to larger model systems. It should be noted that better agreement for double- $\zeta$  basis set (than for triple- $\zeta$  basis set) arises from compensation of errors. However, while the MP2 method is less expensive than CCSD(T), it is still rather computationally demanding. Of the faster DFT techniques, M06-2X is preferable to DFT-D3, since it gave absolute binding energies that better matched the benchmark values. However, when considering the relative magnitudes of the binding energies for the three elements, a different picture emerges. The CCSD(T) benchmark calculations indicate that the binding energy of Pd to benzene is nine times greater than that of Ag and that of Au is two times greater, giving a benchmark Pd:Au:Ag ratio of 9:2:1. The MP2 (10:2:1) and DFT-D3 ratios (7:2:1) matched the benchmark values fairly closely, but the M06-2X results (3:1:1) strongly disfavor Pd. Thus, for comparing the binding energies of different metals, MP2 and DFT-D3 appear to be superior to M06-2X.

Our results strongly contradict the findings of previous studies in which DFT methods were used. For example, DFT/BPW91/TZP calculations<sup>52</sup> on benzene–M (M = Ag and Au) complexes provided binding energies for the (h), (b) and (t) positions of 5.7, 5.3, and 5.3 kcal/mol, respectively, for Ag, and 5.3, 5.1, and 3.9 kcal/mol, respectively, for Au. These findings are clearly incompatible with the benchmark data reported herein, since they suggest that the binding energies for Au are smaller than those for Ag. This is probably due to the neglect of relativistic effects at the DFT/BPW91/TZP level of theory; relativistic effects change the nature of binding in the benzene–Au complexes, as discussed below.

**4.1.3. Nature of the Bonding in Benzene–M Complexes.** The nature of the metal–arene binding in all three complexes differs, as indicated by the differences in the binding energies calculated using different levels of theory. The omission of the correlation energy causes the binding energies to be strongly underestimated. Figures 2–10 show that the HF energy curves for all atoms and all adsorption positions are universally repulsive, i.e., no binding occurs. This indicates that the stabilization of all benzene–M complexes originates from correlation effects. However, while correlation effects are important in the binding of all three of the investigated metals, relativistic effects are only important in the Au complexes. This conclusion is supported by

the calculated one-electron properties shown in Table 1. The calculated ionization potential and electron affinity of Au change dramatically when relativistic effects are included; these in turn affect the benzene–Au binding energies, which are significantly reduced by the omission of relativistic effects. The relativistic CCSD(T)/Pol-DK binding energies are 3.7, 3.7, and 3.1 kcal/mol for the (t), (b), and (h) positions, respectively; the corresponding nonrelativistic binding energies are significantly smaller (2.0, 2.1, and 2.5 kcal/mol, respectively). For the sake of comparison, we also determined the relativistic vs nonrelativistic binding energies for the (t), (b), and (h) positions of the benzene–Ag complex, which were 2.1, 2.2, and 2.6 vs 1.9, 2.0, and 2.4 kcal/mol, respectively.

One of the most reliable ways of obtaining information on the nature of the bonding is to compare the electronic structure of the bound species to that of the isolated atom. In the case of Ag, such comparisons indicate that the stabilization of the benzene–Ag complex is almost entirely due to the London dispersion energy. This is consistent with the high polarizability of Ag and the relatively large distance between the Ag nucleus and the benzene ring, which means that there is very little overlap of the orbitals of the metal and the arene. Indeed, the orbitals of the complex are almost identical to those of its separated constituents. Analysis of the charge transfer in the Ag complexes (Mulliken charges, determined using the MP2/ANO-RCC-VDZP method) revealed that Ag carries a negative charge of  $-0.05 e$  in all of the structures examined, i.e., it acts as an electron acceptor, while benzene is an electron donor. Because of the low electron affinity of Ag and the large separation of the metal atom and the arene, there is relatively little charge transfer from benzene to the Ag atom.

Compared to Ag, Au is significantly less polarizable and has a higher electron affinity (Table 1). The lower polarizability of Au implies that dispersion interactions will be less important in its complexes, and the higher electron affinity is likely to increase the importance of charge-transfer interactions. Mulliken population analyses indicated that the magnitude of the charge transfer in the benzene–Au complexes was approximately twice that in the benzene–Ag complexes, with the Au atom carrying negative charges of  $-0.11$  and  $-0.12 e$  for the (t) and (b) positions, respectively. This enhanced charge transfer is attributable to relativistic effects because their omission halves the electron affinity of the gold atom (Table 1). The stabilization of the benzene–Au complex by charge-transfer interactions is demonstrated by the fact that their binding energies are more than twice as large as those for the corresponding benzene–Ag complexes and by the considerably shorter (by more than 0.5 Å) distances between the benzene ring and the metal atom in the gold complexes. These shorter distances reflect a greater overlap between the orbitals of the two systems. Specifically, the formation of new bonding and antibonding orbitals from the doubly occupied  $5d_0$  orbital of Au and the benzene  $p_z$  orbitals was observed. This interaction model, which highlights the importance of charge transfer, has been presented in previous works.<sup>39–41</sup> The dramatic increase in stability for complexes of Au is due to relativistic effects, which increase the metal's electron affinity and thus favor the transfer of charge from the ligand to the metal. While charge transfer plays a key role for gold atoms in the (b) and (t) positions, it is less pronounced in the (h) position; here, the dispersion energy provides a larger contribution to the binding energy. It is worth noting that for  $Au^+$  and  $Ag^+$  ion–arene complexes, the bonding becomes mainly electrostatic, and binding energies are almost an order of magnitude higher.<sup>64–66</sup>



**Table 2.** Extrapolated Interaction Energies  $\Delta E$  [kcal/mol] and Optimal Bond Lengths  $R$  (in terms of the shortest distance between the metal atom and the benzene plane) [Å] for Benzene– $M$  ( $M = \text{Ag, Au, Pd}$ ) Complexes Calculated at the Various DFT with Dispersion Correction and DK Relativistic and Nonrelativistic WFT Levels

	benzene–Pd			benzene–Ag			benzene–Au		
	(t)	(b)	(h)	(t)	(b)	(h)	(t)	(b)	(h)
DFT-D3/TPSS/def2-QZVP									
$\Delta E$	–28.3	–29.4	–22.1	–3.7	–3.7	–4.0	–7.5	–7.2	–4.6
$R$	2.10	2.07	1.97	3.07	3.10	3.28	2.51	2.56	3.17
M06-2X/lanl2dz									
$\Delta E$	–15.1	–15.2	–12.3	–4.3	–4.6	–5.5	–5.8	–5.9	–6.3
$R$	2.36	2.37	2.45	3.09	3.10	3.12	2.97	2.99	3.10
DK rel. MP2/ANO-RCC-VDZP									
$\Delta E$	–18.5	–19.6	–12.3	–1.5	–1.6	–1.9	–4.2	–4.2	–3.6
$R$	2.11	2.08	1.97	3.34	3.33	3.34	2.66	2.69	3.07
DK rel. MP2/ANO-RCC-VTZP									
$\Delta E$	–28.0	–30.2	–27.5	–2.7	–2.9	–3.3	–8.1	–8.3	–6.1
$R$	2.05	2.01	1.83	3.01	3.01	3.11	2.41	2.39	2.83
DK rel. CCSD(T)/ANO-RCC-VTZP									
$\Delta E$	–18.8	–19.7	–12.8	–1.9	–2.0	–2.3	–4.2	–4.1	–3.2
$R$	2.13	2.11	2.04	3.18	3.18	3.24	2.63	2.67	3.09
DK rel. CCSD(T)/Pol-DK									
$\Delta E$	–	–	–	–2.1	–2.2	–2.6	–3.7	–3.7	–3.1
$R$	–	–	–	3.19	3.19	3.24	2.73	2.79	3.17
nonrel. CCSD(T)/Pol									
$\Delta E$	–	–	–	–1.9	–2.0	–2.4	–2.0	–2.1	–2.5
$R$	–	–	–	3.29	3.29	3.29	3.36	3.37	3.39
GGA PBE									
$\Delta E$	–26.3	–27.3	–19.0	–1.3	–1.2	–1.0	–6.1	–5.6	–1.63
$R$	2.10	2.07	2.01	3.05	3.10	3.39	2.44	2.46	3.09
PBE+vdW									
$\Delta E$	–21.5	–21.8	–13.3	–2.7	–2.7	–2.6	–5.9	–5.5	–3.6
$R$	2.17	2.18	2.16	3.17	3.23	3.41	2.70	2.79	3.21
EE+vdW									
$\Delta E$	–17.2	–18.7	–10.6	–2.4	–2.3	–2.5	–5.1 <sup>a</sup>	–4.8	–3.4
$R$	2.18	2.15	2.16	3.22	3.32	3.41	2.64	2.74	3.22

<sup>a</sup>EE + vdW + spin–orbit coupling (soc) –5.7 kcal/mol.

The metal–ligand bonding in the benzene–Pd complexes differs significantly from that in the Ag and Au complexes due to the different electronic structure of Pd. In the ground state, the valence d-orbitals of palladium are fully occupied, and the first virtual orbital is the 5s. In all of the benzene–Pd complexes examined in this work, the Pd atom carried a small positive charge, indicating that it was acting as an electron donor. Detailed analyses indicated a significant loss of electron density from the Pd valence d-orbitals (relative to the situation in the free atom) and a simultaneous significant increase in electron density in the virtual 5s orbital. This is consistent with the formation of a so-called dative bond, in which charge is transferred from Pd to benzene, leading to an increase in the electron density of the benzene ring and a decrease in that of the Pd atom. This polar complex is then stabilized by back donation of charge from the carbon atom to the valence 5s orbital of Pd. A dative bond of this

kind would account for the high binding energies observed for the benzene–Pd complex.

**4.1.4. Coronene– $X$  Complexes.** Coronene is a more complex model of graphene than benzene. The central aromatic ring of coronene (Figure 1) is surrounded only by other aromatic rings, and all its carbon atoms bind exclusively to other carbon atoms. We investigated the binding of Ag, Au, and Pd atoms to coronene using the MP2/ANO-RCC-VDZP, DFT-D3/def2-QZVP, and M06-2X/lanl2dz methods, as discussed in the preceding section. The size of the coronene complexes meant that it would have been impractical to perform CCSD(T) calculations on them to obtain benchmark binding energies. Therefore, binding energies calculated using the MP2 method were used as reference values for the coronene complex, since this level of theory provided absolute and relative binding energies that were reasonably close to the benchmark CCSD(T) values for all of the benzene–metal

**Table 3.** DK rel. MP2, DFT-D3, and M06-2X Extrapolated Interaction Energies  $\Delta E$  [kcal/mol] and Metal Atom Charges [e] for Coronene–M (M = Ag, Au, Pd) Complexes with an Optimized Bond Length R [Å]

	coronene–Pd			coronene–Ag			coronene–Au		
	(t)	(b)	(h)	(t)	(b)	(h)	(t)	(b)	(h)
MP2									
$\Delta E$	–17.7	–17.9	–13.7	–3.9	–4.0	–4.1	–6.9	–7.0	–6.7
R	2.11	2.09	1.99	~3.17	~3.13	~3.19	2.83	2.82	2.92
charge	0.051	0.045	0.032	–0.052	–0.052	–0.051	–0.068	–0.067	–0.063
DFT-D3									
$\Delta E$	–26.3	–26.9	–24.6	–5.7	–5.8	–6.0	–7.3	–7.3	–6.7
R	2.12	2.08	1.99	3.16	3.14	3.21	2.80	2.84	3.09
M06-2X									
$\Delta E$	–14.0	–14.1	–12.8	–6.3	–6.3	–6.0	–7.2	–7.2	–7.0
R	2.46	2.45	2.47	3.12	3.12	3.12	3.06	3.06	3.09
charge	0.073	0.075	0.067	–0.010	–0.009	–0.007	–0.028	–0.027	–0.027

complexes discussed in the preceding section. The M06-2X method was used to optimize the geometries of the coronene–M complexes and to estimate the changes in the electronic structure of the coronene following adatom adsorption.

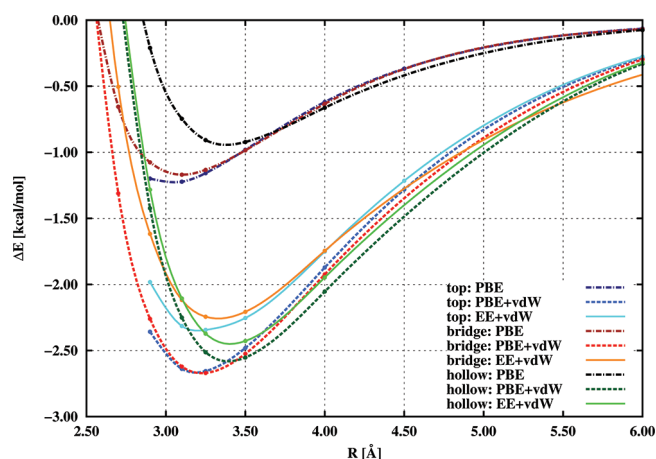
The MP2, DFT-D3, and M06-2X binding energies and equilibrium distances for all of the coronene complexes considered are summarized in Table 3. It is apparent that the binding energies for the coronene complexes differ from their benzene counterparts. At the MP2 level, it was found that the binding energies for Au and Ag increased on going from benzene to coronene, by around 50% in the case of Au and around 100% in the case of Ag. Conversely, going from benzene to coronene reduced the binding energy of Pd by around 10%, although binding in the (h) position was slightly stronger in the coronene complex than in the corresponding benzene species. However, the relative strength of binding to Pd, Au, and Ag remained as it had been in the case of benzene, as did the relative binding energies for adsorption at different positions around the ring. For the Au complexes, the internuclear distances between the metal and the plane containing the arene increased on going from benzene to coronene; for the Ag complexes, the corresponding internuclear distances decreased. However, in both cases, the differences between the distances in the benzene and coronene complexes were small. No significant difference in distance was observed in the Pd complexes. It appears that the nature of the metal–arene bond in the coronene–Pd complexes is very similar to that in the benzene–Pd complexes; a “covalent” bond is formed between the carbon atoms and Pd by the overlap of the d-orbitals of Pd with the  $\pi$  orbitals of the coronene. Silver atoms bind exclusively via dispersion forces; while the polarizability of coronene is greater than that of benzene, this is outweighed by the fact that the coronene complexes have a greater number of carbon atoms and therefore experience more exchange repulsion than their benzene counterparts. This is the cause of the greater carbon–Ag distances in coronene complexes of silver. The situation with the gold complexes is more complicated, because both the dispersion energy and the charge transfer are important in their stabilization. As with the Ag complexes, the exchange repulsion is greater in the coronene complexes of Au than in the benzene species, and so the distances between the Au atom and the plane containing the arene are somewhat greater in the

coronene complexes, although the difference is relatively modest. As with the benzene complexes, it is possible to obtain insights into the bonding and charge transfer in coronene–metal complexes by analyzing the Mulliken charges on the adatoms. Both gold and silver atoms in the coronene complexes carry partial negative charges, indicating that both function as electron acceptors. The MP2 charges, which were used as reference values, were greater than the M06-2X charges and can be compared to those calculated for the benzene complexes. While the extent of charge transfer in the silver complexes of benzene and coronene was very similar, the magnitude of the charge transfer in the coronene–Au complexes was approximately 40% smaller than that in the corresponding benzene complexes.

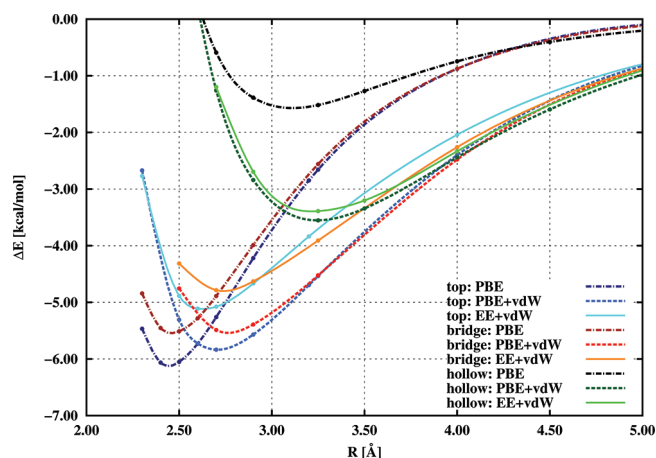
Both the DFT-D3 and M06-2X calculations exhibited trends similar to those observed in the MP2 data, and the relative stabilities of all of the coronene–metal complexes considered were well reproduced. The DFT-D3 interaction energies for the gold complexes were very similar close to those obtained at the MP2 level. However, the DFT-D3 binding energies for the Ag and Pd complexes exceeded the MP2 values by 50% or more. The M06-2X binding energies for the Pd and Au complexes agreed well with the MP2 values, but those for the Ag complex were overestimated by about 60%.

All three methods considered (i.e., MP2, DFT-D3, and M06-2X) indicate that the adsorption of Pd is significantly more favorable than that of Au or Ag, but the extent to which this is the case depends on the method used (MP2, 4:2:1; DFT-D3, 4:1:1; M06-2X, 2:1:1). In all cases, however, the difference between the binding energies for Pd and Ag was smaller than that observed with the corresponding benzene complexes.

Figure 1 shows the (t), (b), and (h) positions for adsorption on coronene and also the M06-2X overlap populations in the C–C bonds that are affected by adsorption. In the case of adsorption of an Ag adatom, there is no significant change in the overlap populations relative to those in the isolated coronene, and the total overlap between Ag and the nearest C is also negligible (–0.001). This is not the case in the corresponding Au complexes, in which all the C–C bonds in coronene are weakened relative to those in the isolated molecule (having electron populations of 0.325, 0.325, and 0.370), but the overlap population of the Au–C bond remains negative (–0.014).



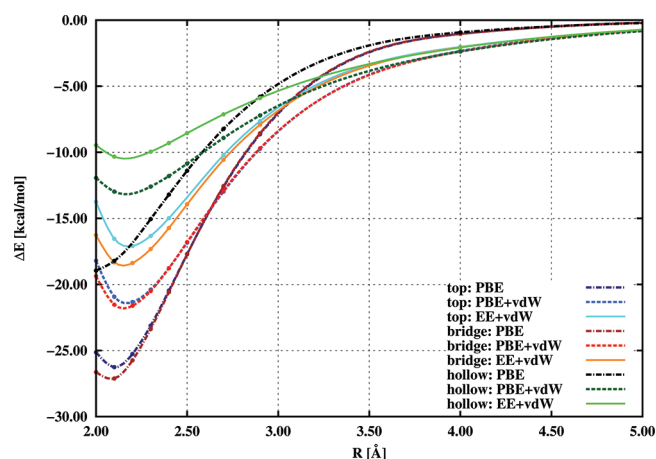
**Figure 11.** Periodic plane-wave DFT/PBE, DFT/PBE+vdW, and DFT/EE + vdW potential curves for the benzene–Ag complex with the metal adsorbed at the (t), (b), and (h) positions.



**Figure 12.** Periodic plane-wave DFT/PBE, DFT/PBE+vdW, and DFT/EE+vdW potential curves for the benzene–Au complex with the metal adsorbed at the (t), (b), and (h) positions.

The C–C bonds in coronene are weakened due to their relatively strong interaction with the Au adatom. Even more dramatic changes occur upon the adsorption of Pd. The Pd–C bond is significantly populated (0.154), and the overlap populations of the C–C bonds are significantly reduced (0.214, 0.214 and 0.243) relative to those in the free coronene. These numbers clearly show that the binding of Ag to coronene (and to some extent, also that of Au) is noncovalent, occurring primarily via dispersion forces, whereas Pd binds covalently. The overlap populations between the Pd and C atoms are comparable to those between carbon atoms in the vicinity of the adsorption site, demonstrating that the adsorption of Pd significantly weakens the covalent C–C bonds in the vicinity of the site of adsorption and results in the formation of a partly covalent bond between the Pd and C atoms.

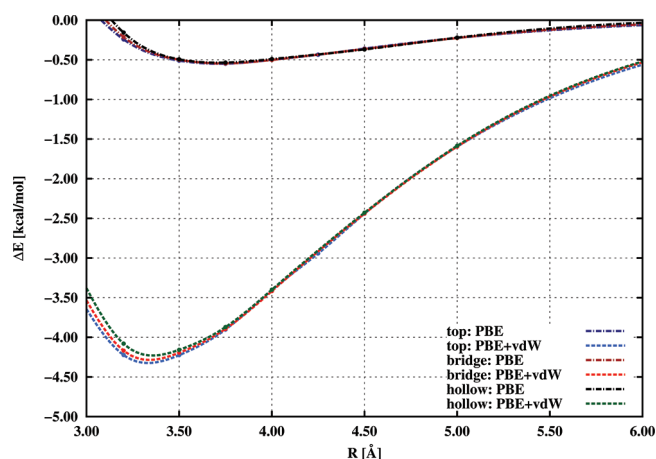
**4.2. Periodic Plane-Wave DFT Calculations.** *4.2.1. Benzene–M Complexes.* Figures 11–13 and Table 2 show the binding energies calculated using the plane-wave approach. The main differences between the investigated elements can be seen even in the results of the simple PBE/GGA calculations, although this method is rather unsatisfactory in quantitative terms. Compared to the CCSD(T) benchmark results, the benzene–Pd



**Figure 13.** Periodic plane-wave DFT/PBE, DFT/PBE+vdW, and DFT/EE+vdW potential curves for the benzene–Pd complex with the metal adsorbed at the (t), (b), and (h) positions.

and –Au complexes are significantly overbound, whereas the benzene–Ag complex is underbound. On examining the PBE+vdW energy curves, it is apparent that this disagreement is primarily due to the neglect of dispersion forces. The inclusion of dispersion forces affords greatly improved agreement with the benchmark values, as discussed in more detail below. It should be noted that the LDA approximation, which is also used in studies of adsorption on graphene, overestimates the binding energy by more than 100% in all cases examined (data not shown) and yields unreasonably short bond distances as well.

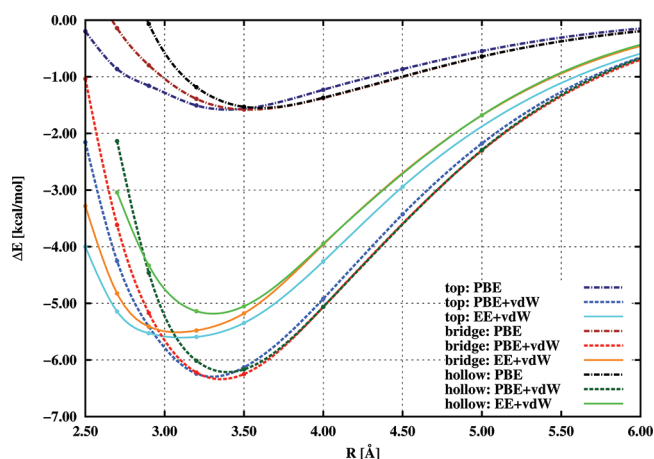
The benzene–Au complex has a total spin moment of  $1 \mu_B$  due to the single valence electron of the Au atom. The spin moment does not change substantially as a function of the distance between the Au atom and the benzene ring, indicating that there is negligible charge transfer between the Au atom and the C atoms of the benzene ring. As suggested by the WFT methods, the (t) position is preferred to the (b) position, although the binding energies for these two spots are very similar and are both significantly greater than that for the hollow (h) position. The relative order of energies is the same for all methods investigated, but the calculated energetic differences are reduced when the nonlocal vdW term is incorporated into the calculations. Inspection of the interaction energy curves in Figure 12 indicates that the vdW term is actually repulsive, i.e., the PBE+vdW equilibrium energies and distances are higher than those given by the PBE calculation. The inclusion of one-quarter of exact exchange in the calculation further reduces the binding energies and yields the best agreement with the benchmark CCSD(T) calculations. The EE+vdW binding energies for the (t), (b), and (h) positions are 5.1, 4.8, and 3.4 kcal/mol, respectively. This means that both the values of binding energies and the differences between the binding energies for the (t), (b), and (h) positions are within 1 kcal of the benchmark CCSD(T) values. As gold is known to display significant relativistic effects, we tested the influence of spin–orbit coupling (soc) on the interaction energy for the (t) position. It was found that soc has a slight effect on the total PBE energy but has little impact on the charge density distribution within the complex, which determines the nonlocal vdW contribution (see eq 1). The binding energy for Au in the (t) position as calculated using the EE+vdW+soc method is 5.7 kcal/mol.



**Figure 14.** Periodic plane-wave DFT/PBE and DFT/PBE+vdW potential curves for the graphene–Ag complex with the metal adsorbed at the (t), (b), and (h) positions.

The potential curves for the benzene–Ag complex are shown in Figure 11 and clearly illustrate the importance of the vdW dispersion term. Using PBE alone, the calculated binding energies for the (t), (b), and (h) positions were 1.3, 1.2, and 1.0 kcal/mol, respectively. Obviously, these values are unrealistically low, especially for the hollow position, which was found to be the preferred site in the CCSD(T) calculations. By including the vdW term, identical binding energies of 2.7 kcal/mol were obtained for the (t) and (b) positions, while the binding energy of 2.6 kcal/mol for the hollow position (h) was slightly lower. While these values are already in very good agreement with the benchmark values, they were further improved upon by adding a fraction of the exact exchange; this made the hollow (h) position the preferred site for adsorption, as predicted by CCSD(T). The EE+vdW binding energies for the (t), (b) and (h) positions were 2.4, 2.3, and 2.5 kcal/mol, respectively. As was the case for the benzene–Au complex, these binding energies are slightly greater than the benchmark values. The spin moment remains constant at  $1 \mu_B$  for all internuclear distances, which is consistent with a negligible electrostatic interaction between the Ag atom and the carbon atoms of the benzene ring.

A different situation obtains for the benzene–Pd complex. Here, the covalent interaction between the metal and the arene means the binding energy is large; using the PBE method, it is predicted to be 26.3, 27.3, and 19.0 kcal/mol for the (t), (b), and (h) positions, respectively. These values are significantly higher than the benchmark CCSD(T) values. In contrast to the situation with the Ag complex, the inclusion of the vdW term substantially reduces the predicted binding energies. While this may be surprising at first sight, the kernel  $\Phi(r, r')$  used to describe the interactions between electron densities (eq 1) becomes repulsive at small distances.<sup>49</sup> Thus, the PBE+vdW calculation corrects the overbinding predicted by PBE alone, giving binding energies of 21.5, 21.8, and 13.3 kcal/mol. It should be noted that such repulsive corrections are impossible in the various empirical DFT+D2 (or D3)<sup>26</sup> approaches, because D2 and D3 terms are always attractive, i.e., they provide a nonzero and positive (in terms of the definition of binding energy used in this paper) contribution to the binding energy. Incorporating a fraction of the exact exchange energy further reduced the calculated binding energy, as was the case for the benzene–Au

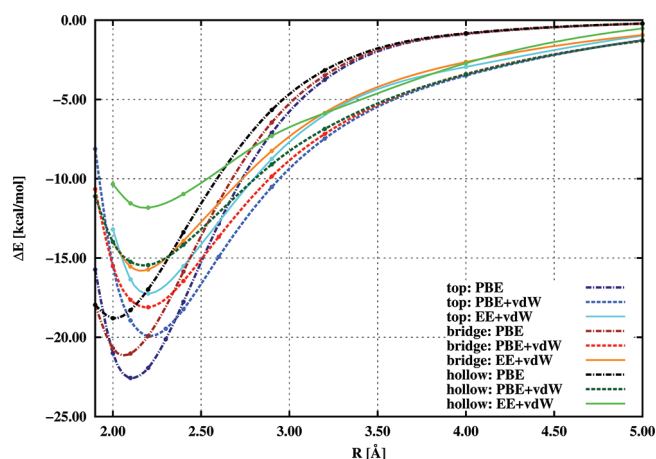


**Figure 15.** Periodic plane-wave DFT/PBE, DFT/PBE+vdW, and DFT/EE+vdW potential curves for the graphene–Au complex with the metal adsorbed at the (t), (b), and (h) positions.

complex. The EE+vdW binding energies for the (t), (b), and (h) positions were thus reduced to 17.2, 18.7, and 10.6 kcal/mol, respectively. As before, the inclusion of one-quarter of the exact exchange yielded DFT results that were very close to the benchmark value (although in this case, the DFT binding energies were slightly lower than the reference values), demonstrating that methods for improving on the treatment of long-range correlation effects (the vdW term) should be used in conjunction with methods that treat midrange exchange properly.

**Graphene–M Complexes.** The main advantage of calculations that use periodic plane-wave basis sets is that they can be applied to the study of extended systems. Our studies on benzene–M complexes demonstrated that the PBE functional can yield binding energies that agree very well with reference CCSD(T) values when augmented with a nonlocal vdW correction and one-quarter of the exact exchange (EE+vdW). We therefore used this method to obtain DFT benchmark energies for the interactions of metal atoms with a graphene sheet. In this context, it should be noted that PBE+vdW interaction energies for Cu, Ag, and Au atoms on graphene have been published very recently.<sup>53</sup> Our calculations differ from those reported in that publication, however, since (i) we included the contribution of the exact HF exchange in order to obtain more reliable interaction energies, and (ii) our calculations used carbon atoms that were fixed in place (i.e., no geometrical relaxation of the graphene sheet was allowed) in order to facilitate comparisons of the bonding of metals adsorbed on graphene with that in benzene and coronene complexes. The role of geometrical relaxation of the graphene surface is thoroughly discussed by Amft et al.<sup>53</sup>

Figures 14–16 and Table 4 summarize the calculated interaction energies for the graphene–M complexes. On examining the binding energy of the graphene–Au complex, it is apparent the bonding is dominated by vdW term, which stands in stark contrast to the situation in the benzene–Au complex. The binding energy calculated using the GGA/PBE approximation alone is very weak ( $\sim 1.6$  kcal/mol), which is consistent with the GGA values of 2.2 kcal/mol for the most favorable (t) position reported in previous works.<sup>12,13</sup> The difference in the GGA binding energies can be attributed to the fact that in those previous works, the geometry of the graphene was allowed to



**Figure 16.** Periodic plane-wave DFT/PBE, DFT/PBE+vdW, and DFT/EE+vdW potential curves for the graphene–Pd complex with the metal adsorbed at the (t), (b), and (h) positions.

**Table 4.** Interaction Energies  $\Delta E$  [kcal/mol] for Graphene–M (M = Ag, Au, Pd) Complexes with an Optimized Bond Length  $R$  [Å]

	graphene–Pd			graphene–Ag			graphene–Au		
	(t)	(b)	(h)	(t)	(b)	(h)	(t)	(b)	(h)
PBE									
$\Delta E$	–22.8	–21.3	–19.0	–0.6	–0.6	–0.6	–1.6	–1.6	–1.6
$R$	2.11	2.07	2.02	3.72	3.72	3.73	3.41	3.54	3.63
PBE+vdW									
$\Delta E$	–20.1	–18.3	–15.6	–4.3	–4.3	–4.2	–6.3	–6.4	–6.2
$R$	2.22	2.20	2.18	3.35	3.35	3.39	3.30	3.36	3.42
EE+vdW									
$\Delta E$	–17.4	–15.9	–12.0	–4.3	–4.3	–4.2	–5.6	–5.5	–5.2
$R$	2.21	2.17	2.18	3.35	3.35	3.39	3.14	3.07	3.33

relax (i.e., was optimized). It should also be noted that our test calculations using the B3LYP hybrid functional (data not shown) predicted no binding at all for gold in the (t) position on graphene, which would appear to support the hypothesis that gold binds only very weakly to graphene surfaces.

On the other hand, LDA calculations gave binding energies of 12.6, 12.2, and 10.3 kcal/mol for the (t), (b), and (h) positions, respectively (data not shown). These energies are twice as high as the benchmark values calculated for the coronene–Au complex, indicating unphysical overbinding by the LDA method; the electronic structures of coronene and graphite are certainly not sufficiently dissimilar to account for this discrepancy. These results clearly show that the LDA is inadequate for modeling the interactions of graphene with gold atoms or surfaces.

The PBE+vdW method gives rather uniform binding energies of 6.3, 6.4, and 6.2 kcal/mol for the (t), (b), and (h) positions, respectively. These values and the differences between them are in good agreement with the MP2 values calculated for the coronene–Au complex. In the original paper by Dion et al.<sup>49</sup> the authors suggested to replace the PBE exchange energy by its revPBE counterpart to obtain more accurate binding energies. We tested this scheme, which is becoming more and more

popular, for the graphene–M (M = Pd, Au) complexes. The revPBE+vdW binding energies of 3.8, 3.9, and 3.9 kcal/mol for the (t), (b), and (h) positions in graphene–Au complex, respectively, are approximately two-times lower than the corresponding benchmark energies for coronene–Au complexes. The same applies also for graphene–Pd complexes (see the following paragraph). The small differences of the electronic structure between graphene and coronene, discussed in the previous paragraph, imply that the revPBE+vdW binding energies are significantly underestimated and that the revPBE+vdW method cannot be recommended for such type of calculations.

The carbon–metal bonding distances are longer than those in the benzene–Au complex because of the greater exchange repulsion between the Au atom and the carbon atoms in the graphene sheet. The elongation of bonding distances with respect to benzene complex is consistent with the elongation of the metal–carbon bond observed in the coronene–Au complex. As was the case with the benzene–M complexes, the incorporation of a fraction of the exact exchange energy slightly reduced the calculated binding energies. The EE+vdW binding energies for the (t), (b), and (h) positions were 5.6, 5.5, and 5.4 kcal/mol, respectively. It should be noted that the preferred (t) position of gold on graphene surface agrees with recent experimental data.<sup>67</sup> The incorporation of exact exchange reduces the distances between the metal atom and the graphene sheet, which are 3.14, 3.07, and 3.33 Å for the (t), (b), and (h) positions, respectively. The distance between the metal atom and the plane containing the arene increases consistently on going from benzene to coronene to graphene, and the calculated binding energies for the adsorption of gold atoms on graphene are somewhat lower than those for the coronene–Au complex. This is largely due to the underestimation of the charge-transfer contribution in the pure PBE GGA calculation, which is highlighted when one compares the results for the graphene and benzene complexes.

The energies of the graphene–Pd complex shown in Figure 16 continue the trend observed on going from benzene to coronene. The interaction energies for the top (t) and bond (b) positions are slightly lowered in comparison with benzene, whereas the energy of the hollow (h) site is higher. The PBE+vdW energies are 20.1, 18.3, and 15.6 kcal/mol and agree very well with those for the coronene–Pd complex. The only difference is that the (t) position is predicted to be the most stable for graphene, whereas the MP2 results for coronene predict that the above-bond position (b) is the most stable. The adsorption of Pd at the (b) position results in the formation of a partial covalent bond with neighboring carbon atoms, as was demonstrated by means of an overlap population analysis in the preceding section. The EE+vdW binding energies for the (t), (b), and (h) positions were 17.4, 15.9, and 12.0 kcal/mol, respectively. For the sake of completeness, the revPBE+vdW binding energies for the (t), (b), and (h) positions were 12.8, 10.9, and 8.2 kcal/mol, respectively.

Finally, examination of the energy profiles for the graphene–Ag complex reveals that silver atoms bind a little more strongly to graphene than to benzene, primarily because of stronger vdW (dispersion) interactions. In this case, the pure GGA predicts only very weak bonding of  $\sim 0.6$  kcal/mol, at large equilibrium distances of around 3.5 Å. Because of the interaction between the silver atom and the graphene sheet is dominated by dispersion forces, the energetic differences between the three adsorption sites examined were negligible. Adsorbed silver atoms can thus easily slide over a graphene surface; the barriers to their diffusion

relate primarily to the buckling of the graphene sheet, which is most pronounced at the hollow site.<sup>53</sup>

The interaction energies for the graphene–Ag complex calculated using the PBE+vdW method were 4.3, 4.3, and 4.2 kcal/mol for the (t), (b), and (h) positions, respectively. The vdW term is obviously dominant and is essential for accurately modeling the adsorption of Ag on graphene. Our results agree very well with PBE+vdW values published by Amft et al. (Ag: 4.5, 4.5, and 4.4 kcal/mol for the (t), (b), and (h) positions, respectively)<sup>53</sup> and also with the MP2 values calculated for the coronene–Ag complex. In accord with the negligible electrostatic interaction between silver and graphene, the inclusion of an exact exchange correction has little impact on the calculated interaction energies, changing them by less than 0.1 kcal/mol for all positions. As such, the PBE+vdW values can effectively be regarded as the benchmark in this case.

On comparing the results for coronene and graphene, it is apparent that MP2 and EE+vdW strongly favor the adsorption of Pd over Au or Ag. Moreover, these two methods both yield similar ratios for the binding energy of Pd relative to Au and Ag; the ratio for MP2 is 9:5:2, while that for EE+vdW is 9:3:2. The two methods also predict similar behavior for the binding energy on switching from benzene to coronene or graphene, in terms of both overall trends and absolute values.

## 5. CONCLUSIONS

WFT and DFT calculations performed for the benzene–M and coronene–M complexes (M = Ag, Au, Pd) indicate that Pd is bound most strongly, followed by Au and Ag. The difference in binding energy between the strongest and weakest complexes is, however, reduced on going from benzene to coronene. The nature of the adsorption of these three elements is different. While silver binds primarily via dispersion forces in both cases, the binding of gold is primarily attributable to charge-transfer interactions between the electron donor (benzene or coronene) and the electron acceptor (the gold atom). Relativistic effects are important in the binding of gold, and their neglect leads to dramatic underestimation of the binding energy. The binding of Pd is quite different again; it forms a (partial) covalent bond with the arene.

The CCSD(T) benchmark binding energies for the benzene–M (M = Pd, Au, Ag) complexes were 19.7, 4.2, 2.4 kcal/mol, respectively; the MP2 binding energies for the coronene–M (M = Pd, Au, Ag) complexes were 17.7, 7.0, 4.1 kcal/mol, respectively. These numbers indicate that the nature of the binding of the metal atoms does not change dramatically on going from benzene to coronene and that the values obtained at the benchmark CCSD(T) level can thus be used to characterize the adsorption of metals on a carbon surface.

Comparison between the reference CCSD(T) and plane-wave DFT calculations demonstrates that neither LDA nor GGA provide reliable binding energies. On the other hand, PBE+vdW performs well, but surprisingly, the revPBE+vdW underbinds studied complexes. The most accurate plane-wave DFT method identified was PBE+vdW with an exact exchange correction; referred here as EE+vdW. Using this method, the binding energies calculated for the benzene–M and graphene–M (M = Pd, Au, Ag) complexes were 18.7, 5.1, and 2.5 kcal/mol and 17.4, 5.6, and 4.3 kcal/mol, respectively. The values obtained for the benzene complexes agree with the benchmark CCSD(T) energies to within chemical accuracy

(~1 kcal/mol). Moreover, calculations using this method accurately reproduced the trends in binding energy observed on switching from benzene to coronene or graphene as well as the corresponding absolute reference values. By comparing the pure GGA binding energies to those calculated using the nonlocal vdW correlation, it was demonstrated that the vdW corrections are purely attractive only in Ag complexes; in Pd complexes, they are repulsive and serve to correct the overbinding predicted by the PBE method. This implies that using empirical corrections to simulate dispersion interactions can be counterproductive when studying graphene–metal systems, since corrections of this kind will always favor binding.

The good agreement obtained with two rather different computational methods (specifically, wave function-based CCSD(T) and MP2 with a local basis set and the density functional-based EE+vdW method, with a plane-wave basis set) indicates that the calculated graphene binding energies reported in this paper can be used as reliable benchmark values and that EE+vdW is a useful and practical method for accurate computational studies of extended systems. Moreover, it also demonstrates that coronene complexes are useful model systems for modeling adsorption on graphene with chemical accuracy.

## AUTHOR INFORMATION

### Corresponding Author

\*E-mail: michal.otyepka@upol.cz; pavel.hobza@uochb.cas.cz.

### Author Contributions

<sup>||</sup>These authors contributed equally to this work.

## ACKNOWLEDGMENT

This work was a part of research project no. Z40550506 of the Institute of Organic Chemistry and Biochemistry, Academy of Sciences of the Czech Republic. It was also supported by the Korea Science and Engineering Foundation (World Class University program R32-2008-000-10180-0), by grants no. LC512 and MSM6198959216 from the Ministry of Education, Youth and Sports of the Czech Republic and by grant No. P208/10/1742 from the Grant Agency of the Czech Republic. It was also supported by the operational program Research and Development for Innovations of European Regional Development Fund (CZ.1.05/2.1.00/03.0058) and the Operational Program Education for Competitiveness of European Social Fund (CZ.1.07/2.3.00/20.0017). The support of Praemium Academiae, Academy of Sciences of the Czech Republic, awarded to P.H. in 2007 is also acknowledged.

## REFERENCES

- (1) Sundaram, R. S.; Steiner, M.; Chiu, H.-Y.; Engel, M.; Bol, A. A.; Krupke, R.; Burghard, M.; Kern, K.; Avouris, P. *Nano Lett.* **2011**, *11*, 3833.
- (2) Baby, T. T.; Aravind, S. S. J.; Arockiadoss, T.; Rakhi, R. B.; Ramaprabhu, S. *Sens. Actuators, B* **2010**, *145*, 71.
- (3) Hong, W. J.; Bai, H.; Xu, Y. X.; Yao, Z. Y.; Gu, Z. Z.; Shi, G. Q. *J. Phys. Chem. C* **2010**, *114*, 1822.
- (4) Li, Y.; Fan, X. B.; Qi, J. J.; Ji, J. Y.; Wang, S. L.; Zhang, G. L.; Zhang, F. B. *Mater. Res. Bull.* **2010**, *45*, 1413.
- (5) Li, Y.; Fan, X. B.; Qi, J. J.; Ji, J. Y.; Wang, S. L.; Zhang, G. L.; Zhang, F. B. *Nano Res.* **2010**, *3*, 429.
- (6) Scheuermann, G. M.; Rumi, L.; Steurer, P.; Bannwarth, W.; Mühlaupt, R. *J. Am. Chem. Soc.* **2009**, *131*, 8262.

- (7) Shan, C. S.; Yang, H. F.; Han, D. X.; Zhang, Q. X.; Ivaska, A.; Niu, L. *Biosens. Bioelectron.* **2010**, *25*, 1070.
- (8) Xiong, Z. G.; Zhang, L. L.; Ma, J. Z.; Zhao, X. S. *Chem. Commun.* **2010**, *46*, 6099.
- (9) Jensen, P.; Blase, X.; Ordejón, P. *Surf. Sci.* **2004**, *564*, 173.
- (10) Wang, G. M.; BelBruno, J. J.; Kenny, S. D.; Smith, R. *Phys. Rev. B* **2004**, *69*, 195412.
- (11) Akola, J.; Häkkinen, H. *Phys. Rev. B* **2006**, *74*, 165404.
- (12) Amft, M.; Sanyal, B.; Eriksson, O.; Skorodumova, N. V. *J. Phys.: Condens. Matter* **2011**, *23*, 205301.
- (13) Chan, K. T.; Neaton, J. B.; Cohen, M. L. *Phys. Rev. B* **2008**, *77*, 235430.
- (14) Jalkanen, J. P.; Halonen, M.; Fernandez-Torre, D.; Laasonen, K.; Halonen, L. *J. Phys. Chem. A* **2007**, *111*, 12317.
- (15) Varns, R.; Strange, P. J. *Phys.: Condens. Matter* **2008**, *20*, 225005.
- (16) Arthur, J. R.; Cho, A. Y. *Surf. Sci.* **1973**, *36*, 641.
- (17) Da Silva, A. J. R.; Carrijo-Faria, J.; da Silva, E. Z.; Fazzio, A. *Nanotechnology* **2003** vol 3, 2003.
- (18) Wang, G. M.; BelBruno, J. J.; Kenny, S. D.; Smith, R. *Surf. Sci.* **2003**, *541*, 91.
- (19) Yagi, Y.; Briere, T. M.; Sluiter, M. H. F.; Kumar, V.; Farajian, A. A.; Kawazoe, Y. *Phys. Rev. B* **2004**, *69*, 075414.
- (20) López-Corral, I.; Germán, E.; Juan, A.; Volpe, M. A.; Brizuela, G. P. *J. Phys. Chem. C* **2011**, *115*, 4315.
- (21) Riley, K. E.; Pitoňák, M.; Jurečka, P.; Hobza, P. *Chem. Rev.* **2010**, *110*, 5023.
- (22) Neogrady, P.; Urban, M. *Int. J. Quantum Chem.* **1995**, *55*, 187.
- (23) Neogrady, P.; Urban, M.; Hubač, I. In *Recent Advances in Coupled-Cluster Methods*; Bartlett, R. J., Ed.; World Scientific: Singapore, 1997, p 275.
- (24) Watts, J. D.; Gauss, J.; Bartlett, R. J. *J. Chem. Phys.* **1993**, *98*, 8718.
- (25) Møller, C.; Plesset, M. S. *Phys. Rev.* **1934**, *46*, 0618.
- (26) Grimme, S.; Antony, J.; Ehrlich, S.; Krieg, H. *J. Chem. Phys.* **2010**, *132*, 154104.
- (27) Zhao, Y.; Truhlar, D. G. *J. Chem. Phys.* **2006**, *125*, 194101.
- (28) Zhao, Y.; Truhlar, D. G. *J. Chem. Theory Comput.* **2007**, *3*, 289.
- (29) Zhao, Y.; Truhlar, D. G. *Theor. Chem. Acc.* **2008**, *120*, 215.
- (30) Heckert, M.; Heun, O.; Gauss, J.; Szalay, P. G. *J. Chem. Phys.* **2006**, *124*, 124105.
- (31) Iliáš, M.; Kellö, V.; Urban, M. *Acta Phys. Slovaca* **2010**, *60*, 259.
- (32) Douglas, M.; Kroll, N. M. *Ann. Phys.* **1974**, *82*, 89.
- (33) Hess, B. A.; Chandra, P. *Phys. Scr.* **1987**, *36*, 412.
- (34) Roos, B. O.; Lindh, R.; Malmqvist, P. A.; Veryazov, V.; Widmark, P. O. *J. Phys. Chem. A* **2004**, *108*, 2851.
- (35) Roos, B. O.; Lindh, R.; Malmqvist, P. A.; Veryazov, V.; Widmark, P. O. *J. Phys. Chem. A* **2005**, *109*, 6575.
- (36) Kellö, V.; Sadlej, A. J. *Theor. Chim. Acta* **1996**, *94*, 93.
- (37) Sadlej, A. J. *Collect. Czech. Chem. Commun.* **1988**, *53*, 1995.
- (38) VanLenthe, J. H.; VanDuijneveldt-van de Rijdt, J. G. C. M.; VanDuijneveldt, F. B. *Adv. Chem. Phys.* **1987**, *69*, 521.
- (39) Antušek, A.; Urban, M.; Sadlej, A. J. *J. Chem. Phys.* **2003**, *119*, 7247.
- (40) Granatier, J.; Urban, M.; Sadlej, A. J. *J. Phys. Chem. A* **2007**, *111*, 13238.
- (41) Granatier, J.; Urban, M.; Sadlej, A. J. *Chem. Phys. Lett.* **2010**, *484*, 154.
- (42) Boys, S. F.; Bernardi, F. *Mol. Phys.* **1970**, *19*, 553.
- (43) Karlström, G.; Lindh, R.; Malmqvist, P. A.; Roos, B. O.; Ryde, U.; Veryazov, V.; Widmark, P. O.; Cossi, M.; Schimmelpfennig, B.; Neogrady, P.; Seijo, L. *Comput. Mater. Sci.* **2003**, *28*, 222.
- (44) Ahlrichs, R.; Bär, M.; Häser, M.; Horn, H.; Kölmel, C. *Chem. Phys. Lett.* **1989**, *162*, 165.
- (45) Frisch, M. J.; Trucks, G. W.; Schlegel, H. B.; Scuseria, G. E.; Robb, M. A.; Cheeseman, J. R.; Scalmani, G.; Barone, V.; Mennucci, B.; Petersson, G. A.; Nakatsuji, H.; Caricato, M.; Li, X.; Hratchian, H. P.; Izmaylov, A. F.; Bloino, J.; Zheng, G.; Sonnenberg, J. L.; Hada, M.; Ehara, M.; Toyota, K.; Fukuda, R.; Hasegawa, J.; Ishida, M.; Nakajima, T.; Honda, Y.; Kitao, O.; Nakai, H.; Vreven, T.; Montgomery, J. A., Jr.; Peralta, J. E.; Ogliaro, F.; Bearpark, M.; Heyd, J. J.; Brothers, E.; Kudin, K. N.; Staroverov, V. N.; Kobayashi, R.; Normand, J.; Raghavachari, K.; Rendell, A.; Burant, J. C.; Iyengar, S. S.; Tomasi, J.; Cossi, M.; Rega, N.; Millam, N. J.; Klene, M.; Knox, J. E.; Cross, J. B.; Bakken, V.; Adamo, C.; Jaramillo, J.; Gomperts, R.; Stratmann, R. E.; Yazyev, O.; Austin, A. J.; Cammi, R.; Pomelli, C.; Ochterski, J. W.; Martin, R. L.; Morokuma, K.; Zakrzewski, V. G.; Voth, G. A.; Salvador, P.; Dannenberg, J. J.; Dapprich, S.; Daniels, A. D.; Farkas, Ö.; Foresman, J. B.; Ortiz, J. V.; Cioslowski, J.; Fox, D. J. *Gaussian 09*, revision A.02; Gaussian, Inc.: Wallingford, CT, 2009.
- (46) Blöchl, P. E. *Phys. Rev. B* **1994**, *50*, 17953.
- (47) Kresse, G.; Joubert, D. *Phys. Rev. B* **1999**, *59*, 1758.
- (48) Perdew, J. P.; Burke, K.; Ernzerhof, M. *Phys. Rev. Lett.* **1996**, *77*, 3865.
- (49) Dion, M.; Rydberg, H.; Schröder, E.; Langreth, D. C.; Lundqvist, B. I. *Phys. Rev. Lett.* **2005**, *95*, 109902.
- (50) Lazić, P.; Atodiresei, N.; Alaei, M.; Caciuc, V.; Blügel, S.; Brako, R. *Comput. Phys. Commun.* **2010**, *181*, 371.
- (51) Klimeš, J.; Bowler, D. R.; Michaelides, A. *J. Phys.: Condens. Matter* **2010**, *22*, 022201 and references therein.
- (52) BelBruno, J. J. *Surf. Sci.* **2005**, *577*, 167.
- (53) Amft, M.; Lebègue, S.; Eriksson, O.; Skorodumova, N. V. *J. Phys.: Condens. Matter* **2011**, *23*, 395001.
- (54) Peterson, K. A.; Puzzarini, C. *Theor. Chem. Acc.* **2005**, *114*, 283.
- (55) Peterson, K. A.; Figgen, D.; Dolg, M.; Stoll, H. *J. Chem. Phys.* **2007**, *126*, 124101.
- (56) Ishikawa, T. *Jpn. J. Appl. Phys.* **1993**, *32*, 4779.
- (57) Scheer, M.; Brodie, C. A.; Bilodeau, R. C.; Haugen, H. K. *Phys. Rev. A* **1998**, *58*, 2051.
- (58) Neogrady, P.; Kellö, V.; Urban, M.; Sadlej, A. J. *Int. J. Quantum Chem.* **1997**, *63*, 557.
- (59) Looock, H. P.; Beaty, L. M.; Simard, B. *Phys. Rev. A* **1999**, *59*, 873.
- (60) Bilodeau, R. C.; Scheer, M.; Haugen, H. K. *J. Phys. B: At. Mol. Opt. Phys.* **1998**, *31*, 3885.
- (61) Brown, C. M.; Ginter, M. L. *J. Opt. Soc. Am.* **1978**, *68*, 243.
- (62) Hotop, H.; Lineberger, W. C. *J. Phys. Chem. Ref. Data* **1985**, *14*, 731.
- (63) Lee, K.; Murray, E. D.; Kong, L.; Lundqvist, B. I.; Langreth, D. C. *Phys. Rev. B* **2010**, *82*, 081101.
- (64) Dargel, T. K.; Hertwig, R. H.; Koch, W. *Mol. Phys.* **1996**, *96*, 583.
- (65) Yi, H.-B.; Lee, H. M.; Kim, K. S. *J. Chem. Theory Comput.* **2009**, *5*, 1709.
- (66) Yi, H.-B.; Diefenbach, M.; Choi, Y. C.; Lee, E. C.; Lee, H. M.; Hong, B. H.; Kim, K. S. *Chem.—Eur. J.* **2006**, *12*, 4885.
- (67) Zan, R.; Bangert, U.; Ramasse, Q.; Novoselov, K. S. *Nano Lett.* **2011**, *11*, 1087.

# Development of Force Field Parameters for Molecular Simulation of Polylactide

James H. McAliley and David A. Bruce\*

Department of Chemical and Biomolecular Engineering, Clemson University, Clemson, South Carolina 29634-0909, United States

**S** Supporting Information

**ABSTRACT:** Polylactide is a biodegradable polymer that is widely used for biomedical applications, and it is a replacement for some petroleum-based polymers in applications that range from packaging to carpeting. Efforts to characterize and further enhance polylactide-based systems using molecular simulations have to this point been hindered by the lack of accurate atomistic models for the polymer. Thus, we present force field parameters specifically suited for molecular modeling of PLA. The model, which we refer to as PLAFF3, is based on a combination of the OPLS and CHARMM force fields, with modifications to bonded and nonbonded parameters. Dihedral angle parameters were adjusted to reproduce DFT data using newly developed CMAP dihedral cross terms, and the model was further adjusted to reproduce experimentally resolved crystal structure conformations, melt density, volume expansivity, and the glass transition temperature of PLA. We recommend the use of PLAFF3 in modeling PLA in its crystalline or amorphous states and have provided the necessary input files required for the publicly available molecular dynamics code GROMACS.

## INTRODUCTION

Polylactide, also called polylactic acid (PLA), is an important polymer for biomedical applications, because it is compatible with living cells and is biodegradable.<sup>1,2</sup> Further, PLA is of interest as a commodity polymer, and is used especially in single-use packaging applications.<sup>3</sup> PLA is an  $\alpha$ -polyester, and the primary structure of its repeat unit is shown in Figure 1.

Classical molecular force fields, such as CHARMM<sup>4</sup> and OPLS,<sup>5</sup> have been widely used in recent decades for simulating organic molecules, by and large with good success. However, neither force field has been parametrized specifically for the dihedral angles present in  $\alpha$ -polyesters such as PLA. In particular, dihedral interaction parameters for the  $O^S-C-C^\alpha-O^S$ ,  $C-C^\alpha-O^S-C$ , and  $C^\alpha-O^S-C-C^\alpha$  motifs, all of which are unique to  $\alpha$ -polyesters, are not found in the parameter databases for these force fields. In lieu of these specific four-atom interaction parameters, one would typically use the so-called *wildcard* parameters included in the force field (these are general parameters represented by  $X-C-C^\alpha-Y$ ,  $X-C^\alpha-O^S-Y$ , and  $X-O^S-C-Y$ , where  $X$  and  $Y$  may be any atom type). Though wildcard parameters provide a reasonable guess for the dihedral interactions in cases where more accurate parameters are unavailable, it has been shown that use of the wildcard parameters for  $\alpha$ -polyesters results in poor accuracy when modeling PLA because the wildcard rotational energy barriers are centered at the wrong dihedral angles and do not describe the barrier heights predicted via quantum (DFT and MP2 level) models well.<sup>6,7</sup>

In this work, we develop a classical force field model specifically suited for polylactides, based on the OPLS and CHARMM forms. The present force field follows the work of O'Brien,<sup>7</sup> in which the PLAFF model was developed and validated extensively for crystalline PLA, and of McAliley,<sup>6</sup> in which the model was further developed for accuracy in modeling amorphous PLA (the resulting force field was referred to as PLAFF2). Our present

model is PLAFF3, and it differs from previous versions in its use of the CMAP cross-term dihedral potential originally developed for CHARMM.<sup>8</sup> This potential function provides more flexibility in fitting barriers to bond rotation and has allowed us to fit the glass transition temperature of PLA with far more accuracy than could be obtained from linear combinations of individual dihedral potentials. We demonstrate that PLAFF3 is better suited than its predecessors for modeling the amorphous dihedral angle distributions in PLA, and that it retains the accuracy of PLAFF in simulating the crystalline phase. We believe that these parameters will be of value to the biological science community in studying PLA. Further, with the growing interest in using renewable polymers for commodity packaging applications, this model will likely be of use to the materials science community in exploring new PLA-based materials. Though other work on PLA force field development has appeared in the literature,<sup>9–12</sup> to our knowledge, the PLAFF3 parameters represent the first noncommercial molecular model validated against electronic structure calculations and experimental data for PLA in melt, glassy, and crystalline phases. As such, we hope this work will allow a larger number of researchers to study the material through simulation than was previously feasible.

## METHODS

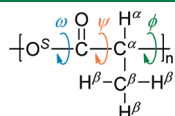
The fitting procedure used in this work is shown schematically in Figure 2. The procedure begins with assembling target data and providing an initial guess for the force field parameters. As a first step, the torsional parameters are adjusted to match DFT data obtained in previous work.<sup>13</sup> Next, the model is tested against experimental crystal structure data for PLA. Dihedral parameters are then adjusted accordingly, until reasonable

**Received:** April 14, 2011

**Published:** September 19, 2011



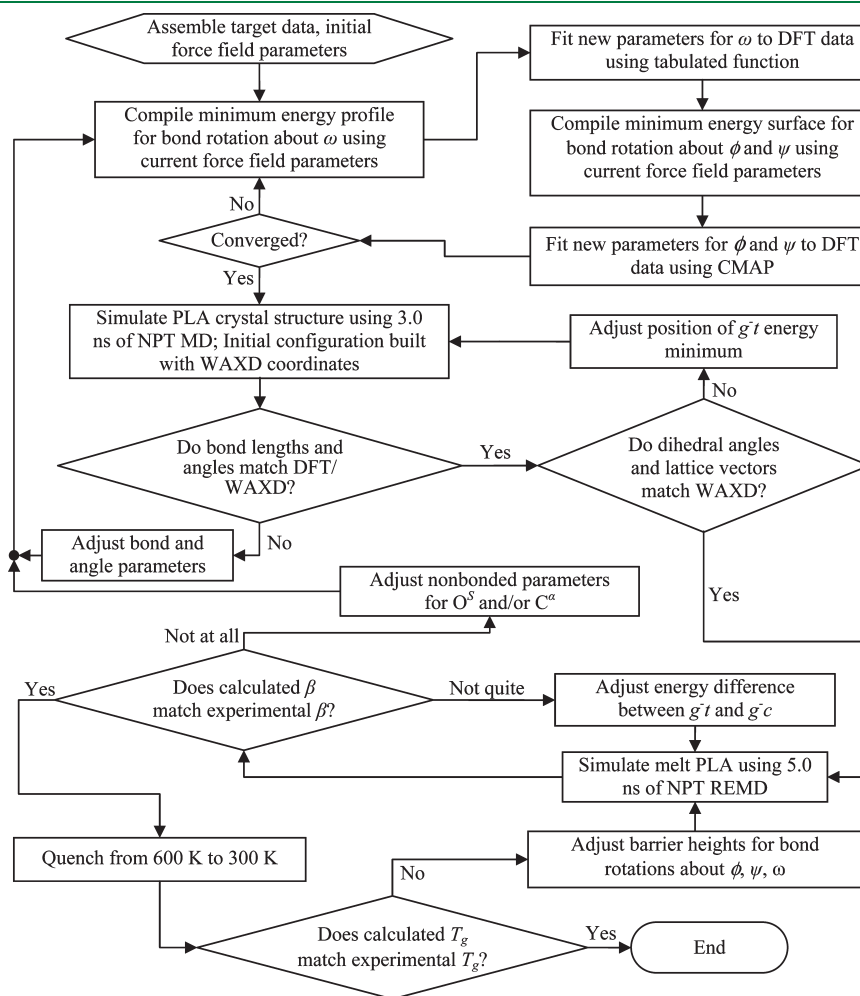
agreement is obtained with the experimental data for crystalline PLA. Following this step, the model is used to simulate the polymer in its melt state. The volume expansivity,  $\beta$ , is estimated from these simulations and compared with experimental dilatometric measurements. In some cases, an adjustment of the relative energies between energy minima can affect a change in  $\beta$ , by establishing a different temperature dependence of the polymer's rotational isomeric states. However, if a more drastic change is required, the nonbonded parameters are adjusted for those atom types that are unique to  $\alpha$ -polyesters, until the density and volume expansivity are near experimental values. After such adjustments, the entire fitting procedure must be repeated to ensure the agreement with DFT and that crystal structure data is maintained. Finally, the model is used in quench simulations, where the polymer is rapidly cooled from the melt



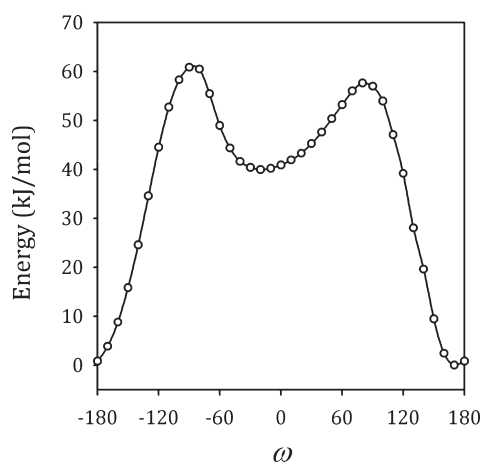
**Figure 1.** Chemical structure of PLA. Superscript labels on atoms are used for reference in the text. The three main chain bond rotations are labeled according to the convention for polypeptides.

state into the glassy state. Using the Williams–Landel–Ferry (WLF) equation, the resulting glass transition temperature,  $T_g$ , may be compared to experimental measurements. If necessary, the energy barriers are then adjusted for rotation about each main chain dihedral angle until agreement is reached with experimental  $T_g$  values. Each of these steps will be discussed in greater detail in the following sections.

**Initial Force Field Parameters.** We considered two force fields as a starting point for the PLAFF models: The Optimized Potentials for Liquid Simulations<sup>5</sup> (OPLS) and the force field from the Chemistry at Harvard Molecular Mechanics (CHARMM) package.<sup>4</sup> The OPLS parameters were taken from the OPLS-AA parameter files as distributed with GROMACS<sup>14</sup> version 3.3.3. For the CHARMM force field, parameters were taken from the CHARMM27 protein–lipid parameter files distributed with CHARMM version c32b2. Atom types were assigned on the basis of chemical functionality (see Supporting Information). Partial atomic charges were unaltered in each force field, with the exception of main-chain atoms and the carbonyl oxygen, which were adjusted slightly to achieve charge neutrality in the lactyl residue and to improve agreement with DFT results. The needed CHARMM27 parameters for PLA were ported into GROMACS, and all further molecular mechanics calculations reported for the CHARMM force field were performed in GROMACS



**Figure 2.** Flow diagram showing the procedure for fitting PLA force field parameters. The dihedral angles for  $\phi$  and  $\psi$  that define the  $g^{-t}$  and  $g^{-c}$  energy minima are shown later in Figure 4.



**Figure 3.** DFT potential energy profile for rotation of the  $O^S-C$  bond through all values of dihedral angle  $\omega$ . Calculations performed on a PLA trimer *in vacuo*.<sup>13</sup>

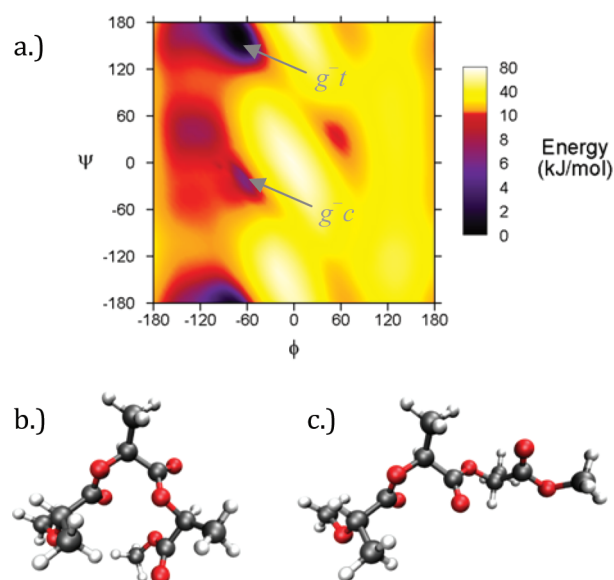
version 3.3.3. It is important to note that the parameters, especially the nonbonded parameters, from the OPLS and CHARMM force fields are not fully compatible. This mixed set of parameters was only used as a starting point for further parameter optimization.

In the fitting procedure, bond stretching and angle bending parameters were optimized to more accurately reproduce energy potentials predicted using DFT (see the Supporting Information). In general, DFT methods are known to give accurate geometries, whereas they are less accurate at predicting vibrational frequencies. For this reason, geometric parameters for bonds and angles (the  $b_0$  parameter for bonds and the  $\theta_0$  parameter for angles) were fit to DFT data, but the bond and angle force constants ( $k_b$  and  $k_\theta$ , respectively) were unaltered from their original OPLS values. In this way, we deviated as little as possible from the OPLS model.

Three dihedral interactions were also adjusted to achieve better agreement with the bond rotational potential energy surfaces calculated from DFT. These correspond to the backbone dihedrals labeled as  $\phi$ ,  $\psi$ , and  $\omega$  in Figure 1 and are defined by the IUPAC convention using the main chain atoms ( $O^S$ ,  $C$ , and  $C^\alpha$ ). In PLAFF3, the potential energies for these three dihedral interactions were represented by tabulated functions. For rotation about  $\omega$ , a one-dimensional tabulated function was used (see the GROMACS User Manual<sup>14</sup> for more information), whereas a two-dimensional tabulated function (also called a correction map or CMAP<sup>15</sup>) was used for each pair of neighboring  $\phi$ ,  $\psi$  dihedrals. The CMAP potential was recently implemented in GROMACS,<sup>8</sup> and all calculations involving such terms were performed with GROMACS version 4.5.1.

**Target Data.** We used several criteria to select target data for parameter fitting. One criterion was that the model should be consistent with results from higher-level molecular simulation methods, namely, our DFT results from previous work.<sup>13</sup> In addition, we aimed to be consistent with experimental results. Because PLA is often used in its semicrystalline form, we desired a model that could reproduce the properties of both the crystalline and amorphous states of the material. Thus, conformational data for the crystalline form of PLA were used, as well as the glass transition temperature and volumetric data for the amorphous polymer.

**DFT Data.** The target potential energy values were taken from *in vacuo* DFT calculations (B3LYP/6-31G\*\*) for a methyl



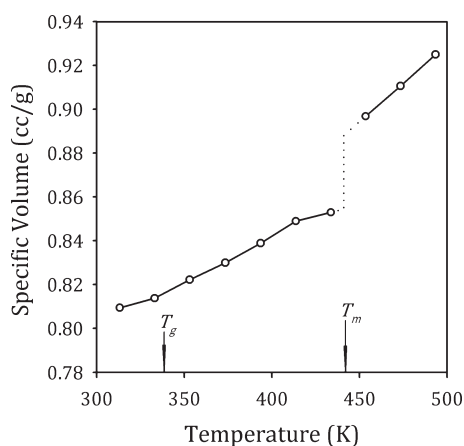
**Figure 4.** (a) DFT potential energy surface, with dihedral angles  $\phi$  and  $\psi$  as independent variables. Calculations performed on a PLA trimer *in vacuo*.<sup>13</sup> The two lowest energy minima, *gauche*<sup>-</sup>/*trans* ( $g^-t$ ) and *gauche*<sup>-</sup>/*cis* ( $g^-c$ ) are indicated. (b) Molecular geometry of  $g^-c$  conformation. (c)  $g^-t$  conformation.

terminated PLA trimer ( $(CH_3(OC(O)CH(CH_3))_3OCH_3)$ , as reported previously.<sup>13</sup> These include an estimate of the potential energy barriers encountered with bond stretching and angle bending (see the Supporting Information), during rotation of the  $O^S-C$  bond described by  $\omega$  (reproduced in Figure 3), and likewise for simultaneous rotation of bonds described by  $\phi$  and  $\psi$  (Figure 4) for the central repeat unit in a PLA trimer.

**Crystal Structure Data.** Several studies on the crystal structure of PLA have appeared in the literature.<sup>16–20</sup> We have chosen to use the structural coordinates from Sasaki and Asakura<sup>20</sup> as our target data. The authors' use of the linked atom full-matrix least-squares (LAFLS) method<sup>21</sup> and the Rietveld whole-fitting method<sup>22</sup> allowed for the positions of individual atoms in the unit cell to be determined with a high degree of accuracy. The authors derived the  $\alpha$ -form of the crystal structure from WAXD data, resulting in a frustrated  $10_3$  helix. The orthorhombic ( $\alpha = \beta = \gamma = 90^\circ$ ) unit cell from that study has  $P2_12_12_1$  symmetry and lattice constants of  $a = 10.66(1)$  Å,  $b = 6.16(1)$  Å, and  $c = 28.88(2)$  Å.<sup>20</sup>

**Volumetric Data for Amorphous/Melt PLA.** To represent the volumetric properties of amorphous PLA, we selected the experimental data from Sato et al.,<sup>23</sup> where the specific volumes of polylactide samples were measured at various temperatures and pressures using metal bellows dilatometry. While the experimental data cover a wide range of temperatures, the only specific volume data used in this study are plotted in Figure 5, and these data correspond to specific volumes measured by heating PLA samples at 1 bar above the melting temperature,  $T_m$ . It can be seen in Figure 5 that an abrupt change in volume occurs upon heating above  $T_m$ , which is attributed to the change in volume that occurs when the crystallites in the sample become amorphous.

In practice, we use molecular simulations to study the melt phase of polymers at temperatures higher than those shown in Figure 5, utilizing the well-known time–temperature superposition principle for polymers.<sup>24</sup> Thus, we look at the volume expansivity (see eq 1 below) to facilitate a comparison. From the data



**Figure 5.** Target volume–temperature data at 1 bar, taken from Sato et al.<sup>23</sup> Arrows indicate the authors' estimate of the glass transition temperature,  $T_g$ , and the melting temperature,  $T_m$ , taken from separate calorimetry data.

in Figure 5, an expansivity of  $\beta = 7.8 \times 10^{-4} \text{ K}^{-1}$  is calculated for the melt state.

**Glass Transition Data.** It is generally accepted that barriers to bond rotation play a major role in determining the glass transition temperature of polymers. Experimentally, the value of the PLA glass transition is dependent on the method used to measure it, and it has also been reported to vary widely with moisture content.<sup>25</sup> Common reported values of  $T_g$  for PLA are in the range of 327 to 345 K (see Table 1), which were obtained using differential scanning calorimetry (DSC) and dielectric relaxation spectroscopy (DRS). The value reported by Auras was measured after extensive drying of the PLA samples.<sup>25</sup> Since water is known to have a plasticizing effect on the material, it follows that this estimate is at the high end of the reported range of  $T_g$  values.

**Fitting Procedure Using DFT Target Data.** In fitting our model to DFT data, we performed energy minimizations using the force field model, each constrained by the same independent variable(s) as the DFT energy minimizations (i.e.,  $\omega$  or  $\phi$  and  $\psi$ ) before comparing energies with the DFT results. This energy minimization step with the force field model adds a high level of nonlinearity to the fitting procedure. Each time the dihedral parameters are adjusted, the minimum energy conformation at each independent variable also changes. Thus, obtaining the optimal torsional potentials according to this prescription requires an iterative scheme.

We began each iteration of the fitting procedure with the dihedral having the largest potential energy barriers. Thus, the rotational energy barrier for the  $\text{O}^{\text{S}}-\text{C}$  bond (dihedral angle  $\omega$ ) was fit to the DFT data shown in Figure 3 first, followed by a simultaneous fit of the  $\phi$  and  $\psi$  potentials to the data shown in Figure 4. A weighted least-squares approach was used in developing PLAFF2, as described elsewhere.<sup>6</sup> However, in PLAFF3, simple adjustments to the tabulated potentials were used to match the DFT data. Weighting of data points was not necessary due to the greater flexibility of the tabulated potentials versus the traditional cosine expansions used in PLAFF2. After each parameter optimization step, the force field minimum energy conformations were re-evaluated using the most current dihedral parameters. This process was repeated until reasonable convergence was achieved with respect to the dihedral parameters and the minimized energies.

**Table 1.** Some Reported Values of the Glass Transition Temperature of PLA

lead author	method	rate	$T_g$ (K)
Dorgan <sup>26</sup>	DSC	10 °C/min	331.6
Sato <sup>23</sup>	DSC		337
Auras <sup>25</sup>	DSC	10 °C/min	344.6
Joziassé <sup>27</sup>	DSC	10 °C/min	336
Kanchanasopa <sup>28</sup>	DRS	$\tau = 100 \text{ s}$	327

**Bounded Adjustment of Dihedral Potentials.** We employ a bounded adjustment procedure for altering the dihedral potentials, to avoid drastically changing the force field parameters and thus the minimum energy conformations. For all iterations, the tabulated potentials were adjusted to be as close as possible to the DFT target data without exceeding a specified change in energy. Otherwise, if unbounded adjustments were allowed, we often observed divergent behavior due to the nonlinear aspects of the iterative scheme. It was found that suitable stability was achieved by limiting the change in energy at each tabulated point to 10 kJ/mol for fitting the  $\omega$  dihedral, and 5 kJ/mol for fitting  $\phi$  and  $\psi$  dimerals. Because these limits gave satisfactory performance, no attempt was made at further tuning the fitting procedure with respect to them.

**Refinement Using Crystal Structure Data.** Returning to our discussion of Figure 2, we proceed, after sufficient convergence is obtained in fitting to DFT bond rotation data, by examining the crystal structure of PLA with the resulting force field parameters. In these simulations, a super cell based on the crystalline unit cell is built according to the WAXD-resolved structure of Sasaki and Asakura.<sup>20</sup> This super cell contains 32 PLA chains, each containing 50 monomers, with the boundary conditions for the polymer being such that monomer number 1 in a chain was bonded to monomer number 50 from the neighboring periodic cell. The selected system size is sufficiently large that no finite size effects are observed with the simulated systems. The system is simulated for 3.0 ns in the NPT ensemble, whereby the lattice or box dimensions are allowed to adjust to their equilibrium values; however, all lattice angles ( $\alpha$ ,  $\beta$ , and  $\gamma$ ) were constrained to 90°, so as to maintain an orthorhombic unit cell that matched experimental observations. Anisotropic pressure coupling was applied with the Berendsen algorithm, such that each box length was adjusted independently.<sup>29</sup> The Nose–Hoover thermostat was used to control temperature at 300 K.<sup>30,31</sup> A cutoff of 1.0 nm was used for van der Waals interactions, while the electrostatics were treated with the Particle-Mesh Ewald (PME) method.<sup>32</sup>

While one of our goals was to have a minimized PLA system that accurately reproduced the experimentally observed dihedral angles from diffraction studies, obtaining these in the crystal structure is difficult without first having accurate bond lengths and valence angles due to packing considerations in the unit cell. Therefore, before earnestly examining the predicted dihedral values, we optimize the bond and angle force field parameters using a series of *in vacuo* DFT calculations that examine the variation in PLA system energy as a function of perturbations to each bond and angle from its minimum energy value (see Supporting Information). This is done in much the same way as fitting the dihedral parameters to DFT bond rotation data, using a self-consistent iteration scheme.

Once the prerequisite of accurate valence geometries is achieved, the dihedral angles ( $\phi$ ,  $\psi$ , and  $\omega$ ) and lattice vectors

are examined over the final 1.0 ns of dynamics of a crystal structure simulation, and their values are compared to those reported in the experimental literature. Should the simulations be inconsistent with the experimental data, the position of the  $g^-t$  energy minimum (see Figure 4) is adjusted with respect to  $\phi$  and  $\psi$ , as is the position of the *trans* energy minimum for  $\omega$ . This process is repeated until the experimental dihedral values are accurately reproduced.

**Refinement Using Melt Phase Target Data.** Once adequate agreement with the crystalline unit cell is obtained, simulations are carried out on amorphous PLA using isothermal–isobaric replica exchange molecular dynamics (NPT-REMD) as implemented in GROMACS. Four separate NPT-REMD runs were performed with unique input configurations, which were initially generated using the Amorphous Cell module in Accelrys' Materials Studio version 4.4.<sup>6</sup>

In our implementation, each replica is comprised of three chains, each containing 500 repeat units (refer to Figure 1), and two lactide molecules. The chain length was chosen to be greater than the experimental entanglement length, which is approximately 125 repeat units. Additionally, limited simulations with other chain length PLA systems showed that the calculated  $T_g$  value decreased with chain length as expected, but no detailed investigation of this effect was attempted. Lactide molecules were also included in the simulated systems because there is always a small percentage of residual lactide monomer in real polylactide samples, and these monomers have a plasticizing effect on the material. With two lactide molecules per simulation cell, our simulated PLA system contains 0.26% residual lactide on a weight basis; the specific amount of lactide present in an industrially produced PLA resin is usually less than 1%,<sup>33</sup> and 0.2 to 0.3 wt %<sup>34</sup> is common.

In each replica exchange simulation, the average volume is calculated for the melt state as a function of temperature. From this, the volume expansivity can be estimated graphically using the relation<sup>35</sup>

$$\beta = \frac{1}{v} \left( \frac{\partial v}{\partial T} \right)_p \quad (1)$$

where  $v$  is the specific volume of the system. Since each replica has the same pressure, a plot can be constructed of  $\ln v$  versus  $T$ , and  $\beta$  can be estimated from its slope. This is compared to experimental measurements of the expansivity of PLA. If satisfactory agreement is not obtained, this indicates that the non-bonded parameters may need further adjustment.

When it was necessary to alter the nonbonded parameters, the atom types for the ester oxygen ( $O^S$ ) and  $\alpha$ -carbon ( $C^\alpha$ ) were chosen for adjustment. These atom types were selected because they are the most likely to deviate from the behavior of normal esters, and no such atom types exist in OPLS or CHARMM for  $\alpha$ -polyesters. Note that, when adjustment of these atoms' non-bonded parameters is necessary, the dihedral angles must again be readjusted to preserve agreement with the crystal structure.

**Refinement Using Glass Transition Target Data.** The glass transition temperature,  $T_g$ , is commonly interpreted for polymers as the temperature below which bond rotations are kinetically trapped. That is, it is the temperature below which torsional energy barriers are crossed at rates much longer than the time scale on which the polymer is observed. As such, the value of  $T_g$  for PLA is influenced by the height of the energy barrier between the various rotational isomeric states.

Many studies have appeared in the literature examining the glass transition temperature via molecular dynamics,<sup>36–38</sup> though relatively few papers address the temporal dependence of the observed glass transition temperature.<sup>39,40</sup> It is well-known, experimentally, that the glass transition will be observed at higher temperatures when a polymer specimen is cooled at a faster rate.<sup>41</sup> This behavior is described very well, over the range of experimental time scales, by the WLF equation:<sup>42</sup>

$$\ln A_T = \frac{\left( \frac{B}{f_0} \right) (T - T_0)}{\frac{f_0}{\alpha_f} + (T - T_0)} \quad (2)$$

where  $A_T$  is the *reduced variables shift factor*,  $B$  is a constant,  $f_0$  is the fractional free volume of the polymer at the reference temperature  $T_0$ , and  $\alpha_f$  is the coefficient of expansion of the free volume. Although the quenching rates accessible to molecular dynamics simulations can differ from experimental cooling rates by 14 orders of magnitude or more, the validity of the WLF equation over such wide a temporal range has been established recently through molecular simulation.<sup>40</sup>

The glass transition temperature of PLA was estimated from our force field model by quenching the amorphous conformations from the NPT-REMD simulations, using a replica at 604.5 K as the starting structure. Simulation conditions were identical to those used in each of the NPT-REMD replicas, except that the set point of the Nose-Hoover thermostat was varied linearly with simulation time over the entire run. Each run lasted until a temperature of 300 K was reached. From each of the four NPT-REMD simulations, six separate quench runs were performed, with quench rates of 15 K/ns, 30 K/ns, 60 K/ns, 150 K/ns, 300 K/ns, and 600 K/ns. The glass transition temperature was estimated for each run by fitting a straight line to a plot of  $\ln v$  versus  $T$ , using all data points below 400 K. A second straight line was drawn through the melt data taken from the NPT-REMD runs, for temperatures above 500 K. The intersection of the two lines was taken as  $T_g$ . Such estimates were then averaged for each quench rate, and then a least-squares fit was performed using the WLF model (eq 2). The reference quench rate was taken to be 10 K/min (normal lab conditions for measuring  $T_g$ ). We found that the so-called *universal* WLF constants ( $B/f_0 = 40.16$  and  $f_0/\alpha_f = 51.6$  K) provided a good fit with the PLA force-field models, reducing the WLF equation to a single adjustable parameter,  $T_0$ . When used in this manner,  $T_0$  corresponds to the  $T_g$  observed at the reference quench rate.

## RESULTS AND DISCUSSION

In implementing the iterative procedure described in Figure 2, many intermediate sets of force field parameters were examined in this study. Here, we compare and discuss the important results from seven different models, ranging from the unaltered CHARMM and OPLS force fields to our intermediate parameter sets derived from those force fields and from the first generation PLAFF force field of O'Brien to the present version of our PLA force field, PLAFF3. Abbreviations for these different models are summarized in Table 2 for ease of reference. These include the models obtained directly after a least-squares fit to the DFT potential energies, referred to as the OPLS' and CHARMM' models. These two models demonstrate that fitting to the DFT energies alone is not sufficient to reproduce experimental data,

Table 2. Description of the Various Classical Models Discussed in the Text<sup>a</sup>

model	description
OPLS	the OPLS force field as developed by Jorgensen and co-workers <sup>5,43</sup> (all-atom version, also known as OPLS-AA)
CHARMM	the CHARMM force field as developed by Brooks and co-workers <sup>4</sup>
OPLS'	OPLS, with backbone torsional potentials refit to DFT data
CHARMM'	CHARMM, with backbone torsional potentials refit to DFT data
OPLS''	OPLS, with CHARMM nonbonded parameters substituted for O1 and C3, selected bond stretching and angle bending terms refit to DFT, and backbone torsional potentials refit to DFT data
PLAFF	The PLA force field developed by O'Brien <sup>7</sup> (all-atom version, also known as PLAFF-AA)
PLAFF2	the second version of PLAFF, with improved accuracy for simulating amorphous PLA; <sup>6</sup> the model is OPLS'', with backbone torsional parameters further adjusted to reproduce crystal structure data and to improve agreement with the experimental glass transition temperature of PLA
PLAFF3	the PLA force field developed in this work; the model is OPLS'', with backbone torsional parameters further adjusted using the CMAP potential to reproduce crystal structure data and to improve agreement with the experimental glass transition temperature of PLA

<sup>a</sup> PLAFF3 is the recommended potential for molecular simulation of PLA.

and further adjustment was required as described in Figure 2. The OPLS'' model shows results in which the parameters have been further adjusted to match the bond lengths and valence angles reported in crystal structure studies of PLA, and in which the nonbonded parameters were adjusted by trial-and-error to better match the melt density and volumetric expansivity of PLA. Finally, the PLAFF2 and PLAFF3 models improve upon OPLS'' by adjusting the dihedral parameters to match crystal structure data and glass transition temperature data. The principal difference between PLAFF2 and PLAFF3 is that PLAFF3 uses CMAP dihedral cross terms, whereas PLAFF2 uses linear combinations of single dihedral potentials.

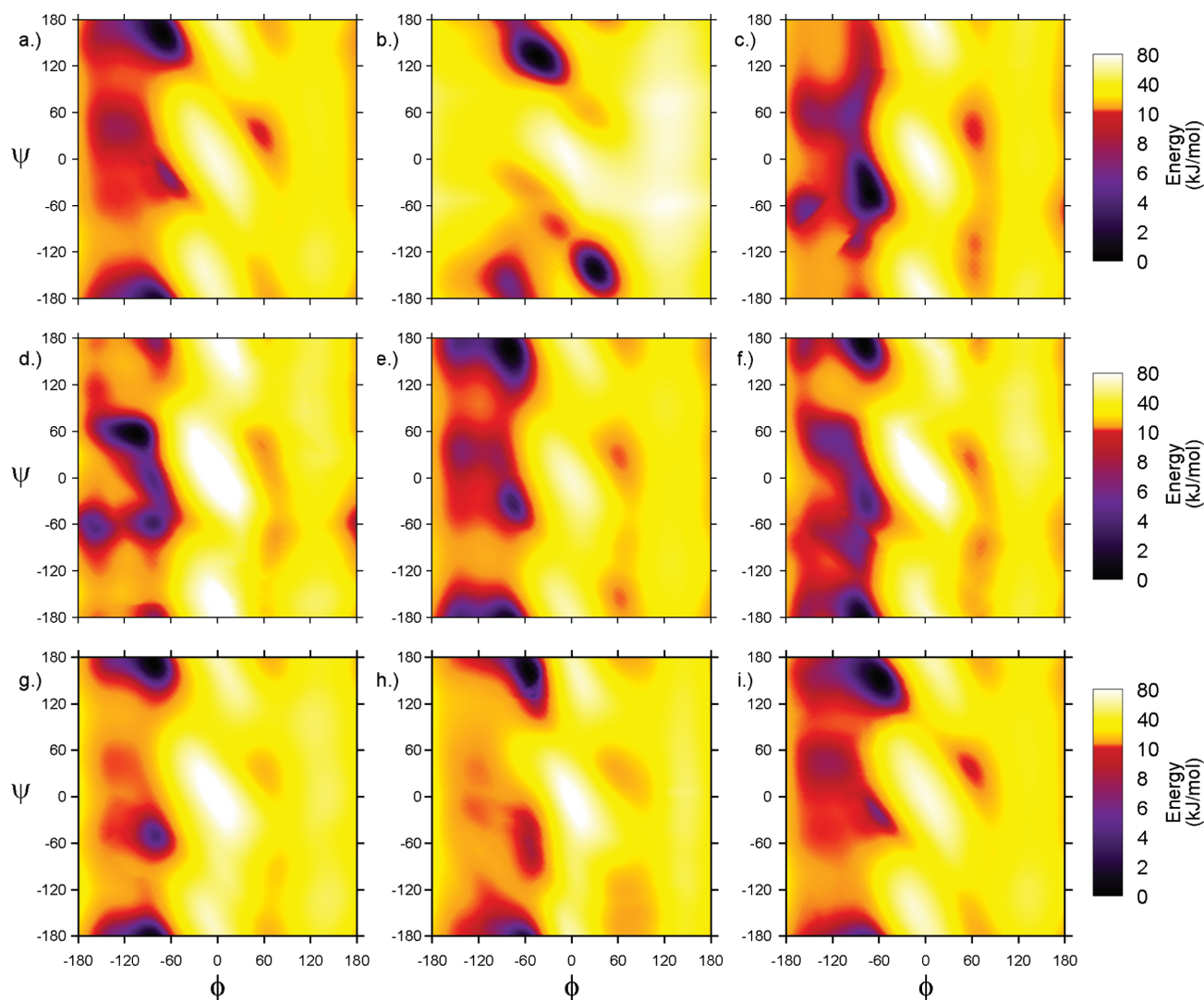
**Comparison of the Classical Models to DFT Data.** The energy landscapes for bond rotation about  $\phi$  and  $\psi$  are shown in Figure 6, calculated using the various models described in Table 2. The figure also shows the DFT target data for comparison. While we do not expect the optimum force field to be in complete agreement with DFT data, we desire the overall shape and location of relative minima/maxima to coincide with the DFT results in order to realistically model the amorphous configuration distribution. Thus, as shown in Figure 6b, the first generation PLAFF force field raises some concern, due to the presence of a low-energy local minimum in the vicinity of  $(\phi, \psi) = (30^\circ, -150^\circ)$  that does not appear in the DFT potential energy surface. Additionally, the  $g^-c$  minimum is predicted by PLAFF to be a much less probable configuration than predicted by DFT. The presence of the extra minimum in Figure 6b is only of concern for applications involving amorphous phases of PLA, in which case the entire dihedral space may be accessed by the simulated polymer chains according to the energetics of the force field model.

While the nonphysical local minimum near  $(\phi, \psi) = (30^\circ, -150^\circ)$  is a striking feature of Figure 6b, it is also obvious from the figure that O'Brien was very successful in fitting the potential energy surface in the vicinity of the global minimum (in the  $g^-t$  position shown in Figure 4). This is evidenced by the remarkable performance of PLAFF in simulating crystalline PLA,<sup>7</sup> and therefore, we feel that the original PLAFF is still very well-suited in modeling the crystalline phase of PLA. When examining the OPLS and CHARMM models in Figure 6c and d, we see that both models lack adequate representation of the global  $g^-t$  minimum predicted by DFT. This observation helps to explain

the superior performance of PLAFF in the crystalline phase as compared with OPLS and CHARMM and suggests that OPLS and CHARMM should not be used for crystalline or amorphous phase simulations without first correcting the backbone torsional potentials.

Figure 6e and f show the results of performing a least-squares fitting procedure to alter the torsional potentials of OPLS and CHARMM, while leaving all other interaction parameters in the models unchanged. This figure demonstrates that there are limitations inherent in each model, preventing a perfect fit to the desired potential energy surface. For example, the CHARMM' potential energy surface in Figure 6f still shows remnants of the local minima, situated in the negative  $\phi$  region between the  $g^-c$  and  $g^-t$  energy minima of the CHARMM model in Figure 6d. The major shortcoming of the models shown in Figure 6b–h is that corrections to the  $(\phi, \psi)$  potential energy surface are limited to linear combinations of separate functions of  $\phi$  and functions of  $\psi$ . Without the use of more sophisticated potential energy functions, e.g., the CMAP dihedral–dihedral cross terms available in recent versions of the CHARMM program,<sup>8,15</sup> accurately reproducing the entire two-dimensional potential energy surface of Figure 6a is highly dependent on the other interactions within the model, such as the bond stretching and angle bending parameters.

Figure 6g and h give the potential energy surfaces after fitting the model to crystal structure data and glass transition data, respectively (see discussion in the following sections). The last plot (Figure 6i) is our currently recommended model, PLAFF3. Note that the agreement between PLAFF2 and the DFT data is diminished when compared to OPLS'', since adjustments to fit one set of target data inevitably alters the performance of the model in reproducing all other target data. The resulting model is a compromise between competing target data. However, the addition of the CMAP dihedral term in PLAFF3 allows for a much more accurate fit of the potential energy surface, which can be altered in very specific local regions without affecting the shape of the surface elsewhere. For example, the  $g^-t$  global minimum of PLAFF3 was shifted by approximately  $+10^\circ$  and  $-15^\circ$  in  $\phi$  and  $\psi$ , respectively, to improve agreement with the crystal structure, yet the remaining portions of the surface are unaffected by this shift.



**Figure 6.** Bond rotational energy profiles for the  $\phi$  and  $\psi$  dihedrals (shown in Figure 1) of a PLA trimer, calculated from (a) B3LYP/6-31G\*\*,<sup>13</sup> (b) PLAFF,<sup>7</sup> (c) OPLS,<sup>5</sup> (d) CHARMM,<sup>4</sup> (e) OPLS', (f) CHARMM', (g) OPLS'', (h) PLAFF2,<sup>6</sup> and (i) PLAFF3. Refer to Table 2 for a description of the models.

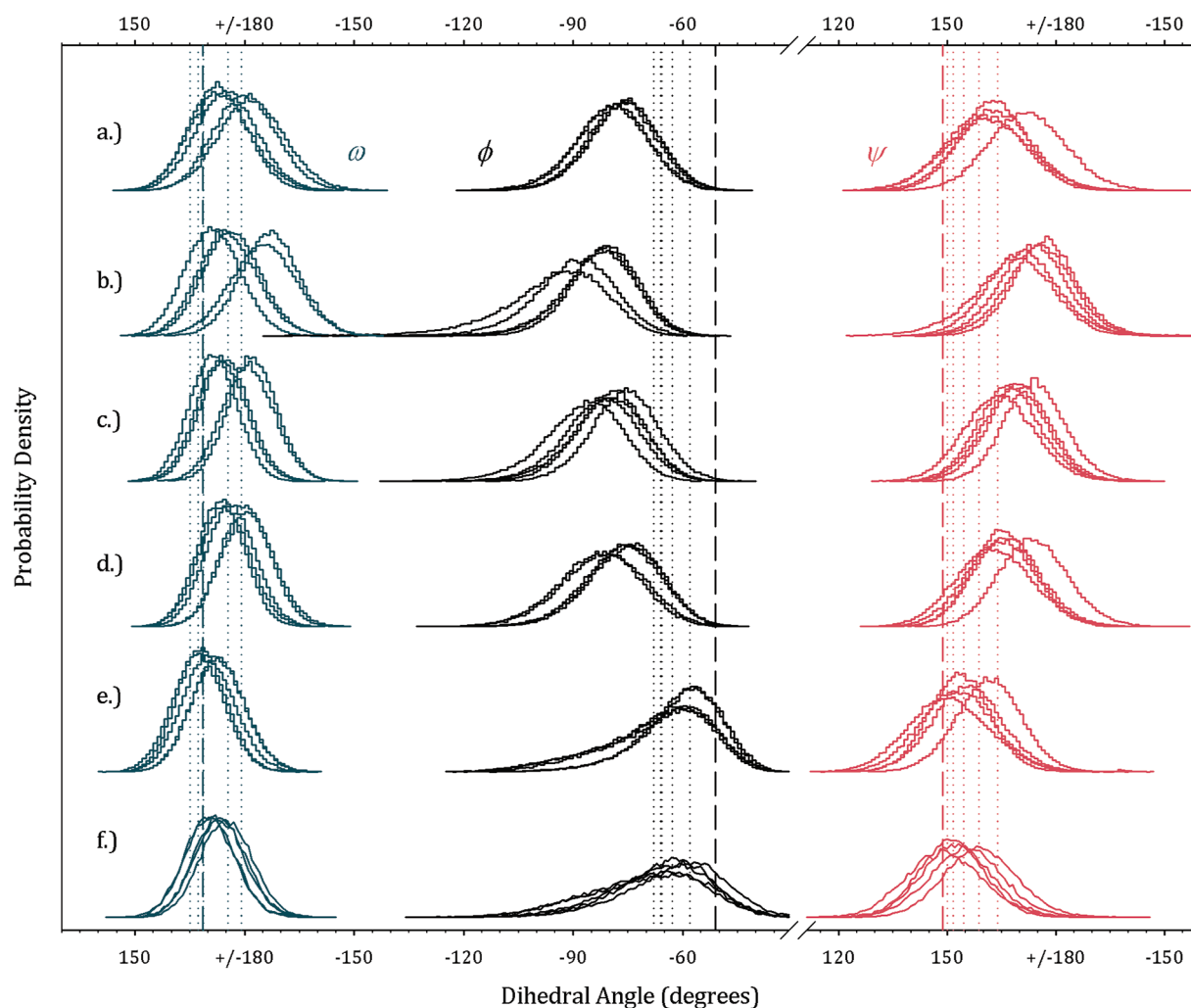
**Table 3. Lattice Dimensions of PLA at 300 K from Published Studies and from Crystal Structure Simulations<sup>a</sup>**

	<i>a</i> (Å)	diff (%)	<i>b</i> (Å)	diff (%)	<i>c</i> (Å)	diff (%)	density (g/cm <sup>3</sup> )	diff (%)
Sasaki <sup>20</sup>	10.66		6.16		28.88		1.261	
Alemán <sup>16</sup>	9.66	−9	5.80	−5	29.01	1	1.472	16.7
Hoogsteen <sup>18</sup>	10.60	−1	6.10	−1	28.80	0	1.285	1.8
de Santis <sup>17</sup>	10.70	0	6.45	5	27.80	−4	1.247	−1.2
OPLS	10.46	−1.9	6.05	−1.8	31.14	7.8	1.214	−3.8
CHARMM	10.72	0.6	5.97	−3.1	31.47	9.0	1.188	−5.8
OPLS'	10.51	−1.4	5.97	−3.1	31.36	8.6	1.216	−3.6
CHARMM'	8.78	−17.6	6.03	−2.1	34.67	20.0	1.303	3.3
OPLS''	10.54	−1.1	6.08	−1.3	30.85	6.8	1.210	−4.1
PLAFF2	10.59	−0.7	6.25	1.5	29.74	3.0	1.215	−3.7
<b>PLAFF3</b>	<b>10.70</b>	<b>0.4</b>	<b>6.14</b>	<b>−0.3</b>	<b>29.98</b>	<b>3.8</b>	<b>1.214</b>	<b>−3.7</b>

<sup>a</sup> Refer to Table 2 for a description of the models. Differences are calculated with respect to the experimental study of Sasaki and Asakura.<sup>20</sup> The recommended model, PLAFF3, is emphasized in bold.

#### Comparison of the Classical Models to Crystal Structure Data. Results from crystal structure simulations using each of

the models are shown in Table 3 and in Figure 7. In each case, the simulation results are compared to reference values from the



**Figure 7.** Dihedral angle distributions for crystalline PLLA at 300 K, simulated with (a) OPLS,<sup>5</sup> (b) CHARMM,<sup>4</sup> (c) OPLS', (d) CHARMM', (e) PLAFF2, (f) PLAFF3. Refer to Table 2 for a description of the models. Vertical dotted lines, values from the WAXD crystal structure analysis of Sasaki and Asakura;<sup>20</sup> vertical dashed lines, averaged values from the PLAFF simulations performed by O'Brien.<sup>7</sup>

experimentally resolved crystal structure(s). We have given priority to ensure PLAFF3 matches the data in Table 3, as the crystalline density and lattice vectors can be measured directly with very few assumptions involved in the experimental analysis. Thus, the dihedral angle distributions in Figure 7 were considered a secondary target.

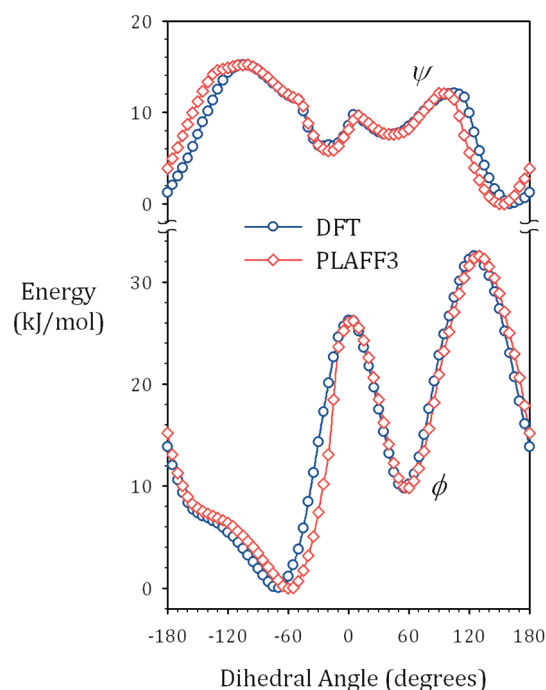
Figure 7 shows the dihedral angle distributions during simulation of crystalline PLA. For each model, five different histograms are accumulated for each backbone dihedral,  $\omega$ ,  $\phi$ , and  $\psi$ . These separate histograms are presented for each of the five unique residues in the frustrated helical structure predicted by Sasaki and Asakura.<sup>20</sup> In each model, it is evident that these five residues take on different dihedral values, according to their orientation inside the unit cell. This supports the existence of a frustrated structure and demonstrates that a helix with perfect screw symmetry is not possible under the crystalline packing conditions of PLA.

From Figure 7a and b, it is apparent that the OPLS and CHARMM models do not predict the same dihedral angle distribution as suggested by the WAXD results.<sup>20</sup> A more surprising result was that refitting the torsional potentials to DFT data had very little effect on the dihedral angle distributions in the crystalline phase, as evident in Figure 7c and d. We found it essential to improve agreement with the

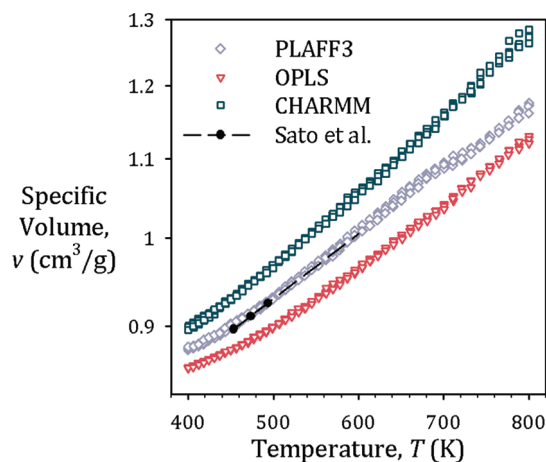
experimental unit cell lattice vectors before adjusting the dihedral parameters. The unit cell dimensions impose constraints on the set of dihedral angles that are probable, given that the crystal structure must be periodic with respect to those dimensions. Further, the set of bond lengths and angles played a vital role in achieving agreement with the crystal structure, as these impose the same sort of constraints on the dihedral angles when a periodic cell is used. Adjustments to bond stretching and angle bending parameters for this purpose are provided in the Supporting Information.

Once the bonded interactions were adjusted and more closely matched those used in the WAXD analysis of Sasaki and Asakura,<sup>20</sup> adjustment of the dihedral angles in the crystalline structure was relatively simple; in practice, we found that all of the backbone dihedral angle distributions could be shifted toward the WAXD values, by altering the potential with respect to the  $\phi$  and  $\psi$  dihedral angles alone. A simple shift in the position of the global minimum was required, as depicted in Figure 8. The minimum was shifted by  $10^\circ$  in the  $\phi$  dihedral angle and  $-15^\circ$  in the  $\psi$  dihedral angle. A similar shift in  $\phi$  was also required by O'Brien in developing PLAFF.<sup>7</sup>

In developing the PLAFF3 model, one of our stated goals was to obtain a force field that is suitable for modeling PLA in its



**Figure 8.** Adjustment of the torsional potential for the  $\phi$  and  $\psi$  dihedral angles, which resulted in improvement of the dihedral angle distributions in crystalline simulations.



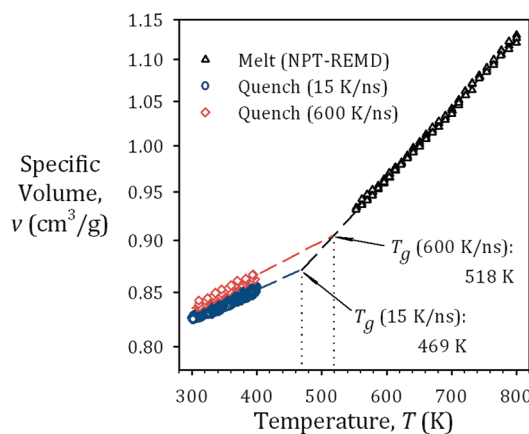
**Figure 9.** Melt phase densities of PLA, plotted from four separate NPT-REMD simulations for each of the CHARMM, OPLS, and PLAFF3 models. The melt phase experimental measurements of Sato et al.<sup>23</sup> are included for comparison and extrapolated toward the higher simulation temperatures.

amorphous state. Simultaneously, we wished to retain the model's accuracy in simulating crystalline PLA, which was a hallmark of O'Brien's original PLAFF.<sup>7</sup> We believe the results presented thus far demonstrate that PLAFF3 does indeed accurately predict the crystalline structure of PLA. In addition, Figure 6i shows the improvement in the topography of the PLAFF3 bond rotational energy landscape, when compared to PLAFF and PLAFF2, and demonstrates that the new model is more likely to have the correct dihedral angle distribution in the melt and amorphous state. In what follows, we show that PLAFF3 is also better suited for simulating PLA in its noncrystalline form,

**Table 4.** Volume Expansivities Estimated for Melt Phase PLA<sup>a</sup>

method/model	$\beta \times 10^4 \text{ (K}^{-1}\text{)}$
OPLS	$7.74 \pm 0.07$
CHARMM	$9.5 \pm 0.1$
PLAFF3	$7.7 \pm 0.4$
Sato et al. <sup>23</sup>	$7.8 \pm 0.4$

<sup>a</sup> Values are calculated from the simulation results shown in Figure 9, by a linear regression (on a log scale plot) of the data points above 550 K. An estimate using the experimental data of Sato et al.<sup>23</sup> is included for comparison. Listed errors are 95% confidence intervals for each slope.



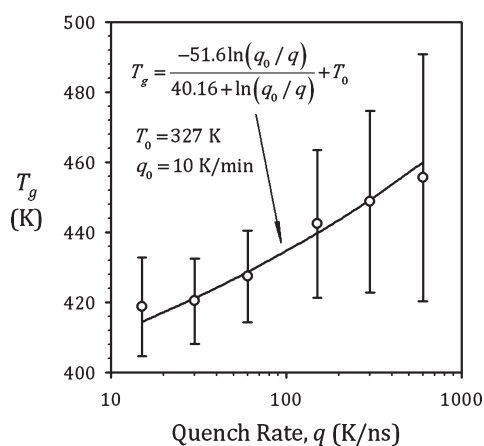
**Figure 10.** Example of the specific volume–temperature ( $v$ – $T$ ) plot used to determine the glass transition temperature. Results are from the OPLS model, using two different quench rates.

up to high temperatures when compared to the other models discussed here.

**Comparison of the Classical Models to Melt Phase Dilatometric Data.** When examining the models' performance in high temperature simulations, we found that the OPLS model underpredicts the specific volume of PLA in the melt phase. This is shown in Figure 9, using results from the NPT-REMD simulations. CHARMM, on the other hand, tends to overestimate the specific volume. Results from the OPLS-based force fields generally reproduced the volume expansivity of PLA, as shown in Table 4, whereas the CHARMM-based models tended to have higher expansivities than indicated in the experimental results of Sato et al.<sup>23</sup> It was found that substituting one or more of the nonbonded parameters (both Lennard-Jones parameters and partial charges) from CHARMM helped to increase the specific volume in the melt, without increasing the expansivity above the desired range. Following this observation, in the PLAFF3 force field, CHARMM nonbonded parameters are used for the O<sup>S</sup> and C<sup>a</sup> atoms. While still slightly lower than the experimental measurements, the melt volumes predicted by PLAFF3 are noticeably closer to the experiment than either OPLS or CHARMM; this result supports our assertion that the model may be used equally well in simulating the melt and/or crystalline states of PLA.

**Comparison of the Classical Models to Glass Transition Data.** The last material property we used in constructing the PLAFF3 set of parameters was the PLA glass transition temperature,  $T_g$ . Figure 10 gives an example of the specific volume intersection method used for determining  $T_g$  at two different quench rates, using the OPLS force field. The results depicted in





**Figure 11.** WLF plot for extrapolating simulation glass transition data to realistic (laboratory scale) quench rates. Results are from the PLAFF3 model. Universal WLF constants are used, with a lab-scale quench rate of  $q_0 = 10$  K/min. Error bars are propagated from 95% confidence intervals on the slopes and intercepts of the melt and glassy  $\nu$ - $T$  plots (see Figure 10). Here,  $T_0 = 327$  K is the glass transition temperature extrapolated to the lab-scale quench rate.

**Table 5. Glass Transition Temperatures Calculated from the Various Models Explored in This Work<sup>a</sup>**

method/model	$T_g$ (K)
simulation data	
OPLS	$388 \pm 14$
CHARMM	$367 \pm 15$
OPLS''	$403 \pm 12$
PLAFF <sup>44</sup>	$408^b$
PLAFF2 <sup>6</sup>	$386 \pm 11$
PLAFF3	$327 \pm 12$
Experimental Data	
Dorgan <sup>26</sup>	331.6
Kanchanasopa <sup>28</sup>	327

<sup>a</sup> Simulation results from previous studies using PLAFF<sup>44</sup> and PLAFF2<sup>6</sup> are included for reference, as well as selected experimental results<sup>25,28</sup> for PLA. <sup>b</sup> Extrapolated to infinite molecular weight limit, not corrected for quench rate dependence.

the figure are generally representative of all intersection plots constructed during this work; the faster quenching rates consistently gave intersection points that are higher up on the melt volumetric curve. Figure 11 demonstrates the extrapolative method used to estimate  $T_g$  for laboratory scale quench rates using the WLF relation. In fitting the WLF equation to the simulation data in Figure 11, the only adjustable parameter used was  $T_0$ , which corresponds to the laboratory-scale glass transition temperature when the universal WLF constants are used (see the Methods section of this paper).

A survey of the glass transition temperatures for some of the models discussed in this work is presented in Table 5. Not all models were tested for the glass transition temperature; following our procedure laid out in Figure 2, we required that our models perform accurately in both the crystalline and melt states before attempting to examine the glass transition temperature. Thus, the OPLS' and CHARMM' models were not examined

with glass transition simulations, as they did not meet the prerequisites in simulating the crystal structure. Similarly, PLAFF was not used because it is believed to give inadequate dihedral angle distributions. We made three exceptions, for demonstration purposes. We chose to estimate  $T_g$  using OPLS, CHARMM, and OPLS'', because these results give some idea of how the glass transition temperature was affected by changes made early on in the fitting procedure.

Most of the simulation-based estimates of  $T_g$  shown in Table 5 are higher than the experimentally observed glass transition temperature, with the PLAFF3 force field being the closest to the experimental value. It is also apparent in Table 5 that the modification of the torsional and other potentials from the OPLS to the OPLS'' model resulted in a worsening of the  $T_g$  estimate using OPLS''. It is obvious that, in adjusting the nonbonded and valence interactions in OPLS to obtain the OPLS'' model, we affected the barrier height of bond rotation about the  $\psi$  dihedral angle. In the PLAFF2 model, we were able to remove this artifact, yet the limitations of using uncorrelated (non-CMAP) dihedral potentials for  $\phi$  and  $\psi$  made further lowering of the barriers difficult without drastically affecting the overall potential energy surface with respect to  $\phi$  and  $\psi$ , as discussed elsewhere.<sup>6</sup> By introducing the CMAP potential in PLAFF3, we were able to remedy this problem.

The considerable freedom entailed in the CMAP model allowed for a nearly exact fit of the DFT potential energy surface in Figure 4, and any adjustments made during crystal structure fitting could be made independently of the bond rotational barrier heights. Whereas such changes in PLAFF2 resulted in an unwanted increase in the barrier to rotation about the  $\psi$  dihedral angle, this barrier height could be preserved using the CMAP potential in PLAFF3. We found that the bond rotational barrier height in PLAFF3 was very close to the DFT-predicted value. PLAFF3 predicted a glass transition temperature within the range of the experimental results, without need for further adjustment of the barrier to rotation about the  $\psi$  dihedral angle. This is a major improvement over previous versions of the force field.

## CONCLUSIONS

In this paper, we presented our work related to the development of an optimized model for the atomistic simulation of polylactide (PLA). The model, PLAFF3, was shown to perform well in simulations of the amorphous and crystalline states of PLA. This model is an update to the previous versions by O'Brien<sup>7</sup> and McAliley,<sup>6</sup> and we have significantly improved the ability of the model to describe the proper dihedral angle distributions in the amorphous states of PLA. On the basis of the results of this work, we recommend the use of the PLAFF3 model under most circumstances.

A major improvement in PLAFF3 over previous models is its ability to predict the glass transition temperature of PLA. This was possible due to the CMAP dihedral cross terms that were used in PLAFF3, in place of the linear combinations of individual dihedral terms used in prior versions of the force field. The inability to reproduce the experimental glass transition temperature was the largest shortcoming of PLAFF2, and PLAFF3 addresses this problem while retaining the accuracy of PLAFF2 in simulating the melt and crystalline states of PLA. We believe the wide range of properties captured by the PLAFF3 model make it well suited for studying a wide range of phenomena, such as crystallization, permeant diffusion, and shear and elongational

flow. Due to the accuracy of PLAFF3 in simulating the pure crystalline form, the model should work equally as well as PLAFF in simulating various surface interactions with the polymer, since the polymer is known to exhibit a high level of crystallization at surface boundaries. We feel that a judicious practitioner of molecular modeling should be able to apply PLAFF3 to successfully simulate any of these phenomena on a molecular level.

## ■ ASSOCIATED CONTENT

**S Supporting Information.** Parameters for the PLAFF3 force field as well as complete simulation input files for PLA. The input files, including atomic coordinate files, are specifically for the GROMACS molecular dynamics package and employ the PLAFF3 force field; however, the parameters in these files can be easily ported to other simulation packages. The atomic coordinate files are for fully equilibrated melt configurations of PLA. Additionally, the DFT derived energy profiles associated with bond stretching and angle bending of an *in vacuo* PLA trimer are provided. This information is available free of charge via the Internet at <http://pubs.acs.org>

## ■ AUTHOR INFORMATION

### Corresponding Author

\*E-mail: [dbruce@clemson.edu](mailto:dbruce@clemson.edu).

## ■ ACKNOWLEDGMENT

The authors acknowledge the support of the Center for Advanced Engineering Fibers and Films and the ERC program of the National Science Foundation (NSF) under Award Number EEC-9731680. Partial support was also provided by the National Institutes of Health (NIH)/National Institute of Biomedical and Bioengineering (NIBIB) under grant number R01 EB006163. The authors would also like to acknowledge the support of the staff from the Cyberinfrastructure Technology Integration group, and the use of the advanced computational resources deployed and maintained by Clemson Computing and Information Technology.

## ■ REFERENCES

- (1) Middleton, J. C.; Tipton, A. J. Synthetic biodegradable polymers as orthopedic devices. *Biomaterials* **2000**, *21* (23), 2335–2346.
- (2) Anderson, J. M.; Shive, M. S. Biodegradation and biocompatibility of PLA and PLGA microspheres. *Adv. Drug Delivery Rev.* **1997**, *28* (1), 5–24.
- (3) Auras, R.; Harte, B.; Selke, S. An overview of polylactides as packaging materials. *Macromol. Biosci.* **2004**, *4* (9), 835–864.
- (4) Brooks, B. R.; Brucoleri, R. E.; Olafson, B. D.; States, D. J.; Swaminathan, S.; Karplus, M. Charmm - a program for macromolecular energy, minimization, and dynamics calculations. *J. Comput. Chem.* **1983**, *4* (2), 187–217.
- (5) Jorgensen, W. L.; Maxwell, D. S.; Tirado-Rives, J. Development and testing of the OPLS all-atom force field on conformational energetics and properties of organic liquids. *J. Am. Chem. Soc.* **1996**, *118* (45), 11225–11236.
- (6) McAliley, J. H. Development of improved torsional potentials in classical force field descriptions of poly (lactic acid). Ph.D. Dissertation, Clemson University, Clemson, SC, 2009. <http://etd.lib.clemson.edu/documents/1252938067/> (accessed Dec 2010)
- (7) O'Brien, C. P. Quantum and molecular modeling of polylactide. Ph.D. Dissertation, Clemson University, Clemson, SC, 2005.

- (8) Bjelkmar, P.; Larsson, P.; Cuendet, M. A.; Hess, B.; Lindahl, E. Implementation of the CHARMM Force Field in GROMACS: Analysis of Protein Stability Effects from Correction Maps, Virtual Interaction Sites, and Water Models. *J. Chem. Theory Comp.* **2010**, *6* (2), 459–466.
- (9) Blomqvist, J. *Ris metropolis monte carlo studies of poly(l-lactic), poly(l,d-lactic) and polyglycolic acids.* *Polymer* **2001**, *42* (8), 3515–3521.
- (10) Blomqvist, J.; Ahjopalo, L.; Mannfors, B.; Pietila, L. O. Studies on aliphatic polyesters i: Ab initio, density functional and force field studies of esters with one carboxyl group. *THEOCHEM* **1999**, *488*, 247–262.
- (11) Blomqvist, J.; Mannfors, B.; Pietila, L. O. Studies on aliphatic polyesters. Part ii. Ab initio, density functional and force field studies of model molecules with two carboxyl groups. *THEOCHEM* **2000**, *531*, 359–374.
- (12) Blomqvist, J.; Pietila, L. O. Amorphous cell studies of polyglycolic, poly(l-lactic), poly(l,d-lactic) and poly(glycolic/l-lactic) acids. *Polymer* **2002**, *43* (17), 4571–4583.
- (13) McAliley, J. H.; O'Brien, C. P.; Bruce, D. A. Continuum electrostatics for electronic structure calculations in bulk amorphous polymers: Application to polylactide. *J. Phys. Chem. A* **2008**, *112* (31), 7244–7249.
- (14) van der Spoel, D.; Lindahl, E.; Hess, B.; Groenhof, G.; Mark, A. E.; Berendsen, H. J. C. GROMACS: Fast, Flexible and Free. *J. Comput. Chem.* **2005**, *26*, 1701–1718.
- (15) Mackerell, A. D.; Feig, M.; Brooks, C. L. Extending the treatment of backbone energetics in protein force fields: Limitations of gas-phase quantum mechanics in reproducing protein conformational distributions in molecular dynamics simulations. *J. Comput. Chem.* **2004**, *25* (11), 1400–1415.
- (16) Alemán, C.; Lotz, B.; Puiggali, J. Crystal structure of the alpha-form of poly(l-lactide). *Macromolecules* **2001**, *34* (14), 4795–4801.
- (17) de Santis, P.; Kovacs, J. Molecular conformation of poly(s-lactic acid). *Biopolymers* **1968**, *6* (3), 299–306.
- (18) Hoogsteen, W.; Postema, A. R.; Pennings, A. J.; Tenbrinke, G.; Zugenmaier, P. Crystal-structure, conformation, and morphology of solution-spun poly(l-lactide) fibers. *Macromolecules* **1990**, *23* (2), 634–642.
- (19) Kobayashi, J.; Asahi, T.; Ichiki, M.; Oikawa, A.; Suzuki, H.; Watanabe, T.; Fukada, E.; Shikami, Y. Structural and optical-properties of poly lactic acids. *J. Appl. Phys.* **1995**, *77* (7), 2957–2973.
- (20) Sasaki, S.; Asakura, T. Helix distortion and crystal structure of the alpha-form of poly(l-lactide). *Macromolecules* **2003**, *36* (22), 8385–8390.
- (21) Arnott, S.; Wonacott, A. J. Atomic coordinates for an alpha-helix refinement of crystal structure of alpha-poly-l-alanine. *J. Mol. Biol.* **1966**, *21* (2), 371–383.
- (22) Rietveld, H. M. A profile refinement method for nuclear and magnetic structures. *J. Appl. Crystallogr.* **1969**, *2* (2), 65–71.
- (23) Sato, Y.; Inohara, K.; Takishima, S.; Masuoka, H.; Imaizumi, M.; Yamamoto, H.; Takasugi, M. Pressure-volume-temperature behavior of polylactide, poly(butylene succinate), and poly(butylene succinate-co-adipate). *Polym. Eng. Sci.* **2000**, *40* (12), 2602–2609.
- (24) Painter, P. C.; Coleman, M. M. *Essentials of polymer science and engineering*; DEStech Publications: Lancaster, PA, 2009; p 464.
- (25) Auras, R.; Harte, B.; Selke, S. An overview of polylactides as packaging materials. *Macromol. Biosci.* **2004**, *4* (9), 835–864.
- (26) Dorgan, J. R.; Lehermeier, H.; Mang, M. Thermal and rheological properties of commercial-grade poly(lactic acid)s. *J. Polym. Environ.* **2000**, *8* (1), 1–9.
- (27) Joziassé, C. A. P.; Veenstra, H.; Grijpma, D. W.; Pennings, A. J. On the chain stiffness of poly(lactide)s. *Macromol. Chem. Phys.* **1996**, *197* (7), 2219–2229.
- (28) Kanchanasopa, M.; Runt, J. Broadband dielectric investigation of amorphous and semicrystalline l-lactide/meso-lactide copolymers. *Macromolecules* **2004**, *37* (3), 863–871.
- (29) Berendsen, H. J. C.; Postma, J. P. M.; DiNola, A.; Haak, J. R. Molecular-Dynamics with Coupling to an External Bath. *J. Chem. Phys.* **1984**, *81* (8), 3684–3690.
- (30) Nose, S. A unified formulation of the constant temperature molecular-dynamics methods. *J. Chem. Phys.* **1984**, *81* (1), 511–519.
- (31) Hoover, W. G. Canonical dynamics: Equilibrium phase-space distributions. *Phys. Rev. A* **1985**, *31* (3), 1695–1697.

- (32) Darden, T.; York, D.; Pedersen, L. Particle Mesh Ewald - An  $N \cdot \log(N)$  Method For Ewald Sums in Large Systems. *J. Chem. Phys.* **1993**, *98* (12), 10089–10092.
- (33) Kolstad, J. J.; Witzke, D. R.; Hartmann, M. H.; Hall, E. S.; Nangeroni, J. Lactic Acid Residue Containing Polymer Composition and Product Having Improved Stability, and Method for Preparation and Use Thereof. U.S. Patent 6,353,086, March 5, 2002.
- (34) Bopp, R. NatureWorks, LLC. Personal communication, 2008.
- (35) O'Connell, J. P.; Haile, J. M. *Thermodynamics: Fundamentals for applications*. Cambridge University Press: New York, 2005; p 82.
- (36) Boyd, R. H. Glass transition temperatures from molecular dynamics simulations. *Trends Polym. Sci.* **1996**, *4* (1), 12–17.
- (37) Han, J.; Gee, R. H.; Boyd, R. H. Glass-transition temperatures of polymers from molecular-dynamics simulations. *Macromolecules* **1994**, *27* (26), 7781–7784.
- (38) Rigby, D.; Roe, R. J. Molecular-dynamics simulation of polymer liquid and glass. I. Glass-transition. *J. Chem. Phys.* **1987**, *87* (12), 7285–7292.
- (39) Buchholz, J.; Paul, W.; Varnik, F.; Binder, K. Cooling rate dependence of the glass transition temperature of polymer melts: Molecular dynamics study. *J. Chem. Phys.* **2002**, *117* (15), 7364–7372.
- (40) Soldera, A.; Metatla, N. Glass transition of polymers: Atomistic simulation versus experiments. *Phys. Rev. E: Stat. Nonlin. Soft Matter Phys.* **2006**, *74* (6–1), 061803/1–061803/6.
- (41) Sperling, L. H. *Introduction to physical polymer science*, 4th ed.; Wiley: Hoboken, NJ, 2006.
- (42) Williams, M. L.; Landel, R. F.; Ferry, J. D. The temperature dependence of relaxation mechanisms in amorphous polymers and other glass-forming liquids. *J. Am. Chem. Soc.* **1955**, *77* (14), 3701–3707.
- (43) Price, M. L. P.; Ostrovsky, D.; Jorgensen, W. L. Gas-phase and liquid-state properties of esters, nitriles, and nitro compounds with the opl-aa force field. *J. Comput. Chem.* **2001**, *22* (13), 1340–1352.
- (44) Zhang, J.; Liang, Y.; Yan, J. Z.; Lou, J. Z. Study of the molecular weight dependence of glass transition temperature for amorphous poly(l-lactide) by molecular dynamics simulation. *Polymer* **2007**, *48* (16), 4900–4905.

# Comparison of the Efficiency of the LIE and MM/GBSA Methods to Calculate Ligand-Binding Energies

Samuel Genheden and Ulf Ryde\*

Department of Theoretical Chemistry, Lund University, Chemical Centre, P.O. Box 124, SE-221 00 Lund, Sweden

Supporting Information

**ABSTRACT:** We have evaluated the efficiency of two popular end-point methods to calculate ligand-binding free energies, LIE (linear interaction energy) and MM/GBSA (molecular mechanics with generalized Born surface-area solvation), i.e. the computational effort needed to obtain estimates of a similar precision. As a test case, we use the binding of seven biotin analogues to avidin. The energy terms used by MM/GBSA and LIE exhibit a similar correlation time ( $\sim 5$  ps), and the equilibration time seems also to be similar, although it varies much between the various ligands. The results show that the LIE method is more effective than MM/GBSA, by a factor of 2–7 for a truncated spherical system with a radius of 26 Å and by a factor of 1.0–2.4 for the full avidin tetramer (radius 47 Å). The reason for this is the cost for the MM/GBSA entropy calculations, which more than compensates for the extra simulation of the free ligand in LIE. On the other hand, LIE requires that the protein is neutralized, whereas MM/GBSA has no such requirements. Our results indicate that both the truncation and neutralization of the proteins may slow the convergence and emphasize small differences in the calculations, e.g., differences between the four subunits in avidin. Moreover, LIE cannot take advantage of the fact that avidin is a tetramer. For this test case, LIE gives poor results with the standard parametrization, but after optimizing the scaling factor of the van der Waals terms, reasonable binding affinities can be obtained, although MM/GBSA still gives a significantly better predictive index and correlation to the experimental affinities.

## INTRODUCTION

One of the most important challenges of computational chemistry is to accurately estimate the free-energy change of a biochemical reaction. For instance, in drug design, one is interested in the binding of small ligands to a biomolecular target, usually a protein. If accurate free energies could be estimated for this reaction by computational methods, billions of dollars could be saved because it would be necessary to synthesize fewer molecules.<sup>1,2</sup>

Many methods are available to estimate free energies, ranging from simple scoring functions that are fast, but not very accurate, to rigorous free energy perturbation (FEP), which is accurate but time-consuming.<sup>3,4</sup> The reason for the latter is that FEP requires extensive sampling using molecular dynamics (MD) or Monte Carlo methods on a series of intermediate, unphysical states. A class of methods that is intermediate in efficiency is the so-called end-point methods, which still are based on physical laws and require sampling, but only of the reactants and the products, not of any intermediate states.<sup>5</sup> However, even with perfect sampling, these methods will not give the exact result, because they are based on several approximations. Therefore, such methods need to be evaluated carefully to identify their strengths and weaknesses.

Two such methods are LIE<sup>6–8</sup> (linear interaction energy) and MM/GBSA<sup>9,10</sup> (molecular mechanics with generalized Born and surface-area solvation). LIE estimates the free energy for the binding of a ligand (L) to its target macromolecule (P) by simulating the free ligand in solution and the ligand–macromolecule complex (PL), using the relation<sup>7</sup>

$$\Delta G = \beta(\langle E_{\text{ele}}^{L-S} \rangle_{\text{PL}} - \langle E_{\text{ele}}^{L-S} \rangle_{\text{L}}) + \alpha(\langle E_{\text{vdW}}^{L-S} \rangle_{\text{PL}} - \langle E_{\text{vdW}}^{L-S} \rangle_{\text{L}}) \quad (1)$$

where  $E_{\text{ele}}^{L-S}$  and  $E_{\text{vdW}}^{L-S}$  are the electrostatic and van der Waals intermolecular interaction energies between the ligand and the

surroundings (S; i.e., protein and solvent),  $\alpha$  and  $\beta$  are two parameters, and the angle brackets indicate ensemble averages from the simulations of either the free ligand or the complex, as indicated by the subscripts.  $\beta$  was originally set to 0.5,<sup>6</sup> because LIE was derived from the linear-response approximation. However, this value was later refined to reflect the chemical nature of the ligand,<sup>7,11,12</sup> based on FEP calculations.  $\alpha$  is usually set to 0.18,<sup>13,14</sup> but this value has been much debated and may be system dependent.<sup>5,8</sup> In several studies, this parameter has been fitted to experimental data for each protein target and ligand type.<sup>5</sup> A third constant term has also been suggested,<sup>15</sup> but it is important only when estimating absolute free energies.<sup>13</sup>

MM/GBSA, on the other hand, estimates the free energy as<sup>9,10</sup>

$$\Delta G = G(\text{PL}) - G(\text{P}) - G(\text{L}) \quad (2)$$

where each free energy is calculated from a sum of six terms

$$G = \langle E_{\text{int}} + E_{\text{ele}} + E_{\text{vdW}} + G_{\text{solv}} + G_{\text{np}} - TS_{\text{MM}} \rangle \quad (3)$$

The three first terms are the molecular-mechanics (MM) internal, electrostatics, and van der Waals energies;  $G_{\text{solv}}$  is the polar solvation energy;  $G_{\text{np}}$  is the nonpolar solvation free energy;  $T$  is the absolute entropy; and  $S_{\text{MM}}$  is an entropy estimate from harmonic frequencies calculated at the MM level. The average in eq 3 should in principle be calculated from three separate simulations PL, P, and L, but for stability reasons,<sup>16</sup> it is more common to simulate only the complex. In that case,  $E_{\text{int}}$  cancels. In the MM/GBSA approach, the polar solvation free energy is calculated by a generalized Born (GB) approach, but it could be calculated by any continuum-solvation method.<sup>17</sup> A common

Received: March 9, 2011

Published: October 04, 2011

choice is the Poisson–Boltzmann method, giving the MM/PBSA approach. The nonpolar solvation free energy is usually estimated by a relation to the solvent-accessible surface area (SASA).<sup>5</sup>

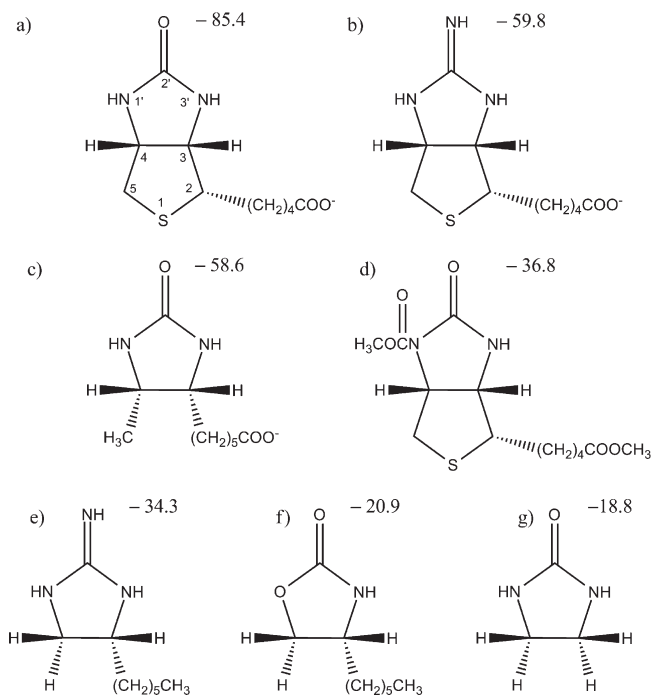
Because LIE and MM/GBSA are two popular end-point methods, it is interesting to compare them. This has been done a few times in the past.<sup>18–21</sup> For the binding of biotin analogues to avidin, MM/PBSA was found to reproduce experimental results more accurately than LIE,<sup>19</sup> although both methods were less accurate than FEP.<sup>22</sup> However, for acetylcholinesterase huprine inhibitors, LIE gave better results than MM/PBSA.<sup>18</sup> For the binding of eight hydroxamate inhibitors to gelatinase A, the two methods showed a similar performance.<sup>20</sup> In all of these studies, the LIE parameters were adjusted to an optimal fit. For the binding of fragment B of protein A to the Fc domain of immunoglobulin G, LIE gave similar results to both MM/PBSA and MM/GBSA, but only one complex was examined, a protein for which it is not clear what  $\alpha$  and  $\beta$  parameters should be used.<sup>21</sup> Apparently, the accuracy of the two methods (i.e., how well they reproduce experimental results) depends strongly on the systems studied, and much larger test sets are needed before any general conclusion can be reached.

In this paper, we will instead focus on the precision (i.e., the statistical uncertainty of the results) and efficiency (i.e., the computer time required to reach a given precision) of the two methods. The statistical precision is important when comparing ligand-affinity methods:<sup>23</sup> Congeneric ligands often have quite similar affinities, and an order of magnitude difference in the binding constant corresponds to only 6 kJ/mol in the free energy of binding. If statistically significant differences should be discerned, a precision of 1–2 kJ/mol is therefore needed. Such precision is also needed to make results obtained by different groups comparable<sup>24</sup> and to avoid the temptation to rerun simulations that gave poor agreement with experiments. On the other hand, we have shown that once such a precision is reached, results obtained by MM/GBSA are reproducible and not sensitive to the setup of the simulations, except for the protonation of residues very close to the ligand.<sup>24</sup>

In a previous paper, we developed a simulation protocol for MM/GBSA that gave a precision of 1 kJ/mol.<sup>23</sup> In particular, we showed that it was more favorable to run several rather short simulations instead of a single long one, as has been concluded also in other studies.<sup>25–27</sup> By running a proper number of independent simulations, any precision can be reached. In this paper, we develop a similar protocol for LIE. This also allows us to discuss the efficiency of the two methods, i.e., to compare the computational effort needed to obtain results of the same statistical precision. If the methods give similar accuracy, of course the more efficient method is preferred. To facilitate the comparison, we use the same test case as for MM/PBSA, viz., the binding of seven biotin analogues to avidin. This test system has been studied before with FEP,<sup>28</sup> MM/PB(GB)SA,<sup>17,19,23,24,29–33</sup> and LIE,<sup>22</sup> and experimental structures<sup>34</sup> as well as affinities are available.<sup>35–37</sup>

## METHODS

**System Preparation.** We have studied the binding of the seven biotin analogues in Figure 1 to avidin. Btn1–Btn3 have a net charge of  $-1$ , whereas the other four ligands are neutral. The structure of avidin was taken from the 1avd crystal structure,<sup>34</sup> which contains a cocrystallized biotin molecule in each subunit of the tetrameric protein. However, in this study, we consider the



**Figure 1.** The seven biotin analogues studied: (a) biotin (Btn1), (b–g) Btn2–Btn7. The numbers shown are experimental affinities in kJ/mol.<sup>36</sup>

binding of only a single ligand to the tetrameric protein. The six biotin analogues were built into the active site to mimic the binding mode of biotin, as has been described previously.<sup>30</sup> In LIE, it is essential that the protein is neutral.<sup>8</sup> Therefore, all titratable residues were neutralized (all of these residues are solvent exposed). This has shown to be the optimal approach to ensure that the complex and free ligand simulations have identical total charge, which is required if we want to ignore long-range effects beyond the simulation sphere.<sup>38</sup> The single histidine residue in each subunit was modeled to be protonated on the NE2 atom.<sup>30</sup> The protein atoms were described by the Amber99SB force field,<sup>39</sup> and parameters for the ligands were taken from the Amber99 force field.<sup>30,40</sup> Ligand charges were calculated with the RESP procedure,<sup>41</sup> using ESP points calculated at the Hartree–Fock 6-31G\* level and sampled with the Merz–Kollman scheme,<sup>42</sup> as has been described before.<sup>30</sup>

Two sets of systems were prepared for each protein–ligand complex, a full system and a truncated system. The full system was prepared by solvating the entire (tetrameric) protein–ligand complex in a sphere with a radius of 47 Å (i.e., extending at least 10 Å outside the protein; in total  $\sim 43\,925$  atoms). The truncated system was prepared by solvating the complex in a 26 Å sphere, centered on the ligand and, thereafter, removing all residues more than 26 Å from the ligand ( $\sim 8325$  atoms). Atoms between 26 and 24 Å were restrained in the simulations, by a harmonic restraint of 41.84 kJ/mol/Å<sup>2</sup>. The truncated system represents a more typical use of LIE.<sup>8</sup> Likewise, two sets of free-ligand systems were created by solvating the ligand in a sphere with a radius of either 47 or 26 Å, because LIE requires that the simulations of the complex and the free ligand have the same size, so that the ignored interactions outside the simulated systems cancel.<sup>8</sup> In these simulations, the geometrical center of the ligand was restrained to the origin using a harmonic restraint of 41.84 kJ/mol/Å<sup>2</sup>. In all simulations, the water model was TIP3P.<sup>43</sup> The systems were

prepared with a combination of the Qprep program of the Q simulation package, the Leap module of the Amber package, and in-house programs.

**Simulations.** All MD simulations were run by the Q simulation package.<sup>44</sup> All bonds involving hydrogen atoms were constrained with the SHAKE approach<sup>45</sup> and a 2 fs time step was employed. The temperature was kept at 300 K using a Berendsen thermostat.<sup>46</sup> The nonbonded cutoff was 10 Å, except for interactions with the ligand, for which an infinite cutoff was applied. Long-range electrostatic interactions were treated with the local reaction-field approximation.<sup>47</sup> Water molecules at the surface of the simulated sphere were subjected to radial and polarization restraints.<sup>44,48</sup>

Prior to the MD simulation, the systems were minimized using the sander module of Amber 10<sup>49</sup> using 100 steps of steepest descent and with a harmonic restraint of 104.6 kJ/mol/Å<sup>2</sup> on all atoms except hydrogen and water atoms. This was followed by starting a number of independent 20 ps MD simulations with the same restraint as in the minimization. Thereafter, an unrestrained MD simulation was carried out for 800 ps (full systems) or 1600 ps (truncated systems) for each of the independent simulations. Snapshots were sampled each picosecond. Twenty independent simulations were employed for the free ligand and for each of the subunits of avidin (i.e., 20 + 80 in total). These independent simulations were initiated by assigning different initial random velocities to all atoms, i.e., the velocity-induced independent-trajectory approach.<sup>24</sup>

**Free Energy Estimates.** The LIE interaction energies in eq 1 (with an infinite cutoff, but without any long-range corrections) were sampled with Q<sup>44</sup> during the simulation and were processed by in-house scripts.  $\beta$  was set to 0.5 for the charged ligands and 0.43 for the other ligands,<sup>7</sup> whereas  $\alpha$  was set to 0.18 as a default for all ligands,<sup>13</sup> although it was also optimized (see below). A new parametrization of  $\beta$  to include more chemical groups has been suggested,<sup>12</sup> but it does not involve thioether and other groups in our ligand set. For the charged ligands, a correction to the neutralization of the charged residues,  $\Delta G_{\text{co}}$ , was estimated by placing a single charge at the position of the CG, CD, NZ, and CZ atoms of Asp, Glu, Lys, and Arg, respectively, and calculating the Coulomb interaction between this charge and all atoms in the ligand for each snapshot, assuming a dielectric constant of 80, as suggested previously.<sup>38,50,51</sup> We tested this correction also for the neutral ligands, but it was found to be negligible,  $\sim 0.1$  kJ/mol.

The MM/GBSA calculations were performed by the Amber 10 package on snapshots from the QMD simulations.<sup>49</sup> The  $E_{\text{ele}}$  and  $E_{\text{vdw}}$  energies in eq 3 were calculated with the same force field as in the simulation and with an infinite cutoff. The polar solvation free energy was estimated by the GB method of Onufriev et al.,<sup>52</sup> model I (OBCI, i.e. with  $\alpha = 0.8$ ,  $\beta = 0$ , and  $\gamma = 2.91$ ). The nonpolar solvation energy was estimated from the SASA according to  $\Delta G_{\text{np}} = \gamma \text{SASA} + b$ , where  $\gamma = 0.0227$  kJ/mol/Å<sup>2</sup> and  $b = 3.85$  kJ/mol.<sup>53</sup> The entropy was estimated by calculating harmonic frequencies at the MM level on a truncated and buffered system (8 + 4 Å from the ligand), as described previously, to improve the statistical precision of the estimate.<sup>32</sup>

**Estimation of the Correlation Time.** The correlation time of the LIE interaction energies was estimated with the statistical inefficiency method.<sup>54,55</sup> In this procedure, the following measure is calculated

$$\Phi = \frac{\tau \sigma^2(Y)_\tau}{\sigma^2(X)} \quad (4)$$

where  $\sigma^2(X)$  is the variance of the distribution  $\{X\}$ , i.e., the variance of the time series of a particular energy, e.g.,  $E_{\text{ele}}^{\text{LIE}}$  in eq 1, and  $\sigma^2(Y)_\tau$  is the variance of the block average of  $\{X\}$ , where the block length is  $\tau$ . This block average is calculated from

$$Y_i = \frac{1}{\tau} \sum_{j=n-i\tau+1}^{n-(i-1)\tau} X_j \quad (5)$$

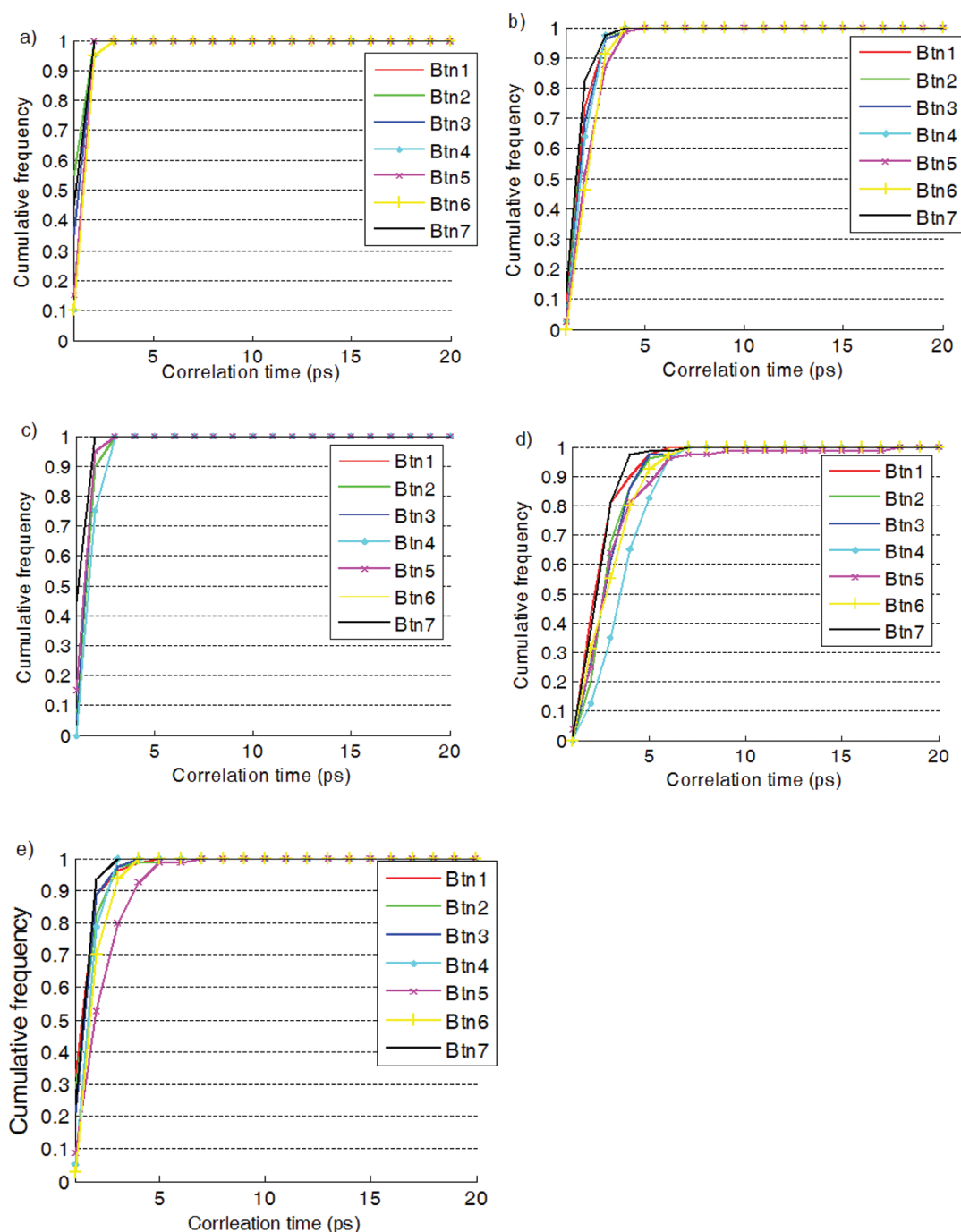
Put in another way,  $\{X\}$  is divided into a number of nonoverlapping segments, each with length  $\tau$ . Once  $\tau$  is so large that the successive values of  $Y_i$  are statistically independent,  $\Phi$  will become a constant and an estimate of the correlation time of  $\{X\}$ . This method is sensitive to equilibration, long-time trends, and bumps in the data (increasing the apparent correlation time). Therefore, we divided the data into segments of 200 ps and calculated the correlation time within each segment separately. The correlation time of the whole simulation was then taken as the median of the calculated correlation time for the segments.

**Error Analysis.** To measure the quality of the free-energy estimates, we use four different estimates: the correlation coefficient between the predicted and experimental data ( $r^2$ ), Pearlman's predictive index (PI),<sup>56</sup> the mean absolute deviation (MAD), and the mean absolute deviation from the best correlation line through the origin (MADtr; i.e., MAD after subtraction of the mean signed deviation). These measures are rather meaningless without an estimate of their statistical uncertainty. They were obtained by a simple parametric bootstrap simulation.<sup>23</sup> Each ligand was assigned a random normal-distributed affinity, with a mean and standard deviation obtained from the free-energy estimates. Then,  $r^2$ , PI, MAD, and MADtr were evaluated, and this procedure was repeated 10 000 times. The standard deviations of these resampled sets are reported as the standard errors of the quality measures. Throughout this paper, all reported statistical uncertainties are standard errors of the mean, i.e., the standard deviation divided by the square root of the number of estimates.

## RESULTS

We have estimated the binding free energy of seven biotin analogues to avidin using the LIE and MM/GBSA approaches. Binding affinities obtained with MM/GBSA have already been published for these ligands.<sup>23</sup> However, LIE requires calculations performed on neutralized spherical systems with only one ligand, and typically also for truncated systems. Therefore, the MM/GBSA calculations were rerun with the same settings as the LIE calculations to make the results completely comparable and also to allow for a comparison of the results obtained with the two setups. Calculations were performed both on the full avidin tetramer and a system truncated to 26 Å around the ligand of interest.

Before the comparison, we must decide how we best can obtain reliable and efficient free energies with a well-defined precision. This has already been done for MM/GBSA,<sup>23</sup> and here we perform a similar analysis for LIE. Following our previous study,<sup>23</sup> as well as several other investigations,<sup>25,26</sup> we will not employ one single long simulation, but instead many shorter independent simulations, generated by using different starting velocities. This gives a more reliable estimate of statistical precision, and we can obtain any desired precision by simply employing a proper number of independent simulations, because the standard deviation of the mean decreases with the square root of the number of independent simulations included in the average. Therefore, we only need to determine the sampling



**Figure 2.** Correlation time of interaction energies of the truncated systems (the results for the full systems are similar): (a)  $\langle E_{\text{vdW}}^{L-S} \rangle_L$ , (b)  $\langle E_{\text{vdW}}^{L-S} \rangle_{\text{PL}}$ , (c)  $\langle E_{\text{ele}}^{L-S} \rangle_L$ , (d)  $\langle E_{\text{ele}}^{L-S} \rangle_{\text{PL}}$  and (e)  $\langle \Delta G_{\text{MM/GBSA}} \rangle$ .

frequency of the energy terms in eqs 1, as well as the length of the equilibration and production parts of the individual simulations.

**Correlation Time.** For all methods that average energies over a series of MD snapshots, it is essential to ensure that the consecutive estimates are independent, i.e., that sampling is not too frequent—otherwise the statistical error will be underestimated. The correlation time of the four time series that are the basis of the LIE estimates (eq 1) were calculated using the method of statistical inefficiency.<sup>54,55</sup> This was done for the whole 800 ps (full system) or 1600 ps (truncated system)

simulations (including equilibration) and for all independent simulations of each ligand.

We soon discovered that the original method is very sensitive to the drift during the equilibration period and also to occasional bumps in the data, which often are seen in long simulations, giving strongly increased correlation times (an example is shown in Figure S1, Supporting Information). As a consequence, the estimated correlation time always increased when the simulation time was increased. Strictly speaking, this shows that there are nanosecond time-scale motions in the structure that may indicate that much longer equilibration and simulation times are needed

to obtain truly uncorrelated data that are independent of the starting structure.<sup>57</sup> However, in standard LIE and MM/GBSA applications, it is assumed that simulations on a nanosecond time-scale started from the crystal structure provide representative structures for binding-energy calculations.<sup>8,10</sup> Therefore, we decided to divide all simulations into segments of 200 ps, for which a separate correlation time was estimated with the original method. This gave correction times of 1–4 ps for most segments, but segments during the equilibration period and segments with bumps gave much higher values. The latter were ignored by taking the median of the segment estimates (cf. Figure S1).

With this approach, we obtained stable results for all systems, which are presented in Figure 2 as the cumulative frequency of the number of simulations (among the 80 or 20 independent simulations) that have a particular correlation time. It can be seen that the correlation time for the  $\langle E_{\text{vdW}}^{\text{L-S}} \rangle_{\text{L}}$  and  $\langle E_{\text{ele}}^{\text{L-S}} \rangle_{\text{L}}$  terms are always 3 ps or less. The  $\langle E_{\text{vdW}}^{\text{L-S}} \rangle_{\text{PL}}$  term has a slightly longer correlation time, up to 5 ps, whereas the  $\langle E_{\text{ele}}^{\text{L-S}} \rangle_{\text{PL}}$  term shows the longest correlation time, up to 18 ps for Btn5 but up to 7 ps for the other systems. The reason for the long correlation time for Btn5 is that this energy term shows a long-term oscillation (see Figure S2, Supporting Information). If segments of 160 or 100 ps are instead used, we obtain correlation times of 4 and 2 ps, respectively.

In our previous study, we showed that the MM/GBSA energies have a correlation time of about 5 ps, at which point 90% of the data were uncorrelated without discarding any data and 98% of the data if the first 100 ps of the simulations were discarded.<sup>23</sup> However, this conclusion was based on only two ligands (Btn1 and Btn2) and with a somewhat different setup (e.g., octahedral systems treated with particle-mesh Ewald summation and no neutralization). Therefore, we repeated the analysis also for MM/GBSA for all ligands. The results in Figure 2e show that the correlation time is up to 7 ps, but 90% of the data are uncorrelated already at 4 ps.

Altogether, these data indicate that the correlation time is shortest for the free-ligand simulations (2–3 ps), slightly longer for the MM/GBSA results (~4 ps), and still somewhat longer for the  $\langle E_{\text{ele}}^{\text{L-S}} \rangle_{\text{PL}}$  term (~6 ps). However, these differences are small and somewhat dependent on the details of the method to calculate the correlation time. Therefore, we decided to use the same correlation time for both methods and also for the two types of LIE simulations, 5 ps.

**Equilibration Time.** The next step is to determine the length of the equilibration period of the simulations, i.e., the part of the simulation that is excluded in the averages. Many methods are available to determine the equilibration time ( $t_{\text{eq}}$ ).<sup>8,23,55,58</sup> We have tested several different variants, e.g., including block averaging or reverse cumulative averaging with two different tests for normal distribution. Unfortunately, all methods to determine  $t_{\text{eq}}$  are sensitive to details and parameters of the algorithms. At the end, we decided to use the following scheme: For each ligand, we calculated block averages of either the MM/GBSA binding energy or the LIE  $\beta \langle E_{\text{ele}}^{\text{L-S}} \rangle + \alpha \langle E_{\text{vdW}}^{\text{L-S}} \rangle$  energy terms for the complex or free-ligand simulations for each 100 ps of the simulations. These averages were compared to the average over the last 400 ps (full system) or 500 ps (truncated system) of the simulation, and if the difference was over 2 kJ/mol, that block was rejected. The equilibration time was taken as the end of the last set of at least two consecutive rejected blocks, but including also isolated rejected blocks if they are one or two blocks away from a set of consecutive rejected blocks. By such a rule, we disregard

occasional isolated rejected blocks late in the simulation, because the aim of the equilibration period is to remove data with a drift at the beginning of the simulation, but not bumps later in the simulation. A minimum equilibration time of 100 ps was assumed for all simulations. Several examples of typical equilibration curves and our selection of  $t_{\text{eq}}$  are shown in Figure S3 (Supporting Information).

For the present comparison between LIE and MM/GBSA, the exact length of the equilibration period is not of prime interest, but rather whether one of the two methods has a longer equilibration time than the other. However, for the present test case, we do not see any clear tendency: The two methods show similar equilibration time for (the complex simulation) of most ligands (eight out of the 7 + 7 simulations with full and truncated protein). When this is not the case, MM/GBSA gives the shorter equilibration time for two of the simulations and LIE a shorter time for four of the simulations. Therefore, we decided to use the same equilibration time for both MM/GBSA and LIE (viz. the largest of the two individual values) to avoid that the comparison is biased by differences in the equilibration.

For the free-ligand simulations (which are relevant only for LIE), the equilibration time is normally shorter than that for the complex. However, for the three charged ligands, it is notable that the free-ligand simulations give a quite large variation in the LIE energies, which quite often give rise to occasional isolated large deviations in block averages (cf. Figure S3e, Supporting Information). However, since the simulations do not show any pronounced trends, only rejected blocks at the beginning of the simulation were omitted (in accordance to the rule given above).

The selected equilibration times are collected in Table S1 (Supporting Information). In variance to our previous investigation of MM/GBSA,<sup>23</sup> we allow for different equilibration times for different ligands, which is more realistic and economic. It can be seen that the equilibration times vary from the requested minimum of 100 ps up to 1000 ps for Btn3. During this process, we decided to increase the simulation time of the truncated systems from 800 to 1600 ps. For the full systems, this could not be afforded (remember that 80 + 20 independent simulations were run for each ligand, giving a total simulation time of 1.12  $\mu\text{s}$ ). It is notable that the equilibration times are longer than in our previous MM/GBSA investigation, in which all simulations were judged to be converged after 100 ps.<sup>23</sup> A typical example of a curve from the previous study is shown in Figure S4 (Supporting Information).

**Length of Production Simulation and Efficiency.** We have now determined the correlation and equilibration times for our simulations of the seven ligands. What remains is to determine the length of the production simulation. This is somewhat involved, because we also run a number of independent simulations for each ligand. Therefore, we can improve the precision either by elongating the production time of each independent simulation or by increasing the number of independent simulations. The latter is more effective because the standard error of the final estimate (average over the independent simulations) will decrease with the square root of the number of independent simulations, whereas the dependence on the production time is less clear, since the results are not fully independent (this is the reason why we use several independent simulations<sup>23</sup>). On the other hand, the independent simulations cost more, because an initial equilibration has to be run. Therefore, to reach an optimum distribution, we need to consider also the computational cost of the simulations and energy calculations (which of course depends on the simulated system and the computer equipment).



We will follow our previous suggestion to optimize the CPU time required to reach a certain precision, e.g.,  $s_{\text{av}} = 1 \text{ kJ/mol}^{23}$  (but the comparison of the two methods will not depend on this limit). We can estimate the standard error of each independent simulation,  $s_{\text{simu}}$ , with a certain number of production snapshots,  $n_{\text{prod}}$ , from our available data. The number of independent simulations ( $n_{\text{av}}$ ) needed to reach the desired precision is then simply

$$n_{\text{av}} = \frac{s_{\text{simu}}^2}{s_{\text{av}}^2} \quad (6)$$

With these data, we can then calculate the required total CPU time: A 50 ps MD simulation takes 8 and 0.6 CPU hours on a single 3.0 GHz Intel Xenon processor for the full and truncated systems, respectively. The MM/GBSA postprocessing energies take  $\sim 0.25$  CPU hours, irrespectively whether full or truncated systems are used (because the time is dominated by the entropy calculations, which are performed on truncated systems in both cases), whereas the LIE energies can be calculated without any overhead. Therefore, the time consumption for LIE is

$$\text{CPU}_{\text{LIE}} = n_{\text{av}}(t_{\text{eq}}c_{\text{MD}} + (n_{\text{prod}} - 1)fc_{\text{MD}}) \quad (7)$$

and for MM/GBSA

$$\text{CPU}_{\text{MM/GBSA}} = n_{\text{av}}(t_{\text{eq}}c_{\text{MD}} + (n_{\text{prod}} - 1)fc_{\text{MD}} + n_{\text{prod}}c_{\text{ene}}) \quad (8)$$

where  $f$  is the sampling frequency (so that  $(n_{\text{prod}} - 1)f$  is the length of the production simulation),  $c_{\text{MD}}$  is the cost of running the MD simulation, and  $c_{\text{ene}}$  is the cost of doing a single MM/GBSA energy calculation. We can now calculate the CPU consumption as a function of  $n_{\text{prod}}$  using the equilibration times and the sampling frequency determined in the previous section, as well as  $s_{\text{simu}}$  obtained from the simulations. As can be seen in Figure S5 (Supporting Information), the CPU shows a minimum when  $n_{\text{prod}}$  is varied, because the first term in eqs 7 and 8 depends on  $n_{\text{av}}$ , which decreases as  $n_{\text{prod}}$  is increased, whereas the other terms depend on  $n_{\text{av}}n_{\text{prod}}$ , which increases with  $n_{\text{prod}}$ . This is the optimum value of  $n_{\text{prod}}$ . Strictly speaking, the results depend on the equilibration time (ligands with long equilibration times prefer somewhat longer production times and therefore fewer independent simulations), but the dependence is rather weak. Moreover, for some ligands, it is also favorable to increase the sampling frequency. Therefore, we have optimized  $n_{\text{prod}}$  and  $f$  separately for each ligand (Table S2, Supporting Information). However, the results are quite similar if we average  $s_{\text{simu}}$  and the CPU time over all seven ligands (and then  $f = 5 \text{ ps}$  is optimal; cf. Figure S5b). These averaged results are given in Table 1.

It can be seen that for the complex simulation of LIE, it is most efficient to run short production simulations,  $\sim 50 \text{ ps}$ , and instead run many independent simulations (91–135). For the free ligand, it is more efficient to run longer simulations ( $\sim 300 \text{ ps}$ ) and fewer independent simulations (7–10). The reason for this is that the standard error of the free-ligand simulations is appreciably smaller for the complex simulations and that it decreases more with the number of snapshots. However, in both cases, the simulation time is shorter than typically is used with LIE, illustrating that long simulations underestimate the statistical uncertainty. On the other hand, the total simulation time, 7.4 + 2.3 ns for the truncated systems and 4.6 + 3.0 ns for the full systems, plus 10–140 ns equilibration, is much longer than normally used with LIE. The total LIE CPU times are 1050 and 3900 CPU hours for the truncated and full systems, respectively.

**Table 1. Optimum Estimates of  $n_{\text{prod}}$  and  $n_{\text{av}}$ , Together with the Corresponding  $s_{\text{simu}}$  (kJ/mol) and CPU (h) Estimates for LIE and MM/GBSA, Following the Procedure Described in the Text (eqs 6–8), Using  $f = 5 \text{ ps}$  and the  $t_{\text{eq}}$  Values Listed in Table S1 (Supporting Information), and Averaging  $s_{\text{simu}}$  and CPU over All Seven Ligands<sup>a</sup>**

system	method	$n_{\text{prod}}$	$n_{\text{av}}$	$s_{\text{simu}}$	CPU
truncated	MM/GBSA	6	366	19.1	3135
	LIE PL	12	135	11.6	1000
	LIE L	67	7	2.6	54
	LIE total		142		1053
full	MM/GBSA	10	168	13.0	5983
	LIE PL	11	91	9.6	3218
	LIE L	60	10	3.1	715
	LIE total		101		3933

<sup>a</sup>  $n_{\text{av}}$  and CPU are calculated for a desired precision of ( $s_{\text{av}} = 1 \text{ kJ/mol}$ ).

Looking at MM/GBSA instead, we see that also for this method, rather short production simulations are most efficient, 25 and 45 ps for the truncated and full systems, respectively. This amounts to a total time of 3135 and 5983 CPU hours, respectively. Thus, this analysis indicates that LIE is more efficient than MM/GBSA, by a factor 3 for the truncated systems and 1.5 for the full systems. Looking at the more detailed data in Table S2 (Supporting Information), it can be seen that for the truncated system, MM/GBSA requires between 1.7 and 6.7 times more CPU than LIE for the truncated system, whereas for the full system the ratios are 1.0–2.4.

From eqs 7 and 8, it can be seen that the CPU consumption depends on three terms, two of which are common to both methods, whereas the last one, the cost of the energy calculations, only applies to MM/GBSA. This inherent difference between the two methods will always favor LIE and also lead to MM/GBSA typically preferring a slightly lower  $n_{\text{prod}}$  than LIE. However, the importance of this difference decreases with the size of the system, because for the truncated system the third term is 4–8 times larger than the second term, whereas for the full system, it is only 30–60% of the second term.

On the other hand, this effect is counteracted by the fact that LIE requires simulations of both the free ligand and the complex, whereas MM/GBSA is based only on the simulations of the complex. If everything else were equal, this would compensate for the extra cost of the energy calculations, and MM/GBSA would always be preferable.

However, a third factor is also important, viz., the standard errors ( $s_{\text{simu}}$ ) of the various energies and their dependence on the number of snapshots, which will affect  $n_{\text{av}}$  and the optimum  $n_{\text{prod}}$ . As we will see below, there is little difference in the standard error of the MM/GBSA and LIE estimates of the binding energy from the simulations of the complex (although the various ligands show a rather large variation). However, the standard error for the free-ligand simulations are appreciably smaller than for the complex simulations and also shows a larger decrease with  $n_{\text{prod}}$ , as can be seen in Figure S6 (Supporting Information). This leads to a lower total number of required snapshots ( $n_{\text{av}}n_{\text{prod}}$ ) for the free ligand than for the complex (469 compared to 1620 for the truncated system and 600 compared to 1001 for the full system). The combination of these three factors gives the net efficiency illustrated in Table 1.

Table 2. Binding Free Energies for the Various Methods on the Truncated Systems in kJ/mol<sup>a</sup>

$\alpha$	LIE	LIE ( $\Delta G_{cc}$ )	LIE (opt)	LIE ( $\Delta G_{cc,opt}$ )	MM/GBSA
	0.18	0.18	1.15	1.15	
with $\Delta G_{cc}$	no	yes	no	yes	
Btn1	-9.7±2.0	-11.7±2.0	-117.1±2.2	-119.0±2.2	-125.2±2.3
Btn2	2.4±2.0	0.3±2.0	-107.9±2.2	-109.9±2.2	-105.3±2.7
Btn3	-5.7±1.6	-7.7±1.6	-105.8±1.7	-107.8±1.7	-111.4±1.8
Btn4	-19.7±1.0	-19.7±1.0	-133.2±1.3	-133.2±1.3	-98.1±1.1
Btn5	-10.9±0.7	-10.9±0.7	-86.6±0.9	-86.6±0.9	-54.2±2.3
Btn6	-16.4±0.5	-16.4±0.5	-86.0±0.7	-86.0±0.7	-58.6±1.3
Btn7	-11.7±0.7	-11.7±0.7	-45.0±0.7	-45.0±0.7	-14.2±0.7
MAD	34.7±0.5	34.7±0.5	52.4±0.6	53.3±0.6	37.4±0.7
MADtr	24.8±0.6	24.8±0.6	16.2±0.5	15.7±0.4	16.2±0.6
$r^2$	-0.27±0.08	-0.27±0.08	0.38±0.02	0.41±0.02	0.76±0.01
PI	-0.70±0.15	-0.70±0.15	0.69±0.02	0.69±0.02	0.95±0.01

<sup>a</sup> A negative  $r^2$  indicates that  $r$  is negative.

Table 3. Binding Free Energies Using Various Methods on the Full Systems in kJ/mol<sup>a</sup>

$\alpha$	LIE	LIE ( $\Delta G_{cc}$ )	LIE (opt)	LIE ( $\Delta G_{cc,opt}$ )	MM/GBSA
	0.18	0.18	1.15	1.15	
with $\Delta G_{cc}$	no	yes	no	yes	
Btn1	-2.5±1.5	-9.5±1.5	-109.3±1.7	-116.4±1.7	-123.5±1.4
Btn2	9.3±1.4	2.1±1.4	-103.6±1.6	-110.7±1.6	-114.4±1.4
Btn3	2.5±1.1	-4.2±1.1	-94.2±1.3	-100.9±1.3	-106.3±1.5
Btn4	-14.2±1.3	-14.2±1.3	-121.0±1.6	-121.0±1.6	-93.3±1.1
Btn5	-9.7±0.9	-9.7±0.9	-82.7±1.1	-82.7±1.1	-56.2±1.2
Btn6	-15.0±0.7	-15.0±0.7	-81.3±0.9	-81.3±0.9	-53.6±0.7
Btn7	-13.0±0.8	-13.0±0.8	-43.0±0.9	-43.0±0.9	-13.5±0.8
MAD	38.9±0.4	35.9±0.4	45.8±0.5	48.8±0.6	36.7±0.5
MADtr	27.6±0.4	24.2±0.4	15.9±0.5	14.1±0.5	16.1±0.4
$r^2$	-0.53±0.06	-0.33±0.08	0.39±0.02	0.50±0.03	0.79±0.01
PI	-0.75±0.04	-0.75±0.11	0.69±0.02	0.69±0.04	1.00±0.01

<sup>a</sup> A negative  $r^2$  indicates that  $r$  is negative.

**Affinity Estimates.** Next, we estimated the binding free energy for the seven biotin analogues using the optimum equilibration time and sampling frequency, but using all available snapshots for production (because these are the best estimates we can obtain with available data). For LIE, we used  $\beta = 0.5$  for the charged ligands (Btn1–Btn3) and 0.43 for the neutral ligands (Btn4–Btn7) and  $\alpha = 0.18$  for all ligands. These values are usually denoted the standard parametrization.<sup>13,14</sup> The binding free energies for the truncated and full systems are shown in Tables 2 and 3, respectively. It can be seen that the results with the standard parametrization are poor with negative predictive indices and negative correlation coefficients (although  $r^2$  is positive by definition). This is because LIE predicts a higher affinity for the neutral ligands than for the charged ones.

Kollman and co-workers have studied a similar set of biotin analogues with the LIE method, but they obtained a reasonable, positive correlation.<sup>19</sup> The reason for this is that they used a special value for  $\alpha$  (1.0), fitted to the experimental data. Using the standard parametrization for their data (available for all our biotin analogues, except Btn3), we obtain a negative correlation

and  $r^2 = 0.01$ . Reported electrostatic and van der Waals energies are rather similar to ours with correlation coefficients ( $r^2$ ) of 0.57 and 0.96 (Table S3, Supporting Information). The rather large difference in the electrostatic energy is probably caused by differences in the simulation setup: They employed smaller systems and neutralized only a minimum amount of titrable residues.<sup>19</sup>

Therefore, we also tried to optimize the  $\alpha$  value, keeping  $\beta$  at the default values. Since we use four different quality estimates (MAD, MADtr,  $r^2$ , and PI), we fitted  $\alpha$  to optimize each of these measures by varying  $\alpha$  from -5 to +5 with increments of 0.05. The results are shown in Table 4 for the truncated system (the full system gave similar results). Optimizing  $\alpha$  against  $r^2$  gave unrealistic binding affinities because the correlation coefficient benefits from large energy differences, which are obtained when the energies are scaled up.  $r^2$  converges asymptotically at  $\alpha > 20$ , but as can be seen in Table 4, both MAD and MADtr are poor already at  $\alpha = 5$ .

Optimizing  $\alpha$  against PI gave a nonsmooth dependence on  $\alpha$ , although it gave similar results as when optimizing against

**Table 4. Results for LIE Obtained after Optimizing the  $\alpha$  Parameter with Respect to the Four Quality Measures MAD, MADtr,  $r^2$ , and PI**

optimized measure	$\alpha$	MAD	MADtr	$r^2$	PI
MAD	0.70	17.9	16.5	0.23	0.64
MADtr	1.15	52.4	16.2	0.38	0.69
$r^2$	5.00	377.0	91.9	0.50	0.61
PI	1.10	35.3	16.2	0.34	0.69
Kollman et al. <sup>22</sup>	1.00	35.3	16.3	0.34	0.65

MADtr. On the other hand, fitting  $\alpha$  against MADtr and MAD gave a smooth, parabolic dependence on  $\alpha$ . The fits gave optimal values of  $\alpha = 0.70$  and  $1.15$ , respectively. For these two  $\alpha$  values, PI and MADtr differ by only 0.05 and 0.3, respectively, which probably are not statistically significant. On the other hand, MAD and  $r^2$  differ significantly. We consider it more important to obtain good relative estimates (high  $r^2$ ) than absolute estimates (low MAD), and we therefore prefer the  $\alpha$  value obtained by fitting to MADtr, 1.15.

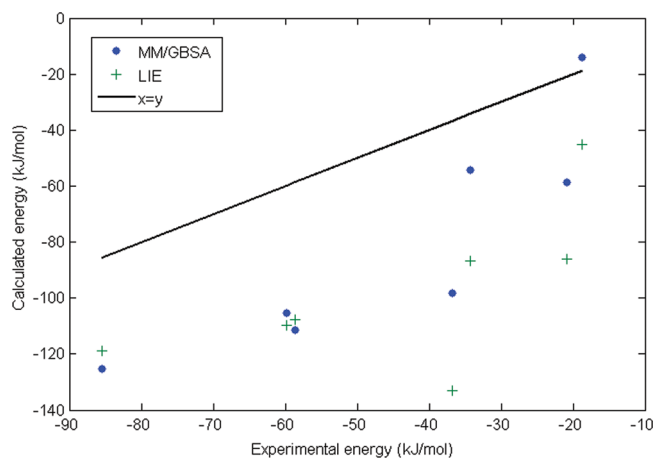
The LIE results using  $\alpha = 1.15$  are included in Tables 2 and 3 (column LIE (opt)) and are plotted in Figure 3. It can be seen that the largest error is found for Btn4. Kollman et al. also had a problem with this ligand. They argued that the error comes from the fact that the protein needs to be reorganized to accommodate the ester group of Btn4 and that such a term is missing in the LIE approach.<sup>19</sup>

Finally, we also considered the correction to the binding affinities of the charged ligands from the omitted surface charges,  $\Delta G_{cc}$ . For the truncated systems, this amounts to  $\sim 2$  kJ/mol, irrespectively of the ligand, and for the full systems it amounts to  $\sim 7$  kJ/mol. The largest contribution from a residue is  $\sim 2$  kJ/mol, showing that neutralization does not affect individual interactions so much. However, since avidin contains almost 100 charged surface residues, the sum is significant for the full systems. The effect of adding  $\Delta G_{cc}$  is shown in Table 2 for the truncated system, and as the correction is small, it hardly affects the result at all. However, for the full system, the correlation coefficient increases to 0.50, and MADtr decreases to 14 kJ/mol (see Table 3).

The MM/GBSA estimates of binding free energies are also included in Tables 2 and 3 and are plotted in Figure 3. For the truncated system, MADtr is 16 kJ/mol,  $r^2 = 0.76$ , and PI = 0.95. The full-system estimates are slightly better (see Table 3). Comparing with LIE, the difference in MADtr is not statistically significant, but the better results for the correlation coefficient and the PI are statistically significant. This is because MM/GBSA does not have any problems with Btn4. It is notable that the LIE and MM/GBSA results are quite similar for the three charged ligands (differences less than 6 kJ/mol), whereas for the neutral systems, the difference is appreciably larger (27–35 kJ/mol), MM/GBSA always giving a less favorable binding.

The binding affinities obtained with the truncated and full systems are quite similar. For LIE, they differ by up to 12 kJ/mol, the results of the full system almost always being more positive. For MM/GBSA, the results differ by 1–10 kJ/mol in a less systematic way.

Compared to our previous MM/GBSA binding affinities,<sup>23</sup> the difference is 1–19 kJ/mol (for the full system), i.e., much larger than the statistical uncertainty. The difference is systematic in that the negative ligands give more negative affinities (by 5–12 kJ/mol) in the present calculations, whereas the neutral ligands



**Figure 3.** Correlation between predicted and experimental free energies of the seven biotin analogues for the truncated systems. The LIE results are obtained with  $\alpha = 1.15$  and with the  $\Delta G_{cc}$  corrections.

give more positive affinities (by 1–19 kJ/mol). As a consequence, the present simulations give a similar MADtr (16 compared to 15 kJ/mol) but improved  $r^2$  (0.79 compared to 0.59) and PI (1.00 compared to 0.85). This shows that the MM/GBSA results depend much more strongly on the water model (TIP3P or TIP4P-Ewald), the treatment of long-range electrostatics (a spherical system with reaction-field corrections or an octahedral system with Ewald summation), and the treatment of surface charges (neutralization or not) than on the placement of the explicit water molecules, the initial velocities, and the protonation and rotation of residues, which gave variations of less than 1 kJ/mol for Btn1 in a previous investigation.<sup>24</sup> We have previously compared the results of spherical vs periodic simulations and neutralized systems with MM/PBSA, but the precision was too low to discern differences of relevant sizes.<sup>30</sup>

**Precision.** The statistical precision of the free-energy estimates is also shown in Tables 2 and 3 (standard deviation of the mean). It can be seen that LIE and MM/GBSA give similar uncertainties, 1–3 kJ/mol. For the truncated system, LIE with fitted  $\alpha$  and charge corrections gives a smaller uncertainty than MM/GBSA for five ligands. In general, the charged ligands show a slightly larger uncertainty than the neutral ligands. For the full systems, the precision is often slightly better (1–2 kJ/mol), and in most cases, MM/GBSA has a lower uncertainty than LIE (not for Btn3 and Btn5). This better precision of the full system is unexpected considering that the simulations are only half as long (0.8 ns compared to 1.6 ns). This shows that the intrinsic standard deviation of the data is larger for the truncated system than for the full system.

It should be noted that the LIE terms are scaled by the parameters  $\alpha$  and  $\beta$ . Without this scaling, the uncertainty of the LIE energies is larger than that of the MM/PBSA energies (and the two methods would become more equal in efficiency). The fact that LIE with  $\alpha = 1.15$  gives only a slightly larger uncertainty than with  $\alpha = 0.18$  (by 0.1–0.3 kJ/mol), although the van der Waals term is scaled up by a factor of  $1.15/0.18 \approx 6$ , shows that the precision of LIE is strongly dominated by the electrostatic term. This is also the reason why the optimization of the LIE procedure does not need to be redone with the optimized  $\alpha$  value.

It is somewhat disappointing that even with 80 + 20 independent simulations of 1.6 ns length, we have not been able to reach a

Table 5. Binding Free Energies (kJ/mol) for Each Subunit of Avidin in the Truncated Simulations<sup>a</sup>

	LIE				MM/GBSA			
	A	B	C	D	A	B	C	D
Btn1	-132.1±2.3	-99.6±3.0	-135.8±3.1	-100.8±2.7	-128.9±3.2	-112.8±3.5	-136.2±4.0	-120.6±2.8
Btn2	-121.9±3.4	-92.5±3.8	-124.6±3.2	-92.5±2.6	-107.0±2.0	-86.0±2.7	-104.4±1.1	-96.9±2.3
Btn3	-116.1±2.7	-95.6±2.4	-115.6±3.0	-96.0±2.3	-89.8±4.7	-115.3±2.4	-97.3±5.7	-110.6±2.0
Btn4	-137.0±2.5	-133.5±1.9	-131.9±2.4	-130.3±2.1	-99.8±1.4	-98.1±1.6	-99.5±1.6	-100.3±1.5
Btn5	-83.5±1.8	-89.3±1.4	-83.0±1.5	-90.6±1.4	-45.5±2.7	-66.6±2.7	-39.1±3.7	-63.4±1.7
Btn6	-86.3±1.1	-86.7±1.2	-82.9±1.6	-88.1±1.1	-58.6±3.0	-58.9±1.3	-52.7±2.0	-56.6±1.6
Btn7	-41.5±0.8	-48.1±0.6	-42.2±1.0	-48.0±0.6	-9.1±1.0	-21.3±0.7	-10.2±0.8	-20.5±0.8

<sup>a</sup> For LIE,  $\alpha = 1.20$  was used.

Table 6. Binding Free Energies (kJ/mol) for Each Subunit of Avidin in the Simulations of the Full System<sup>a</sup>

	LIE				MM/GBSA			
	A	B	C	D	A	B	C	D
Btn1	-108.8±2.9	-109.2±1.8	-109.1±2.9	-110.3±2.4	-116.8±3.4	-122.3±1.2	-121.5±2.7	-118.2±2.8
Btn2	-106.8±3.0	-103.5±1.9	-101.2±2.9	-102.7±2.6	-119.0±1.6	-112.4±0.7	-105.1±1.6	-101.5±1.6
Btn3	-95.0±3.1	-95.9±1.4	-91.2±2.6	-94.7±1.4	-100.4±4.2	-102.6±1.1	-99.5±3.2	-103.2±0.7
Btn4	-117.9±2.6	-124.4±1.9	-118.1±2.6	-124.2±1.9	-92.3±2.2	-91.9±2.0	-93.7±2.0	-93.7±2.1
Btn5	-77.9±1.6	-87.0±1.4	-80.0±1.8	-86.3±1.3	-51.0±2.6	-61.2±1.4	-48.3±1.9	-58.6±1.8
Btn6	-78.9±1.4	-83.1±0.9	-77.6±1.3	-85.1±1.2	-49.9±1.1	-56.0±0.8	-49.4±1.4	-58.2±0.9
Btn7	-39.6±1.8	-46.8±0.5	-37.9±1.5	-47.0±0.4	-9.4±1.2	-17.4±0.8	-7.7±1.3	-19.2±0.4

<sup>a</sup> For LIE,  $\alpha = 1.15$  was used.

precision of 1 kJ/mol for four of the ligands. In fact, for Btn1, which has the poorest precision, this would require 390 simulations or 620 ns simulation time. It is also notable that the precision of the MM/GBSA results is worse than in our previous investigation,<sup>23</sup> for which a precision of 1 kJ/mol was reached for all seven ligands with 20–50 300 ps simulations of the complex. This emphasizes a specific shortcoming of LIE for this tetrameric protein: With LIE, we can only simulate one ligand at a time, whereas for MM/GBSA, we could obtain four affinity estimates from the simulation of the complex with four ligands. We employed that opportunity in the previous study, but in this study, we used the same simulations for both LIE and MM/GBSA. Moreover, it seems that the convergence of the MM/GBSA energy terms is slower for a neutralized and truncated protein.

**Affinities of Individual Subunits.** We have previously shown that MM/GBSA estimates employing several independent simulations gave identical affinities for the four subunits in avidin within statistical uncertainty.<sup>23</sup> It is of interest to see if the same holds also for LIE. The binding free energies for the four subunits are shown in Tables 5 and 6 for the truncated and full systems, respectively, using  $\alpha = 1.15$ . It is evident that the four subunits do not give the same binding affinities for the truncated systems: The four subunits give results that differ by up to 21–36 kJ/mol for the charged ligands and by 5–8 kJ/mol for the neutral ones. This is much more than expected from the standard errors of the estimates, 2–4 and 1–2 kJ/mol, respectively. On the other hand, subunits B and D give binding affinities that are the same within statistical uncertainty (the difference is less than 3 kJ/mol), and the same applies to subunits A and C, although the difference is up to 5 kJ/mol. This indicates that the differences are caused by differences in the subunits, probably the fact that subunits A and

C have one less amino-terminal residue than subunits B and D in the crystal structure.

Surprisingly, for the full systems, the differences between the subunits are smaller, 2–9 kJ/mol, with no difference between the charged and neutral ligands. In this case, the differences are statistically significant only for Btn5–Btn7. Subunits B and D still give very similar results (within 2 kJ/mol), whereas the differences for subunits A and C are larger, up to 6 kJ/mol for Btn2 (but they are not statistically significant at the 95% level).

These large differences between the subunits for LIE led us to check also the MM/GBSA results. From Table 6, it can be seen that MM/GBSA actually gives similar differences between the subunits to LIE, 2–28 kJ/mol differences for the truncated systems and 2–18 kJ/mol for the full systems. This is a surprising difference compared to our previous results,<sup>23</sup> which most likely is caused by differences in the setup of the two sets of calculations, in particular the neutralization of surface charges, which may emphasize the one-residue difference between the subunits. Moreover, it is clear that the differences are amplified by the truncation of the protein. This is an important issue that will be the subject of future investigations of other proteins.

## CONCLUSIONS

In this study, we have designed a simulation protocol for the LIE method to reach a certain level of statistical precision in the predicted affinities, in the same way as in a previous study with MM/GBSA.<sup>23</sup> Our results indicate that for this biotin–avidin system, a sampling frequency of ~5 ps and equilibration times of 0.1–1.0 ns are appropriate. By optimizing the CPU time, we also suggest that rather short simulations should be used for the complex (50 ps after equilibration) but longer for the free ligand

( $\sim 300$  ps). Then, a proper number of independent simulations should be run until the desired precision is obtained. The sampling frequency and equilibration times are probably similar for most systems, whereas the length of the production simulation may depend on the simulated system and the computational equipment. Considering the long equilibration times observed for some systems, it would probably have been better to first run one long equilibration of each avidin–ligand complex (1–5 ps), before the independent simulations were started by using different sets of starting velocities. Then, the equilibration for the individual simulations could most likely have been reduced, as was observed for galectin-3 with MM/GBSA.<sup>23</sup>

In parallel, similar calculations have been performed with the MM/GBSA approach, based on the same simulations. This allows us to compare the efficiency of these two popular endpoint methods. We have reached several interesting conclusions:

- The correlation time of the LIE and MM/GBSA energies is similar.
- The equilibration time varies heavily with the ligand, but the two methods seem to require similar equilibration times.
- In general, LIE seems to be more efficient than MM/PBSA by a factor of 2–7 for the truncated systems, but by a factor of 1.0–2.4 for the full system (i.e., it gives the same statistical precision with a computational effort that is lower by these factors). The lower efficiency of MM/GBSA comes from the extra time required for the entropy calculation, which more than compensates for the fact that LIE requires an extra simulation (of the free ligand).
- On the other hand, in variance to MM/GBSA, LIE contains one empirical parameter,  $\alpha$ . If the standard value ( $\alpha = 0.18$ ) is used, LIE gives very poor results for this test case, with negative correlation and PI. However, if  $\alpha$  is fitted, LIE and MM/GBSA give similar MADtr,  $\sim 16$  kJ/mol, although MM/GBSA still outperforms LIE for  $r^2$  and PI. This is mainly due to LIE problems with a single ligand, Btn4.
- LIE is more restrictive in the setup of the simulation: It requires that the size of the simulated systems is the same for the complex and the free ligand and also that the protein is neutralized in the simulations. Our results indicate that this neutralization may slow the convergence and make the result different for the four subunits in avidin.
- Moreover, LIE simulations are typically performed on truncated systems with a radius of  $\sim 25$  Å.<sup>8</sup> Our results indicate that such a truncation may also slow the convergence and emphasize differences between the subunits.
- The change of the water model, the treatment of long-range electrostatics, and the neutralization of the protein have a quite large effect on the MM/GBSA binding energies (1–19 kJ/mol), much larger than the initial solvation, the starting velocities, as well as the protonation and rotation of residues.<sup>23</sup> It remains to be seen with larger test sets which of these setups is preferred, but for the present systems, the current setup gives a somewhat better correlation and PI compared to the experimental results.

Thus, we can conclude that LIE is inherently more effective than MM/GBSA (giving a certain precision at a smaller expense in computation time), at least for the present test case. Considering that the avidin tetramer is rather large and that LIE is typically run on truncated proteins, it is likely that this conclusion is valid also for other proteins, although more tests are required to confirm this. However, if the entropy term in MM/GBSA is

ignored, as has been done in many studies,<sup>5,59,60</sup> MM/GBSA is expected to become the more effective method. This might be an interesting alternative for MM/GBSA, especially as the entropy term has been criticized<sup>61</sup> and it limits both the precision and the CPU consumption.

On the other hand, we have seen that LIE depends on an empirical parameter and that it has more restrictions on the setup of the calculations. Clearly, MM/GBSA is disfavored by the LIE setup used in this study, giving a slower convergence, and it may be more effective with a more typical MM/GBSA setup. In particular, MM/GBSA may obtain four energy estimates from each snapshot for this tetrameric protein. Even more seriously, it is clear that the setup of the calculations quite strongly affects the results. It remains to be shown on much larger and more diverse test sets which of the setups are more realistic and which of the MM/GBSA and the LIE methods give the more accurate results.

## ■ ASSOCIATED CONTENT

**S Supporting Information.** Equilibration times and optimum  $f$ ,  $n_{\text{prod}}$  and  $n_{\text{av}}$  estimates for the various simulations; comparison with the results in ref 19; times series for two representative simulations; examples of the selection of the equilibration time; as well as the dependence of the CPU time and  $s_{\text{simu}}$  on  $n_{\text{prod}}$ . This material is available free of charge via the Internet at <http://pubs.acs.org>.

## ■ AUTHOR INFORMATION

### Corresponding Author

\*Tel.: +46 – 46 2224502. Fax: +46 – 46 2228648. E-mail: Ulf.Ryde@teokem.lu.se.

## ■ ACKNOWLEDGMENT

We thank Johan Åqvist and co-workers for help with LIE and the Q software package. This investigation has been supported by grants from the Swedish Research Council (project 2010-5025) and from the Research School in Pharmaceutical Science. It has also been supported with the computer resources of Lunarc at Lund University, C3SE at Chalmers University of Technology, and HPC2N at Umeå University.

## ■ REFERENCES

- (1) Michel, J.; Essex, J. W. *J. Comput.-Aided Mol. Des.* **2010**, *24*, 639–658.
- (2) Shirts, M.; Mobley, D. L.; Chodera, J. D. *Annu. Rep. Comput. Chem.* **2007**, *3*, 41–59.
- (3) Gohlke, H.; Klebe, G. *Angew. Chem., Int. Ed.* **2002**, *41*, 2644–2676.
- (4) Steinbrecher, T.; Labahn, A. *Curr. Med. Chem.* **2010**, *17*, 767–785.
- (5) Foloppe, N.; Hubbard, R. *Curr. Med. Chem.* **2006**, *13*, 3583–3608.
- (6) Åqvist, J.; Medina, C.; Samuelsson, J.-E. *Protein Eng.* **1994**, *7*, 385–391.
- (7) Hansson, T.; Marelus, J.; Åqvist, J. *J. Comput.-Aided Mol. Des.* **1998**, *12*, 27–35.
- (8) Brandsdal, B. O.; Östberg, F. M.; Almlöf, M.; Feierberg, I.; Luzhkov, V. B.; Åqvist, J. *Adv. Protein Chem.* **2003**, *66*, 123–158.
- (9) Srinivasan, J.; Cheatham, T. E., III; Cieplak, P.; Kollman, P. A.; Case, D. A. *J. Am. Chem. Soc.* **1998**, *37*, 9401–9409.
- (10) Kollman, P. A.; Massova, I.; Reyes, I.; Kuhn, B.; Huo, S.; Chong, L.; Lee, M.; Lee, T.; Duan, Y.; Wang, W.; Donini, O.; Cieplak, P.; Srinivasan, J.; Case, D. A.; Cheatham, T. E., III. *Acc. Chem. Res.* **2000**, *33*, 889–897.
- (11) Åqvist, J.; Hansson, T. *J. Phys. Chem.* **1996**, *100*, 9512–9521.

- (12) Almlöf, M.; Carlsson, J.; Åqvist, J. *J. Chem. Theory Comput.* **2007**, *3*, 2162–2175.
- (13) Almlöf, M.; Brandsdal, B. O.; Åqvist, J. *J. Comput. Chem.* **2004**, *25*, 1242–1254.
- (14) Carlsson, J.; Boukharta, L.; Åqvist, J. *J. Med. Chem.* **2008**, *51*, 2648–2656.
- (15) Carlson, H. A.; Jorgensen, W. L. *J. Phys. Chem.* **1995**, *99*, 10667–10673.
- (16) Swanson, J. M. J.; Henchman, R. H.; McCammon, J. A. *Biophys. J.* **2004**, *86*, 67–74.
- (17) Genheden, S.; Luchko, T.; Gusarov, S.; Kovalenko, A.; Ryde, U. *J. Phys. Chem. B* **2010**, *114*, 8505–8516.
- (18) Barril, X.; Gelpi, J. L.; Lopez, J. M.; Orozco, M.; Luque, F. J. *Theor. Chem. Acc.* **2001**, *106*, 2–9.
- (19) Kollman, P. A.; Kuhn, B. *J. Med. Chem.* **2002**, *43*, 3786–3791.
- (20) Hou, T.; Guo, S.; Xu, X. *J. Phys. Chem. B* **2002**, *106*, 5527–5535.
- (21) Salvalaglio, M.; Zamolo, L.; Busini, V.; Moscatelli, D.; Cavallotti, C. *J. Chrom. A* **2009**, *1216*, 8678–8686.
- (22) Wang, J.; Dixon, R.; Kollman, P. A. *Proteins: Struct., Funct., Genet.* **1999**, *34*, 69–81.
- (23) Genheden, S.; Ryde, U. *J. Comput. Chem.* **2010**, *31*, 837–846.
- (24) Genheden, S.; Ryde, U. *J. Comput. Chem.* **2011**, *32*, 187–195.
- (25) Lawrenz, M.; Baron, R.; McCammon, J. A. *J. Chem. Theory Comput.* **2009**, *5*, 1106–1116.
- (26) Fujitani, H.; Tanidal, Y.; Ito, M.; Jayachandran, G.; Snow, C. D.; Shirts, M. R.; Sorin, E. J.; Pande, V. S. *J. Chem. Phys.* **2005**, *123*, 084108.
- (27) Genheden, S.; Diehl, C.; Akke, M.; Ryde, U. *J. Chem. Theory Comput.* **2010**, *6*, 2176–2190.
- (28) Miyamoto, S.; Kollman, P. A. *Proteins: Struct., Funct., Bioinf.* **1993**, *16*, 226–245.
- (29) Kuhn, B.; Gerber, P.; Schulz-Gasch, T.; Stahl, M. *J. Med. Chem.* **2005**, *48*, 4040–4048.
- (30) Weis, A.; Katebzadeh, K.; Söderhjelm, P.; Nilsson, I.; Ryde, U. *J. Med. Chem.* **2006**, *49*, 6596–6606.
- (31) Brown, S. P.; Muchmore, S. W. *J. Chem. Inf. Model.* **2006**, *46*, 1493.
- (32) Kongsted, J.; Ryde, U. *J. Comput.-Aided Mol. Des.* **2009**, *23*, 63–71.
- (33) Söderhjelm, P.; Kongsted, J.; Ryde, U. *J. Chem. Theory Comput.* **2010**, *6*, 1726–1737.
- (34) Pugliese, L.; Coda, A.; Malcovati, M.; Bolognesi, M. *J. Mol. Biol.* **1993**, *231*, 698–710.
- (35) Green, N. M. *Biochem. J.* **1966**, *101*, 774.
- (36) Green, N. M. *Adv. Protein Chem.* **1975**, *29*, 85–133.
- (37) Green, N. M. *Methods Enzymol.* **1999**, *184*, 51–67.
- (38) Åqvist, J. *J. Comput. Chem.* **1996**, *17*, 1587–1597.
- (39) Hornak, V.; Abel, R.; Okur, A.; Strockbine, B.; Roitberg, A.; Simmerling, C. *Proteins: Struct., Funct., Bioinf.* **2006**, *65*, 712–725.
- (40) Cornell, W. D.; Cieplak, P.; Bayly, C. I.; Gould, I. R.; Merz, K. M.; Ferguson, D. M.; Spellmeyer, D. C.; Fox, T.; Caldwell, J. W.; Kollman, P. A. *J. Am. Chem. Soc.* **1995**, *117*, 5179–5197.
- (41) Bayly, C. I.; Cieplak, P.; Cornell, W. D.; Kollman, P. A. *J. Phys. Chem.* **1993**, *97*, 10269–10280.
- (42) Besler, B. H.; Merz, K. M.; Kollman, P. A. *J. Comput. Chem.* **1990**, *11*, 431–439.
- (43) Jorgensen, W. L.; Chandrasekhar, J.; Madura, J. D.; Impley, R. W.; Klein, M. L. *J. Chem. Phys.* **1983**, *79*, 926–935.
- (44) Marelus, J.; Kolmodin, K.; Feierberg, I.; Åqvist, J. *J. Mol. Graphics Model.* **1998**, *16*, 213–225.
- (45) Ryckaert, J. P.; Ciccotti, G.; Berendsen, H. J. C. *J. Comput. Phys.* **1977**, *23*, 327–341.
- (46) Berendsen, H. J. C.; Postma, J. P. M.; Van Gunsteren, W. F.; Dinola, A.; Haak, J. R. *J. Chem. Phys.* **1984**, *81*, 3684–3690.
- (47) Lee, F. S.; Warshel, A. *J. Chem. Phys.* **1992**, *97*, 3100–3107.
- (48) King, G.; Warshel, A. *J. Chem. Phys.* **1989**, *91*, 3647–3661.
- (49) Case, D. A.; Darden, T. A.; Cheatham, T. E., III; Simmerling, C. L.; Wang, J.; Duke, R. E.; Luo, R.; Crowley, M.; Walker, R.; Zhang, W.; Merz, K. M.; Wang, B.; Hayik, S.; Roitberg, A.; Seabra, G.; Kolossvary, I.; Wong, K.; Paesani, F.; Vanicek, J.; Wu, X.; Brozell, S. R.; Steinbrecher, T.; Gohlke, H.; Yang, L.; Tan, C.; Mongan, J.; Hornak, V.; Cui, G.; Mathews, D. H.; Seetin, M. G.; Sagui, C.; Babin, V.; Kollman, P. A. *Amber 10*; University of California, San Francisco: San Francisco, CA, 2008.
- (50) Brandsdal, B. O.; Smalås, A. O.; Åqvist, J. *FEBS Lett.* **2001**, *49*, 171–175.
- (51) Schutz, C. N.; Warshel, A. *Proteins* **2001**, *44*, 400–417.
- (52) Onufriev, A.; Bashford, D.; Case, D. A. *Proteins* **2004**, *55*, 383–394.
- (53) Kuhn, B.; Kollman, P. A. *J. Med. Chem.* **2000**, *43*, 3786–3791.
- (54) Bishop, M.; Frinks, S. *J. Chem. Phys.* **1987**, *87*, 3675–3678.
- (55) Yang, W.; Bitetti-Putzer, R.; Karplus, M. *J. Chem. Phys.* **2004**, *120*, 2618–2629.
- (56) Pearlman, D. A.; Charifson, P. S. *J. Med. Chem.* **2001**, *44*, 3417–3423.
- (57) Zuckerman, D. M. *Annu. Rev. Biophys.* **2011**, *40*, 41–62.
- (58) Allen, M. P.; Tildesley, D. J. *Computer Simulation of Liquids*; Clarendon Press: Oxford, U. K., 1991.
- (59) Hayes, J. M.; Skamnaki, V. T.; Archontis, G.; Lamprakis, C.; Sarrou, J.; Bischler, N.; Skaltsounis, A.-L.; Zographos, S. E.; Oikonomakos, N. H. *Proteins* **2011**, *79*, 704–719.
- (60) Hou, T.; Wang, J.; Li, Y.; Wang, W. *J. Chem. Inf. Model.* **2011**, *51*, 69–82.
- (61) Singh, N.; Warshel, A. *Proteins* **2010**, *78*, 1705–1723.

# Molecular Mechanics Investigation of an Adenine–Adenine Non-Canonical Pair Conformational Change

Keith P. Van Nostrand,<sup>†,‡</sup> Scott D. Kennedy,<sup>†</sup> Douglas H. Turner,<sup>‡,§</sup> and David H. Mathews<sup>\*,†,‡,||</sup>

<sup>†</sup>The Department of Biochemistry and Biophysics, University of Rochester Medical Center, 601 Elmwood Avenue, Box 712, Rochester, New York 14642, United States

<sup>‡</sup>Center for RNA Biology, University of Rochester Medical Center, 601 Elmwood Avenue, Box 712, Rochester, New York 14642, United States

<sup>§</sup>Department of Chemistry, University of Rochester, RC Box 270216, Rochester, New York 14627-0216, United States

<sup>||</sup>Department of Biostatistics and Computational Biology, University of Rochester Medical Center, 601 Elmwood Avenue, Box 712, Rochester, New York 14642, United States

## S Supporting Information

**ABSTRACT:** Conformational changes are important in RNA for binding and catalysis, and understanding these changes is important for understanding how RNA functions. Computational techniques using all-atom molecular models can be used to characterize conformational changes in RNA. These techniques were applied to an RNA conformational change involving a single base pair within a nine base pair RNA duplex. The adenine–adenine (AA) noncanonical pair in the sequence 5'GGUGAAGGCU3' paired with 3'PCCGAAGCCG5', where P is purine, undergoes conformational exchange between two conformations on the time scale of tens of microseconds, as demonstrated in a previous NMR solution structure [Chen, G.; et al. *Biochemistry* 2006, 45, 6889–903]. The more populated, major, conformation was estimated to be 0.5 to 1.3 kcal/mol more stable at 30 °C than the less populated, minor, conformation. Both conformations are trans-Hoogsteen/sugar edge pairs, where the interacting edges on the adenines change with the conformational change. Targeted molecular dynamics (TMD) and nudged elastic band (NEB) were used to model the pathway between the major and minor conformations using the AMBER software package. The adenines were predicted to change conformation via intermediates in which they are stacked as opposed to hydrogen-bonded. The predicted pathways can be described by an improper dihedral angle reaction coordinate. Umbrella sampling along the reaction coordinate was performed to model the free energy profile for the conformational change using a total of 1800 ns of sampling. Although the barrier height between the major and minor conformations was reasonable, the free energy difference between the major and minor conformations was the opposite of that expected on the basis of the NMR experiments. Variations in the force field applied did not improve the misrepresentation of the free energies of the major and minor conformations. As an alternative, the molecular mechanics Poisson–Boltzmann surface area (MM-PBSA) approximation was applied to predict free energy differences between the two conformations using a total of 800 ns of sampling. MM-PBSA also incorrectly predicted the major conformation to be higher in free energy than the minor conformation.

## INTRODUCTION

There are many known noncoding RNAs involved in diverse biological processes, such as binding and catalysis, where conformational changes are crucial for function.<sup>1–3</sup> Important biological roles of RNA depend upon single base pairs and bulged bases.<sup>4–7</sup> Diverse conformational changes in structured noncoding RNAs, such as rRNA and other ribozymes, are known to occur.<sup>8,9</sup> Conformational changes can alter binding surfaces to change their specificity, facilitate movement such as translocation of the ribosome,<sup>10–12</sup> or change activity as in riboswitches.<sup>13–15</sup> Thus, it is important to develop methods to model conformational changes to improve the understanding of RNA function.

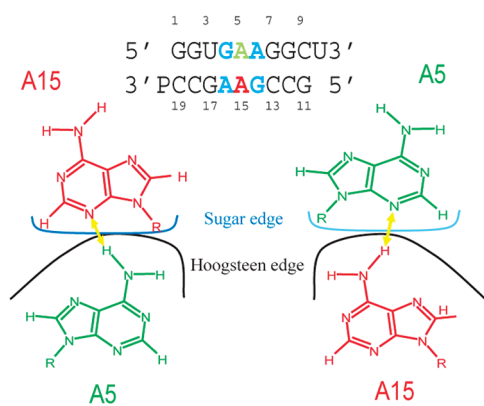
Molecular dynamics and umbrella sampling have been used previously to study conformational changes in both DNA<sup>16–18</sup> and RNA.<sup>19–21</sup> Many of these studies examine base pair opening<sup>17,19</sup> and base flipping,<sup>16,18,22</sup> where a base becomes unpaired, breaks its stacking interactions, and leaves the helix. This is described as a looped out or extra-helical state. A variety of

reaction coordinates are used in these studies, including a pseudodihedral angle known as the center of mass (COM or CPD) dihedral,<sup>18</sup> a projection of the glycosidic bond into a plane which is normal to a local helical axis vector and also includes the C'–C' vector,<sup>17</sup> and an improper dihedral angle defined by four atoms.<sup>19</sup> These studies demonstrate the use of free energy methods to give free energy for conformational changes involving individual bases or base pairs. RNA hairpin stability was also investigated using the end-to-end distance for two hairpin sequences known to be of different stabilities.<sup>21</sup>

The AA noncanonical pair system studied here is an RNA duplex of nine base pairs. The center base pair is a noncanonical pair consisting of adenine 5 (A5) and adenine 15 (A15) (Figure 1) that undergoes a conformational change from a major, i.e., more populated, A15–A5 to a minor, i.e., less populated,

Received: April 1, 2011

Published: October 04, 2011



**Figure 1.** Diagram of the AA noncanonical pair system, including the minor form on the left and major form on the right. The sequence for the AA noncanonical pair system is given at the top, including dangling ends that were removed for all simulations. Flanking GA pairs are in blue, A5 in green, and A15 in red. The yellow arrow indicates the single hydrogen bond stabilizing the noncanonical pair.

A5–A15 trans-Hoogsteen/sugar edge noncanonical base pairing interaction.<sup>23</sup> The NMR data obtained by Chen et al. for the duplex contained NOEs that were inconsistent with a single structure.<sup>24</sup> The NMR data were divided into two sets of NOE and dihedral angle restraints, based on a prior well-determined structure,<sup>25</sup> that were separately used to model structures for the major (Protein Data Bank, PDB, accession #: 2DD2) and minor (PDB: 2DD3) conformations.<sup>24</sup> Estimates of chemical shift differences of the H2 of adenine-5 between major and minor forms suggest an exchange rate of at least  $435 \text{ s}^{-1}$ , while line widths suggest higher rates ranging from  $20\,000$  to  $65\,000 \text{ s}^{-1}$ .<sup>24</sup> The observed chemical shifts of A5 and A15 H2 protons suggested that the major conformation is  $0.5$  to  $1.3 \text{ kcal/mol}$  more stable than the minor conformation at  $30 \text{ }^\circ\text{C}$ .

This study focused on modeling the conformational change pathway to provide insight about the pathway and to test the accuracy of modern molecular mechanics methods. For example, two types of pathways can be imagined for the conformational change. One pathway would have stacked intermediates where the adenines are stacked on each other, and the other would have intermediates where the adenines are hydrogen bonded in a common plane.

In this study, a number of computational methods, targeted molecular dynamics (TMD),<sup>26,27</sup> nudged elastic band (NEB),<sup>20,28–30</sup> umbrella sampling free energy calculations,<sup>31–34</sup> and MM-PBSA,<sup>35–43</sup> were used to model the conformational change pathway and equilibrium for a trans Hoogsteen/sugar edge AA noncanonical pair.<sup>23</sup> TMD applies a biasing potential to drive a starting structure toward a target structure in a molecular dynamics simulation. This can therefore generate plausible pathways for conformational changes. In contrast, NEB uses a string of states attached by virtual springs between fixed end states to generate a low potential energy pathway as determined by the potential energy landscape. These low potential energy conformational change pathways are close to likely pathways but are approximate because they neglect entropic effects.<sup>28–30</sup> True pathways involve thermal fluctuation and thus undergo random motions that do not follow exact minimum energy pathways; though the conformations visited by molecules undergoing conformational transitions are frequently along minimum potential energy pathways. NEB requires that the first and last

images of the string of conformations, i.e., the reactant and product structures, be fixed in conformation. Thus, the major and minor conformations are unchanged by the calculation. The low potential energy pathways determined with NEB are independent of the directionality of the pathway, and therefore the choice of reactant and product structures is arbitrary. TMD and NEB provide complementary information.

TMD and NEB predicted pathways where the adenines stack in the intermediates. These calculations suggested a reaction coordinate<sup>16</sup> described by an improper dihedral angle. Umbrella sampling and the weighted histogram analysis method (WHAM) were applied to predict relative free energies along the reaction coordinate.<sup>31–34</sup> The reaction coordinate observed in both TMD- and NEB-predicted conformational change pathways was used to facilitate umbrella sampling to determine the potential of mean force for the pathway.

The free energy profiles from umbrella sampling calculations with a molecular mechanics force field contradicted the experimental results. The relative free energy change between the major and minor conformations was overestimated and opposite that given the relative populations of these forms in solution. Free energy calculations were repeated using the parmbsc0<sup>44</sup> force field parameters, variations in salt concentration, or the TIP4PEW<sup>45,46</sup> water model. Each of these alternatives also was unable to correctly model the free energy difference between the major and minor conformations.

An alternative method, MM-PBSA/GBSA, for estimating free energies was also applied. This method uses a force field to estimate the molecular mechanics potential energy of the RNA in gas phase and then applies either the Poisson–Boltzmann<sup>35–37,39–42</sup> or generalized Born<sup>38,43</sup> surface area (PBSA or GBSA) methods of implicit solvation to estimate the free energy in solution. In addition, normal-mode analysis is used to predict the conformational entropy. MM-PBSA/GBSA also predicted a higher free energy for the major conformation than for the minor conformation and thus also did not qualitatively agree with experimental results. These results suggest that the AMBER99 force field is unable to adequately represent the conformational free energy change of this RNA molecule.

## METHODS

**Model Structures.** The models of the AA noncanonical pair system were those by Chen et al. for the major (PDB: 2DD2) and minor (PDB: 2DD3) RNA structures.<sup>24</sup> The dangling ends (unpaired U and purine), added to stabilize the structures for NMR, were removed for these calculations. Both structures have the same covalent structure and consist of 587 atoms. The lowest potential energy structure as evaluated by the AMBER99<sup>47,48</sup> force field from the reported set of NMR-guided models was selected as the representative structure for each conformation. This was structure 23 and structure 9 for 2DD2 and 2DD3, respectively.

**Modeling.** The AMBER<sup>48–50</sup> molecular dynamics package (version 9 or 10) was used for all calculations. Unless stated otherwise, the Cornell et al. ff99 (AMBER99) force field was used.<sup>47,48</sup> Trajectories were analyzed using ptraj<sup>51</sup> included with AMBER and the LOOS<sup>52</sup> software package.

**Implicit Solvent MD.** For implicit solvent molecular dynamics, generalized Born (GB) implicit solvation<sup>53–56</sup> was used, and the pairwise interactions were modeled with the HCT method.<sup>56</sup> A salt concentration of  $0.1 \text{ M}$  was simulated using a



Debye–Hückel limiting law for interaction screening.<sup>55</sup> An effective generalized Born radius of 25 Å was used with a nonbonded cutoff of 100 Å. Energy minimization was first performed with steepest descent for 1000 steps, and then for 10 000 steps via the conjugate gradient method. A 2 fs time step was used for dynamics. To start dynamics, the system was heated from 0 to 300 K over 60 ps followed by 100 ps of equilibrating MD. SHAKE<sup>57,58</sup> was applied for bonds to hydrogen. A Langevin thermostat, with a collision frequency of  $1 \text{ ps}^{-1}$ ,<sup>59,60</sup> was used.

**Explicit Solvent MD.** Explicit solvent NPT molecular dynamics utilized the TIP3P water model<sup>61,62</sup> with periodic boundary conditions. The TIP4PEW water model<sup>45,46</sup> was also investigated, where only the water model and ion parameters were changed and all other parameters were the same. Electrostatics were modeled with the particle mesh Ewald (PME) method<sup>63–65</sup> with a direct space sum cutoff of 10 Å. Neutralizing  $\text{Na}^+$  ions were added to counter backbone phosphates. Then, 1 M NaCl was added where the number of  $\text{Na}^+$  and  $\text{Cl}^-$  ions added was calculated by dividing the total number of water molecules added to the simulation box by the 55.5 M concentration of pure water.

The structure models from the PDB were subjected to two rounds of minimization, followed by heatup and equilibration. In the first minimization, the RNA was held fixed by a harmonic potential at  $500 \text{ kcal}/(\text{mol} \times \text{Å}^2)$  for 500 steps of steepest descent minimization followed by 500 steps of conjugate gradient minimization and with the system at constant volume. In the second round, the RNA was freed of restraint, and 1000 steps of steepest descent followed by 1500 steps of conjugate gradient minimization at constant volume were performed. The subsequent MD was run with a pressure relaxation time of 2 ps. A warmup simulation was performed with the RNA fixed in space with a harmonic potential of  $10 \text{ kcal}/(\text{mol} \times \text{Å}^2)$  for 20 ps of simulation. Dynamics were run with a 2 fs time step with SHAKE<sup>57,58</sup> constraining hydrogens and a Langevin thermostat collision frequency<sup>59,60,66</sup> of  $1.0 \text{ ps}^{-1}$ . Equilibration simulations were run at a constant pressure of 1 atm at a temperature of 300 K for 480 ps. Four separate trials of 100 ns of MD were performed for both the major and minor conformations by varying the random number seed used by the Langevin thermostat.

**TMD.** TMD<sup>26,27</sup> calculations were run from an energy-minimized structure to an energy-minimized target structure where the RMSD bias was applied to the entire molecule with the target RMSD set to 0 Å. GB implicit solvation<sup>53–56</sup> was used. A nonbonded and generalized Born radii cutoff of 100 Å was used. A salt concentration of 0.2 M was simulated by Debye–Hückel screening.<sup>55,56</sup> The temperature was set to 300 K, and Langevin dynamics<sup>66</sup> was used as a thermostat with a collision frequency<sup>59,60</sup> of  $1.0 \text{ ps}^{-1}$ . The time step was 2 fs, and SHAKE<sup>57,58</sup> was used.

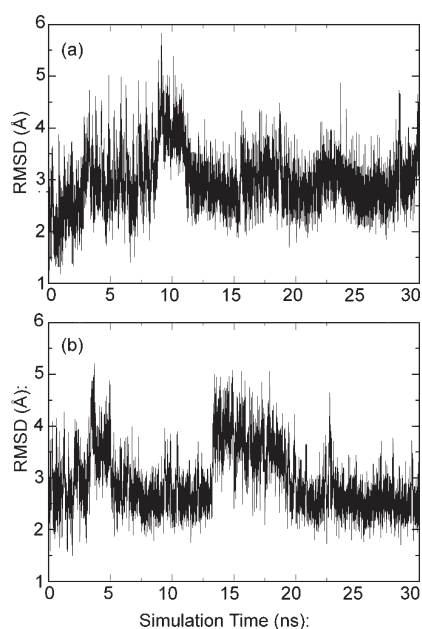
**NEB.** NEB<sup>20,28–30</sup> was done using AMBER with a 12 step simulated annealing<sup>67</sup> protocol that was concluded by quenched dynamics. The implementation uses revised tangents to prevent kinks in the pathway.<sup>30</sup> Trials were varied by changing the random number seed. The simulated annealing protocol applied here was the same as in previous work.<sup>68</sup> GB implicit solvation<sup>53–56</sup> was used in NEB calculations with a nonbonded and GB radii cutoff of 15 Å and no Debye–Hückel screening.<sup>55</sup> Langevin temperature scaling was used with a collision frequency of  $1000 \text{ ps}^{-1}$ .<sup>59,60,66</sup> SHAKE<sup>57,58</sup> was used. Fifteen trials were performed with 30 images, where the first 15 images were started as the major conformation and the last 15 images were started as

the minor conformation. The first image was fixed as the experimental model of the minor conformation and the last image as the major conformation. An additional trial with 60 images using 30 major and 30 minor images was performed to test resolution.

**Umbrella Sampling.** Umbrella sampling<sup>31</sup> calculations were performed along the improper dihedral reaction coordinate from  $-210^\circ$  to  $30^\circ$  with the equilibrium position of the harmonic potential in windows separated by  $10^\circ$ . Initial atomic structures for each window were selected manually from available structures from NEB-determined pathways. Structures were selected that were close to the equilibrium dihedral values for the windows. The improper dihedral restraint was applied with parabolic sides extending  $\pm 40^\circ$  from the window center with a maximum  $100 \text{ kcal}/(\text{mol} \times \text{rad}^2)$  restraining potential outside of the well. The sampling of the improper dihedral angle was verified to remain within the parabolic well for all umbrella sampling windows. The force constant value allowed for adequate overlap between distributions of reaction coordinate values between neighboring windows. NEB images selected for windows were solvated with TIP3P water<sup>61,62</sup> with an isometric box with nearest contact to the RNA at a radius of 10 Å. Neutralizing  $\text{Na}^+$  was added, and then an additional 1 M NaCl ions was added. The number of  $\text{Na}^+$  and  $\text{Cl}^-$  ions added was determined by dividing the number of water molecules in the box by the 55.5 M concentration of pure water. Umbrella sampling was repeated with neutralizing  $\text{Na}^+$  but without additional NaCl ions to the simulation box as well to test whether the presence of salt had an effect on the free energy surface.

Energy minimization was then performed as for explicit MD for the end states above followed by a heatup with the RNA solute fixed. Equilibration was run the same as for the explicit MD above with the addition of the improper dihedral restraint for the umbrella sampling window. The restraint was ramped up for the first 20 ps of simulation time from 0% to 100% strength and then maintained for the duration of equilibration and sampling. Different trials were run by changing the random number seed for all simulations. Following umbrella sampling, the weighted histogram analysis method (WHAM)<sup>33,34</sup> was used to obtain the free energy profile, i.e., the potential of mean force (PMF) along the reaction coordinate. The AMBER99 force field<sup>47–49</sup> with TIP3P water<sup>61,62</sup> was used unless otherwise noted. This procedure was applied similarly for umbrella sampling done with the TIP4PEW<sup>45,46</sup> water model and with the parmbsc0<sup>44</sup> force field.

**Averaging Multiple Free Energy Calculations.** When plotting multiple free energy profiles (free energy as a function of the reaction coordinate) on the same plot or when averaging free energy profiles, a reference profile was chosen, and all other profiles were adjusted in the free energy axis to minimize the sum of squares difference with the reference. This adjustment was calculated by averaging the difference of each point within the free energy profile from the corresponding point in the reference free energy profile. Then, the average of these differences was subtracted from the nonreference profile to obtain the new profile, adjusted to the minimum square difference relative to the reference profile. To average multiple profiles, the reference free energy profile and the profiles with the minimized sum of squares difference to the reference profile were averaged. The RMSD between free energy profiles was calculated using one profile as a reference and the other compared profile adjusted to minimize the sum of squares difference.



**Figure 2.** Mass weighted RMSD of all of the atoms to the solution structure for implicit solvent simulations for the minor (a) and major (b) conformations.

**Molecular Mechanics Poisson–Boltzmann and Generalized Born Surface Area (MM-PBSA/GBSA) Calculations.** MM-PBSA<sup>35–43</sup> calculations were performed on the eight total 100 ns trajectories of the major and minor conformations. The AMBER 10 mmpbsa.pl script was used to automate the procedure.<sup>69</sup> Snapshots were taken at 1 ps intervals from the simulations. The calculation was run including normal-mode analysis with 100 000 maximum cycles of minimization for each image from the MD trajectories to ensure convergence with a relaxed convergence criterion for the energy gradient of 0.001 kcal/mol. For GB implicit solvent calculations, Debye–Hückel screening for 1 M NaCl was used. The free energy for configurations sampled by a trajectory was the sum of the gas phase molecular mechanics energy, the estimated hydrophobic and reaction field energy contributions from either the PBSA or GBSA methods, the hydrophobic contribution to the solvation free energy for either the PB calculation or the GB calculation, and the product of conformational entropy estimated by normal-mode analysis and the temperature of 300 K.

## RESULTS

**Molecular Dynamics.** The models for the major and minor conformations of the AA noncanonical pair system were subjected to molecular dynamics as an initial test of their stability. The dangling unpaired uridine and purine present in the experimental model structures were not included in the calculations performed in this work. Given the estimated rate of exchange between the major and minor conformations from the NMR experiment between 20 000 and 65 000 s<sup>-1</sup>, it was expected that a spontaneous change in conformation was unlikely to be observed on molecular dynamics simulation time scales. A total of 30 ns of generalized Born implicit solvent molecular dynamics<sup>53–56</sup> was applied to both the major and minor structures. The structures have an RMSD to the starting experimental structure of less than 5.9 Å throughout the trajectory (Figure 2) with means of 2.93 Å

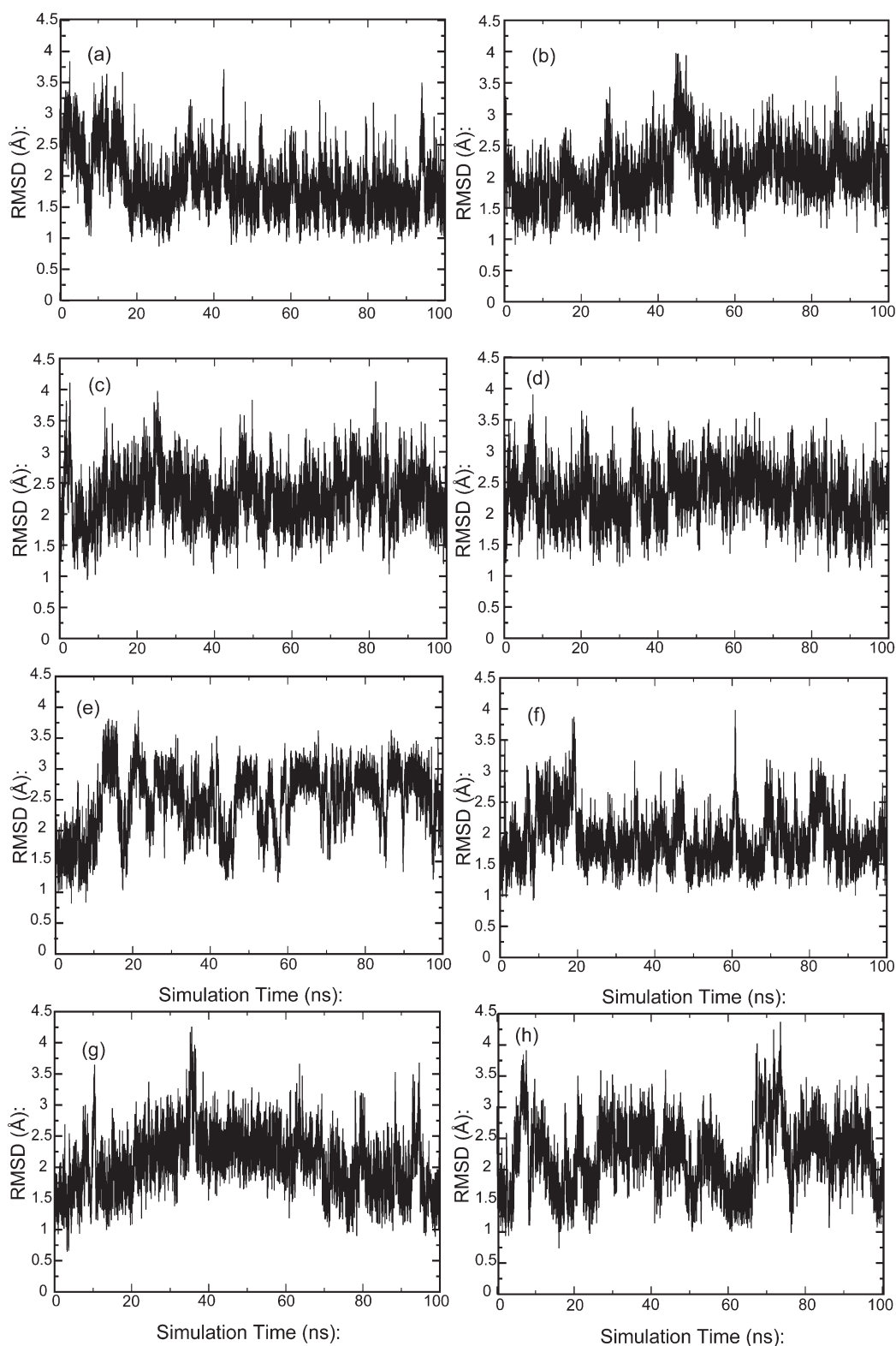
and 2.98 Å for the major and minor structures, respectively. In the major structure, A15 moves out of the helix away from A5 at 8.1 ns and makes temporary van der Waals and hydrogen bonded interactions with the edges of paired bases in the stem of the helix between the A5–A15 noncanonical pair and the G2–C17 base pair. In the minor structure, A15 leaves the helix and moves away from A5 at 5 ns and forms a hydrogen bond with the 2' hydroxyl of G8 for the remainder of the simulation. Thus, the unrestrained AA noncanonical pair was not stable during implicit solvent molecular dynamics for both the major and minor conformations.

Four 100 ns explicit solvent simulations were performed on each the major and minor conformations. All simulations had consistently lower RMSDs than the implicit solvent simulations (Figure 3, Table 1). The AA noncanonical pair retained its native (trans-Hoogsteen/sugar edge) pairing throughout each of the simulations. Over the four simulations, the mean RMSD was 2.17 and 2.13 Å, for the major and minor conformations, respectively.

Sugar puckers were also determined for all of the nucleotides in the major and minor conformations of the AA noncanonical pair structure from the TOCSY and NOESY spectra.<sup>24</sup> The Altona and Sunderlingam<sup>70</sup> convention was used to measure the sugar pucker along the MD trajectories of the major and minor conformations. The NMR results indicate that the nucleotides have C3' endo sugar puckers with angles ranging from 0 to 36°. Exceptions to this are for A5, which goes from C3' endo in the minor conformation to C2' endo with an angle range of 144–180° in the major conformation, and for A15 which goes from C2' endo in the minor conformation to C3' endo in the major conformation. The sugar puckers for A5 and A15 are plotted in Figure S1 for the minor conformation and Figure S2 for the major conformation (see Supporting Information). The time courses for sugar pucker indicate that A15 in the minor conformation usually shifts from the experimental C2' endo conformation fluctuating near 180° to the C3' endo conformation at 0° before 50 ns of simulation time. Otherwise, the sugar puckers of A5 in the minor conformation and A5 and A15 in the major conformation reasonably agree with the NMR experiment.

**Conformational Change Pathway Hypotheses.** The conformational change pathway translates A5 and A15 as shown in Figure 1 from the major configuration on the left to the minor configuration on the right. This can be imagined to occur in three possible ways. One possibility is that the bases move around the edges of each other in an edge-on-edge conformational change pathway, involving hydrogen-bonded intermediates. The other possibility is that the adenine bases of the noncanonical pair may slide one over the other through a stacked intermediate. The sliding pathway can occur in two alternative directions where the faces of the bases in the stacked intermediate are different. These two stacked pathways can be defined by which faces of the adenines were toward each other (Figure 4), thus giving a 5'-facing and 3'-facing stacked pathway.

**TMD Modeling.** TMD calculations were performed using a number of different force constants to give an initial estimate of conformational change pathways. The ideal range for the biasing potential force constant was determined to be 0.14–0.19 kcal/mol. Force constants below 0.14 kcal/mol resulted in trajectories where the RMSD with the target structure as a reference remained at 1.4 Å for the complete 1.5 ns of TMD simulation, indicating that the change to the target structure was not achieved (Figure S3, Supporting



**Figure 3.** Mass weighted RMSD of all solute atoms to the solution structure for 100 ns of explicit solvent simulations in TIP3P water with neutralizing  $\text{Na}^+$  plus 1 M NaCl for the minor (a, b, c, d) and major (e, f, g, h) conformations. Four simulations were performed for each by changing the random number seed.

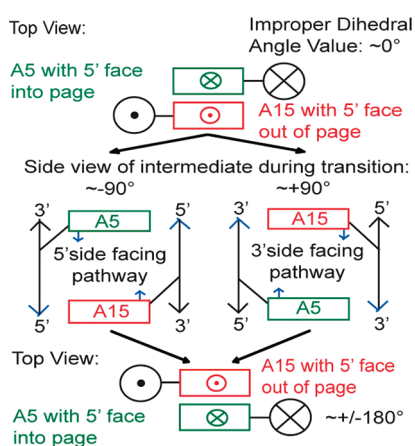
Information). Force constants greater than 0.19 kcal/mol were too high because the conformational change occurred within the first 30 ps of the trajectory.

Figure S4 (Supporting Information) shows RMSD as a function of time for four force constants. For each simulation, the RMSD started at the RMSD between the starting and target

**Table 1. RMSD Information for Implicit and Explicit Solvent Simulations of Major and Minor Forms of the AA Non-Canonical Pair System<sup>a</sup>**

simulation	avg. RMSD (Å)	min. RMSD (Å)	min. at time (ns)	max RMSD (Å)	max. at time (ns)	std. dev. (Å)
implicit minor	2.98	1.10	0.086	5.88	9.162	0.58
implicit major	2.93	1.34	0.880	5.32	3.713	0.60
explicit minor 1	1.87	0.83	88.856	3.83	2.586	0.46
explicit major 1	2.48	0.82	4.185	3.95	21.466	0.52
explicit major 2	2.05	0.86	2.827	3.98	44.645	0.43
explicit major 2	1.89	0.82	8.571	3.98	44.645	0.42
explicit minor 3	2.28	0.94	7.160	4.25	81.683	0.43
explicit major 3	2.08	0.61	3.266	4.26	35.586	0.48
explicit minor 4	2.31	1.06	92.692	3.93	7.467	0.40
explicit major 4	2.22	0.69	16.111	4.40	73.560	0.54

<sup>a</sup>Implicit simulations were both 30 ns in length, and all four independent explicit solvent simulations were 100 ns in length for both conformations.



**Figure 4.** Definition of 5'- and 3'-side facing pathways. The circle with a dot indicates the backbone with the 5' end coming out of the page. The circle with an X is the backbone with the 5' end going into the page.

structure (2.2 Å) and then dropped to an RMSD of approximately 1.4 Å within the first picosecond of TMD. The plot of RMSD to the target structure for the 0.14 kcal/mol force constant TMD calculation (Figure S4b) indicated a conformation in which a 1.4 Å RMSD was maintained for 0.6 ns of simulation before the conformation changed to the target structure. Adenine-5 and adenine-15 were stacked on each other within the helix during this time period, with the 5' faces of the adenines toward each other. Thus, TMD results showed a stacked intermediate structure along the pathway between minor and major configurations. Thirteen TMD trajectories with the biasing force constants from 0.01 to 0.13 kcal/mol each moved to a 5'-facing stacked configuration but did not continue to approach the target conformation. Seven TMD trajectories with biasing force constants from 0.14 to 0.2 kcal/mol underwent a 5'-facing pathway. A separate set of 11 TMD calculations with the major conformation as the initial state and the minor conformation as the target gave similar results.

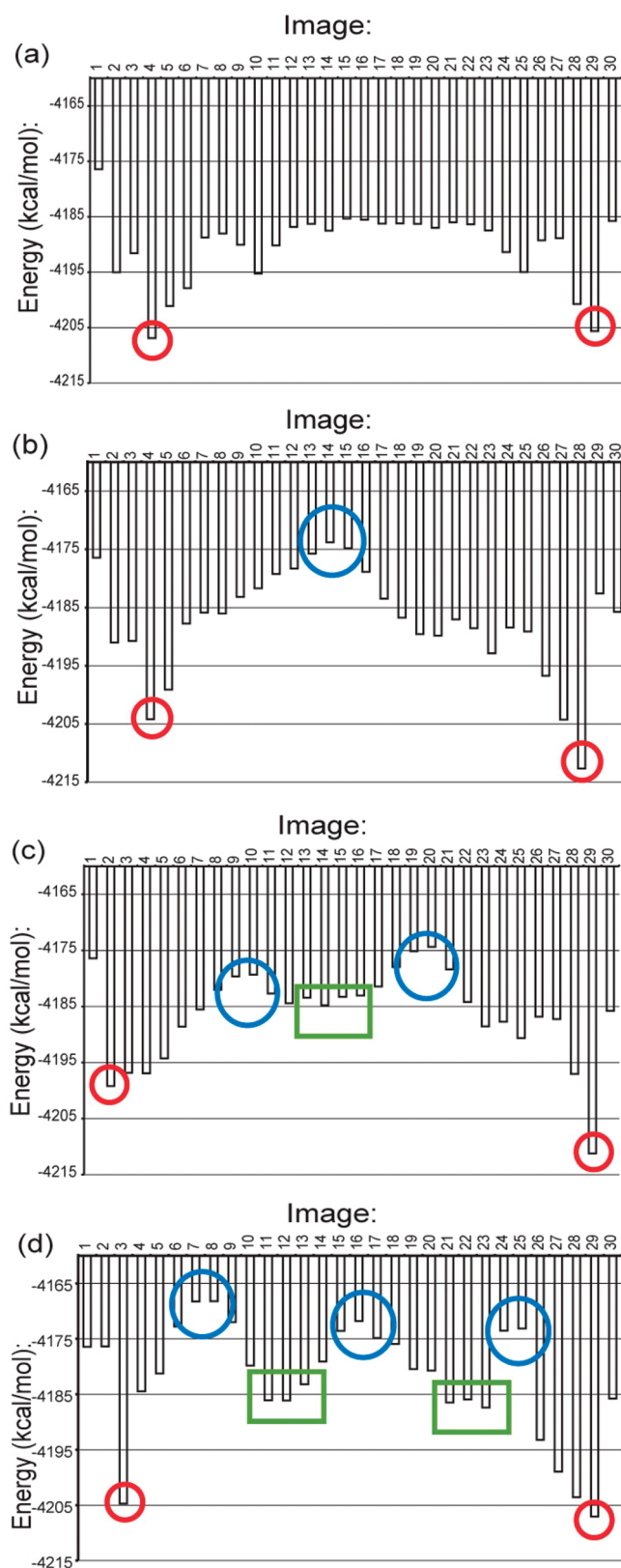
To address whether either face of the adenines can stack during the conformational change, TMD calculations were used to generate the alternative 3'-facing pathway within the same force constant range by using a modified starting structure. The modified structure was generated by first applying three distance restraints, found empirically, to an MD simulation of the minor

conformation to move adenine-5 to the 3' side of adenine-15. The distance restraints were all flat bottom harmonic potential wells with force constants of 30 kcal/mol. Outside of the well, the potential was at this constant. The three restraints were from the C8 of A5 to C2 of A15 with the well centered from 4.425 Å to 5.925 Å with outside bounds of 4.175 Å and 6.175 Å, N9 of G8 to C2 of A15 with the well centered from 9.278 Å to 10.278 Å with outside bounds of 8.278 Å and 11.278 Å, and C8 of A5 to N9 of G17 with the well centered from 4.53 Å to 5.53 Å with outside bounds of 3.53 Å and 6.53 Å. The starting minor experimental structure measured 3.88 Å for C8 of A5 to C2 of A15, 11.70 Å for N9 of G8 to C2 of A15, and 7.30 Å for C8 of A5 to N9 of G17. A structure was selected from the MD trajectory with the restraints applied as a starting configuration for TMD. With the altered starting structure, TMD followed the 3'-side pathway rather than the 5'-side pathway.

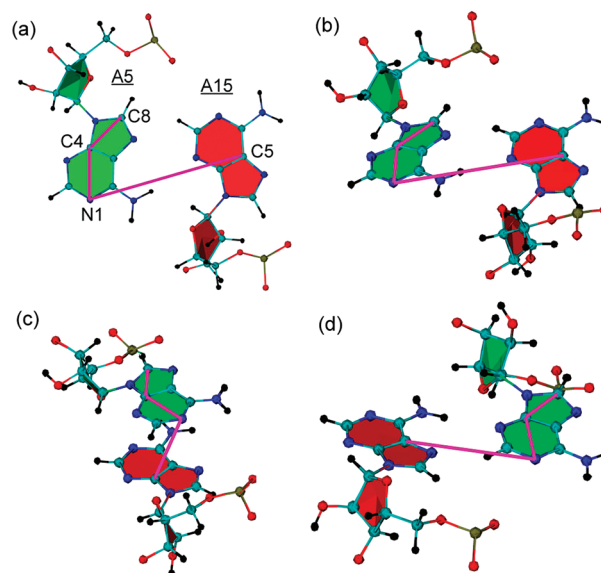
**NEB.** Fifteen NEB calculations with different random number seeds were performed. The end points are held fixed for NEB, thus reducing the error associated with their instability observed in implicit solvent calculations. Like the TMD simulations, the initial NEB pathways predicted stacked intermediates in the pathways with the adenines sliding by the 5' face (Figure 4). Although each NEB pathway slides by the 5' face, potential energy profiles of different trials revealed variability in the numbers of transition states and intermediates (Figure 5).

For the fifteen different NEB calculations, the images along the pathways showed varying degrees of progress along the conformational change from the major to minor forms. These NEB pathways provide a model for the conformational change. The initial images of the NEB pathways involve breaking of the hydrogen bonding interaction between the noncanonical pair adenines and the stacking interactions with the neighboring sheared GA pairs of the minor conformation, as shown in Figure 4. The stacking interactions broken include those between A5 and both A6 and A16, as well as those between A15 and both G14 and G4. The noncanonical pair adenines move to form the stacking interactions of the major configuration, as shown in Figure 6, by moving through configurations where the two adenines are stacked with each other through the 5'-facing pathway described in Figure 4. In the major form, the stacking interactions include those between A5 and both G4 and G14 as well as A15 with both A6 and A16.

It is possible that the bases could also traverse a pathway where the 3' faces are toward each other during sliding (Figure 4).



**Figure 5.** Variation in potential energy profiles for NEB trials. Red circles indicate images where the AMBER99 force field finds a more stable configuration than the experimental structure, an artifact of the force field. Potential energy profiles appear with no well-defined features (a), with single transition states (blue circles) (b), two transition states and one intermediate (green box) (c), and three transition states and two intermediates (d).

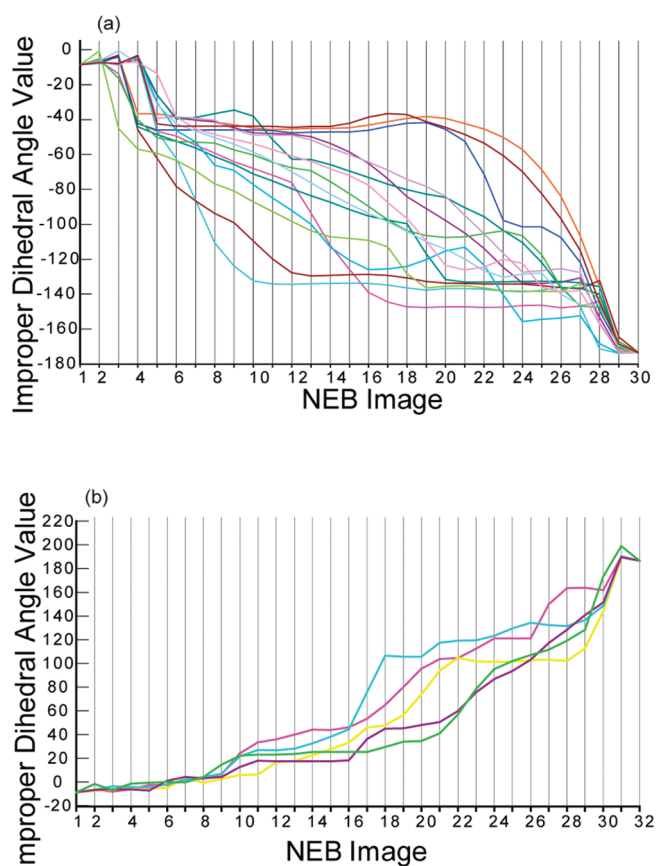


**Figure 6.** Definition of improper dihedral coordinate for the AA noncanonical pair. (a) Atoms for the improper dihedral are labeled. (b) The minor or reactant state starts at  $-8.7^\circ$ . (c) The intermediate state for the  $5'$ -side pathway occurs at  $-93.5^\circ$ . (d) The major or product state is at  $-173.5^\circ$ .

The  $3'$ -facing pathway was explored by NEB. Snapshots from the  $3'$ -stacking TMD trajectories were used as starting images for five NEB calculations. The same NEB protocol was used for these calculations, but with 32 images total. These NEB pathways each converged on  $3'$ -stacked intermediates. The potential energies of the intermediate structures in the  $3'$ -side-biased NEB trials were similar to those of the  $5'$ -side NEB pathways described above. The average, maximum, minimum, and standard deviation of the potential energies for the images are plotted in Figures S5 and S6 in the Supporting Information. The NEB pathways appear to depend on the starting configurations, which has been noted previously for systems with periodic reaction coordinates.<sup>20</sup> Starting with the experimental major and minor conformations poises the NEB method to produce the  $5'$ -side pathway. Moving the adenines as described above using restraints biases both TMD and NEB to produce the  $3'$ -side pathway.

**Reaction Coordinate.** Given the pathways observed by TMD and NEB, an improper dihedral reaction coordinate defined by C8, C4, and N1 on adenine-5 and C5 on adenine-15 described the conformational change (Figure 6). Other candidate coordinates were tested, such as the glycosidic sugar angle ( $\chi$ ) and sugar pucker, but were found to be unsuitable reaction coordinates (results not shown). The glycosidic angle of A15 did not differ for the major and minor conformations. The glycosidic angle for A5 did have distinct values for the major and minor conformation; however, the value fluctuated up and down along the pathways and thus did not follow reaction progress. The sugar pucker was found to change suddenly, within two or three images at different points along the pathway, for both A5 and A15, and thus did not adequately follow reaction progress.

The A5 base plane was defined by choosing three atoms, namely, C8, N4, and N1. Defining the final atom as C5 of base A15 resulted in an improper dihedral angle that brought one mismatch based over the face of the other. The  $5'$ - and  $3'$ -facing pathways were distinguishable by whether the angle changed

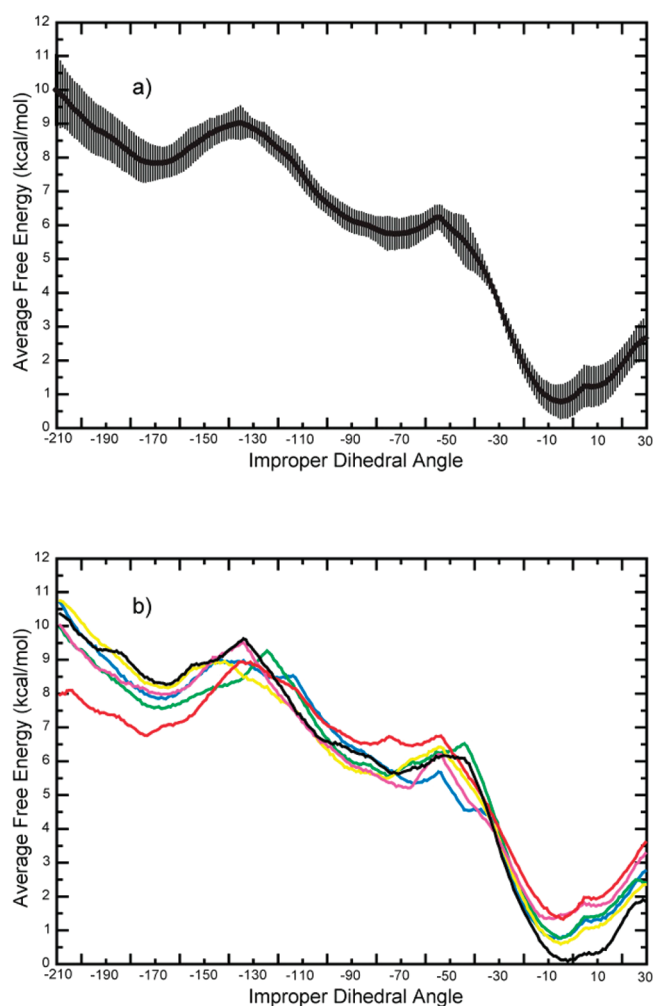


**Figure 7.** Plot of improper dihedral angle for initial NEB trials. (a) Plot of 15 NEB trials that follow the 5'-side pathway. Each trial is plotted in a distinct color. (b) Plot of improper dihedral angle for the five NEB trials that follow the 3'-side pathway.

negatively from  $-8.7^\circ$  to  $-173.5^\circ$  (5'-facing) or positively from  $-8.7^\circ$  to  $+186.5^\circ$  (3'-facing; Figure 6). The improper dihedral was plotted together for the images from the 15 NEB pathways (Figure 7a). The improper dihedral angle values were also plotted for the five NEB trails undergoing the 3'-facing pathway (Figure 7b). The improper dihedral value went from a definitive reactant value to a product value in a smooth continuous manner for each of the trials, showing that this behaved well as a reaction coordinate.

The improper dihedral angle was measured in the 100 ns molecular dynamics trajectories of the major and minor conformations (Figure S7, Supporting Information). The improper dihedral angle maintained expected values for all simulations, with the minor conformation simulations averaging  $-4.2^\circ$  and the major conformation remaining near  $-170.6^\circ$ . The average standard deviation was  $\pm 7.9^\circ$  for the minor conformation and  $\pm 14.0^\circ$  for the major conformation. This demonstrates that the improper dihedral angle fluctuates more during unrestrained molecular dynamics of the major conformation than for the minor conformation.

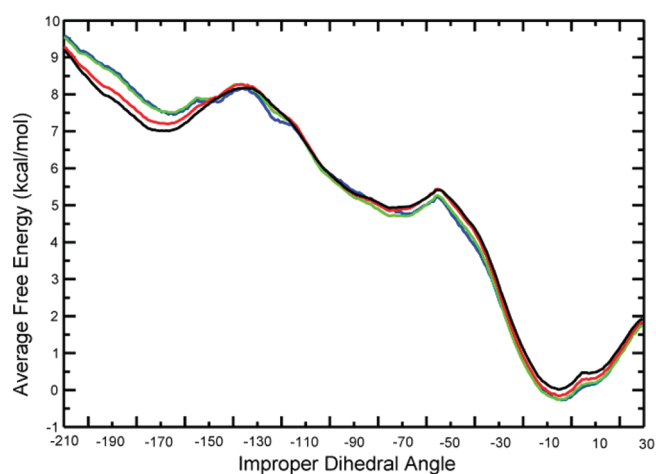
**Free Energy Calculations.** Umbrella sampling<sup>31</sup> and WHAM<sup>32–34</sup> were applied to predict the conformational free energy change along the reaction coordinate. Six independent calculations were run by varying the random number seed used for Langevin dynamics. Twenty-five windows for equilibrium sampling were used for each calculation. A total of 12 ns of



**Figure 8.** Free energy profiles from sampling with the AMBER99 force field<sup>47</sup> and 1 M NaCl. (a) A 12 ns sampling of 25 windows for six random number seeds was combined to produce the free energy profile with WHAM, and the plotted error is the standard deviation error between trials. In total, 1800 ns of sampling was used to generate the free energy profile. (b) The free energy profiles for the six random number seeds are plotted separately.

sampling was done for each window in each trial, yielding a total of 1800 ns of sampling. The average of the free energy profiles for each trial is plotted in Figure 8a. The standard deviations were calculated using the six trials, and these ranged from 0.19 to 1.04 kcal/mol along the profile. The individual free energy curves were also plotted separately (Figure 8b). Overall, the location of minima and maxima and shape of free energy barriers remained consistent between the trials, although some differences in barrier heights and the relative free energy between maxima and minima occur. The greatest free energy difference of 2.42 kcal/mol from the reference profile occurred at the bin centered at  $-209.5^\circ$  for one of the trials. This profile also had the highest RMSD difference from the reference profile of 1.16 kcal/mol. The greatest free energy difference of 2.81 kcal/mol between two trials also occurred at the bin centered at  $-209.5^\circ$ .

Separate umbrella sampling calculations were attempted for the 3'-facing pathway. Windows with the improper dihedral angle value centered at  $90^\circ$  to  $140^\circ$  consistently changed to an alternative conformation that was not supported by the NMR



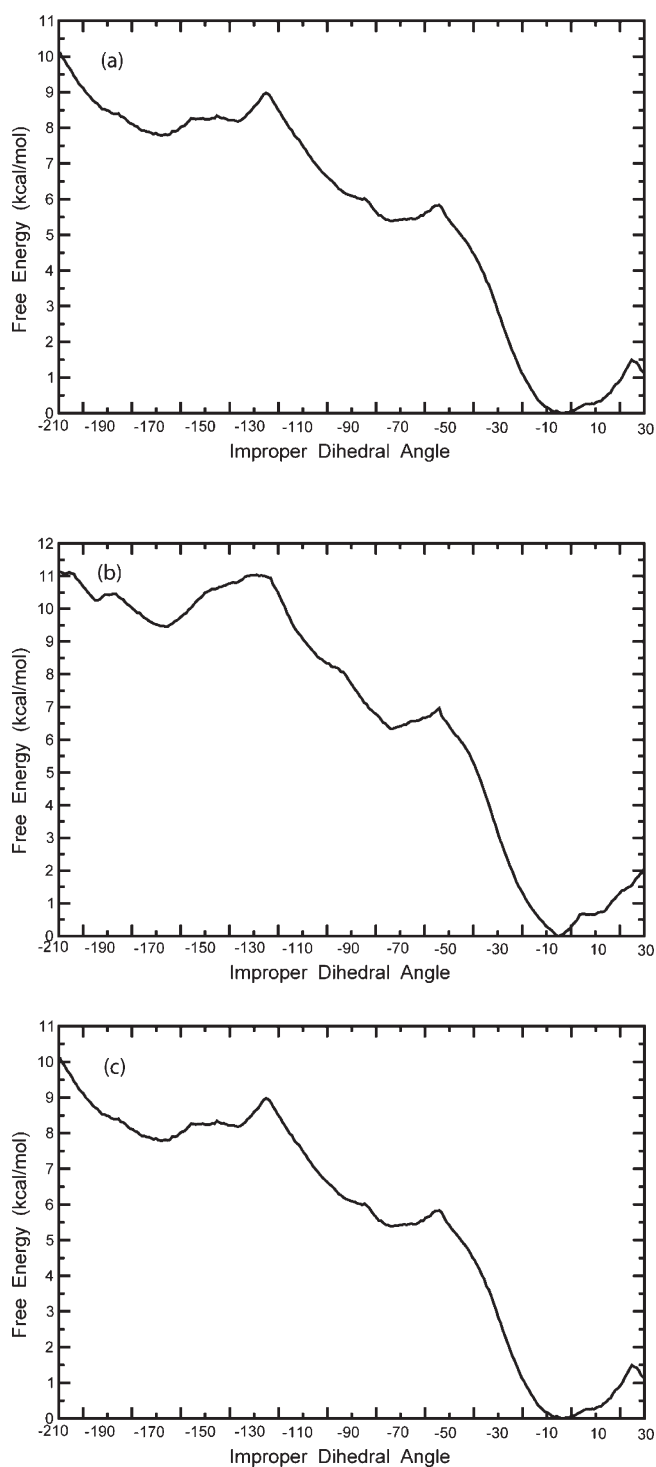
**Figure 9.** Convergence in time for the umbrella sampling calculations. Free energy profiles were generated by combining six random number seeds where the colored lines correspond to 2 (blue), 4 (green), 8 (red), and 12 (black) ns of sampling for all windows.

spectra. For these calculations, A5 moves parallel to the direction of the helix, no longer stacking or hydrogen bonding with A15. A5 also lost its stacking interactions with the flanking bases. A5 remained outside the helix in a configuration where the A5 edges were at a right angle to the edges of bases G4 and G14, possibly partly stabilized by transient electrostatic interactions between heteroatoms. Therefore, these calculations were stopped. This also implied that the 5'-facing pathway is a more viable model for the conformational change.

To test convergence, free energy profiles were generated combining the data of all six random number seeds for 2, 4, 8, and 12 ns of sampling at each equilibrium position (Figure 9). The greatest free energy difference of 0.81 kcal/mol between the 2 and 12 ns of sampling occurred at the bin centered at  $-186.5^\circ$ . The RMSD difference between the free energy profiles for 2 and 12 ns of sampling was 0.36 kcal/mol. Between 6 and 12 ns of sampling, this reduced to 0.22 kcal/mol. Free energy differences between the major and minor conformations decreased as sampling time was increased. Although the free energy difference changed, the locations of energy maxima and minima remained the same as the sampling time was increased. Therefore, 2 ns of sampling for each window was enough to establish the overall features of the free energy profile. The RMS difference of 0.13 kcal/mol between the free energy profile from 8 and 12 ns of sampling suggests that 12 ns sampling was adequately converged.

The relative free energy change between the major and minor conformations was estimated from the umbrella sampling, and the minor conformation was predicted to be more stable. From NMR data, the population of major to minor conformations was estimated from 70:30% to 90:10%, giving a favorable free energy for the major conformation of  $-0.51$  to  $-1.31$  kcal/mol.<sup>24</sup> The calculated free energy surface, using AMBER99, gave a free energy difference of 7.04 kcal/mol in favor of the minor conformation. This indicated the force field misrepresented the stability of the major and minor conformations.

Although evidence suggests that the umbrella sampling simulation time was adequate, one additional concern is that the adjacent windows are not sampling adjacent conformations in degrees of freedom aside from the chosen reaction coordinate. To test this, both the stacking interactions of the bases and the



**Figure 10.** Free energy profiles from calculations using modified force field parameters. (a) The free energy profile produced using the Barcelona parameters (parmbsc0)<sup>44</sup> with 12 ns of sampling for all 25 windows with TIP3P water, neutralizing  $\text{Na}^+$ , and 1 M NaCl. (b) The free energy profile produced using the AMBER99 force field with 12 ns of sampling and the TIP4PEW water model with neutralizing  $\text{Na}^+$  and 1 M NaCl. (c) Free energy profile from 12 ns of sampling of the AA noncanonical pair in neutralizing  $\text{Na}^+$  only and no additional salt. Overall features of the free energy profiles are similar to the final free energy curve with 1 M NaCl (Figure 8a), with major and minor states having similar energies and improper dihedral angle values.

Table 2. MM-PBSA Results from the Four 100 ns Explicit Solvent Trajectories for Both the Major and Minor Conformations<sup>a</sup>

	Nmode	MM-PBSA		MM-GBSA	
	$T\Delta S^{ob}$ (kcal/mol)	energy <sup>c</sup> (kcal/mol)	$\Delta G^o$ (kcal/mol) <sup>d</sup>	energy <sup>e</sup> (kcal/mol)	$\Delta G^o$ (kcal/mol) <sup>f</sup>
major structure mean	492.42 ± 1.12	−3776.07 ± 20.07	−4268.49 ± 19.66	−3690.05 ± 9.91	−4182.47 ± 9.60
simulation 1	493.27	−3746.19	−4239.46	−3675.62	−4168.89
simulation 2	490.86	−3782.55	−4273.41	−3691.67	−4182.53
simulation 3	493.18	−3788.33	−4281.51	−3695.59	−4188.77
simulation 4	492.37	−3787.20	−4279.57	−3697.33	−4189.70
minor structure mean	493.57 ± 0.07	−3791.76 ± 2.23	−4285.33 ± 2.18	−3696.78 ± 1.08	−4190.35 ± 1.05
simulation 1	493.60	−3788.76	−4282.36	−3695.19	−4188.79
simulation 2	493.62	−3791.44	−4285.06	−3697.04	−4190.66
simulation 3	493.46	−3793.75	−4287.21	−3697.48	−4190.94
simulation 4	493.57	−3791.76	−4285.33	−3696.78	−4190.35
total $\Delta$ (major − minor)	−1.15 ± 1.12	15.70 ± 20.20	16.84 ± 19.78	6.73 ± 9.97	7.88 ± 9.66

<sup>a</sup>Total  $\Delta$  values are determined by subtracting the mean of the minor conformation simulation from the mean of the major conformation and propagating the errors from the means. Mean and standard deviation values for major and minor conformations are calculated on the four trajectories.

<sup>b</sup>The conformational entropy from normal model analysis times the temperature of 300 K. <sup>c</sup>The free energy value reported here is the sum of PBSUR, the hydrophobic contribution to the solvation free energy from the Poisson–Boltzmann calculation, and PBCAL, the reaction field energy from the PB calculation. <sup>d</sup>The MM-PBSA free energy minus the normal-mode analysis entropy term. <sup>e</sup>Free energy value reported here is the sum of GBSUR, the hydrophobic contribution to the solvation free energy from the generalized Born calculation, and GB, the reaction field energy from the GB calculation.

<sup>f</sup>The MM-GBSA free energy minus the normal-mode analysis entropy term.

backbone dihedral angles were examined for all windows. There was no evidence that the sampling was structurally disconnected in adjacent windows.

The three base pairs of both stem regions were consistently stacked in helical form on both sides of the three base pair GAA internal loop in all umbrella sampling windows. Stacking interactions of the AA noncanonical pair and the flanking GA pairs were observed to change smoothly from one window to the next, where the minor conformation's stacking interactions were gradually lost for intermediate windows at and around  $-90^\circ$  and progressively reformed in windows closer to the major conformation's improper dihedral angle value of  $-170^\circ$ .

Backbone dihedral angles were measured for all of umbrella sampling windows and all residues. Histograms of representative residue backbone dihedrals are shown in Figure S8 (Supporting Information). These dihedrals were chosen because they show different distributions in the major and minor conformations, so that intermediate umbrella sampling windows can be checked for the degree with which they transition from one distribution to the other. The  $\delta$  dihedral angle for A5 was known from experiment to be  $160 \pm 30^\circ$  in the major conformation and  $122 \pm 67.5^\circ$  in the minor conformation.<sup>24</sup> The distribution of A5  $\delta$  angles in the umbrella sampling windows for the major and minor conformation agree with this result, showing major A5  $\delta$  to have an average value of  $141.9^\circ$  and minor A5  $\delta$  to have an average value of  $83.8^\circ$  (Figure S8b). Observation of the distributions of all umbrella sampling windows between the major and minor conformation shows a smooth change as the histograms of neighboring umbrella sampling windows have significant overlap. Three other backbone dihedral angles also exhibit a similar pattern with different values for the major and minor conformations and an overlapping distribution of neighboring windows along the conformational change pathway between the two conformations, including A6  $\alpha$  (Figure S8c,d), G14  $\delta$  (Figures S8e,f), and G4  $\zeta$  (Figure S8g,h). The average values of these four backbone dihedrals in the major and minor conformation umbrella sampling windows agree with the backbone dihedral

angles in the experimental model of the major and minor conformations. These plots support that sampling orthogonal to the improper dihedral angle, such as these backbone dihedrals, was related between umbrella window simulations along the conformational change pathway.

**Modified Molecular Mechanics Models.** Variation in potentials, salt concentration, water model, and periodic box size were tested to check whether the predicted relative free energies of the major and minor conformations would be closer to the expected experimental result. An update to the force field parameters that improved upon the  $\alpha$  and  $\gamma$  dihedrals known as the Barcelona parameters or parmbsc0<sup>44</sup> was tested and resulted in little change in the free energy profile as plotted in Figure 10a. The RMSD between the free energy profile produced with the original AMBER99 force field and the parmbsc0 force field was 0.54 kcal/mol with a maximum difference of 1.23 kcal/mol occurring at the bin centered at  $-165.5^\circ$ .

The TIP4PEW water model<sup>45,46</sup> and ion parameters were also tested with 1 M NaCl, and the resulting free energy profile is plotted in Figure 10b. The free energy profile has barriers in the same location as the TIP3P free energy profile, but the free energy difference between the major and minor conformations increased to 9.45 kcal/mol. The RMSD between the TIP3P and TIP4PEW free energy profiles was 1.04 kcal/mol with a maximum difference of 1.72 kcal/mol occurring at the bin centered at  $24.5^\circ$ . Thus, using the TIP4PEW water model instead of TIP3P did not improve the results produced by the AMBER99 force field.

Running the free energy calculations at 1 M NaCl was intended to mimic the electrolyte strength under physiological conditions. The NMR experiments<sup>24</sup> were performed in a solution containing 80 mM NaCl, 10 mM sodium phosphate, and 0.5 mM disodium EDTA. To investigate the effect of salt, the free energy calculations were repeated with only neutralizing  $\text{Na}^+$  ions and no additional NaCl with 12 ns of sampling per window (Figure 10c). The free energy profiles with and without 1 M NaCl have barriers at roughly the same locations. The profiles



have an RMSD of 0.39 kcal/mol and a maximum difference of 1.17 kcal/mol at the bin centered at 29.5°. Therefore, the 1 M NaCl simulations were considered adequate because running with minimal salt did not change the free energy profile significantly.

The size of the isometric solvent box for umbrella sampling calculations was determined using a distance of nearest contact to the solute. This was set to 10 Å as a compromise between having a box large enough to be physically reasonable but not so large as to become computationally prohibitive. To ensure that the box size does not affect the free energy result, a larger solvent box using a radius of nearest contact to the RNA of 20 Å was used in a separate trial of simulations. The TIP3P water model and neutralizing Na<sup>+</sup> ions plus 1 M NaCl ions were added as for the 10 Å box. A total of 2 ns of sampling for each of the 25 windows was performed and the resulting free energy profile generated with WHAM and plotted in Figure S9 (Supporting Information). The free energy of the major conformation was still estimated to be 8.85 kcal/mol greater than the minor conformation. Thus, increasing the explicit solvent periodic box size did not affect the resulting free energy.

**Molecular Mechanics Poisson–Boltzmann and Generalized Born Surface Area (MM-PBSA/GBSA) Calculations.** As an alternative to umbrella sampling, the MM-PBSA and MM-GBSA methods were used to predict the free energy changes between the major and minor conformations. Calculations of the free energy for the major and minor conformations were used to predict the free energy change for the conformational transition by subtracting the free energy of the minor conformation from the major conformation (Table 2). Eight total calculations were performed using the four independent 100 ns trajectories that were run for each conformation. This facilitated an estimation of the standard deviations of means. Both PB and GB solvation methods incorrectly predicted that the minor conformation was more stable than the major. MM-PBSA gave a larger positive free energy change,  $16.84 \pm 19.78$  kcal/mol, than MM-GBSA,  $7.88 \pm 9.66$  kcal/mol. MM-PBSA/GBSA gave an incorrect positive free energy change similar to the umbrella sampling result.

The MM-PBSA/GBSA results demonstrate a potential problem with the method. The first 100 ns major structure molecular dynamics simulation MM-PBSA free energy deviates from the other three simulations by about 40 kcal/mol. The major contributor to this energy difference results from a much lower gas phase electrostatic potential energy from the force field. The RMSD and structure of the simulation appears similar to the other simulations, indicating that the configurations visited by this simulation happen to produce a much lower electrostatic energy. Because electrostatic interactions can be strong and highly sensitive to atomic coordinates, it is plausible that a small change in conformation could yield a large change in electrostatic energy. For MM-PBSA/GBSA, the free energy difference between the major and minor conformations is calculated by taking the difference between two large numbers and is thus inherently prone to error.

## DISCUSSION

This study reports the modeling of the conformational change of an AA noncanonical pair. It highlights the capability of modern methods for modeling conformational change and predicting conformational free energy changes. This study also

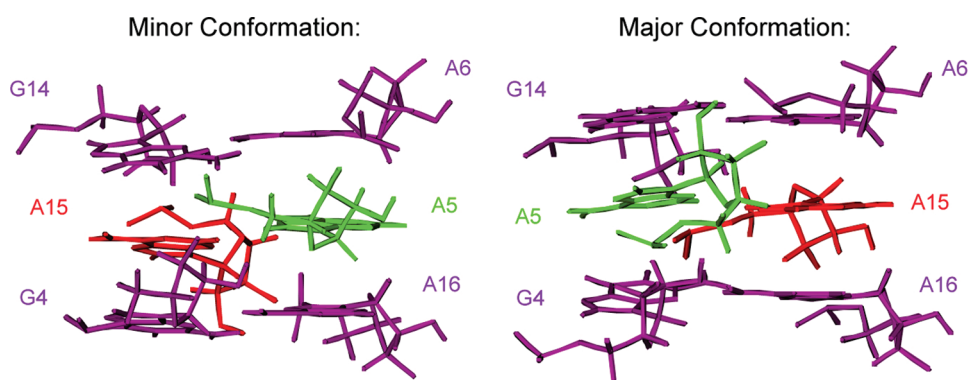
reveals limitations in the AMBER force field, a commonly used force field for modeling RNA dynamics.

TMD and NEB provided complementary information for the modeling of the conformational change. Each method predicted a pathway that involved stacked intermediates, and interestingly, the RNA structure showed enough flexibility for the AA to change conformation without the adjacent base pairs being broken. This was also previously observed in modeling the conformational change of a GG noncanonical pair.<sup>20</sup>

The TMD and NEB pathways suggested a dihedral reaction coordinate that could be used to follow the conformational change. This was then used to predict the free energy change along the pathway using umbrella sampling and WHAM. The rate of convergence and the magnitude of uncertainty for umbrella sampling were explored by repeating the umbrella sampling calculations with six independent trials. A total of 1800 ns of sampling was performed. The maxima, minima, and barriers of the free energy profile appear in the first 2 ns of sampling, but the independent trials have profiles differing with RMSDs as high as 0.56 kcal/mol even out to 12 ns. It can thus be concluded that the true rate of convergence is slower than the 12 ns of sampling performed per window. The need for extensive sampling is consistent with prior observations in peptide systems.<sup>71</sup> Over six trials, the error in free energies can be estimated as the standard deviations of the separate trials. These errors are a result of incomplete sampling of conformational space available to the molecule and are estimated to be less than 1.04 kcal/mol. Experimental errors are typically on the same order.

Both TMD and NEB predicted only a pathway in which the 5' sides of the bases face each other in the intermediate structure. To model a pathway with the 3' bases facing each other, the starting structures needed to be altered so that the bases were poised to follow that pathway. The reaction coordinate that describes the conformational change is periodic, where the dihedral angle is defined by atoms C8, C4, and N1 on adenine 5 and C5 on adenine 15. It has been noted previously that sampling pathways with NEB can be incomplete for a periodic system. Here, it is shown that TMD can also miss pathways around a periodic coordinate. In this case, the 3'-facing pathway was ruled out subsequently as a viable pathway because the intermediate dihedrals were unstable in that direction during umbrella sampling.

Free energy changes along a conformational change reaction coordinate provide both an estimate of the relative stability of the ends and the barrier(s) between states. For this AA system, the AMBER ff99 force field incorrectly predicted the minor conformation as the more stable conformation by 7.04 kcal/mol. All six independent calculations provided the same conclusion. Furthermore, variations in the force field, including the Barcelona dihedral parameters, using TIP4PEW water, altering the concentration of additional salt in the box, and using a larger water box did not change this incorrect prediction. Salt concentration and the ions are important for determining the RNA structure.<sup>72</sup> The simulations appear relatively insensitive to salt concentration, as the calculated RMSD between the free energy profile from 1 M NaCl and from only neutralizing Na<sup>+</sup> with TIP3P water was only 0.39 kcal/mol. Umbrella sampling with TIP4PEW and the latest ion parameters for the water model did not improve upon the relative free energies of the major and minor conformations and, in fact, increased the free energy of the major conformation to 9.45 kcal/mol.



**Figure 11.** Minor and major conformation loop region with residues labeled. The minor conformation is shown on the left and the major conformation on the right. This illustration shows that A15 is stacked on both G4 and G14 in the major conformation, but stacked on A6 and A16 in the minor conformation. A5, however, is stacked on A6 and A16 in the major conformation and stacked between G4 and G14 in the minor conformation.

The parmbsc0 modification of the AMBER99 force field gave a free energy profile with an RMSD of only 0.54 kcal/mol from the original AMBER99 profile and thus did not improve upon the result.

To test the possibilities that the umbrella sampling was insufficient or that the force field inaccuracies are for structures in the transition state, the MM/PBSA and MM/GBSA methods were applied to predict free energy changes between the end points. These methods also incorrectly favored the minor conformation above the major. This supports the hypothesis that the force field is also inaccurate for native structures, not just for structures along conformational change pathways.

## CONCLUSION

In this study, the conformational change pathway of an AA noncanonical pair in RNA was modeled using both TMD and NEB. Free energy calculations were then performed along a reaction coordinate. This simple system provides an excellent test case of available computational methods because of the availability of NMR data.<sup>24</sup> Direct observation of the actual conformational change pathways and the configurations visited during the pathways, however, are not possible with current technology. Several conclusions can be drawn. Both NEB and TMD predict a conformational change pathway where the adenines slide one over the other within the helix with the 5'-side facing the opposing adenine. The pathways involve breaking of the hydrogen bond between the adenines of the noncanonical pair as well as the stacking interactions between the adenines and the flanking GA pairs as defined by the minor conformation as shown in Figure 11 with subsequent movement through intermediates with the adenines stacked next to each other. The pathways move through this stacked intermediate and reform the hydrogen bond and stacking that defines the major structure, as shown in Figure 11.

Predicted free energy changes using both umbrella sampling with 1800 ns of total sampling and MM-PBSA/GBSA incorrectly favor the minor conformation, suggesting that the AMBER99 force field can be improved. This system, because of its size and the available solution structures, provides a useful benchmark for testing force fields.

## ASSOCIATED CONTENT

**S Supporting Information.** Figures of the sugar pucker dihedrals of the A5 and A15 backbone sugars from the MD

trajectories used in the MMPBSA calculations; plots of the final RMSD to target for TMD trajectories as a function of force constant; RMSD to target for selected TMD trajectories as a function of time; maximum, minimum, and average potential energies for the NEB calculations of the 5'- and 3'-facing pathways; improper dihedral angle plots for the four 100 ns simulations of the major and minor conformations; histogram plots of selected backbone dihedrals for all umbrella sampling windows and the major and minor conformations; and the resulting free energy profile from sampling with a larger solvent box. This material is available free of charge via the Internet at <http://pubs.acs.org>.

## AUTHOR INFORMATION

### Corresponding Author

\*E-mail: David\_Mathews@urmc.rochester.edu.

## ACKNOWLEDGMENT

The authors thank Alan Grossfield, Tod Romo, Matthew Seetin, and Harry Stern for many helpful suggestions. This study was funded by a National Institutes of Health grants R01HG004002, to D.H.M., and R01GM22939 to D.H.T.

## REFERENCES

- (1) Williamson, J. R. Induced fit in RNA-protein recognition. *Nat. Struct. Mol. Biol.* **2000**, *7*, 834–837.
- (2) Patel, D. J. Adaptive recognition in RNA complexes with peptides and protein modules. *Curr. Opin. Struct. Biol.* **1999**, *9*, 74–87.
- (3) Leulliot, N.; Varani, G. Current topics in RNA-protein recognition: Control of specificity and biological function through induced fit and conformational capture. *Biochemistry* **2001**, *40*, 7947–7956.
- (4) Nagan, M. C.; Beuning, P.; Musier-Forsyth, K.; Cramer, C. J. Importance of discriminator base stacking interactions: molecular dynamics analysis of A73 microhelix(Ala) variants. *Nucleic Acids Res.* **2000**, *28*, 2527–2534.
- (5) Sarzynska, J.; Kulinski, T.; Nilsson, L. Conformational dynamics of a 5S rRNA hairpin domain containing loop D and a single nucleotide bulge. *Biophys. J.* **2000**, *79*, 1213–1227.
- (6) Paillart, J. C.; Westhof, E.; Ehresmann, C.; Ehresmann, B.; Marquet, R. Non-canonical interactions in a kissing loop complex: The dimerization initiation site of HIV-1 genomic RNA. *J. Mol. Biol.* **1997**, *270*, 36–49.
- (7) Berglund, J. A.; Rosbash, M.; Schultz, S. C. Crystal structure of a model branchpoint-U2 snRNA duplex containing bulged adenosines. *RNA* **2001**, *7*, 682–691.

- (8) Rázga, F.; Koča, J.; Šponer, J.; Leontis, N. B. Hinge-like motions in RNA kink-turns: The role of the second A-minor motif and nominally unpaired bases. *Biophys. J.* **2005**, *88*, 3466–3485.
- (9) Lilley, D. M. Analysis of global conformational transitions in ribozymes. *Methods Mol. Biol.* **2004**, *252*, 77–108.
- (10) Peske, F.; Savelsbergh, A.; Katunin, V. I.; Rodnina, M. V.; Wintermeyer, W. Conformational changes of the small ribosomal subunit during elongation factor G-dependent tRNA-mRNA translocation. *J. Mol. Biol.* **2004**, *343*, 1183–1194.
- (11) Matassova, N. B.; Rodnina, M. V.; Wintermeyer, W. Elongation factor G-induced structural change in helix 34 of 16S rRNA related to translocation on the ribosome. *RNA* **2001**, *7*, 1879–1885.
- (12) Rodnina, M. V.; Savelsbergh, A.; Wintermeyer, W. Dynamics of translation on the ribosome: molecular mechanics of translocation. *FEMS Microbiol. Rev.* **1999**, *23*, 317–333.
- (13) Avihoo, A.; Gabdank, I.; Shapria, M.; Barash, D. In silico design of small RNA switches. *IEEE Trans. Nanobiosci.* **2007**, *6*, 4–11.
- (14) Noeske, J.; Buck, J.; Furtig, B.; Nasiri, H. R.; Schwalbe, H.; Wöhnert, J. Interplay of 'induced fit' and preorganization in the ligand induced folding of the aptamer domain of the guanine binding riboswitch. *Nucleic Acids Res.* **2007**, *35*, 572–583.
- (15) Schwalbe, H.; Buck, J.; Furtig, B.; Noeske, J.; Wöhnert, J. Structures of RNA switches: Insight into molecular recognition and tertiary structure. *Angew. Chem., Int. Ed.* **2007**, *46*, 1212–1219.
- (16) Song, K.; Campbell, A. J.; Bergonzo, C.; de los Santos, C.; Grollman, A. P.; Simmerling, C. An Improved Reaction Coordinate for Nucleic Acid Base Flipping Studies. *J. Chem. Theory Comput.* **2009**, *5*, 3105–3113.
- (17) Hart, K.; Nyström, B.; Öhman, M.; Nilsson, L. Molecular dynamics simulations and free energy calculations of base flipping in dsRNA. *RNA* **2005**, *11*, 609–618.
- (18) Banavali, N. K.; MacKerell, A. D. Free energy and structural pathways of base flipping in a DNA GCGC containing sequence. *J. Mol. Biol.* **2002**, *319*, 141–160.
- (19) Giudice, E.; Várnai, P.; Lavery, R. Base pair opening within B-DNA: free energy pathways for GC and AT pairs from umbrella sampling simulations. *Nucleic Acids Res.* **2003**, *31*, 1434–1443.
- (20) Mathews, D. H.; Case, D. A. Nudged elastic band calculation of minimal energy paths for the conformational change of a GG non-canonical pair. *J. Mol. Biol.* **2006**, *357*, 1683–1693.
- (21) Deng, N. J.; Cieplak, P. Free Energy Profile of RNA Hairpins: A Molecular Dynamics Simulation Study. *Biophys. J.* **2010**, *98*, 627–636.
- (22) Giudice, E.; Várnai, P.; Lavery, R. Energetic and conformational aspects of A: T base-pair opening within the DNA double helix. *ChemPhysChem* **2001**, *2*, 673–677.
- (23) Leontis, N. B.; Westhof, E. Geometric nomenclature and classification of RNA base pairs. *RNA* **2001**, *7*, 499–512.
- (24) Chen, G.; Kennedy, S. D.; Qiao, J.; Krugh, T. R.; Turner, D. H. An alternating sheared AA pair and elements of stability for a single sheared purine-purine pair flanked by sheared GA pairs in RNA. *Biochemistry* **2006**, *45*, 6889–903.
- (25) Chen, G.; Znosko, B. M.; Kennedy, S. D.; Krugh, T. R.; Turner, D. H. Solution structure of an RNA internal loop with three consecutive sheared GA pair. *Biochemistry* **2005**, *44*, 2845–2856.
- (26) Schlitter, J.; Engels, M.; Kruger, P. Targeted Molecular-Dynamics - a New Approach for Searching Pathways of Conformational Transitions. *J. Mol. Graphics* **1994**, *12*, 84–89.
- (27) Schlitter, J.; Engels, M.; Kruger, P.; Jacoby, E.; Wollmer, A. Targeted Molecular-Dynamics Simulation of Conformational Change - Application to the T[ $\rightarrow$ ]R Transition in Insulin. *Mol. Simul.* **1993**, *10*, 291–308.
- (28) Alfonso, D. R.; Jordan, K. D. A flexible nudged elastic band program for optimization of minimum energy pathways using ab initio electronic structure methods. *J. Comput. Chem.* **2003**, *24*, 990–996.
- (29) Jónsson, H.; Mills, G.; Jacobsen, K. W. Nudged elastic band method for finding minimum energy paths of transitions. In *Classical and Quantum Dynamics in Condensed Phase Simulations*; Berne, B. J., Ciccotti, G., Coker, D. F., Eds.; World Scientific: Singapore, 1998; pp 384–404.
- (30) Henkelman, G.; Jónsson, H. Improved tangent estimate in the nudged elastic band method for finding minimum energy paths and saddle points. *J. Chem. Phys.* **2000**, *113*, 9978–9985.
- (31) Torrie, G. M.; Valleau, J. P. Non-Physical Sampling Distributions in Monte-Carlo Free-Energy Estimation - Umbrella Sampling. *J. Comput. Phys.* **1977**, *23*, 187–199.
- (32) Souaille, M.; Roux, B. Extension to the weighted histogram analysis method: combining umbrella sampling with free energy calculations. *Comput. Phys. Commun.* **2001**, *135*, 40–57.
- (33) Kumar, S.; Bouzida, D.; Swendsen, R. H.; Kollman, P. A.; Rosenberg, J. M. The Weighted Histogram Analysis Method for Free-Energy Calculations on Biomolecules. I. The Method. *J. Comput. Chem.* **1992**, *13*, 1011–1021.
- (34) Roux, B. The Calculation of the Potential of Mean Force Using Computer-Simulations. *Comput. Phys. Commun.* **1995**, *91*, 275–282.
- (35) Floris, F.; Tomasi, J. Evaluation of the Dispersion Contribution to the Solvation Energy - A Simple Computational Model in the Continuum Approximation. *J. Comput. Chem.* **1989**, *10*, 616–627.
- (36) Gallicchio, E.; Kubo, M. M.; Levy, R. M. Enthalpy-entropy and cavity decomposition of alkane hydration free energies: Numerical results and implications for theories of hydrophobic solvation. *J. Phys. Chem. B* **2000**, *104*, 6271–6285.
- (37) Honig, B.; Nicholls, A. Classical Electrostatics in Biology and Chemistry. *Science* **1995**, *268*, 1144–1149.
- (38) Kollman, P. A.; Massova, I.; Reyes, C.; Kuhn, B.; Huo, S. H.; Chong, L.; Lee, M.; Lee, T.; Duan, Y.; Wang, W.; Donini, O.; Cieplak, P.; Srinivasan, J.; Case, D. A.; Cheatham, T. E. Calculating structures and free energies of complex molecules: Combining molecular mechanics and continuum models. *Acc. Chem. Res.* **2000**, *33*, 889–897.
- (39) Lu, Q.; Luo, R. A Poisson-Boltzmann dynamics method with nonperiodic boundary condition. *J. Chem. Phys.* **2003**, *119*, 11035–11047.
- (40) Luo, R.; David, L.; Gilson, M. K. Accelerated Poisson-Boltzmann calculations for static and dynamic systems. *J. Comput. Chem.* **2002**, *23*, 1244–1253.
- (41) Sigalov, G.; Scheffel, P.; Onufriev, A. Incorporating variable dielectric environments into the generalized Born model. *J. Chem. Phys.* **2005**, *122*, 094511–15.
- (42) Sitkoff, D.; Sharp, K. A.; Honig, B. Accurate Calculation of Hydration Free-Energies Using Macroscopic Solvent Models. *J. Phys. Chem.* **1994**, *98*, 1978–1988.
- (43) Srinivasan, J.; Cheatham, T. E.; Cieplak, P.; Kollman, P. A.; Case, D. A. Continuum solvent studies of the stability of DNA, RNA, and phosphoramidate - DNA helices. *J. Am. Chem. Soc.* **1998**, *120*, 9401–9409.
- (44) Peréz, A.; Marchán, I.; Svozil, D.; Šponer, J.; Cheatham, T. E.; Laughton, C. A.; Orozco, M. Refinement of the AMBER force field for nucleic acids: Improving the description of alpha/gamma conformers. *Biophys. J.* **2007**, *92*, 3817–3829.
- (45) Horn, H. W.; Swope, W. C.; Pitera, J. W. Characterization of the TIP4P-Ew water model: Vapor pressure and boiling point. *J. Chem. Phys.* **2005**, *123*, 194504.
- (46) Horn, H. W.; Swope, W. C.; Pitera, J. W.; Madura, J. D.; Dick, T. J.; Hura, G. L.; Head-Gordon, T. Development of an improved four-site water model for biomolecular simulations: TIP4P-Ew. *J. Chem. Phys.* **2004**, *120*, 9665–9678.
- (47) Cornell, W. D.; Cieplak, P.; Bayly, C. I.; Gould, I. R.; Merz, K. M.; Ferguson, D. M.; Spellmeyer, D. C.; Fox, T.; Caldwell, J. W.; Kollman, P. A. A second generation force field for the simulation of proteins, nucleic acids, and organic molecules (1995, vol 117, p 5179). *J. Am. Chem. Soc.* **1996**, *118*, 2309–2309.
- (48) Case, D. A.; Cheatham, T. E., 3rd; Darden, T.; Gohlke, H.; Luo, R.; Merz, K. M., Jr.; Onufriev, A.; Simmerling, C.; Wang, B.; Woods, R. J. The Amber biomolecular simulation programs. *J. Comput. Chem.* **2005**, *26*, 1668–88.

- (49) Case, D. A.; Darden, T. A.; Cheatham, T. E., III; Simmerling, C. L.; Wang, J.; Duke, R. E.; Luo, R.; Merz, K. M.; Pearlman, D. A.; Crowley, M.; Walker, R. C.; Zhang, W.; Wang, B.; Hayik, S.; Roitberg, A.; Seabra, G.; Wong, K. F.; Paesani, F.; Wu, X.; Brozell, S.; Tsui, V.; Gohlke, H.; Yang, L.; Tan, C.; Mongan, J.; Hornak, V.; Cui, G.; Beroza, P.; Mathews, D. H.; Schafmeister, C.; Ross, W. S.; Kollman, P. A. *AMBER 9*; University of California: San Francisco, CA, 2006.
- (50) Pearlman, D. A.; Case, D. A.; Caldwell, J. W.; Ross, W. S.; Cheatham, T. E.; Debolt, S.; Ferguson, D.; Seibel, G.; Kollman, P. *AMBER, A Package of Computer-Programs for Applying Molecular Mechanics, Normal-Mode Analysis, Molecular-Dynamics and Free-Energy Calculations to Simulate the Structural and Energetic Properties of Molecules. Comput. Phys. Commun.* **1995**, *91*, 1–41.
- (51) Cheatham, T. E., III; Crowley, M.; Tsui, V.; Pitera, J.; Case, D. A.; Gohlke, H.; Tanner, S.; Absgarten, E.; Roe, D.; Frybarger, P.; Walker, R. C. *ptraj*, version 1.3; University of Utah: Salt Lake City, UT, 2010.
- (52) Romo, T. D.; Grossfield, A. LOOS: An extensible platform for the structural analysis of simulations. *Conf. Proc. IEEE Eng. Med. Biol. Soc.* **2009**, *1*, 2332–5.
- (53) Tsui, V.; Case, D. A. Theory and applications of the generalized Born solvation model in macromolecular Simulations. *Biopolymers* **2000**, *56*, 275–291.
- (54) Weiser, J.; Shenkin, P. S.; Still, W. C. Approximate atomic surfaces from linear combinations of pairwise overlaps (LCPO). *J. Comput. Chem.* **1999**, *20*, 217–230.
- (55) Srinivasan, J.; Trevathan, M. W.; Beroza, P.; Case, D. A. Application of a pairwise generalized Born model to proteins and nucleic acids: inclusion of salt effects. *Theor. Chem. Acc.* **1999**, *101*, 426–434.
- (56) Hawkins, G. D.; Cramer, C. J.; Truhlar, D. G. Pairwise Solute Descreening of Solute Charges from a Dielectric Medium. *Chem. Phys. Lett.* **1995**, *246*, 122–129.
- (57) Ryckaert, J. P.; Ciccotti, G.; Berendsen, H. J. C. Numerical-Integration of Cartesian Equations of Motion of a System with Constraints - Molecular-Dynamics of N-Alkanes. *J. Comput. Phys.* **1977**, *23*, 327–341.
- (58) Miyamoto, S.; Kollman, P. A. Settle - an Analytical Version of the Shake and Rattle Algorithm for Rigid Water Models. *J. Comput. Chem.* **1992**, *13*, 952–962.
- (59) Uberuaga, B. P.; Anghel, M.; Voter, A. F. Synchronization of trajectories in canonical molecular-dynamics simulations: Observation, explanation, and exploitation. *J. Chem. Phys.* **2004**, *120*, 6363–6374.
- (60) Izaguirre, J. A.; Catarello, D. P.; Wozniak, J. M.; Skeel, R. D. Langevin stabilization of molecular dynamics. *J. Chem. Phys.* **2001**, *114*, 2090–2098.
- (61) Jorgensen, W. L. Revised Tips for Simulations of Liquid Water and Aqueous-Solutions. *J. Chem. Phys.* **1982**, *77*, 4156–4163.
- (62) Jorgensen, W. L.; Chandrasekhar, J.; Madura, J. D.; Impey, R. W.; Klein, M. L. Comparison of Simple Potential Functions for Simulating Liquid Water. *J. Chem. Phys.* **1983**, *79*, 926–935.
- (63) Cheatham, T. E.; Miller, J. L.; Fox, T.; Darden, T. A.; Kollman, P. A. Molecular-Dynamics Simulations on Solvated Biomolecular Systems - the Particle Mesh Ewald Method Leads to Stable Trajectories of DNA, RNA, and Proteins. *J. Am. Chem. Soc.* **1995**, *117*, 4193–4194.
- (64) Darden, T.; York, D.; Pedersen, L. Particle Mesh Ewald - an N. Log(N) Method for Ewald Sums in Large Systems. *J. Chem. Phys.* **1993**, *98*, 10089–10092.
- (65) Essmann, U.; Perera, L.; Berkowitz, M. L.; Darden, T.; Lee, H.; Pedersen, L. G.; Smooth, A Particle Mesh Ewald Method. *J. Chem. Phys.* **1995**, *103*, 8577–8593.
- (66) Wu, X. W.; Brooks, B. R. Self-guided Langevin dynamics simulation method. *Chem. Phys. Lett.* **2003**, *381*, 512–518.
- (67) Kirkpatrick, S.; Gelatt, C. D.; Vecchi, M. P. Optimization by Simulated Annealing. *Science* **1983**, *220*, 671–680.
- (68) Réblová, K.; Štřelcová, Z.; Kulhánek, P.; Beššeová, I.; Mathews, D. H.; Van Nostrand, K.; Yildirim, I.; Turner, D. H.; Šponer, J. An RNA Molecular Switch: Intrinsic Flexibility of 23S rRNA Helices 40 and 68 5'-UAA/5'-GAN Internal Loops Studied by Molecular Dynamics Methods. *J. Chem. Theory Comput.* **2010**, *6*, 910–929.
- (69) Case, D. A.; Darden, T. A.; Cheatham, T. E., III; Simmerling, C. L.; Wang, J.; Duke, R. E.; Luo, R.; Crowley, M.; Walker, R. C.; Zhang, W.; Merz, K. M.; Wang, B.; Hayik, S.; Roitberg, A.; Seabra, G.; Kolossváry, I.; Wong, K. F.; Paesani, F.; Vanecek, V.; Wu, X.; Brozell, S. R.; Steinbrecher, T.; Gohlke, H.; Yang, L.; Tan, C.; Mongan, J.; Hornak, V.; Cui, G.; Mathews, D. H.; Seetin, M. G.; Sagui, C.; Babin, V.; Kollman, P. A. *AMBER 10*; University of California: San Francisco, CA, 2008.
- (70) Altona, C.; Sundaralingam, M. Conformational-Analysis of Sugar Ring in Nucleosides and Nucleotides - Improved Method for Interpretation of Proton Magnetic-Resonance Coupling-Constants. *J. Am. Chem. Soc.* **1973**, *95*, 2333–2344.
- (71) Grossfield, A.; Feller, S. E.; Pittman, M. C. Convergence of molecular dynamics simulations of membrane proteins. *Proteins: Struct., Funct., Bioinf.* **2007**, *67*, 31–40.
- (72) Draper, D. E. RNA Folding: Thermodynamic and Molecular Descriptions of the Roles of Ions. *Biophys. J.* **2008**, *95*, 5489–5495.

# A New Coarse-Grained Force Field for Membrane–Peptide Simulations

Zhe Wu, Qiang Cui,\* and Arun Yethiraj\*

Theoretical Chemistry Institute and Department of Chemistry, University of Wisconsin, Madison, 1101 University Avenue, Madison, Wisconsin 53706, United States

**S** Supporting Information

**ABSTRACT:** We present a new coarse-grained (CG) model for simulations of lipids and peptides. The model follows the same topology and parametrization strategy as the MARTINI force field but is based on our recently developed big multipole water (BMW) model for water (*J. Phys. Chem. B* **2010**, *114*, 10524–10529). The new BMW-MARTINI force field reproduces many fundamental membrane properties and also yields improved energetics (when compared to the original MARTINI force-field) for the interactions between charged amino acids with lipid membranes, especially at the membrane–water interface. A stable attachment of cationic peptides (e.g., Arg<sub>8</sub>) to the membrane surface is predicted, consistent with experiment and in contrast to the MARTINI model. The model predicts electroporation when there is a charge imbalance across the lipid bilayer, an improvement over the original MARTINI. Moreover, the pore formed during electroporation is toroidal in nature, similar to the prediction of atomistic simulations but distinct from results of polarizable MARTINI for small charge imbalances. The simulations emphasize the importance of a reasonable description of the electrostatic properties of water in CG simulations. The BMW-MARTINI model is particularly suitable for describing interactions between highly charged peptides with lipid membranes, which is crucial to the study of antimicrobial peptides, cell penetrating peptides, and other proteins/peptides involved in the remodeling of biomembranes.

## I. INTRODUCTION

Many biological processes that occur at the cellular membrane involve lipid membrane deformations at many length scales,<sup>1,2</sup> which are either triggered or facilitated by small peptides<sup>3</sup> or complex protein machineries.<sup>4,5</sup> To effectively complement experimental studies of these processes, it is important to develop computational models that are capable of describing membrane deformations as well as interactions between the membrane and peptides/proteins. The latter requirement highlights the importance of developing particle-based coarse-grained (CG) models for membrane systems, which are particularly useful for phenomena that occur on length and time scales too large for atomistic simulations but where continuum mechanical models<sup>1,6–8</sup> are not appropriate. In this paper we report a new CG model for lipids and peptides that is based on an accurate CG model for water developed in our groups.

The past decade has seen a flurry of activity in the development of CG models for biomolecules and lipids.<sup>9–13</sup> By grouping several atoms into a single unit, thus decreasing computational cost, CG models have proven valuable in many simulation studies of biomembranes and their interactions with peptides and proteins.<sup>12</sup> For example, the MARTINI force field<sup>14,15</sup> has been successfully applied to study lipid vesicle formation and fusion,<sup>16–19</sup> lipid phase transformation,<sup>20</sup> structure and dynamics of lipid bilayers and monolayers,<sup>14,15,21</sup> and effects of various molecules (e.g., cholesterol, proteins) on the shape and phase behaviors of complex membranes.<sup>22–25</sup> Solvent-free CG models for membranes have also been proposed and found useful in a number of studies,<sup>26–28</sup> although transferable protein models that are compatible with these membrane models have not yet been reported to our knowledge.

An important aspect of computational biophysics is the treatment of electrostatic interactions. Driven by a desire for computational efficiency, many CG models choose to remove and/or simplify the treatment of electrostatic interactions. In MARTINI, for example, four water molecules are grouped into a single uncharged unit (bead), and electrostatic interactions between charged beads, which represent either lipid head groups or charged amino acid side chains, are treated with a cutoff scheme and a fairly large dielectric constant ( $\epsilon = 15–20$ ). Although this model can be effective for describing interactions between lipids and nonpolar groups, we expect it to be less appropriate for describing the interaction between lipid membrane and highly charged species, such as cell penetrating peptides and antimicrobial peptides, for which a proper treatment of electrostatics is likely crucial. This is supported by the observation that although MARTINI gives satisfactory results (when compared to atomistic simulations) for the potential of mean force (PMF) for the penetration of nonpolar and polar (neutral or with a small dipole moment) amino acids into lipid bilayers, it incurs large errors for the PMF of charged amino acids, especially positively charged residues.<sup>23</sup>

A major source of the error in MARTINI can be attributed to the treatment of water. It has been well established that water molecules near the lipid–water interface make a major contribution to the electrostatic potential profile near that interface.<sup>29</sup> Since water molecules are treated as uncharged beads in MARTINI, they do not contribute directly to the electrostatic potential, which is the reason that the calculated interfacial potential at the lipid bilayer–water interface with MARTINI is grossly incorrect

Received: May 1, 2011

Published: September 20, 2011

(MARTINI predicts a value of  $-0.4$  V which may be compared to the experimental estimate<sup>30,31</sup> of  $+0.22$  to  $0.28$  V and results from atomistic simulations<sup>32,33</sup> of  $+0.4$  to  $1.0$  V). A natural remedy, therefore, is to include electrostatics for the description of water, which has been recently pursued by several CG models.

In the CG force field of Essex and co-workers,<sup>34</sup> water is treated using the soft sticky dipole model of Ichiye and co-workers.<sup>35</sup> Although the model has been shown to provide impressive results for the mechanical and electrostatic properties of lipid bilayers, it offers limited computational advantage over atomistic models of water. Moreover, the orientation of lipid head groups in their model is different from atomistic simulations, thus the contributions from different groups to the interfacial potential do not match atomistic results. An improved version of MARTINI, called the polarizable MARTINI force field,<sup>36</sup> has been proposed and features a water model with two charged-sites; recent test calculations suggest that the polarizable MARTINI model leads to much improved results for the insertion PMF of charged pentapeptides into a lipid bilayer.<sup>37</sup> However, this model still gives qualitatively incorrect results for the interfacial dipole potential (reported as  $-2$  V in ref 36), which suggests that the improved PMF is due in part to error cancellation; it has also been discussed in the literature that the insertion PMF is not a simple function of the interfacial dipole potential.<sup>38</sup> More recently, a four-site CG water model (WAT FOUR)<sup>39</sup> has been reported. Although the model was constructed to map 5 water molecules into 1 CG unit, the calculated properties (e.g., density) suggested that 1 CG unit effectively reflects 11 water molecules; moreover, the model overestimates the dielectric constant and underestimates the surface tension, thus hampering its applicability to bilayer systems. Finally, a dipolar CG model that represents five water molecules has been proposed very recently;<sup>40</sup> it reproduces the key properties of water well but has not yet been used to parametrize CG models for other biomolecules.

We have recently developed and reported a new CG model for water, termed as the big multipole water (BMW) model;<sup>41</sup> it features the same four to one mapping as the original MARTINI but includes three explicitly charged sites for each CG unit. Our basic approach starts with using atomistic simulations to characterize the electrostatic properties (multipole moments) and the nonbonded interactions of four water clusters. The results inform us of the appropriate functional forms of electrostatic and nonpolar components of the model; for example, we found it was necessary to describe the nonpolar component with a much softer potential than the commonly used Lennard-Jones (LJ) form. The parameters in the model are then fitted based on comparing experimental and computed properties of water, including bulk density, isothermal compressibility, dielectric permittivity, surface tension, and air–water interface potential. A preliminary combination of the BMW model with MARTINI lipids resulted in a membrane dipole potential that is in good agreement with experimental estimates. Finally, in a recent study,<sup>42</sup> we have shown that the electrostatic features of the water model also appear important to a proper description of the hydrophobic effect. For the association of two hydrophobic peptides, BMW simulations predict the process as entropy driven, in agreement with atomistic studies,<sup>43,44</sup> while several nonelectrostatic/dipolar CG water models (MARTINI,<sup>15</sup> polarizable MARTINI<sup>36</sup> and the model of Shinoda et al.)<sup>45</sup> predict the process as enthalpy driven with very small entropic contributions.

In this work, we report our continuing efforts in developing a CG force field for lipids and amino acids with BMW as the basis.

Since the basic topology and parametrization strategy of the force field follow the MARTINI convention, we refer the model as BMW-MARTINI. With a proper treatment of electrostatics but similar limitations in secondary structure descriptions as the original MARTINI, our model in current form is particularly useful for describing the interaction between lipid membrane and highly charged peptide or protein motifs that are either disordered or rigid (i.e., the model is not capable of describing coupled binding/folding processes).

In the following, we first describe how the model is constructed and how the parameters in the model are determined. Next, we present results that illustrate the performance of the BMW-MARTINI model for a fairly broad range of properties concerning lipid membrane and its interaction with amino acids; these include: (i) thermodynamics (free energy and enthalpy–entropy components) of hydration and partitioning between oil–water for basic bead types; (ii) mechanical properties of lipid bilayers and self-assembly of lipids; (iii) potential of mean force for the penetration of amino acid side chains into a lipid bilayer; (iv) behavior of a highly charged peptide (Arg<sub>8</sub>) on the surface of a lipid bilayer; and (v) electroporation. Finally, we draw a few conclusions and comment on possible directions for future developments.

## II. MODEL AND METHODS

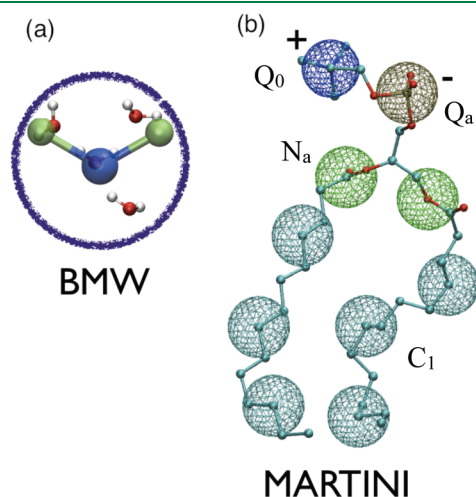
**A. Big Multipole Water (BMW).** Since the BMW water model forms the basis of our CG force field, we briefly summarize its key features. Four water molecules are mapped into one CG unit with three charged sites (see Figure 1a). Since our atomistic simulations<sup>41</sup> indicated that the distributions of dipole moment and quadrupole moment tensor of four water clusters are similar in bulk water, at the air–water interface, and in salt solutions, we chose the geometry and charges of these sites to reproduce the most probable dipole and quadrupole tensors of four water clusters from atomistic simulations. Nonpolar interactions between sites are represented by a modified Born–Mayer–Huggins (BMH)<sup>46,47</sup> potential, which features a softer interaction at short distance than the commonly used LJ potential; the soft-core interaction is crucial for avoiding spurious long-range correlations between water molecules. The BMW model is capable of reproducing key properties of bulk water and air–water interface, most notably bulk permittivity, surface tension, and air–water interfacial potential. The BMW is more computationally intensive (by a factor of 6) than the original MARTINI model due to the larger number of sites and the use of particle mesh Ewald (PME) for electrostatics but is nevertheless about two orders of magnitude more computationally efficient than atomistic simulations.

**B. Parameterization Strategies.** The new force field is parametrized with a strategy that maintains the self-consistency among models for lipids and amino acids. The parametrization is done in a multistage fashion based on carefully monitoring a broad set of properties related to hydration/transfer free energies, lipid bilayer properties, lipid self-assembly, and interaction between amino acids and a lipid bilayer.

First, partition free energies of uncharged bead types between water–air or water–hexadecane are tuned to give the approximate sets of scaling factors for nonpolar interactions. Then, lipid bilayer structural properties and lipid self-assembly phase behaviors are used to fine-tune these scaling factors as well as parameters for charged group (Q)-BMW interactions and the

angle bending force constant for lipid tails. Next, potentials of mean force (PMFs) for the penetration of different amino acids into a dioleoylphosphatidylcholine (DOPC) bilayer are calculated, and the values in the interfacial region are compared with atomistic results<sup>48</sup> to determine interactions between charged (lipid head groups) and neutral (neutral amino acids) groups. Finally, the insertion PMFs of charged amino acid side chains into a bilayer are used to further refine the interaction between CG beads and BMW water. The initial parametrization of the models is described in this section, and a fine-tuning of the parameters is presented in the Results and Discussion section.

**C. Initial Parametrization of Lipid and Amino Acid Potentials.** For lipids and amino acids, we follow the same mapping scheme as the MARTINI force field [illustrated in Figure 1b for a dimyristoylphosphatidylcholine (DMPC) molecule].<sup>15,23</sup> We adopt the same classification scheme for beads: charged (Q), polar (P), nonpolar (N), and apolar (C) for representing similar chemical structures. We also use the same sets of subtypes (with a few changes, see below) to label different levels of effective nonelectrostatic interactions between beads. On the other hand, since the MARTINI force field employs hydration and transfer free energies as the guiding properties for parametrization, changing the underlying water model to BMW requires an extensive reparameterization.



**Figure 1.** Mapping between the chemical structure and the CG model for water and DMPC lipid. Topologies for nonwater components are taken directly from MARTINI.

Several modifications are made to the CG particle (or bead) types in MARTINI. First, water specific types are added. The modified BMH potential is used for BMW-BMW (only between the charge-negative sites), while LJ is used for all other interactions, including those between BMW and other bead types. The antifreezing particle type BP4 is deleted, because with the soft interaction the BMW water does not (unphysically) freeze, and antifreeze particles are therefore not required. Second, a super repulsive (with  $\sigma = 0.62$  nm) interaction between charged (Q) and apolar (C) types is no longer needed, because the interaction between charged beads is now characterized by Coulombic interactions with a small amount of screening (the screening dielectric constant in the BMW model is 1.3 instead of 15–20 in the original MARTINI). Subtypes AC1 and AC2 are therefore deleted in the new force field, and all Q–C interactions are assigned with  $\sigma = 0.47$  nm ( $\sigma = 0.43$  nm between bead types in rings). Finally, new subtypes are added for amino acids to introduce additional flexibility:  $RQ_d$  for guanidinium group in arginine and  $AQ_a$  for aspartate and glutamate. For ions and peptide terminal groups, the  $Q_d$  subgroups are the same, while all  $Q_a$  subgroups are replaced by  $AQ_a$  (e.g., for  $Cl^-$ ).

For the nonelectrostatic interaction levels among bead types, those between groups of P, N, and C are mostly inherited from MARTINI, except for the following:  $P_5$  and  $P_4$  interact with  $N_{da}$ ,  $N_a$ , and  $N_d$  with  $\epsilon = 5.6$  kJ/mol; all ring beads (label starting with S) interact with C1 with a scaling factor of 90% for the original  $\epsilon$  in MARTINI. By contrast, bead–water interactions have to be modified to reproduce relevant hydration and transfer free energies. Similar to the polarizable MARTINI model,<sup>36</sup> scaling factors for the well-depths,  $\epsilon$ , are introduced to reduce the strength of interaction (relative to the original MARTINI) between uncharged bead types and water (BMW); the factor is 71% for levels with  $\epsilon < 4.5$  kJ/mol in the original MARTINI and 75% otherwise. Further more, nonelectrostatic interactions between charged groups (type Q) and all other beads are modified, since electrostatic interactions are treated differently in the new model. As shown in Table 1, besides Q–Q and Q–BMW interactions, levels for Q to other uncharged groups, especially apolar types (C), are also tuned to ensure reasonable partitioning free energies. This is required because hydration free energies of charged groups are altered upon using the BMW model for water. Meanwhile, levels for Q–P and P–N interactions remain very similar to the original MARTINI because the nonpolar interactions between these bead types still implicitly represent both van der Waals and dipolar/hydrogen-bonding contributions.

The bonded parameters (e.g., bond, angle, and torsional angle force constants) for the new CG model are largely the same as the

**Table 1.** Levels of Nonpolar Interactions among Charged Groups, BMW Water, and Uncharged Groups in the BMW-MARTINI Model<sup>a</sup>

	BMW	Q						P					N				C				
		$Q_{da}$	$Q_d$	$RQ_d$	$Q_a$	$Q_{da}$	$Q_0$	$P_5$	$P_4$	$P_3$	$P_2$	$P_1$	$N_{da}$	$N_d$	$N_a$	$N_0$	$C_5$	$C_4$	$C_3$	$C_2$	$C_1$
$Q_{da}$	I	O	O	O	O	O	II	O	O	O	I	I	O	O	O	IV	III	IV	IV	IV	IV
$Q_d$	I	O	I	I	O	O	II	O	O	O	I	I	O	III	O	IV	III	IV	IV	IV	IV
$RQ_d$	IV	O	I	I	O	O	II	O	O	O	I	I	O	III	O	IV	III	IV	IV	IV	IV
$Q_a$	I	O	O	O	I	I	II	O	O	O	I	I	O	O	III	IV	III	IV	IV	IV	IV
$AQ_a$	I	O	O	O	I	I	II	O	O	O	I	I	O	O	III	IV	III	I	I	I	I
$Q_0$	I	II	II	II	II	II	IV	I	O	I	II	III	III	III	III	IV	III	IV	IV	IV	IV

<sup>a</sup> Level of interaction indicates the well depth in the LJ potential: O,  $\epsilon = 5.6$  kJ/mol; I,  $\epsilon = 5.0$  kJ/mol; II,  $\epsilon = 4.5$  kJ/mol; III,  $\epsilon = 4.0$  kJ/mol; IV,  $\epsilon = 3.5$  kJ/mol; V,  $\epsilon = 3.1$  kJ/mol; VI,  $\epsilon = 2.7$  kJ/mol; VII,  $\epsilon = 2.3$  kJ/mol; and VIII,  $\epsilon = 2.0$  kJ/mol. The LJ parameter  $\sigma = 0.47$  nm ( $\sigma = 0.43$  nm for rings) is used for all interaction levels. The same grouping criteria (including subgroups) are applied as in the original MARTINI scheme.<sup>15</sup>

original MARTINI. Only the force constant for angle bending in hydrocarbons (lipid tails) is modified from 25 to 10 kJ/(mol·degree<sup>2</sup>), because it is reported that the angle distribution in MARTINI is narrower than the atomistic counterpart.<sup>49</sup> Also, the scheme of restraining secondary structure elements in MARTINI is adopted, and relieving such restraints<sup>50</sup> will be an interesting direction for further developments.

**D. Simulation Protocols.** The simulation protocols used for the new CG model are largely the same as for BMW water simulations.<sup>41</sup> A time step of 20 fs is used with GROMACS 4.0.5.<sup>51</sup> Temperature and pressure are kept constant by using the Berendsen scheme,<sup>52</sup> with coupling times of  $\tau_T = 1$  ps and  $\tau_P = 5$  ps. The SETTLE algorithm<sup>53</sup> is used to constrain “bonds” in CG water, and LINCS<sup>54</sup> is used for bonds in ring structures in several amino acids. PME with a spacing of 0.2 nm and  $\epsilon_r = 1.3$  are applied for electrostatics. Similar to the original MARTINI, LJ interactions are excluded between bonded beads but not for second nearest neighbors. The shift cutoff scheme ( $r_{\text{shift}} = 0.9$  nm and  $r_{\text{cut}} = 1.2$  nm) is applied to all LJ interactions, while the switch scheme ( $r_{\text{switch}} = 1.2$  nm and  $r_{\text{cut}} = 1.4$  nm) is used for water–water BMH interactions (see Supporting Information for details). With these protocols, simulations with the new force field are slower than the original MARTINI but about 2–3 orders of magnitude faster than atomistic simulations. As discussed in MARTINI applications,<sup>15</sup> since the lipid lateral diffusion rate is about 4 times larger than experimental measurement, time scales in all simulation are interpreted as 4 times the actual simulation lengths.

**E. Properties Calculated.** The free energies of hydration and partitioning between water and octanol for a number of n-alkanes are calculated using thermodynamic integration (TI),<sup>55,56</sup> which is carried out at 300 K with 21 evenly spaced  $\lambda$  windows, where  $\lambda$  is the coupling parameter in the TI method, each sampled for 160 ns; the target bead is decoupled from its surrounding solvents with a soft core potential.

The area per lipid of lipid bilayers is calculated from  $NP_{xy}P_zT$  simulations, where  $N$  is the number of molecules,  $P_{xy}$  and  $P_z$  are the transverse and normal components of the pressure tensor, and  $T$  is the temperature. Patches of 512 lipids are simulated with a hydration level of  $\sim 60$  water molecules per lipid; the results are averaged over 240 ns production run after equilibration.

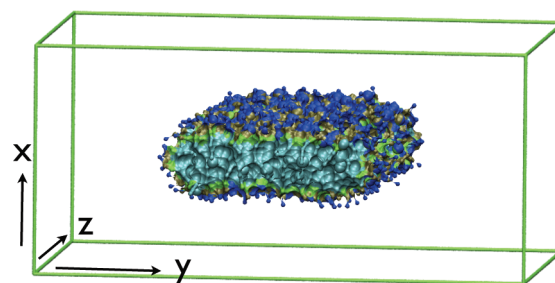
The area compressibility modulus,  $K_A$ , is calculated from the relation between membrane area per lipid  $A$ , tension-free equilibrium area per lipid  $A_0$  and surface tension  $\gamma$  (eq 1):<sup>57</sup>

$$K_A = 2A_0 \left( \frac{\partial \gamma}{\partial A} \right)_T \quad (1)$$

$NA_{xy}P_zT$  simulations are performed on dipalmitoylphosphatidylcholine (DPPC) lipids with five different restrained membrane areas (evenly from 60 to 68 Å<sup>2</sup>/lipid in the  $xy$  plane), and the corresponding surface tension per leaflet  $\gamma$  is calculated. Test calculations indicate only a small finite size effect on the calculated  $K_A$ , with the difference between a bilayer patch of 128 and 512 lipids being within statistical uncertainties.

The line tension is computed by constructing a ribbon structure of bilayers, continued in the  $z$  direction as shown in Figure 2. This structure is simulated in the  $NP_{xy}L_zT$  (325 K) ensemble, and line tension is calculated from the edge along  $z$  from eq 2:<sup>58</sup>

$$\Lambda = \frac{1}{2} \left\langle L_x L_y \left[ \frac{P_{xx} + P_{yy}}{2} - P_{zz} \right] \right\rangle \quad (2)$$



**Figure 2.** Ribbon structure of 512 DPPC lipids and the simulation box (shown with VMD)<sup>59</sup> in the line tension calculation. For clarity, all waters (16 096 BMW water) are omitted, and the ribbon is intercepted in the  $x$ – $y$  plane as shown, with tails in cyan, glycerol in green, and head groups in tan and blue.

### III. RESULTS AND DISCUSSION

**A. Free Energies of Hydration and Partitioning.** Calibration of free energies of hydration and partitioning between water and hexadecane plays an important role in the development of GROMOS and MARTINI force fields. We also use these quantities to calibrate LJ interaction parameters for uncharged bead types in the BMW-MARTINI model. For charged groups, due to the lack of relevant experimental data, they are parametrized based on lipid properties and amino acid PMFs, as described in later sections.

Due to the CG nature of the model, it is only meaningful to compare calculated free energies with experimental values for a range of similar compounds, as was done with the original MARTINI model. As shown from Table 2, results with the new CG model are overall in good agreement with experimental and original MARTINI results. Thus simply using two scaling factors for the LJ well-depth between uncharged groups and water is sufficient for the current purpose. Due to the soft-core nature of the BMW model, which allows smaller energy variation upon solute insertion, the scaling factors are smaller than the one (95%) introduced in the polarizable MARTINI model.<sup>36</sup> Similar to the original MARTINI model,<sup>15</sup> our model systematically underestimates the hydration free energies as compared to experimental values. The magnitude of the underestimation is larger for more polar beads (i.e., P groups), indicating that improving the description of polar groups (e.g., by including explicit dipoles) can be a future direction for development. Since our current model mainly focuses on the partitioning of beads between polar and apolar environments, as in the original MARTINI, interactions between P and C groups are underestimated to compensate for the underestimated hydration of P groups. For charged particles, because electrostatics are explicitly included in the new model, the solvation free energies are comparable to atomistic results (also see Supporting Information for discussions of ions).

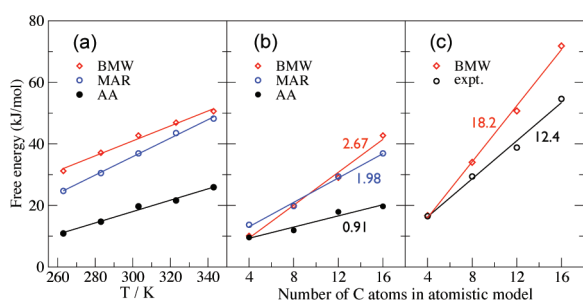
**B. Hydration Thermodynamics.** In addition to comparing hydration free energies, it is of interest to compare the enthalpic–entropic components of hydration at the CG and atomistic levels to ensure that the CG model captures the proper physics of solvation, especially that of hydrophobic groups. For this purpose, we study the hydration of N-hexadecane (four bonded C1 beads in the CG model), which forms the tail of DPPC. Solvation free energies computed at different temperatures (263–343 K) are decomposed into solute–solvent interaction energy  $U_{uv}$  and



**Table 2.** Free Energy of Hydration and Partition between Water and Hexadecane for Neutral Groups Is Compared to Experimental Values for Target Compounds and MARTINI<sup>a</sup>

scaling factor	type	hydration			partition		
		expt.	MARTINI	BMW	expt.	MARTINI	BMW
75%	P <sub>5</sub>	−40	−25	−21	−27	−28	−27
	P <sub>4</sub>	−27~−35	−18	−15	−21~−25	−23	−22
	P <sub>3</sub>	−29	−18	−15	−19	−21	−19
	P <sub>2</sub>	−21	−14	−11	−13	−17	−14
	P <sub>1</sub>	−20~−21	−14	−11	−9~−10	−11	−10
71%	N <sub>d</sub> /N <sub>a</sub> /N <sub>da</sub>	−12~−20	−9	−5	−4~−6	−7	−4
	N <sub>0</sub>	−8	−2	−1	−1	−2	0
	C <sub>5</sub>	−6	1	2	7	5	7
	C <sub>4</sub>	−4~−2	5	5	7~11	9	10
	C <sub>3</sub>	−1~−2	5	5	12	13	14
	C <sub>2</sub>	8	10	8		16	17
	C <sub>1</sub>	9~10	14	10	18	18	19

<sup>a</sup> The estimated experimental values are taken directly from ref 15 and compound names are not listed here. The scale in the first column is the scaling factor for MARTINI potential levels between the groups and water. All units are in kJ/mol, and the new model is labeled as BMW.



**Figure 3.** Hydration free energy in water for (a) n-hexadecane as a function of temperature and (b) n-alkanes at 303 K. Panel (c) shows the partition free energy for n-alkanes between water and octanol at 303 K. Data for both atomistic (AA) and MARTINI (MAR)<sup>14</sup> are taken directly from ref 61, while the partitioning data are compared to experimental values (expt.).<sup>62</sup> The corresponding linear regression curves and their slopes are also displayed in the same color.

entropy  $S_{uv}$ ,<sup>60</sup> and the results are compared to both atomistic and MARTINI results.<sup>61</sup>

In all cases, the hydration free energies predicted by our new CG model are similar to the original MARTINI model, and systematically overestimate the atomistic value for the free energy at all temperatures, as can be seen in Figure 3. Decomposition of the solvation free energy into enthalpic and entropic terms (not shown) indicates that both models overestimate both the solute–solvent interaction energy and the entropy by significant amounts (of the order of 50 kJ/mol and 100–200 kJ/(mol·K), respectively), i.e., this does not simply originate from an overestimation of the water–oil repulsion as previously reported.<sup>61</sup> As the chain length increases (Figure 3b), the enthalpic contribution in the CG models decreases more slowly compared to atomistic results, while contributions from entropy ( $-TS$ ) change at a similar pace in the three models. In other words, with a CG model, change in the solute–solvent enthalpy as the chain elongates is not strong enough to compensate for the change in solute–solvent entropy.

For the partition free energy of solutes between water and octanol (see Figure 3c), the discrepancy between CG and atomistic

**Table 3.** Area Per Lipid for Common Saturated and Unsaturated PC, PE, and PS Phospholipids in the CG Model<sup>a</sup>

systems	expt.	MARTINI	BMW
DPPC (325 K)	63 <sup>63</sup>	64	64
DPPC (338 K)	64–67 <sup>65,66</sup>	66	65
DOPC (300 K)	67 <sup>64</sup>	67	64
DOPE (273 K)	65 <sup>67</sup>	61	60
DOPS (303 K)	65 <sup>68</sup>	67	62

<sup>a</sup> All units are in  $\text{\AA}^2$ . Typical uncertainties in simulation results and experiments are 1 and 2  $\text{\AA}^2$ , respectively. For MARTINI, all values are shown as reported in ref<sup>14</sup> except for DOPS, which is calculated from this work.

results is less sensitive to the chain length, due likely to error cancellation. This is encouraging because partition free energies are more important for the development of both original and our improved MARTINI model.

**C. Lipid Properties.** 1. *Bilayer Structural, Elastic, and Dynamic Properties.* The BMW-MARTINI model values for the area per lipid in bilayers are similar (somewhat lower) to that obtained from experiment or from the MARTINI model. Table 3 compares the area per lipid obtained from various models for several common saturated, unsaturated, and charged lipid bilayers. Given the uncertainty in experimental measurements, the results from the new model are satisfying, especially for charged lipids, such as DOPS. Similar agreement is found for the thickness of lipid bilayers. For DPPC bilayers (at 325 K), the experimentally measured thickness is 3.8 nm,<sup>63</sup> MARTINI gives 4.0 nm, and BMW-MARTINI gives 3.9 nm. For DOPC bilayers (at 300 K), both the original MARTINI (4.5 nm) and BMW-MARTINI (4.6 nm) values are larger than the experimental result (3.7 nm).<sup>64</sup> Considering the CG nature of the lipid tails (one bead representing  $\sim 4$  CH<sub>2</sub> groups), the agreement can be considered satisfactory. Further more, the density profiles of the CG bilayer are in good agreement with experiment (see Supporting Information).

The BMW-MARTINI results for the area compressibility modulus,  $K_A$ , and the line tension are significantly higher than the experimental values or those obtained from the MARTINI

model. For DPPC bilayers, the BMW-MARTINI prediction for  $K_A$  is  $585 \pm 16$  dyn/cm, which is higher than the experimental value<sup>69</sup> of 234 dyn/cm and the MARTINI result of 292 dyn/cm (also calculated using eq 1). The model predicts a line tension of  $118 \pm 11$  pN, which is higher than the experimental estimate of 10–30 pN for similar lipids<sup>70–72</sup> and the MARTINI value of 64 pN.<sup>73</sup>

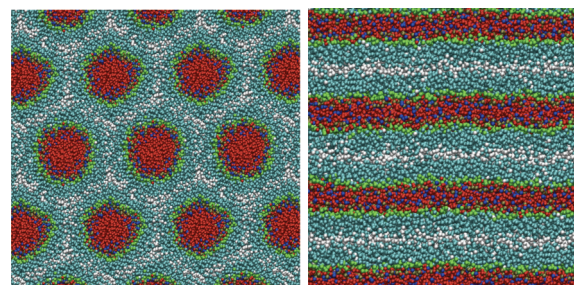
The high line tension in the BMW-MARTINI model will likely result in the model significantly overestimating the barrier to pore formation in lipid bilayers. The energy  $E$  required to form a pore inside the bilayers can be approximated by<sup>74</sup>  $E(r) = 2\pi\Lambda - \pi r^2\gamma$ , where  $r$  is the pore radius and  $\gamma$  the tension on membrane. At the critical tension  $\gamma^*$ , the edge energy from  $\Lambda$  is overcome, and a pore is stable at  $r = \gamma^*/\Lambda$ . Therefore, overestimating  $\Lambda$  means that the critical tension required to form a pore is likely to be overestimated. The overestimated line tension is likely related to the CG nature of the solvent, where the partitioning of individual water molecules at the bilayer edge is not captured. By grouping four waters together into a single site, the energy penalty for bringing a CG water into the hydrophobic region becomes larger. This energy penalty is large in the BMW-MARTINI model, compared to the MARTINI model, because of the presence of electrostatic interactions.

The lipid (DPPC) lateral diffusion constant is  $0.8 \pm 0.1 \times 10^{-7}$  cm<sup>2</sup>/s (using the effective scaling factor of 4 due to coarse-graining),<sup>15</sup> which is similar to the experimental value of  $1.0 \times 10^{-7}$  cm<sup>2</sup>/s.<sup>75</sup>

**2. Lipids Self-Assembly.** The microphase morphology of lipids is an important benchmark for interaction parameters of charged groups. Experimental studies indicate that DOPE lipids assemble into an inverted hexagonal phase at temperatures above 280–300 K with 1:16 lipid/water hydration level,<sup>76</sup> while DOPC lipids assemble into the lamellar phase under similar conditions.<sup>77</sup> This phase behavior reflects the spontaneous curvature of the lipid layers, which are negative for DOPE and about zero for DOPC.<sup>78,79</sup> The interactions between lipid head groups (Q–Q) and between the head groups and water (Q–BMW) play a key role, since they modulate the effective shape of the lipid molecules (e.g., cone vs cylinder).

The BMW-MARTINI model gives the correct phases for DOPE and DOPC, which suggests that the spontaneous curvatures for bilayers of both lipids are properly captured by the current model. To study the phase behavior of lipids, a system of 1000 lipids mixed randomly with 4000 CG solvents (corresponding to 16 water molecules per lipid) is simulated with completely anisotropic pressure coupling at 318 K. Because the lamellar phase has a faster water exchange rate, it is prepared as the initial configuration: A lamellar phase is first assembled by artificially setting the lipid head groups to interact with water with the highest interaction strength in the MARTINI force field (level O), then the proper level of interaction is used for subsequent simulations. For DOPE, “stalks” of lipids are gradually formed between the lamellar layers, and the inverted hexagonal phase is formed within  $\sim 5 \mu\text{s}$  (shown in Figure 4). The resulting hexagonal spacing (distance between the central axes of water channels) is 6.8 nm, which agrees well with the estimate from SAXS data as 7.1 nm.<sup>80</sup> For DOPC, no stalk formation is observed, and the lamellar phase remains stable for the subsequent 5  $\mu\text{s}$  simulations, and the resulting lamellar repeat spacing is 5.8 nm; the experimental value for fully hydrated lipids is 6.3 nm.<sup>63</sup>

**3. Membrane Interface Electrostatic Properties.** The dipole potential at the membrane-water interface is significantly improved



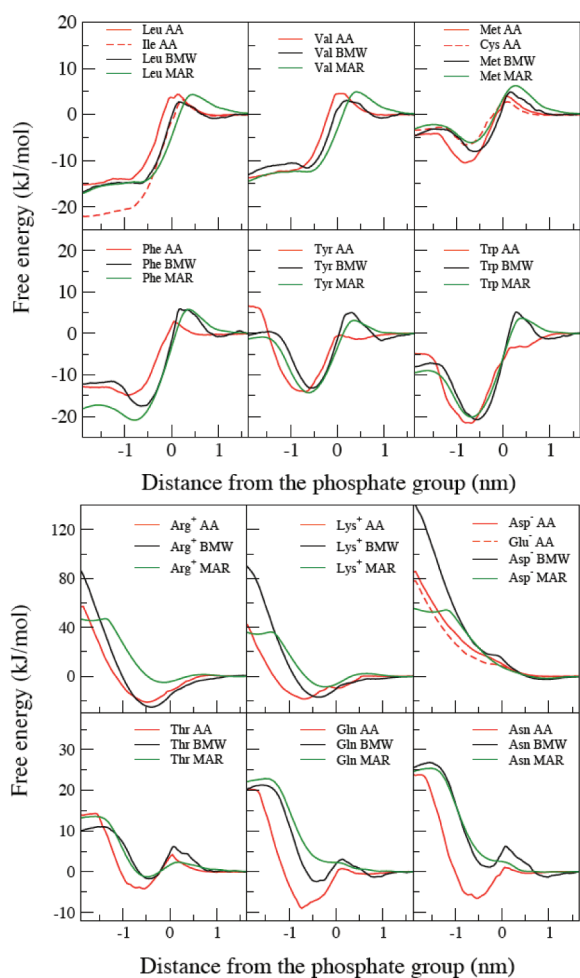
**Figure 4.** BMW-MARTINI results for the inverted hexagonal phase of DOPE (left) and lamellar phase of DOPC (right); the beads are color coded: water in red, lipid tails in cyan, glycerol in green, lipid head groups in blue and tan, and lipid tail terminals (last bead) in white.

with the BMW model over previous CG models, such as MARTINI.<sup>41</sup> With refined parameters in this work, this feature is maintained. The calculated value of the interfacial potential is +0.23 V (the value was +0.30 V with a preliminary combination<sup>41</sup> of BMW and MARTINI lipid models), in good agreement with experimental estimates of 0.22–0.28 V for DPPC bilayers.<sup>30</sup>

**D. Interaction between Amino Acids and a Lipid Bilayer.** *1. Partition of Amino Acid Side Chains.* The PMF is calculated for each amino acid side chain as a function of the distance from the center of a DOPC lipid bilayer. Results are compared to both atomistic<sup>48</sup> and MARTINI<sup>23</sup> models to further fine-tune interaction levels in our CG model, especially for charged groups. For each simulation, two side chain analogues are placed at a distance of 4.5 nm from each other, one in the center of bilayer and the other in the bulk, in a system of 96 DOPC lipids and 1300 BMW waters; the PMF is then calculated with the standard umbrella sampling protocol, with 46 windows and 80 ns for each window, and force constants of  $\sim 1000$  kJ/(mol·nm<sup>2</sup>). The PMF is averaged over the symmetric halves of the bilayer. For charged side chains, two ions ( $\text{Na}^+$  or  $\text{Cl}^-$ ) are added to maintain charge neutrality.

Some of the parameters are fine-tuned after a comparison with atomistic results, so this is strictly not a test of the model. New bead types are introduced:  $\text{RQ}_d$  for the arginine guanidinium group and  $\text{AQ}_d$  for aspartate and glutamate. The other side chains are represented with the same topology as MARTINI. For neutral amino acids the parameters are not changed beyond what was used for the partition free energy between water and hexadecane (Table 2). For the aromatic side chains, a scaling factor of 90% for all the ring groups (start with S) to C1 is applied. New parameters are fit for the charged side groups. For highly polar but neutral side chains, Gln ( $\text{P}_4$ ) and Asn ( $\text{P}_5$ ), their interactions to  $\text{N}_{da}$ ,  $\text{N}_d$ , and  $\text{N}_a$  are changed as  $\epsilon = 5.6$  kJ/mol.

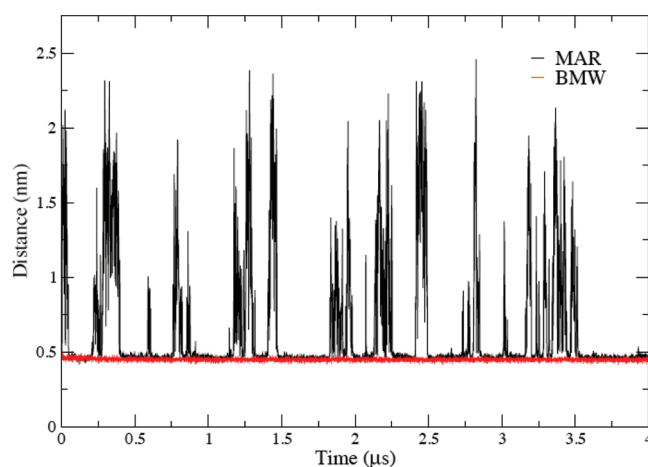
The new force field is more accurate for the side chain PMF than the MARTINI force field, when compared to results from atomistic simulations. In Figure 5, representative PMF profiles for hydrophobic, polar, aromatic, and charged amino acid side chain analogues are shown. For neutral amino acids, the PMFs have the proper values in both solution and membrane tail regions. PMF in the interfacial region is determined by interactions between charged groups (lipid head groups) and the side chain. Take leucine for example, modeled as  $\text{C}_1$ , the scaling factor for the nonpolar interactions determines its insertion PMF between water and the bilayer center, and the barrier at the interface is defined by Q– $\text{C}_1$  interactions. Both the original MARTINI and the BMW-MARTINI are accurate for these PMFs. For the polar amino acids the performance of the new force field is slightly superior to the MARTINI force field, but



**Figure 5.** Insertion PMF for amino acids into a DOPC bilayer (negative distances indicate the interior of the bilayer). Note that although the position of phosphate groups is strictly defined with electron density maximum for each model, the position of the amino acid relative to the phosphate ( $x$ -axis) in the MARTINI and BMW-MARTINI models is subject to some uncertainty (up to 0.2 nm) due to the CG nature of these models. Both atomistic (AA)<sup>48</sup> and MARTINI (MAR)<sup>23</sup> data are obtained from Tieleman and co-workers (private communication).

neither force field is in quantitative agreement with the atomistic simulation results.

For the charged amino acids, their PMFs are closely coupled to membrane–water defects induced by the penetration of the side chains. This makes the parametrization more complicated than for neutral groups because all Q–Q, Q–BMW, Q–neutral interactions should, in principle, be considered. However, since our model gives the correct electrostatic profile at the interface, a deep minimum is always observed for the cationic side chains although the depth of the minimum depends on the parametrization. Compared with the original MARTINI (Arg<sup>+</sup>:  $-6$  kJ/mol and Lys<sup>+</sup>:  $-9$  kJ/mol), our new model gives a substantially deeper minimum:  $-22$  kJ/mol for arginine and  $-20$  kJ/mol for lysine; the latter values are close to the atomistic results. In the center of the bilayer, the calculated PMF with BMW-MARTINI is too high for all charged amino acids compared to atomistic results. However, even different atomistic force fields give rather different values in this region with an uncertainty up to 25 kJ/mol.<sup>81</sup> Therefore, we do not consider the discrepancy in this region as a

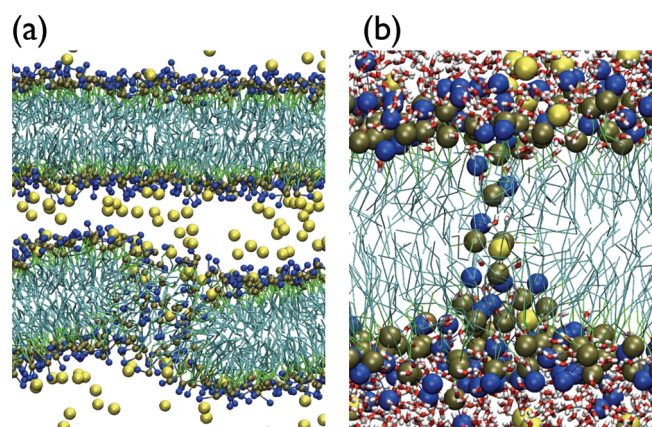


**Figure 6.** The minimum distance between Poly-Arg<sub>8</sub> and a DOPC bilayer as a function of time from MARTINI (MAR) and BMW-MARTINI (BMW) simulations. The van der Waals diameter of the relevant bead types is 0.47 nm, thus the peptide is considered attached to the surface if the minimum distance is below 0.47 nm. Somewhat similar differences are observed also for Poly-Arg<sub>8</sub> at the surface of an anionic membrane (70% DOPC and 30% DOPS).

significant limitation of the model, although this issue deserves further exploration especially in the context of studying pore formation, where the line tension of the bilayers will play an important part.

**2. Poly-Arg<sub>8</sub> Attachment on Membrane.** We use the new force field to investigate the attachment of peptides to a membrane surface. Simulations of a poly-Arg<sub>8</sub> cationic peptide (analogue to HIV-TAT peptide)<sup>82</sup> on the surface of a 128 DOPC lipid bilayer are performed with both the new model and original MARTINI; the simulation cells are charge neutral with counterions (0.15 M NaCl) (see details for ion in Supporting Information). According to recent experimental observations,<sup>83</sup> the peptide attaches to the membrane surface with a long residence time of up to a second; atomistic simulations<sup>82,84</sup> up to 400 ns also showed stable attachment. With the original MARTINI, the peptide does not absorb in a stable fashion and desorbs frequently with the longest residence time on the membrane surface of  $\sim 350$  ns (Figure 6a).

The BMW-MARTINI model predicts a stable attachment of poly-Arg<sub>8</sub> to a DOPC bilayer. Figure 6 depicts the minimum distance between the peptide and the bilayer as a function of time and shows that with the BMW-MARTINI model the peptide stays attached up to tens of microseconds (throughout the entire simulation) unlike the original MARTINI model. On the surface of an anionic bilayer (70% DOPC and 30% DOPS), the new model gives very stable attachment of poly-Arg<sub>8</sub>. With MARTINI, however, the peptides are still observed to desorb from the membrane surface although the residence time ( $\sim \mu$ s) is longer than that at the surface of a zwitterionic (DOPC) bilayer. This difference is anticipated based on the difference in the insertion PMF for arginine side chains at the membrane surface (see Figure 5 on DOPC), for which the MARTINI model predicts a much shallower minimum than the BMW-MARTINI model. This qualitative difference between BMW-MARTINI and MARTINI highlights the importance of properly describing electrostatics for the analysis of highly charged peptides or protein motifs near membrane surface.



**Figure 7.** Snapshots from simulations of electroporation (with an initial charge imbalance of  $26 e^-$ ). (a) The initial structure of the pore, formed within 1 ns (the water molecules are omitted for clarity). (b) The final structure of the (toroidal) water defect. Color code: yellow:  $\text{Na}^+$  or  $\text{Cl}^-$ ; blue: choline; and tan: phosphate.

In a recent study<sup>37</sup> of insertion PMF for pentapeptides that feature both charged and hydrophobic groups, Singh and Tieleman found that the polarizable MARTINI model leads to a substantial improvement over the original MARTINI, and the results were in good agreement with atomistic simulations. This is interesting because the interfacial potential calculated by polarizable MARTINI still has the incorrect sign.<sup>36</sup> On the other hand, as analyzed by Allen and co-workers,<sup>38</sup> the membrane permeation energetics for charged groups are not a simple function of the interfacial potential because the dipole potential is not fully sensed at the locally deformed bilayer interface. Moreover, the presence of both hydrophobic and charged groups in the pentapeptides may have further helped attenuate the errors associated with the charged group. In the near future, as the amino acids parameters for the polarizable MARTINI model<sup>36</sup> become available, it is valuable to compare the BMW-MARTINI and polarizable MARTINI models in a systematic fashion.

Simulations with multiple copies of poly-Arg<sub>8</sub> (in the presence of  $\text{Na}^+/\text{Cl}^-$  counterions) on a DOPC bilayer show no membrane penetration during up to  $20 \mu\text{s}$ , in agreement with more recent atomistic simulations.<sup>84</sup> With 8 peptides, the area per lipid increases slightly from the peptide-free value of  $0.64\text{--}0.66 \text{ nm}^2$ . Order parameter calculations indicate that the head groups become more ordered and perpendicular to the membrane normal in the presence of multiple cationic peptides; at the same time, water clusters near the membrane surface (below phosphate groups) become slightly less ordered.

Finally, interesting phase behaviors have been observed experimentally for the mixture of cationic peptides (e.g., poly-Arg vs poly-Lys) and lipids.<sup>85</sup> Simulations with a reliable CG model are expected to be very effective at complementing experiments to better understand the connections between peptide sequence and phase behavior.<sup>86</sup> Such studies are in progress and will be reported separately.

**3. Electroporation.** When two bilayers are constructed with an imbalance of ions on the two sides, the chemical potential difference and local electric field drive the transfer of ions through the bilayer via the formation of a water pore. The process is referred to as electroporation and has been used to test both atomistic<sup>87,88</sup> and CG<sup>36</sup> lipid models. Here we investigate the same systems studied previously by the polarizable MARTINI model.<sup>36</sup> A typical

system consists of two DPPC bilayers, which contain 512 lipid molecules, 5632 CG BMW water molecules, 52 evenly distributed (over the two water compartments)  $\text{Cl}^-$  ions and 52  $\text{Na}^+$  in one water compartment; this charge imbalance of  $26 e^-$  results in an electric field of  $0.7 \text{ V/nm}$ . Simulations are performed in the  $NP_{xy}P_zT$  ensemble, at 325 K, and for  $6 \mu\text{s}$ .

With this charge imbalance, the simulations show the opening of a water pore, transportation of both  $\text{Na}^+$  and  $\text{Cl}^-$  ions, and then closing of the pore. Typically, a water pore is formed within 1 ns, and ions diffuse from one water compartment to the other through the pore (shown in Figure 7a);  $\text{Na}^+$  and  $\text{Cl}^-$  ions translocate in opposite directions during the process. The water pore grows in size until approximately 20–40 ns, when it reaches its maximum diameter ( $\sim 4 \text{ nm}$ ) and then starts to shrink in size. After approximately 100 ns, only a water defect remains with about 5 CG water in the membrane interior. With the small charge imbalances at this stage, the ions still translocate through a toroidal pore (shown in Figure 7b), in agreement with previous atomistic simulations.<sup>88</sup> By contrast, the defect observed at this stage with the polarizable MARTINI model does not involve significantly displaced lipid headgroups and thus closer to the barrel stave model.<sup>36</sup> Eventually, after approximately 400 ns, the water defect completely seals, and the two compartments have different concentrations of ions (containing about  $36 \text{ Na}^+/34 \text{ Cl}^-$  and  $16 \text{ Na}^+/18 \text{ Cl}^-$ , respectively) with a negligible charge imbalance of  $1\text{--}2 e^-$ .

Spontaneous electroporation also occurs for smaller ( $20 e^-$ ) and larger ( $52 e^-$ ) charge imbalances. For large charge imbalances, multiple pores located on different bilayers are observed. With a smaller charge imbalance ( $20 e^-$ ), pore formation takes a longer time, which could be several  $\mu\text{s}$ . The total number of  $\text{Na}^+$  and  $\text{Cl}^-$  transferred to the opposite water compartment is similar if the initial charge imbalance is small, indicating little membrane selectivity toward anions ( $\text{Cl}^-$ ) over cations ( $\text{Na}^+$ ), in agreement with atomistic studies.<sup>88</sup>

## IV. CONCLUSIONS

We report a new CG force field, called BMW-MARTINI, for simulations of lipids and peptides in water. The model follows the same strategy as the original MARTINI force field but is based on the BMW water model, which includes electrostatic interactions. The interactions between almost all the CG sites are reparameterized.

The new force field provides a reasonably accurate description of the hydration of the CG sites, the transfer free energy of sites between hexadecane and water, lipid phase behavior, membrane electrostatic properties, and insertion potential of mean force for amino acids into a bilayer. For most of these properties, the BMW-MARTINI model gives similar results as the original MARTINI.

For membrane electrostatic properties and the insertion potential of mean force of charged amino acids into a bilayer, the new model is superior to both the original and the polarizable MARTINI models and predicts a much deeper free energy minimum at the membrane–water interface. As a consequence, the new model predicts a stable attachment of cationic peptides to both zwitterionic and negatively charged membrane surfaces, as observed in experiment and atomistic simulation, while frequent desorptions are observed with the original MARTINI force field. The model also predicts electroporation when there is a charge imbalance across the lipid bilayer, in contrast to the original MARTINI. The pore formed during electroporation is

toroidal in nature, similar to the prediction of atomistic simulations but distinct from results of polarizable MARTINI for small charge imbalances.

In terms of efficiency, the new model is more computationally intensive than the nonelectrostatic MARTINI by a factor of less than 6, but still more than 2 orders of magnitude more efficient than atomistic models. Therefore, the BMW-MARTINI model is a useful alternative to the (polarizable) MARTINI model, and future studies are required to systematically compare the merits and limitations of these models in realistic applications.

In the current form, the BMW-MARTINI model is not as accurate for the mechanical properties of the membrane, with values for the area compressibility modulus and line tension being significantly higher than experiment or the original MARTINI model. We attribute this to the mapping of four water molecules to one CG site which makes the process of transferring a water from the aqueous phase to hydrophobic region unrealistic. The electrostatic interactions in the BMW water model exacerbate this problem with the coarse-graining procedure. This overestimation of mechanical properties is likely to result in a high barrier to pore formation. A possible solution to this problem is to use an adaptive resolution scheme for water,<sup>89,90</sup> where CG sites can be transformed into atomistic water molecules in the interface and pore region. Alternatively, soft interactions, which are limited to water–water interactions in the current model, can be introduced between water and other components of the system.

## ■ ASSOCIATED CONTENT

**S Supporting Information.** Additional benchmark results and discussions are included. This material is available free of charge via the Internet at <http://pubs.acs.org>.

## ■ AUTHOR INFORMATION

### Corresponding Author

\*E-mail: [cui@chem.wisc.edu](mailto:cui@chem.wisc.edu); [yethiraj@chem.wisc.edu](mailto:yethiraj@chem.wisc.edu).

## ■ ACKNOWLEDGMENT

We thank Prof. P. Tieleman and Dr. L. Monticelli for sharing the insertion PMF data for amino acids into a lipid bilayer, for sending us the preprint of ref 37, and for a critical reading of the manuscript. The research has been supported by the National Science Foundation (CHE-0957285 to QC and CHE-0717569 and CHE-1111835 to A.Y.). Computational resources from the National Center for Supercomputing Applications at the University of Illinois and the Centre for High Throughput Computing (CHTC) at UW-Madison are greatly appreciated. Computations are also supported in part by National Science Foundation through a major instrumentation grant (CHE-0840494).

## ■ REFERENCES

- (1) Phillips, R.; Ursell, T.; Wiggins, P.; Sens, P. *Nature* **2009**, *459*, 379–385.
- (2) Andersen, O. S.; Koeppe, R. E., II *Annu. Rev. Biophys. Biomol. Struct.* **2007**, *36*, 107–130.
- (3) Huang, H. W. *Biochem.* **2000**, *39*, 8347–8352.
- (4) Chapman, E. R. *Annu. Rev. Biochem.* **2008**, *77*, 615–641.
- (5) Doherty, G. J.; McMahon, H. T. *Annu. Rev. Biochem.* **2008**, *37*, 65–95.
- (6) Tang, Y.; Cao, G.; Chen, X.; Yoo, J.; Yethiraj, A.; Cui, Q. *Biophys. J.* **2006**, *91*, 1248–1263.
- (7) Chen, X.; Cui, Q.; Tang, Y. Y.; Yoo, J.; Yethiraj, A. *Biophys. J.* **2008**, *95*, 563–580.
- (8) Ma, L.; Yethiraj, A.; Chen, X.; Cui, Q. *Biophys. J.* **2009**, *96*, 3543–3554.
- (9) Tozzini, V. *Curr. Opin. Struct. Biol.* **2005**, *15*, 144–150.
- (10) Clementi, C. *Curr. Opin. Struct. Biol.* **2008**, *18*, 10–15.
- (11) Ayton, G. S.; Noid, W. G.; Voth, G. A. *Curr. Opin. Struct. Biol.* **2007**, *17*, 192–198.
- (12) Marrink, S. J.; de Vries, A. H.; Tieleman, D. P. *Biochim. Biophys. Acta, Biomembr.* **2009**, *1788*, 149–168.
- (13) Shinoda, W.; DeVane, R.; Klein, M. L. *J. Phys. Chem. B* **2010**, *114*, 6836–6849.
- (14) Marrink, S. J.; de Vries, A. H.; Mark, A. E. *J. Phys. Chem. B* **2004**, *108*, 750–760.
- (15) Marrink, S. J.; Risselada, H. J.; Yefimov, S.; Tieleman, D. P.; de Vries, A. H. *J. Phys. Chem. B* **2007**, *111*, 7812–7824.
- (16) Marrink, S. J.; Mark, A. E. *J. Am. Chem. Soc.* **2003**, *125*, 15233–15242.
- (17) Marrink, S. J.; Mark, A. E. *J. Am. Chem. Soc.* **2003**, *125*, 11144–11145.
- (18) Kasson, P. M.; Kelley, N. W.; Singhal, N.; Vrljic, M.; Brunger, A. T.; Pande, V. S. *Proc. Natl. Acad. Sci. U.S.A.* **2006**, *103*, 11916–11921.
- (19) Smirnova, Y. G.; Marrink, S. J.; Lipowsky, R.; Knecht, V. *J. Am. Chem. Soc.* **2010**, *132*, 6710–6718.
- (20) Marrink, S. J.; Mark, A. E. *Biophys. J.* **2004**, *87*, 6710–6718.
- (21) Baoukina, S.; Monticelli, L.; Risselada, H. J.; Marrink, S. J.; Tieleman, D. P. *Proc. Natl. Acad. Sci. U.S.A.* **2008**, *105*, 10803–10808.
- (22) Risselada, H. J.; Marrink, S. J. *Proc. Natl. Acad. Sci. U.S.A.* **2008**, *105*, 17367–17372.
- (23) Monticelli, L.; Kandasamy, S. K.; Periole, X.; Larson, R. G.; Tieleman, D. P.; Marrink, S. J. *J. Chem. Theory Comput.* **2008**, *4*, 819.
- (24) Yefimov, S.; van der Giessen, E.; Onck, P. R.; Marrink, S. J. *Biophys. J.* **2008**, *94*, 2994–3002.
- (25) Yoo, J.; Cui, Q. *Biophys. J.* **2009**, *97*, 2267–2276.
- (26) Lu, L. Y.; Voth, G. A. *J. Phys. Chem. B* **2009**, *113*, 1501–1510.
- (27) Deserno, M. *Macromol. Rapid Commun.* **2009**, *30*, 752–771.
- (28) Wang, Z. J.; Deserno, M. *J. Phys. Chem. B* **2010**, *114*, 11207–11220.
- (29) Lin, J. H.; Baker, N. A.; McCammon, J. A. *Biophys. J.* **2002**, *83*, 1374–1379.
- (30) Clarke, R. J. *Adv. Colloid Interface Sci.* **2001**, *89*, 263–281.
- (31) Shapovalov, V. L.; Kotova, E. A.; Rokitskaya, T. I.; Antonenko, Y. N. *Biophys. J.* **1999**, *77*, 299–305.
- (32) Siu, S. W. I.; Vácha, R.; Jungwirth, P.; Böckmann, R. A. *J. Chem. Phys.* **2008**, *128*, 125103.
- (33) Harder, E.; MacKerell, A. D.; Roux, B. *J. Am. Chem. Soc.* **2009**, *131*, 2760–2761.
- (34) Orsi, M.; Haubertin, D. Y.; Sanderson, W. E.; Essex, J. W. *J. Phys. Chem. B* **2008**, *112*, 802–815.
- (35) Liu, Y.; Ichiye, T. *J. Phys. Chem.* **1996**, *100*, 2723–2730.
- (36) Yesylevskyy, S. O.; Schäfer, L. V.; Sengupta, D.; Marrink, S. J. *PLoS Comput. Biol.* **2010**, *6*, e1000810.
- (37) Singh, G.; Tieleman, D. P. *J. Chem. Theory Comput.* **2011**, *7*, 2316–2324.
- (38) Vorobyov, I.; Bekker, B.; Allen, T. W. *Biophys. J.* **2010**, *98*, 2904–2913.
- (39) Darre, L.; Machado, M. R.; Dans, P. D.; Herrera, F. E.; Pantano, S. *J. Chem. Theory Comput.* **2010**, *6*, 3793–3807.
- (40) Riniker, S.; van Gunsteren, W. F. *J. Chem. Phys.* **2011**, *134*, 084110.
- (41) Wu, Z.; Cui, Q.; Yethiraj, A. *J. Phys. Chem. B* **2010**, *114*, 10524–10529.
- (42) Wu, Z.; Cui, Q.; Yethiraj, A. *J. Phys. Chem. Lett.* **2011**, *2*, 1794–1798.
- (43) MacCallum, J. L.; Moghaddam, M. S.; Chan, H. S.; Tieleman, D. P. *Proc. Natl. Acad. Sci. U.S.A.* **2007**, *104*, 6206–6210.
- (44) Berne, B. J.; Weeks, J. D.; Zhou, R. H. *Annu. Rev. Phys. Chem.* **2009**, *60*, 85–103.

- (45) Shinoda, W.; Devane, R.; Klein, M. L. *Mol. Simul.* **2007**, *33*, 27–36.
- (46) Fumi, F. G.; Tosi, M. P. *J. Phys. Chem. Solids* **1964**, *25*, 31–43.
- (47) Fumi, F. G.; Tosi, M. P. *J. Phys. Chem. Solids* **1964**, *25*, 45–52.
- (48) MacCallum, J. L.; Bennett, W. F. B.; Tieleman, D. P. *Biophys. J.* **2008**, *94*, 3393–3404.
- (49) Baron, R.; de Vries, A. H.; Hünenberger, P. H.; van Gunsteren, W. F. *J. Phys. Chem. B* **2006**, *110*, 8464–8473.
- (50) Alemani, D.; Collu, F.; Cascella, M.; Peraro, M. D. *J. Chem. Theory Comput.* **2010**, *6*, 315–324.
- (51) Hess, B.; Kutzner, C.; van der Spoel, D.; Lindahl, E. *J. Chem. Theory Comput.* **2008**, *4*, 435–447.
- (52) Berendsen, H. J. C.; Postma, J. P. M.; van Gunsteren, W. F.; Dinola, A.; Haak, J. R. *J. Chem. Phys.* **1984**, *81*, 3684–3690.
- (53) Miyamoto, S.; Kollman, P. A. *J. Comput. Chem.* **1992**, *13*, 952–962.
- (54) Hess, B.; Bekker, H.; Berendsen, H. J. C.; Fraaije, J. G. E. M. *J. Comput. Chem.* **1997**, *18*, 1463–1472.
- (55) Beveridge, D. L.; DiCapua, F. M. *Annu. Rev. Biophys. Biophys. Chem.* **1989**, *18*, 431–492.
- (56) van Gunsteren, W. F.; Daura, X.; Mark, A. E. *Helv. Chim. Acta* **2002**, *85*, 3113–3129.
- (57) Venable, R. M.; Skibinsky, A.; Pastor, R. W. *Mol. Simul.* **2006**, *32*, 849–855.
- (58) Jiang, F. Y.; Bouret, Y.; Kindt, J. T. *Biophys. J.* **2004**, *87*, 182–192.
- (59) Humphrey, W.; Dalke, A.; Schulten, K. *J. Mol. Graphics* **1996**, *14.1*, 33–38.
- (60) Yu, H.; Karplus, M. *J. Chem. Phys.* **1988**, *89*, 2366.
- (61) Baron, R.; Trzesniak, D.; de Vries, A. H.; Elsener, A.; Marrink, S. J.; van Gunsteren, W. F. *Chem. Phys. Chem.* **2007**, *8*, 452–461.
- (62) Khadikar, P. V.; Mandloi, D.; Bajaj, A. V.; Joshi, S. *Bioorg. Med. Chem. Lett.* **2003**, *13*, 419–422.
- (63) Nagle, J. F.; Tristram-Nagle, S. *Biochim. Biophys. Acta* **2000**, *1469*, 159–195.
- (64) Kučerka, N.; Nagle, J. F.; Sachs, J. N.; Feller, S. E.; Pencser, J.; Jackson, A.; Katsaras, J. *Biophys. J.* **2008**, *95*, 2356–2367.
- (65) Balgavý, P.; Dubnicková, M.; Kucerka, N.; Kiselev, M. A.; Yaradaikin, S. P.; Uhríková, D. *Biochim. Biophys. Acta* **2001**, *1*, 40–52.
- (66) Petrache, H. I.; Dodd, S.; Brown, M. *Biophys. J.* **2000**, *79*, 3172–3192.
- (67) Rand, R. P.; Parsegian, V. A. *Biochim. Biophys. Acta* **1989**, *988*, 351–376.
- (68) Petrache, H. I.; Stephanie, T. N.; Gawrisch, K.; Harries, D.; Parsegian, V. A.; Nagle, J. F. *Biophys. J.* **2004**, *86*, 1574–1586.
- (69) Rawicz, W.; Olbrich, K. C.; McIntosh, T.; Needham, D.; Evans, E. *Biophys. J.* **2000**, *79*, 328–339.
- (70) Genco, I.; Gliozzi, A.; Relini, A.; Robello, M.; Scallan, E. *Biophys. J.* **1993**, *1149*, 10–18.
- (71) Zhelev, D.; Needham, D. *Biochim. Biophys. Acta* **1993**, *1147*, 89–104.
- (72) Moroz, J. D.; Nelson, P. *Biophys. J.* **1997**, *72*, 2211–2216.
- (73) de Joannis, J.; Jiang, F. Y.; Kindt, J. T. *Langmuir* **2006**, *22*, 998–1005.
- (74) Glaser, R. W.; Leikin, S. L.; Chernomordik, L. V.; Pastushenko, V. F.; Sokirko, A. I. *Biochim. Biophys. Acta* **1988**, *940*, 275–287.
- (75) Kuo, A.-L.; Wade, C. G. *Biochemistry* **1979**, *18*, 2300–2308.
- (76) Rand, R. P.; Fuller, N. L. *Biophys. J.* **1994**, *66*, 2127–2138.
- (77) Zimmerberg, J.; Kozlov, M. M. *Nat. Rev. Mol. Cell Biol.* **2006**, *7*, 9–19.
- (78) Kirk, G. L.; Gruner, S. M.; Stein, D. L. *Biochemistry* **1984**, *23*, 1093–1102.
- (79) Gruner, S. M. *Proc. Natl. Acad. Sci. U.S.A.* **1985**, *82*, 3665–3669.
- (80) Turner, D. C.; Gruner, S. M. *Biochemistry* **1992**, *31*, 1340–1355.
- (81) Li, L.; Vorobyov, I.; Allen, T. W. *J. Phys. Chem. B* **2008**, *112*, 9574–9587.
- (82) Herce, H. D.; Garcia, A. E. *Proc. Natl. Acad. Sci. U.S.A.* **2007**, *104*, 20805–20810.
- (83) Lee, H. L.; Dubikovskaya, E. A.; Hwang, H.; Semyonov, A. N.; Wang, H.; Jones, L. R.; Twieg, R. J.; Moerner, W. E.; Wender, P. A. *J. Am. Chem. Soc.* **2008**, *130*, 9364–9370.
- (84) Yesylevskyy, S.; Marrink, S. J.; Mark, A. E. *Biophys. J.* **2009**, *97*, 40–49.
- (85) Schmidt, N.; Mishra, A.; Lai, G. H.; Wong, G. C. L. *FEBS Lett.* **2010**, *584*, 1806–1813.
- (86) Schmidt, N. W.; Mishra, A.; Lai, G. H.; Davis, M.; Sanders, L. K.; Tran, D.; Garcia, A.; Tai, K. P.; McCray, P. B., Jr.; Ouellette, A. J.; Selsted, M. E.; Wong, G. C. L. *J. Am. Chem. Soc.* **2011**, *133*, 6720–6727.
- (87) Gurtovenko, A. A.; Vattulainen, I. *J. Am. Chem. Soc.* **2005**, *127*, 17570–17571.
- (88) Gurtovenko, A. A.; Vattulainen, I. *Biophys. J.* **2007**, *92*, 1878–1890.
- (89) Praprotnik, M.; Matysiak, S.; Site, L. D.; Kremer, K.; Clementi, C. *J. Phys.: Condens. Matter* **2007**, *19*, 292201.
- (90) Matysiak, S.; Clementi, C.; Praprotnik, M.; Kremer, K.; Site, L. D. *J. Chem. Phys.* **2008**, *128*, 024503.

# A New Approach for Investigating the Molecular Recognition of Protein: Toward Structure-Based Drug Design Based on the 3D-RISM Theory

Yasuomi Kiyota,<sup>†</sup> Norio Yoshida,<sup>†,‡</sup> and Fumio Hirata<sup>\*,†,‡</sup>

<sup>†</sup>Department of Theoretical Molecular Science, Institute for Molecular Science, Okazaki 444-8585, Japan

<sup>‡</sup>Department of Functional Molecular Science, The Graduate University for Advanced Studies, Okazaki 444-8585, Japan

**ABSTRACT:** A new approach to investigate a molecular recognition process of protein is presented based on the three-dimensional reference interaction site model (3D-RISM) theory, a statistical mechanics theory of molecular liquids. Numerical procedure for solving the conventional 3D-RISM equation consists of two steps. In step 1, we solve ordinary RISM (or 1D-RISM) equations for a solvent mixture including target ligands in order to obtain the density pair correlation functions (PCF) among molecules in the solution. Then, we solve the 3D-RISM equation for a solute–solvent system to find three-dimensional density distribution functions (3D-DDF) of solvent species around a protein, using PCF obtained in the first step. A key to the success of the method was to regard a target ligand as one of “solvent” species. However, the success is limited due to a difficulty of solving the 1D-RISM equation for a solvent mixture, including large ligand molecules. In the present paper, we propose a method which eases the limitation concerning solute size in the conventional method. In this approach, we solve a solute–solute 3D-RISM equations for a protein–ligand system in which both proteins and ligands are regarded as “solutes” at infinite dilution. The 3D- and 1D-RISM equations are solved for protein–solvent and ligand–solvent systems, respectively, in order to obtain the 3D- and 1D-DDF of solvent around the solutes, which are required for solving the solute–solute 3D-RISM equation. The method is applied to two practical and noteworthy examples concerning pharmaceutical design. One is an odorant binding protein in the *Drosophila melanogaster*, which binds an ethanol molecule. The other is phospholipase A2, which is known as a receptor of acetylsalicylic acid or aspirin. The result indicates that the method successfully reproduces the binding mode of the ligand molecules in the binding sites measured by the experiments.

## 1. INTRODUCTION

The molecular recognition (MR) in living systems is a crucial elementary process for biomolecules to perform their functions as, for example, enzymes or ion channels. The MR process can be defined as a molecular process in which one or few guest molecules are bound in high probability at a particular site, a cleft or a cavity, of a host molecule in a particular orientation. The process is governed essentially by the two physicochemical properties: (1) difference in the thermodynamic stability (or free energy) between the bound and unbound states of host and guest molecules, and (2) structural fluctuation of molecules. In this article, we propose a new approach to describe the molecular recognition process based on the statistical mechanics of molecular liquids.

In the last three decades, many computational methodologies for investigating a MR process have been proposed.<sup>1–18</sup> As is mentioned above, the focus of the MR in silico is concerned with the prediction of ligands or drugs that would be strongly bound to key regions of a receptor or an enzyme. The popular “docking simulation” for drug design uses essentially a trial and error scheme to find a “best-fit complex” of host and guest molecules based on geometrical and/or energetic criteria.<sup>3,4</sup> However, the best-fit complex in a geometrical sense is not necessarily the most stable one in terms of the thermodynamics because it cannot account for the solvent; so neither the dehydration penalty nor the entropy barrier is taken into account. By the “dehydration

penalty,” we mean a free energy penalty concerning a molecular process in which a water molecule detaches from a binding site.

The so-called implicit solvent models, the generalized Born<sup>5</sup> and the Poisson–Boltzmann equations,<sup>6</sup> which have been used most popularly for evaluating the solvation thermodynamics of biomolecules, are not accurate and insightful for this problem under concern, because by definition they do not have a molecular view for solvent. It is impossible to define a dielectric constant of solvent inside a host cavity, and therefore, it cannot account for the dehydration penalty, especially that from the host cavity. At best, those quantities can be calculated by fitting the empirical parameters, such as the boundary conditions and the dielectric constants, with experimental data, but then it loses credibility as a theory to predict the phenomena.

The molecular simulation, on the other hand, can provide the most detailed molecular view for the process. The simulation methods, molecular dynamics (MD) and Monte Carlo (MC), sample the configuration space of a ligand–receptor system in solvent using the numerical integration of the Newtonian equation (MD) or the probabilistic search along the Markov path (MC) in order to evaluate the free energy difference between the bound and unbound states of the host–guest system. However, this type of simulation does not work for the

Received: May 29, 2011

Published: September 15, 2011

problem well, because a MR process is usually slow as well as rare events. A common strategy adopted by the simulation community to overcome the difficulty is a non-Boltzmann-type sampling which defines a “reaction coordinate” or an “order parameter” onto which all other degrees of freedoms are projected. The best example is the “umbrella” sampling to realize the potential of mean force or the free energy along a conduction path of an ion in an ion channel.<sup>7</sup> The method is quite powerful for sampling the configuration space around an order parameter if the parameter is unique and if the configuration space to be projected on the parameter is sufficiently small. Unfortunately, the problems in the biochemical processes are not so simple as can be described by a unique order parameter. So, it is often the case that the results of the simulation depend on choices of order parameters and on “scheduling” of the sampling. The other methodology employed to accelerate the sampling is to apply an artificial external force on the system. That kind of simulation should verify that the configuration of water satisfies the Boltzmann distribution.

A different approach to the MR process has been developed based on the three-dimensional reference interaction site model (3D-RISM) theory during the last five years.<sup>8–18</sup> The method integrates *analytically* the configuration space of a ligand–receptor system in solvent by means of the statistical mechanics. The analytical integration which extends over the infinitely large configuration space is the advantage which distinguishes the method from the molecular simulation. Due to this advantage, the 3D-RISM theory is free from the difficulties which the simulation methods are facing.

The 3D-RISM equation was derived from the molecular Ornstein–Zernike (MOZ) equation, the most fundamental equation to describe the density pair correlation of liquids, for a solute–solvent system in the infinite dilution by taking a statistical average over the orientation of solvent molecules.<sup>19–21,35–37</sup> By solving the combined 3D-RISM with RISM equations, the latter providing the solvent structure in terms of the site–site density pair correlation functions, one can get the “solvation structure” or the solvent distributions around a solute. The high peak of the solvent distributions indicates that the solvent affinity of target protein or receptor at that point is high. Therefore, the MR process can be probed by the solvent or ligand distribution. The method produces naturally all the solvation thermodynamics as well, including energy, entropy, free energy, and their derivatives, such as the partial molar volume and compressibility. Unlike the molecular simulation, there is no necessity for concern about size of the system and “sampling” of the configuration space, because the method treats essentially the infinite number of molecules and integrates over the entire configuration space of a system.<sup>22</sup>

By the way, in all previous studies of MR by 3D-RISM theory, the receptor protein and ligand molecule were regarded as solute and solvent, respectively. In those cases, the MR process is analyzed in terms of solvent distribution around a solute molecule, which is called a sol’ute–sol’vent density distribution function (uv-DDF). The RISM equation for solvent system should be solved before a 3D-RISM calculation. There, “solvent” consists of ligand molecules, water, and other components of solution. Because the RISM and 3D-RISM equations coupled with closure relation (i.e., hypernetted chain or Kovalenko–Hirata (KH) closure) are nonlinear integral equations, those are solved in an iterative manner.<sup>23,24</sup>

Although many theoretical and methodological efforts have been devoted to solve the RISM and 3D-RISM equations, the

system has been largely limited to that including relatively small solvent molecules, such as water, ions, carbon monoxide, and the largest being glycerol.<sup>14,25,26</sup> A reason why is because numerical solution of the RISM equation for solution including large ligands becomes increasingly unstable as the size of ligands increases. However, many ligands of biological interests, including ordinary drug molecules, are not so small. Therefore, we propose a new approach to tackle MR of large ligand molecules by protein based on the 3D-RISM and RISM theories. The strategy of the method is to regard a ligand molecule as a solute in addition to a receptor protein, which are immersed in solvent in the infinite dilution. The distribution of ligand molecules around a receptor protein is described by the sol’ute–sol’ute density distribution function (uu-DDF), instead of uv-DDF. In this sense, the new method is named as “uu-3D-RISM.” Under the treatment of this method, interactions between ligand molecules can be ignored completely from the consideration, because the density of ligand molecule is vanishingly small at the limit. Therefore, it is not necessary to solve the ligand–ligand RISM equation, the most unstable equation, anymore. This assumption stabilizes the numerical solutions of a set of the 3D-RISM and RISM equations dramatically.

An approach to uu-DDF has already been proposed by Kovalenko and Hirata to investigate the potential of mean force between two molecular ions in a polar molecular solvent.<sup>21</sup> In their method, uu-DDF is a function of position and orientation of two solute molecules. Therefore, in order to obtain the uu-DDF for all possible positions and orientations of a ligand molecule, the method requires sampling in the entire coordinate space. On the contrary, the uu-DDF can be obtained by a single calculation in the present method, because one of the solute molecules, usually a ligand molecule, is treated in terms of an interaction-site model.

This paper is organized as follows. In Section 2, we briefly review the RISM and 3D-RISM theories in order to identify each member of the RISM family and derive the uu-3D-RISM equation. We also clarify how the molecular recognition process of biological system is treated with the uu-3D-RISM method. Section 3 is devoted to applications of the uu-3D-RISM method to two practical and noteworthy examples, odorant binding protein<sup>27</sup> and phospholipase A2.<sup>28–31</sup> Section 4 concludes the paper.

## 2. METHOD

**2.1. Outline of 3D-RISM Theory.** The MOZ integral equation for a multicomponent system is written as

$$h^{ij}(1,2) = c^{ij}(1,2) + \sum_l \int c^{ij}(1,3) \rho^l h^{lj}(3,2) d(3) \quad (1)$$

where  $\rho^l$  is the number density of species  $l$ ,  $h^{ij}(1,2)$  and  $c^{ij}(1,2)$  denote the total and direct correlation functions between a pair of molecular species  $i$  and  $j$  in a solution, respectively.<sup>32</sup> The numbers in the parentheses represent the coordinates of molecules in the liquid system, including both the position  $\mathbf{r}$  and the orientation  $\Omega$ . The total correlation function  $h^{ij}(1,2)$  is related to the density pair correlation functions  $g^{ij}(1,2)$  by  $h^{ij}(1,2) = g^{ij}(1,2) - 1$ . The summation in the right-hand side runs over species in a mixture.

The eq 1 depends essentially on six coordinates in the Cartesian space, and they include a six-fold integral. This integral is the one which had prevented the theory from applying to



polyatomic molecules. It is the interaction-site model and the RISM approximation proposed by Chandler and Andersen<sup>33</sup> that enabled one to solve the equations. The idea behind the model is to project the functions onto the one-dimensional space along the distance between a pair of interaction sites, usually placed on the center of atoms, by taking the statistical average over the angular coordinates of molecules, fixing the separation between two interaction sites. The projection can be accomplished by the following equation:

$$h_{\alpha\gamma}(r) = \frac{1}{\Omega^2} \int \delta(|\mathbf{r}_1 + \mathbf{I}_1^\alpha|) \delta(|\mathbf{r}_2 + \mathbf{I}_2^\gamma| - r) h(1,2) d(1) d(2) \quad (2)$$

where  $\mathbf{r}_1$  and  $\mathbf{I}_1^\alpha$  indicate the position of molecule 1 in laboratory frame and the position of site  $\alpha$  of molecule 1 in molecular frame, respectively.

Now, we classify molecular species in the system into two categories, “solute” and “solvent,” respectively, as previous works.<sup>19–21</sup> After this, the superscripts “u” and “v” denote solute and solvent, respectively. For example,  $h^{vv}$  is the total correlation function between different molecular species  $v$  and  $v'$  in solvent, and  $\rho^v$  and  $\rho^u$  also are elements of diagonal matrices which denote density of each species in solvent and solute. The summations concerning  $v$  and  $v'$  in the equations run over solvent species, while  $u$  and  $u'$  run over solute species. The most interesting case to investigate “solvation” of a biomolecule can be realized by taking the infinite dilution limit for all the solute species, namely,  $\rho^u \rightarrow 0$ . Then, eqs 1 and 2 are constructed from solvent–solvent, solute–solvent, and solute–solute systems, and you note that these equations can be solved sequentially, because the former equation is independent from the later.

The RISM equation can be derived from eqs 1 and 2 with a super position approximation for the direct correlation function, which reads

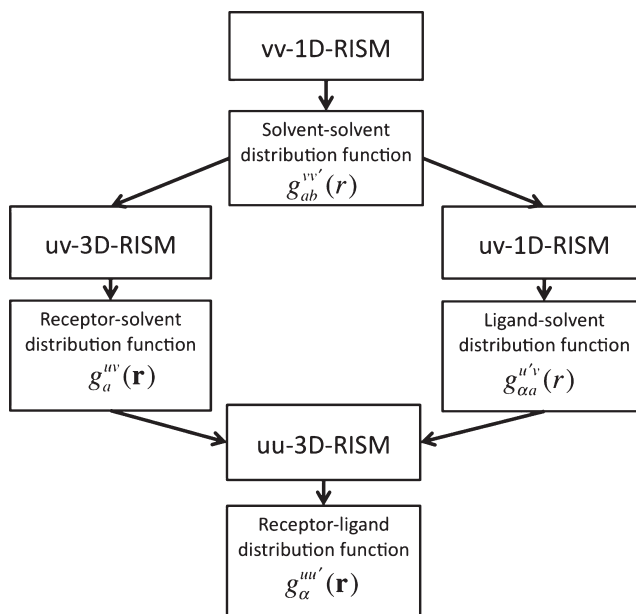
$$h_{\alpha\lambda}^{vv'}(r) = \sum_{\substack{\eta \in v' \\ \gamma \in v}} \omega_{\alpha\gamma}^v * c_{\gamma\eta}^{vv'} * \omega_{\eta\lambda}^{v'}(r) + \sum_{v''} \rho^{v''} \sum_{\substack{\eta \in v'' \\ \gamma \in v}} \omega_{\alpha\gamma}^v * c_{\gamma\eta}^{vv''} * h_{\eta\lambda}^{v''v'}(r) \quad (3)$$

where  $\omega$  is an intramolecular correlation function, the asterisk denotes the convolution integrals, and  $\rho^v$  denotes the number density of solvent species  $v$ . For clarity, we refer to eq 3 as vv-1D-RISM equation hereafter. A similar equation can be derived from eqs 1 and 2 for a solute–solvent system in the infinite dilution limit as follows:

$$h_{\alpha\lambda}^{uv'}(r) = \sum_{\substack{\eta \in v \\ \gamma \in u}} \omega_{\alpha\gamma}^u * c_{\gamma\eta}^{uv} * \omega_{\eta\lambda}^v(r) + \sum_{v'} \rho^{v'} \sum_{\substack{\eta \in v' \\ \gamma \in v}} \omega_{\alpha\gamma}^u * c_{\gamma\eta}^{uv'} * h_{\eta\lambda}^{v'v'}(r) \quad (4)$$

The theory has been proven to be so successful for describing structure and thermodynamics of liquid and liquid mixtures, including a variety of aqueous solutions.<sup>20</sup> However, the theory has exhibited serious breakdown, especially when it was applied to solutions including macromolecules, such as protein as a solute.<sup>34</sup> Then, in order to avoid the problem, alternative

Scheme 1. Scheme of uu-RISM Method



approaches have been developed during two decades.<sup>19,35</sup> We have derived from eq 1 by taking an average over angular coordinates only for solvent coordinates, not for solute coordinates.<sup>36,37</sup>

$$h_{\alpha}^{uv}(\mathbf{r}) = \frac{1}{\Omega} \int \delta(\mathbf{r}_2 + \mathbf{I}_2^\alpha - \mathbf{r}) h^{uv}(1,2) d(2) = \sum_{v' \in \text{solvent}} \sum_{\gamma \in v'} c_{\alpha}^{uv'} * [\omega_{\gamma\alpha}^{v'} + \rho^{v'} h_{\gamma\alpha}^{v'v'}](\mathbf{r}) \quad (5)$$

where  $v'$  runs over the all solvent species. Where we also employed a super position approximation for the solute–solvent direct correlation function:

$$c_{\alpha}^{uv}(1,2) \equiv \sum_{a \in v} c_{\alpha}^{ua}(1,2) \quad (6)$$

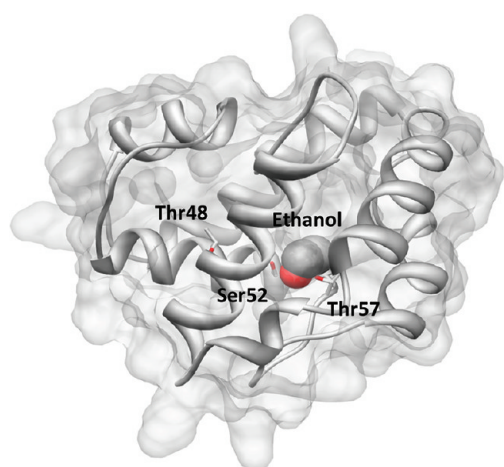
This is the basic assumption of the 3D-RISM theory. The solvent–solvent (vv) total correlation function appeared in the right-hand side of eq 5 is evaluated from the vv-1D-RISM equation, or eq 4, in advance. The  $\rho h(\mathbf{r})$  is essentially the “second moment of the density fluctuation” of two spatial points, which can be identified as a “mean excess density” due to the method devised by Percus.<sup>38</sup>

The MR phenomena can be described by solving the vv-1D-RISM and 3D-RISM equations sequentially, considering that a receptor protein is immersed in solvent–ligand mixture in the infinite dilution limit of the receptor. In other words, ligand molecules are treated as one of components of a solvent mixture. MR is realized in terms of the  $\rho h(\mathbf{r})$  of ligand atoms at a binding site, relative to bulk solutions; if  $\rho h(\mathbf{r})$  is greater than zero, then we conclude that the ligand is “recognized” by the site. So, the procedure of realizing MR by 3D-RISM is quite straightforward. There is no necessity to define order parameters, such as “reaction coordinates” and “umbrella”, for exploring the configuration space of ligands, which is the case in the molecular simulations.

Table 1. Summary of Performed Calculations<sup>a</sup>

species	odorant binding protein			phospholipase A2	
		uu-RISM	uv-RISM		uu-RISM
protein	1OOF	3D representation (solute)	3D representation (solute)	1OXR	3D representation (solute)
ligand	ethanol	site representation (solute)	site representation (solvent)	aspirin	site representation (solute)
solvent	water	site representation (solvent)	site representation (solvent)	water	site representation (solvent)

<sup>a</sup>Note that the type of ligand is different between uu- and uv-3D-RISM.

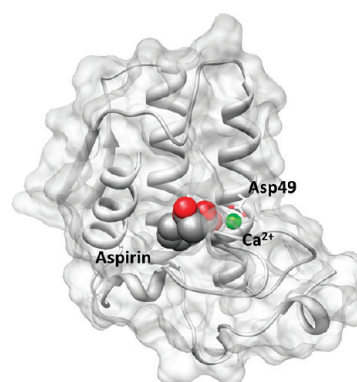


**Figure 1.** X-ray crystal structure of the odorant binding protein LUSH from *Drosophila melanogaster*. It has a specific alcohol binding site which can bind a series of short-chain n-alcohols. The structure taken from PDB (PDB ID: 1OOF) includes one ethanol molecule which is represented by the VDW surface. The protein surfaces are represented as a gray transparent surface. The binding site is constructed by a group of amino acids, Thr57, Ser52, and Thr48. These amino acids form a network of concerted hydrogen bonds between the protein and the alcohol.

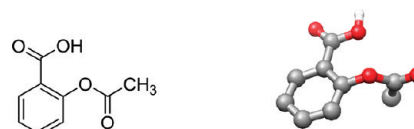
It is the inclusion of ligand species in solvent that gave 3D-RISM/vv-1D-RISM a great advantage. In fact, many applications of 3D-RISM to MR processes so far have been so successful as long as small ligands including water, CO, NH<sub>3</sub>, and metal ions, etc. are concerned. However, the advantage turns into disadvantage when large molecules including most of the drug compounds are involved. The problem originates in the vv-1D-RISM equation, not in the 3D-RISM equation. To describe a MR process with 3D-RISM, we have to solve the vv-1D-RISM equation for solvent mixture, including ligand species in advance. However, according to our experience, numerical solutions of the vv-1D-RISM equation for a mixture including a large compound are quite unstable due to inherent nonlinearity of the integral equation, which increases with increasing size and complexity of molecules.

In the following subsection, we develop a new approach which is derived from the solute–solute MOZ equation in which both receptor and ligand molecules are dissolved in a solvent mixture at the infinite dilution.

**2.2. uu-3D-RISM Equation.** The strategy of the new approach is to regard a ligand molecule as a solute molecule, in addition to a receptor protein, which is immersed in solvent in the infinite dilution limit. By this assumption, eq 1 can be regarded as a protein–ligand uu-MOZ equation. In the present approach,



**Figure 2.** X-ray crystal structure of the complex formed between phospholipase A2 (PLA2) and 2-acetoxybenzoic acid, aspirin. It can be taken from PDB (PDB ID: 1OXR). Phospholipase A2 can bind aspirin, which is represented by VDW surface, for anti-inflammatory effects in its specific binding site. The protein surfaces are represented as a gray transparent surface. The aromatic ring of aspirin is embedded in the hydrophobic environment, and other substituted groups form several important attractive interactions with calcium ion, His48, and Asp49. Calcium ion is shown as a green sphere.



**Figure 3.** Structure aspirin depicted with two different presentations.

since a ligand molecule can be assumed to be reasonably small, we apply the interaction site model to a ligand molecule in a manner similar to solvent. The uu-3D-RISM can be derived from eq 1 by taking an average over angular coordinates for ligand(solute) coordinates, not for protein(solute) coordinate. The solute–solute total correlation functions can be written as

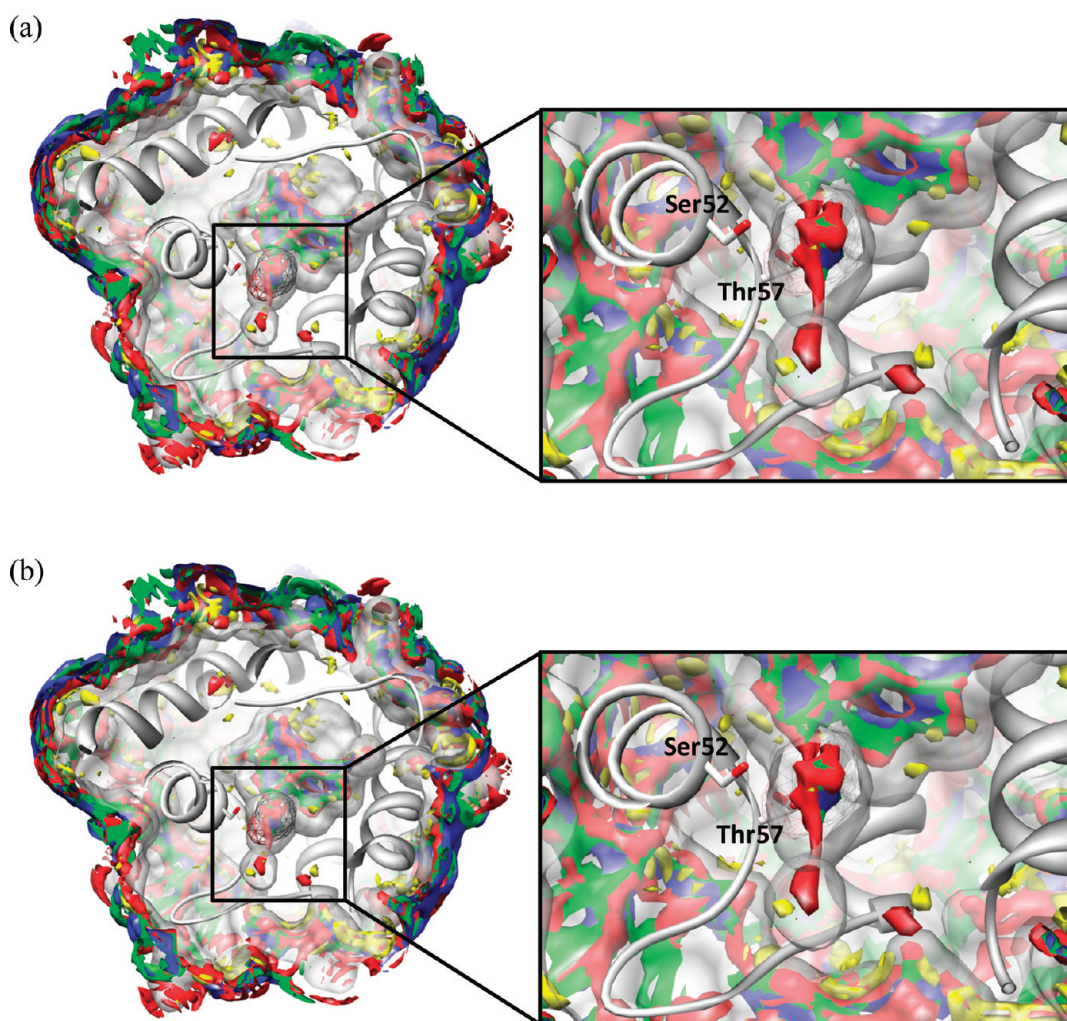
$$h_{\alpha}^{uu'}(\mathbf{r}) = \frac{1}{\Omega} \int \delta(\mathbf{r}_2 + \mathbf{I}_2^{\alpha} - \mathbf{r}) h^{uu'}(1, 2) d(2) \quad (7)$$

Accordingly, the Fourier transform of total correlation functions is obtained by orientational reduction of Fourier transform of the molecular total correlation functions:

$$\tilde{h}_{\alpha}^{uu'}(\mathbf{k}) = \frac{1}{\Omega} \int d\Omega_2 e^{i\mathbf{k}\cdot\mathbf{r}_2} \tilde{h}^{uu'}(k, \Omega_1, \Omega_2) \quad (8)$$

We employ a super position approximation for the solute–solute direct correlation function:

$$c_{\alpha}^{uu'}(1, 2) \equiv \sum_{\alpha \in v} c_{\alpha}^{uu'}(\mathbf{r}) \quad (9)$$



**Figure 4.** The 3D-DDF of ethanol around and inside odorant binding protein, LUSH, obtained by (a) uu-3D- and (b) uv-3D-RISM with the threshold  $g_{\gamma}(\mathbf{r}) > 2$ : blue, CH<sub>3</sub>; green, CH<sub>2</sub>; red, oxygen atom of hydroxyl group; yellow, hydrogen atom of hydroxyl group. The protein surfaces are represented as a gray transparent surface. The location of ethanol in X-ray structure is depicted with a wire frame.

The Fourier transform of solute–solute direct correlation functions can be defined as

$$\begin{aligned}\tilde{c}^{uu'}(k, \Omega_1, \Omega_2) &= \int d\mathbf{r}_{12} e^{i\mathbf{k}\mathbf{r}_{12}} c^{uu'}(r_{12}, \Omega_1, \Omega_2) \\ &= \sum_{\alpha} e^{-i\mathbf{k}\mathbf{r}_{\alpha 2}} c_{\alpha}^{uu'}(\mathbf{k})\end{aligned}\quad (10)$$

The solute–solute MOZ eq 1 in Fourier space can be written as

$$\begin{aligned}\tilde{h}^{uu'}(k, \Omega_1, \Omega_2) &= \int d\mathbf{r}_{12} e^{i\mathbf{k}\mathbf{r}_{12}} h^{uu'}(r_{12}, \Omega_1, \Omega_2) \\ &= \int d\mathbf{r}_{12} e^{i\mathbf{k}\mathbf{r}_{12}} c^{uu'}(r_{12}, \Omega_1, \Omega_2) \\ &+ \frac{1}{\Omega} \sum_{\nu} \int d\mathbf{r}_{12} \int d(3) e^{i\mathbf{k}\mathbf{r}_{12}} c^{uv}(r_{12}, \Omega_1, \Omega_2) \rho^{\nu} h^{\nu u'}(r_{12}, \Omega_1, \Omega_2)\end{aligned}\quad (11)$$

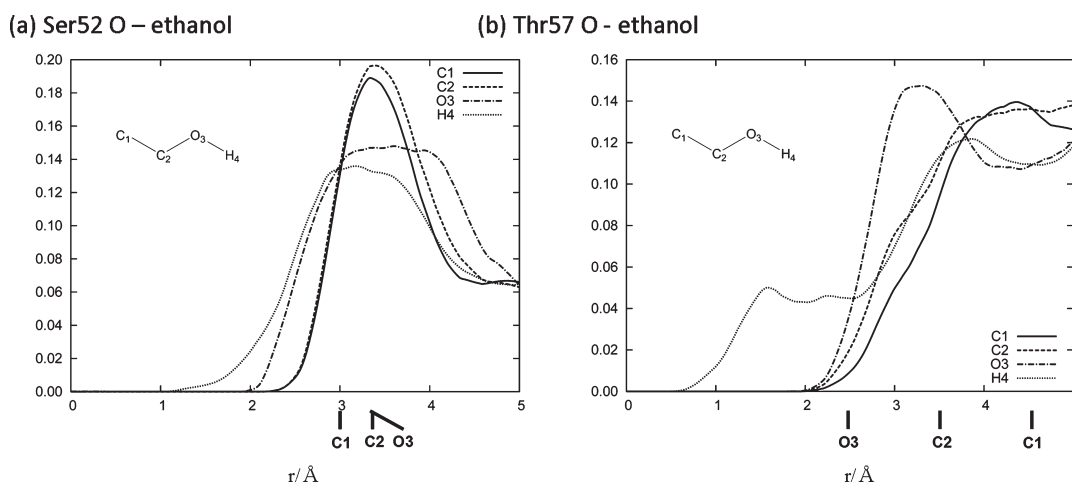
From eqs 2 and 8–11, the uu-3D-RISM equation is obtained as

$$h_{\alpha}^{uu'}(\mathbf{r}) = \sum_{\gamma} c_{\gamma}^{uu'} * \omega_{\gamma\alpha}^{u'}(\mathbf{r}) + \sum_{\nu} \sum_{\gamma} c_{\gamma}^{uv} * \rho^{\nu} h_{\gamma\alpha}^{\nu u'}(\mathbf{r})\quad (12)$$

In order to solve those 1D- and 3D-RISM equations obtained above, we need another equation which complements or “closes” the equations. Here, we employ the KH closure which reads<sup>39</sup>

$$\begin{aligned}g_{\alpha}^{uu'}(\mathbf{r}) &= \begin{cases} \exp(d_{\alpha}^{uu'}(\mathbf{r})) & \text{for } d_{\alpha}^{uu'}(\mathbf{r}) \leq 0 \\ 1 + d_{\alpha}^{uu'}(\mathbf{r}) & \text{for } d_{\alpha}^{uu'}(\mathbf{r}) > 0 \end{cases} \\ d_{\alpha}^{uu'}(\mathbf{r}) &= -\beta u_{\alpha}^{uu'}(\mathbf{r}) + h_{\alpha}^{uu'}(\mathbf{r}) - c_{\alpha}^{uu'}(\mathbf{r})\end{aligned}\quad (13)$$

The procedure to obtain the receptor–ligand distribution function is shown in Scheme 1. First, vv-DDF is evaluated by vv-1D-RISM, where solvent includes water, electrolyte, organic solvent, and so on. The vv-DDF is used in both uv-3D- and uv-1D-RISM calculations. The uv-3D- and uv-1D-RISM calculations are carried out to obtain receptor–solvent and ligand–solvent DDF, respectively. By inserting these two DDFs, uu-3D-RISM can be solved to get receptor–ligand DDF.



**Figure 5.** RDFs of ethanol around hydroxyl groups of (a) Ser52 and (b) Thr57, respectively. Oxygen atoms in each hydroxyl groups were chosen as the averaging center.

**Table 2.** Distance (Å) Matrix between the Specific Sites of LUSH and Ethanol

sites	Ser52-O	Thr57-O
C1	3.0	4.6
C2	3.3	3.5
O3	3.3	2.3

From receptor–ligand DDF, a MR analysis, such as a binding pocket search, can be performed.

**2.3. Extracting Binding Mode of Ligand Inside Protein from 3D-DDF.** At this point, we would like to make a comment on general difficulty to extract “binding mode,” or position and orientation, of a ligand inside protein from 3D-DDF. The problem is common to the experimental methodologies, such as X-ray and neutron diffractions, since both the theory and the experiments are observing essentially the same property, namely, the density distribution of ligand atoms, which are atomic positions statistically averaged over thermal motion. For the purpose of comparing the results from 3D-DDF with that from the experiment, it will be best if we can do it directly by defining a measure for “distance” between the two distributions. However, such a methodology has not been well developed yet. Moreover, general measure for 3D-DDF is available neither in literature nor in the Protein Data Bank. Instead, the most probable binding mode of a ligand is presented, which of course depends on a way of analyzing the 3D-DDF data.

Considering such a situation, we are developing two methods which provide us information concerning binding mode of ligand inside a protein. One of those methods is based on the radial distribution functions (RDF) of ligand atoms from atoms in amino acid residues of protein. A position of the first peak in RDF corresponds to an average distance between a pair of atoms of a ligand and a protein. By making the analyses of distances among the several atoms in ligand and protein, we can extract the binding mode of a ligand inside of protein. When the size of ligand is small enough, such as carbon monoxide and ethanol, the method works quite well to determine the binding mode.

The other method we are developing to abstract the binding mode of a ligand inside a protein from 3D-DDF is similar

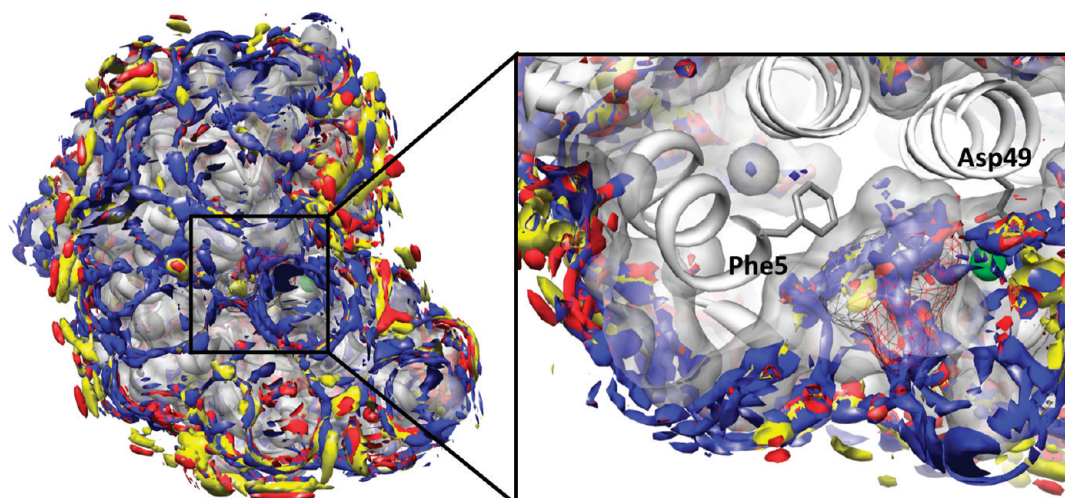
essentially to that used in the analysis of 3D-DDF from the diffraction measurement. The method defines two “score” functions, one corresponding to the position of ligand and the other to its orientation, in terms of 3D-DDF and of trial geometry of a ligand. The level of agreement between a trial geometry and 3D-DDF is ranked according to the score functions.

**2.4. Computational Detail.** The approach proposed in this article overcomes the difficulty associated with the uv-3D-RISM approach, and it will provide a new tool for the rational drug design. Here, we demonstrate robustness and capability of the uu-3D-RISM theory by applying the approach to two systems which are of great interest in biochemistry and pharmacology. One is an odorant binding protein known as LUSH, and the other is phospholipase A2. Table 1 shows the outline of these works.

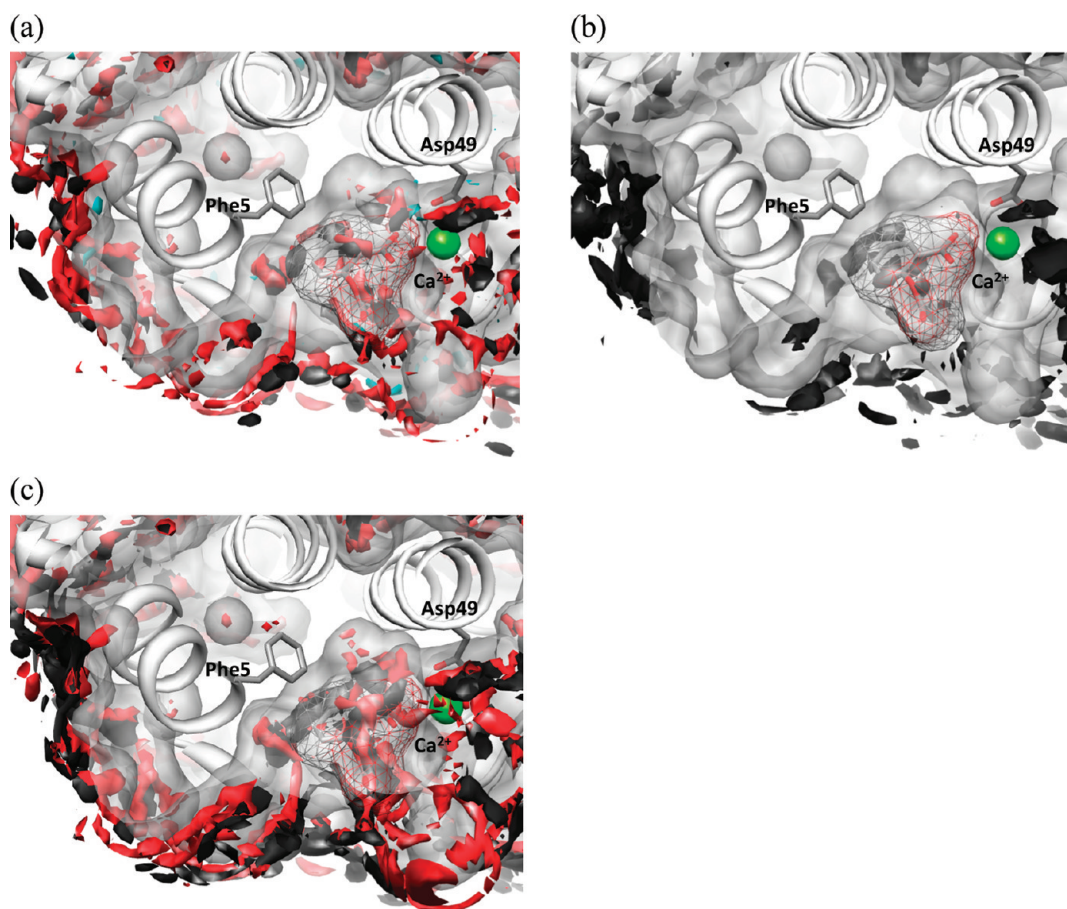
**2.4.1. Odorant Binding Protein (PDB ID: 1OOF).** In order to examine robustness of the new approach, we consider binding of an ethanol molecule to an odorant binding protein (LUSH), see Figure 1. The ligand is small enough to be treated with the uv-3D-RISM, so that one can compare the results from the methods with that from the uu-3D-RISM theory. In the case of uv-3D-RISM, the ligand is regarded as a component of solvent, while in the case of uu-3D-RISM, it is considered as a solute. In both cases, the odorant binding protein is treated as a solute receptor protein.

The Amber-99 parameter set<sup>40</sup> was employed for the protein, and the general amber force field (GAFF)<sup>41</sup> was employed for the ligand ethanol and for the acetic acids, which is part of receptor protein. TIP3P water<sup>42</sup> was chosen as solvent at 298 K and 0.9979 g/cm<sup>3</sup>. The uu-3D- and uv-3D-RISM equations were solved on a grid of 160<sup>3</sup> points in a cubic supercell of 80 Å<sup>3</sup>. The density of ethanol was so chosen that the volume ratio of water to ethanol becomes 99:1%.

**2.4.2. Phospholipase A2 (PDB ID: 1OXR).** In order to demonstrate the capability of the new approach, we examine MR of aspirin by PLA2 (Figure 2). Since aspirin is a rather large ligand, having 14 specific interaction sites (Figure 3) in the neutral state, it may not be treated readily with the ordinary uv-3D-RISM due to the difficulty stated above. So, this is a good example to demonstrate capability of the new method. The Amber-99 parameter set was employed for the protein, while GAFF was



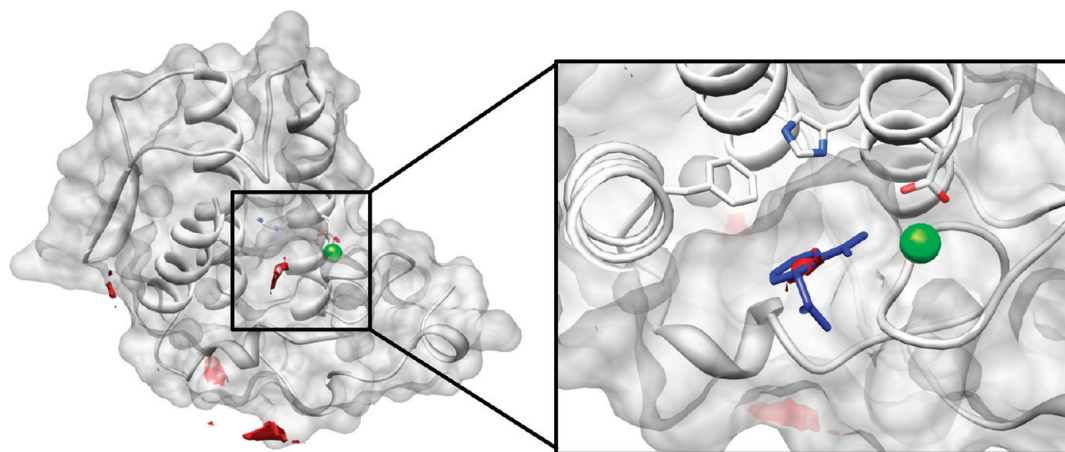
**Figure 6.** The 3D-DDF of protonated (neutral) aspirin around and inside phospholipase A2, obtained by uu-3D-RISM with the threshold  $g_{\gamma}(\mathbf{r}) > 2$ : red, COOH; yellow, aromatic ring; and blue, OCOCH<sub>3</sub>. The protein surfaces are represented as a gray transparent surface. The location of aspirin in X-ray structure is depicted with a wire frame.



**Figure 7.** The 3D-DDF of neutral aspirin around and inside phospholipase A2, obtained by uu-3D-RISM with the threshold  $g_{\gamma}(\mathbf{r}) > 2$ . (a) Carboxyl group, COOH; (b) aromatic ring; and (c) acetoxy group, OCOCH<sub>3</sub>. The color code is assigned to oxygen (red), carbon (black), and hydrogen (cyan). The protein surfaces are represented as a gray transparent surface.

employed for aspirin with a united-atom modification concerning hydrogen atoms. TIP3P water was chosen as solvent at 298 K

and 0.9979 g/cm<sup>3</sup>. The 3D-RISM equation was solved on a grid of 200<sup>3</sup> points in a cubic supercell of 100 Å<sup>3</sup>.



**Figure 8.** Affinity of an aspirin molecule to binding site in the phospholipase A2 estimated by the function of DC based on local PMF with the threshold  $f_{DC}(\mathbf{x}) > 1.45$ . The protein surfaces are represented as a gray transparent surface. In the top view, the location of aspirin in X-ray structure is depicted with blue sticks.

### 3. RESULTS AND DISCUSSION

**3.1. Odorant binding protein, LUSH.** *3.1.1. 3D-Distribution Functions of Ethanol around and inside LUSH.* In Figure 4, the 3D-DDFs of ethanol obtained by uu-3D- and uv-3D-RISM are compared. The 3D-DDFs are depicted by isosurface representation with the threshold  $g_\gamma(\mathbf{r}) > 2$ . This threshold implies that the probability of finding site  $\gamma$  at the position  $\mathbf{r}$  is twice as large as that in the bulk. The gray surface represents the protein, whereas the blue, green, red, and yellow surfaces depict the distribution of  $\text{CH}_3$ ,  $\text{CH}_2$ , O, and H sites of alcohol, respectively. At a glance, DDF from uu-3D-RISM, which is depicted in Figure 4a, shows good agreement with that from uv-3D-RISM shown in Figure 4b, especially nearby the binding site. Both 3D-DDFs are also in accord with the results from X-ray crystallography.

It is clear from the formulation described in the previous section that the uu-3D-DDF is equivalent to uv-3D-DDF in the low density limit of ligand concentration. The difference of these two 3D-DDFs are measured by the root-mean-square deviation (rmsd)  $d_{\text{rmsd}} = (\sum_{i=1}^n (g_\gamma^{\text{uu}}(x_i) - g_\gamma^{\text{uv}}(x_i))/n)^{1/2}$ ,  $x_i$  denotes each grid point);  $\text{CH}_3$  site,  $d_{\text{max}} = 0.0764$  and  $d_{\text{rmsd}} = 0.0055$ ;  $\text{CH}_2$  site,  $d_{\text{max}} = 0.0433$  and  $d_{\text{rmsd}} = 0.0066$ ; O site,  $d_{\text{max}} = 0.1015$  and  $d_{\text{rmsd}} = 0.0099$ ; and H site,  $d_{\text{max}} = 0.1362$  and  $d_{\text{rmsd}} = 0.0079$ . These differences are small enough for an actual application to systems in which the ligand concentration is low. In reality, the difference is well within a thermal fluctuation of the atomic position of ligands at a binding site in protein. Therefore, the uu-3D-DDF evaluated by uu-3D-RISM can be employed to molecular recognition problem instead of uv-3D-DDF for ligands in a low-concentration region.

*3.1.2. Radial Distribution Functions of Ethanol from Hydroxyl Groups in a Binding Site.* The RDFs between atoms in a ligand and those belonging to amino acids in protein allow us to investigate binding modes, position and orientation, of ligands inside protein. The RDFs can be obtained by averaging the 3D-DDF over the direction around a specified center:

$$g_a^{\text{1D}}(r, \mathbf{r}_0) = \frac{1}{4\pi} \int g_a(\mathbf{r}_0 + \mathbf{r}) d\hat{\mathbf{r}} \quad (14)$$

where  $\hat{\mathbf{r}}$  is the direction of  $\mathbf{r}$ , and  $\mathbf{r}_0$  indicates a center for averaging. The averaging centers were selected near the binding

**Table 3.** Top Three Peaks List of Affinity Based on the Overlap between the Ligand and Its Distribution Function

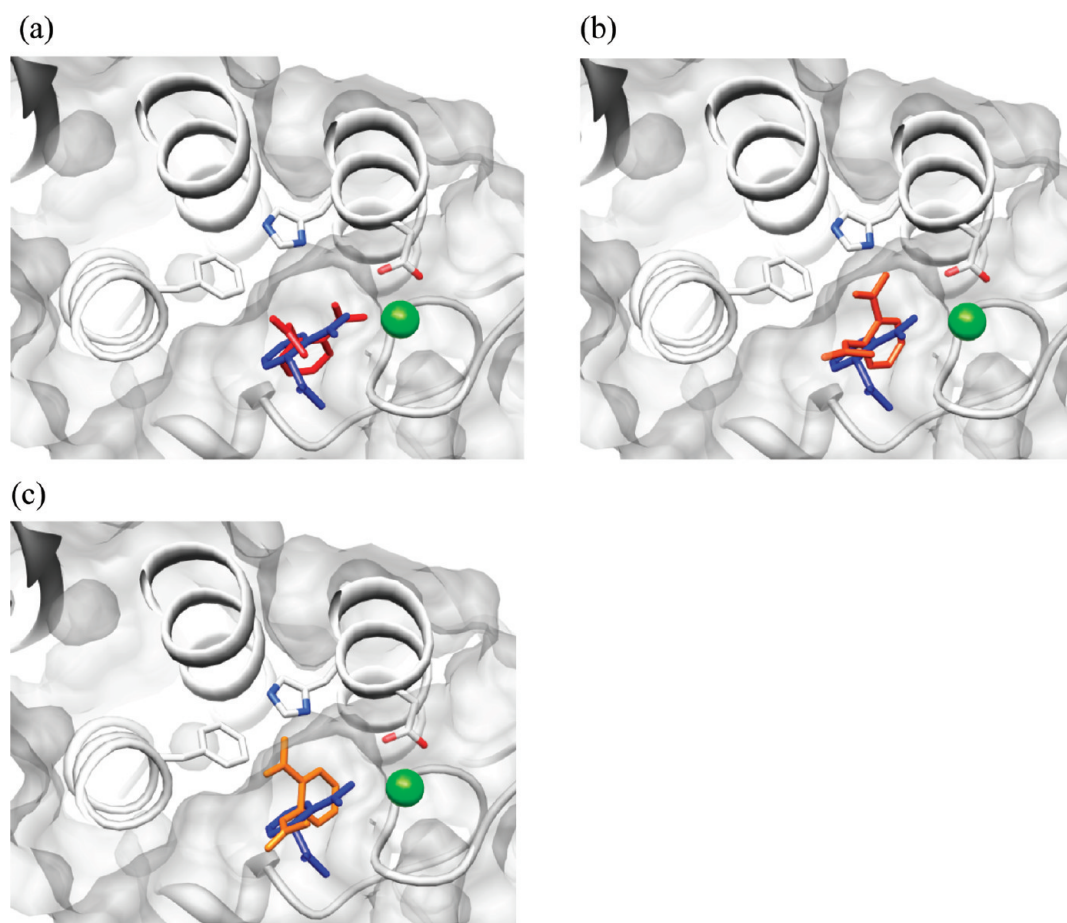
	order	$\alpha$	$\beta$	$\gamma$	overlap
(a)	1	10	330	50	$2.33 \times 10^5$
(b)	2	40	260	40	$2.02 \times 10^5$
(c)	3	60	220	50	$1.83 \times 10^5$

sites. In the present case, we choose the oxygen atoms of hydroxyl groups in Ser52 and Thr57 as an averaging center. The RDFs around these sites are shown in Figure 5. For comparing the peak positions of RDFs with those from the X-ray structure, we refer to the distance between these specific sites and each site of an ethanol molecule obtained from the experiment (Table 2). The distances experimentally determined are also marked by bars in the  $x$ -axis of the Figure 5a and 5b. Although each peak of RDF is not very sharp reflecting thermal fluctuation of the ligand inside the binding site, the positions of peaks in RDFs are consistent with those deduced from the X-ray crystallography.

The RDFs of hydrogen atoms, which are not treated by the X-ray diffraction measurement, are also depicted in Figure 5. The results may provide an orientation of hydrogen atom of the ligand molecule. In the case of Ser52, any peak indicating such orientation does not appear between two oxygen atoms, while a discernible peak appears at  $r \sim 1.6 \text{ \AA}$  in the case of Thr57. It is not clear at this moment whether the peak indicates the existence of a hydrogen bond or not. However, it is clear that the hydroxyl group is oriented toward Thr57, which is also consistent with the result from the X-ray crystallography.

**3.2. Phospholipase A2, PLA2.** It is not a straightforward task to apply the uv-3D-RISM method to the problem due to the reason described in detail in the previous section: The ligand molecule or aspirin is too large to get a convergent result of the vv-RISM equation for the solvent mixture, including the ligand molecules. Here, we only apply the uu-3D-RISM to the molecular recognition of aspirin to PLA2.

Aspirin, acetylsalicylic acid, is a weak acid in aqueous solution. We employed a neutral state, which was shown in Figure 3, because the affinity of the neutral state to binding site is much higher than a charged state.



**Figure 9.** Predicted structures of ligand from the top three peaks of affinity based on the overlap between the ligand and its distribution function (Table 3). The protein surfaces are represented as a gray transparent surface. The location of aspirin in X-ray structure is depicted with blue sticks.

**3.2.1. 3D-DDF of Aspirin Around and Inside Phospholipase A2.** The 3D-DDF of aspirin at the nonprotonated state is shown in Figure 6 with the threshold  $g_{\gamma}(\mathbf{r}) > 2$ . The red, yellow, and blue surfaces are the distributions of carboxyl group, COOH (number of sites is 4), the aromatic ring (number of sites is 6), and the acetoxy group, OCOCH<sub>3</sub> (number of sites is 4), respectively. We can easily observe these distributions not only inside the binding site but also around the protein. Since the distributions shown in Figure 6 are jumbled inside the binding site, the contribution from each site group, COOH, aromatic ring, and OCOCH<sub>3</sub>, is separately depicted in Figure 7. The color code assigned to each atom in the figures is as follows: red, oxygen; black, carbon; and blue, hydrogen.

As you can see in Figure 7a and c, the carboxyl and acetoxy groups are widely distributed inside and around the binding site region. On the other hand, the distribution of the aromatic ring is seen only inside the binding site and is apparently accommodated well within the pocket. Indeed, the binding pocket is formed by hydrophobic amino acid residues, like leucine, phenylalanine, and so on. In that sense, this result represents the case in which the hydrophobic effect makes essential contributions to the molecular recognition for the aromatic ligand.

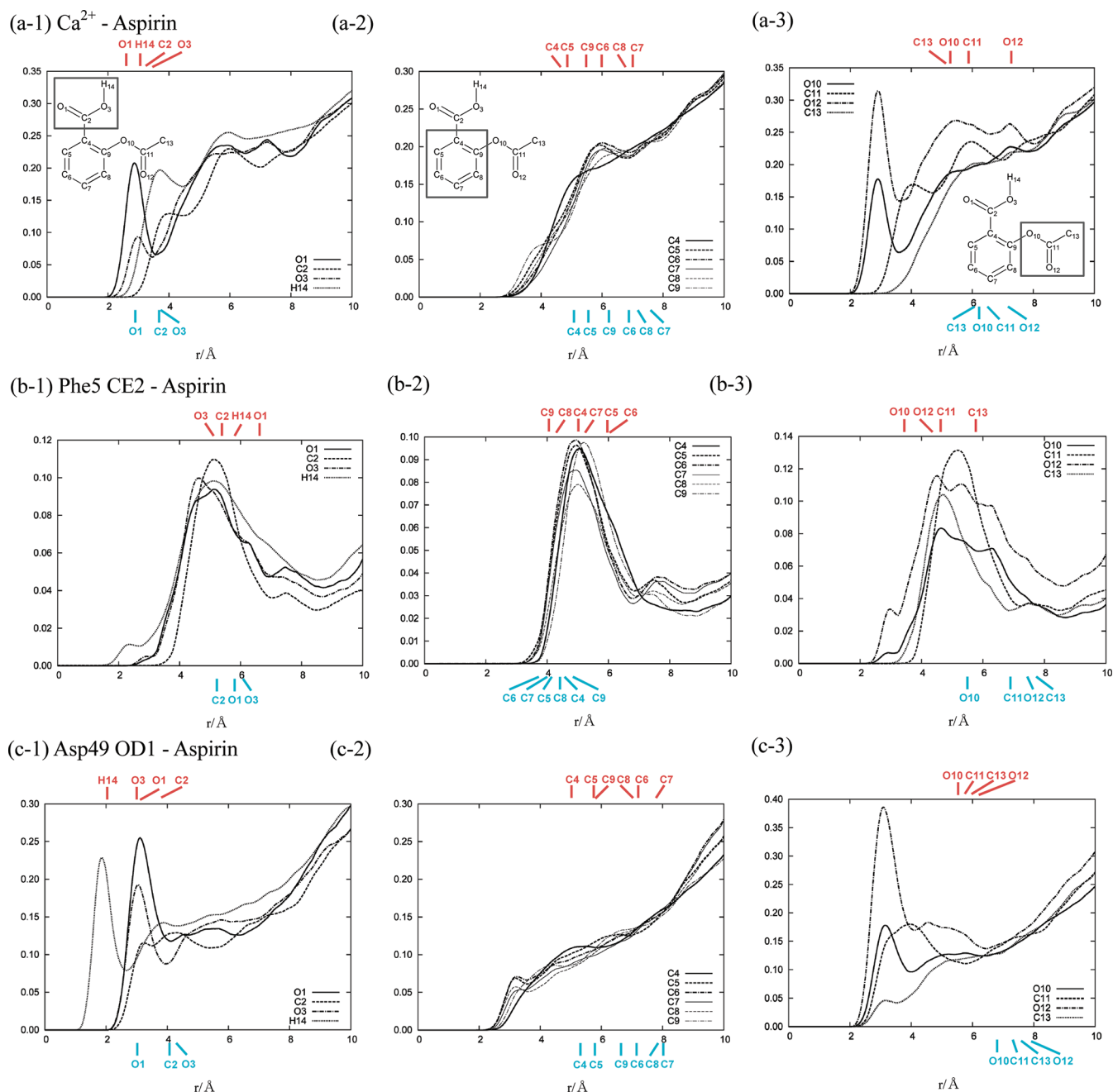
**3.2.2. Binding Mode of an Aspirin Molecule to Binding Site in the Phospholipase A2.** Since the final goal of our study is to establish a method to probe a large ligand molecule recognized by protein using 3D-RISM, we must determine the location of

binding site which has the highest affinity of ligand in target protein. The 3D-DDFs or the potential of mean force (PMF) are a good indicator to evaluate the affinity. Those have been successfully applied to measure the affinity or selectivity of solvent in protein.<sup>12–14</sup> However, since 3D-DDF is the distribution function of an individual site consisting a ligand molecule, it is difficult to evaluate the affinity of whole ligand molecule directly.

In this section, we introduce a function for distribution center (DC) of ligand to measure the affinity of ligand molecule. The function of DC is defined as

$$f_{\text{DC}}(\mathbf{x}) = \begin{cases} \left( \frac{N}{\prod_{\gamma} V_{\text{box}} - V_{\text{protein}}(\mathbf{x})} \int_{V_{\text{box}}(\mathbf{x})} g_{\gamma}(\mathbf{r}) d\mathbf{r} \right)^{1/N} & \text{for } V_{\text{box}} - V_{\text{protein}}(\mathbf{x}) \geq V_{\text{ligand}} \\ 0 & \text{for } V_{\text{box}} - V_{\text{protein}}(\mathbf{x}) < V_{\text{ligand}} \end{cases} \quad (15)$$

where  $\mathbf{x}$  denotes the center of box,  $N$  is the total number of sites of ligand molecule for normalization,  $V_{\text{box}}$  is the volume of the box, and  $V_{\text{protein}}(\mathbf{x})$  is the excluding volume of the solute protein in the box. Therefore,  $V_{\text{box}} - V_{\text{protein}}(\mathbf{x})$  denotes the space where ligand can be distributed. Note that the integrations in right-hand side of eq 15 are only performed inside  $V_{\text{box}}$  centered at  $\mathbf{x}$ . The size of box is adjusted to the length of a ligand molecule. The uu-3D-DDF is integrated in the box, and the result is projected to



**Figure 10.** RDFs of aspirin around (a) calcium ion, (b) Phe5, and (c) Asp49, respectively. In part b, CE2 atom was chosen as the averaging center. In part c, OD1 atom was chosen as the averaging center. The distances between each specific site and the atom of predicted aspirin as top peak are marked by upper bars in the  $x$ -axis (red indices). The distances between the sites and the atom of aspirin in X-ray structure are also marked by lower bars (blue indices).

the center of box. If the value is larger than one, the probability of finding “an aspirin molecule” at the position is higher than bulk. Although the DC only gives us the rough estimate of the location of binding site, it is helpful to guide a further analysis concerning the binding mode in more detail using information of RDFs, which will be discussed later.

We preformed the calculation of DC based on eq 15 in order to estimate the affinity of ligand to the binding site. The result of DC is shown in Figure 8 with the threshold  $f > 1.45$ . The maximum value is 1.56. The DC function does not take quite a high value because it is averaged over the sites and volume.

Note that the result of DC is projected onto the center of the calculated box. In Figure 8, we observe the highest peak at the center of the binding site, which is determined by the X-ray crystallography. The results demonstrate that the new method is capable of locating the binding site in the protein and the affinity of a ligand to the site properly.

In this paper, although we focused the highest peak on the binding site, there is another important finding in Figure 8. One may identify a long distribution stretched from the binding site to the bulk. This peak may imply the pathway through which the ligand is entering and escaping. Other peaks observed around the



surface of protein are of interest, because those peaks may be related to the escaping and entering pathways of ligand in diffusive motion. Further analysis of ligand distributions around protein surface will be presented in future.

In order to understand a mechanism of molecular recognition process, it is important to determine an explicit structure of the ligand inside the specific binding site, because we need the distinct structure to calculate some physiological properties, like free energy. Actually, we can also investigate the orientation of the ligand by calculating the overlap between the structure of the ligand and the 3D distribution functions which are obtained by uu-3D-RISM. We define the target function for the orientation by the following equation:

$$f_{\text{ori}}(\mathbf{x}, \Omega) = \prod_{\gamma}^N g_{\gamma}(\mathbf{x} + \mathbf{l}_{\gamma} \cdot \hat{\mathbf{R}}(\Omega)) \quad (16)$$

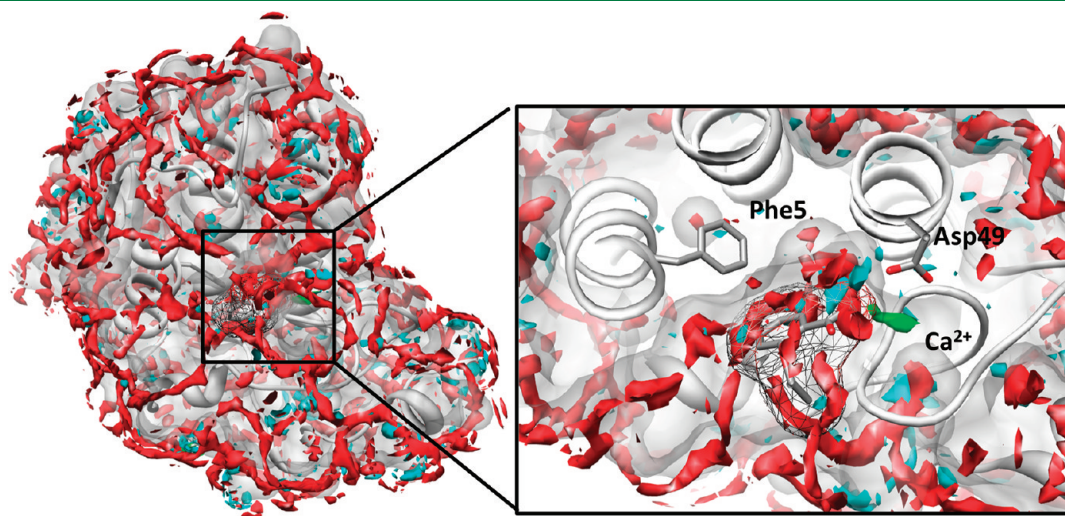
where  $\mathbf{x}$  denotes center of box which is obtained by eq 15,  $\mathbf{l}_{\gamma}$  denotes the internal coordinate of site  $\gamma$ , and  $\hat{\mathbf{R}}(\Omega)$  denotes

**Table 4. Distance (Å) between the Specific Sites of PLA2 and the Atoms of Aspirin**

sites	Ca <sup>2+</sup>		Phe5–CE2		Asp49–OD1	
	theory	exptl.	theory	exptl.	theory	exptl.
O1	2.4	2.8	6.7	5.8	3.1	3.0
C2	3.3	3.7	5.5	5.3	3.5	4.0
O3	3.4	3.7	5.2	6.0	3.0	4.2
C4	4.4	5.0	5.1	4.5	5.0	5.4
C5	4.9	5.5	6.0	4.1	5.9	5.8
C6	6.2	6.9	6.0	3.9	7.2	7.1
C7	7.0	7.7	5.3	4.0	7.8	8.0
C8	6.8	7.4	4.2	4.3	7.2	7.9
C9	5.6	6.2	4.0	4.5	5.9	6.6
O10	5.6	6.2	3.5	5.4	5.5	6.7
C11	5.9	6.3	4.4	6.9	5.8	7.3
O12	6.8	7.1	4.3	7.6	6.3	8.0
C13	5.6	6.1	5.9	7.8	6.1	7.6

rotational matrix with the Euler angles to search the entire orientational space. As we mentioned above, the ligand molecule just fits the size of box. It means that the center of the box coincides with the center of the molecule, which is the origin of the rotational matrix at the same spatial point. The advantage of this approach is that we can measure the affinity quantitatively as the degree of overlap. The results of this searching are summarized in Table 3 and Figure 9. In Table 3, the top three orientations, ranked based on the degree of overlap, are listed in terms of the Euler angles,  $\alpha$ ,  $\beta$ , and  $\gamma$ . Figure 9 shows the explicit ligand structures corresponding to the top three orientations listed in Table 3. The location and orientation of aspirin in the X-ray structure are depicted with blue sticks. It is worthwhile to note that the structure corresponding to one with the greatest overlaps has the same orientation with the X-ray structure, concerning the molecular axis aligning C<sub>2</sub>, C<sub>4</sub>, and C<sub>7</sub> atoms, although the rotational angle around the axis is somewhat different from each other. This orientation seems to be induced by the calcium ion located at the binding site, since the carboxylic group of aspirin faces to the calcium ion. In case of the other two structures with lower score, the carboxylic groups are facing toward His48, which is positively polarized as well.

**3.2.3. RDFs of Aspirin around Specific Sites Inside or around Binding Sites of PLA2.** In order to find the orientation of aspirin in the binding pocket, we examine the RDFs of each site of aspirin using eq 14. Three specific sites of the residues around the binding pocket are chosen as the averaging centers in order to calculate RDFs. These are the calcium ion, the CE2 atom in Phe5, and the OD2 atom in Asp49. The reasons why those atoms are chosen as the averaging centers are because the calcium ion and the OD2 atom in Asp49 help aspirin to bind in the pocket through the carboxyl group and because the CE2 atom exists in the side opposite to the calcium ion across the pocket. The RDFs are shown in Figure 10. For the purpose of comparing the peaks of RDFs with the corresponding information from the X-ray structure, the distances between these specific sites of the amino acid residue and each atomic site of the aspirin molecule, determined by the X-ray crystallography, are summarized in Table 4. The distances are marked by bars (blue) in the  $x$ -axis of the Figure 10. The distances corresponding to the structure



**Figure 11.** The 3D-DDFs of water around and inside phospholipase A2 are obtained by uv-3D-RISM with the threshold  $g_{\gamma}(\mathbf{r}) > 3$ : red, oxygen atom of water; and cyan, hydrogen atom of water. The 3D-DDF of calcium ion is also obtained with the threshold  $g_{\text{Ca}^{2+}}(\mathbf{r}) > 40$  as a green spot. The protein surfaces are represented as a gray transparent surface. The location of aspirin in X-ray structure is depicted with a wire frame.

(position and orientation) deduced from the spatial distribution functions due to uu-3D-RISM, with the highest score (Figure 9a), are also marked with red in the same figures. The position of each peak in RDFs is in general consistent with corresponding distance from the X-ray diffraction as well as from the 3D-RISM, except for the distinct peaks in those corresponding to O12 and O10 in Figure 10 (a-3) and to O10 and O12 in (c-3). Those peaks are assigned to distributions of the corresponding atoms existing outside the binding site, and they are irrelevant to the ligand bound in the active site.

Especially interesting among RDFs is the carboxyl group around the calcium ion Figure 10 (a-1) and the aromatic ring around Phe5 (b-2). The peak positions of the RDFs coincide well with those deduced from the orientation of ligands determined both by the experiment and from the analysis of our spatial distribution functions. These suggest strongly the importance of a role played by the calcium ion in recognizing aspirin in the active site. We have confirmed the distinct binding of a calcium ion ( $g_{Ca^{2+}}(r) > 40$ ) at the binding site by means of the 3D-RISM calculation in Figure 11.

The RDFs of carboxyl group around Asp49, shown in Figure 10 (c-1), are worthwhile to draw special attention, since they are suggestive of a mechanism concerning the recognition of aspirin by the protein. According to the results, the ligand molecule is forming a hydrogen bond with one of the carboxylic oxygen atoms of Asp49 through its carboxylic hydrogen atom; note the sharp peak around  $r = 1.8 \text{ \AA}$  in the RDF of hydrogen (H14). Apparently, the carboxylic oxygen of Asp49 was supposed to make a hydrogen bond with solvent water, if the position was not invaded by the ligand.

Figure 11 shows the 3D-DDFs of water molecule at the binding site without the ligand. The region of binding site is constructed by hydrophobic residues, like phenylalanine, leucine, isoleucine, and so on. However, water molecules can be bound with main chain or with a charged residue, such as Asp49, through hydrogen bonds. Water molecules are apparently making a hydrogen-bond network or train inside the binding site. It suggests that the dehydration penalty will be extremely high when a ligand replaces those water molecules, and ordinary docking algorithms might not be able to find the binding site.<sup>31</sup>

So, we can draw a hypothetical scheme concerning the recognition mechanism of aspirin to PLA2. The recognition process is largely motivated by the calcium ion, which was already bound at the active site before any aspirin is put in the solution. The recognition process is initiated first by the coulomb interaction between the calcium ion and the carbonyl–oxygen of the carboxyl group, which is followed by formation of the hydrogen bond between the carboxyl–oxygen of Asp49 and the carboxyl–hydrogen of the ligand. In the latter process, a water molecule which was hydrogen bonded to the carboxyl–oxygen, prior to the ligand invasion, is excluded from the binding site. The ligand is further stabilized by the hydrophobic interaction between the phenyl group of the ligand and the hydrophobic residues consisting of the other side of the binding site.

#### 4. CONCLUDING REMARK

We proposed a new approach, the uu-3D-RISM theory, to investigate the molecular recognition in biological system. A motivation to develop the new approach was that the ordinary RISM/3D-RISM approach has difficulty in solving the solvent–solvent RISM equation involving large ligand molecules, which of course have vital importance in the rational drug design. The uu-

3D-RISM is formulated from the general equation of molecular Ornstein–Zernike by considering both a receptor and a ligand as “solutes” immersed in solvent at the infinite dilution limit.

In order to confirm the robustness of the new approach, we calculated the spatial distribution of ethanol at the active site of an odorant binding protein, LUSH, based on the two methods, the ordinary RISM/3D-RISM theory and the uu-3D-RISM, since an ethanol molecule is small enough to be handled with the old method. The new approach reproduced the results from the old method, with subtle difference expected from the discrepancy in the concentration of ligand: one in a finite concentration and the other in the infinite dilution. The analysis based on the radial distribution function (RDF) indicates that the position and the orientation of the ligand inside the binding pocket are consistent with those from the experimental results due to the X-ray crystallography. Robustness of the new approach was thus verified.

We then applied the new approach to an aspirin binding protein, phospholipase A2 (PLA2), with aspirin as a ligand. The process may not be tractable by the old method due to the reason stated above. Since the size of aspirin is much larger and more complex than previous application, or ethanol, analyzing the spatial distribution (uu-DDF) of the ligand inside the binding site, obtained from uu-3D-RISM, is not a trivial problem anymore. So, we developed a new approach to analyze uu-DF, defining a new function referred to as “distribution center (DC),” which locates the center of the most probable distribution of ligand. The position and orientation of aspirin inside the binding site of PLA2 were determined from DC and RDFs of atomic sites of the ligand around particular residues consisting the binding pocket. The binding configuration of the ligand inside the pocket was in fair agreement with that determined from the X-ray crystallography. We will report details about analyses for the origin of binding affinity in a following paper.

The second application of the uu-3D-RISM method clearly demonstrates that the theory is a prospective tool for discovering or designing a new drug, because aspirin itself is already one of the most popular drugs in the market.

#### ■ AUTHOR INFORMATION

##### Corresponding Author

\*E-mail: hirata@ims.ac.jp. Telephone: +81(Japan)-564-55-7314.

#### ■ ACKNOWLEDGMENT

Authors are also grateful to Prof. Sato in Kyoto U. for his helpful discussion. This work is supported by Grants-in-Aid for Scientific Research on Innovate Areas of Molecular Science of Fluctuations toward Biological Functions from the MEXT in Japan. Authors are also supported by the Next Generation Super Computing Project, Nanoscience Program, a Grant-in-Aid for Scientific Research. N.Y. is grateful to a Grant-in-Aid for Young Scientists. Y.K. is grateful to the support by the grant from the Japan Society for the Promotion of Science (JSPS) and the fellow. Molecular graphics images were produced using the UCSF Chimera<sup>43</sup> package and gOpenMol.<sup>44,45</sup> A part of this work was carried out by super computer in the Research Center for Computational Science (RCCS), Okazaki Research Facilities, National Institutes of Natural Sciences (NINS). The uu-3D-RISM program has been implemented to RISM/3D-RISM solver developed by our group in IMS.<sup>46,47</sup>

## REFERENCES

- (1) Jorgensen, W. L. *Science* **2004**, *303*, 1813–1818.
- (2) Kitchen, D. B.; Decornez, H.; Furr, J. R.; Bajorath, J. *Nature Reviews* **2004**, *3*, 935–949.
- (3) Karplus, M.; McCammon, A. *Nat. Struct. Biol.* **2002**, *9*, 646–652.
- (4) Zwier, M. C.; Chong, L. T. *Curr. Opin Pharmacol.* **2010**, *10*, 1–8.
- (5) Essex, J. W.; Severance, D. L.; Tirado-Rives, J.; Jorgensen, W. L. *J. Phys. Chem. B* **1997**, *101*, 9663–9666.
- (6) Chang, M. *Monte Carlo simulation for the pharmaceutical industry: concepts, algorithms, and case studies*; CRC Press: Boca Raton, FL, 2010.
- (7) Raha, K.; Peters, M. B.; Wang, B.; Yu, N.; Wollacott, A. M.; Westerhoff, L. M.; Merz, K. M., Jr. *Drug Discovery Today*. **2007**, *12*, 725–731.
- (8) Imai, T.; Hiraoka, R.; Kovalenko, A.; Hirata, F. *J. Am. Chem. Soc. Commun.* **2005**, *127*, 15334–15335.
- (9) Imai, T.; Hiraoka, R.; Seto, T.; Kovalenko, A.; Hirata, F. *J. Phys. Chem. B* **2007**, *111*, 11585–11591.
- (10) Yoshida, N.; Phongphanphanee, S.; Maruyama, Y.; Imai, T.; Hirata, F. *J. Am. Chem. Soc. Commun.* **2006**, *273*, 12042–12043.
- (11) Yoshida, N.; Phongphanphanee, S.; Hirata, F. *J. Phys. Chem. B* **2007**, *111*, 4588–4595.
- (12) Phongphanphanee, S.; Yoshida, N.; Hirata, F. *Chem. Phys. Lett.* **2007**, *449*, 196–201.
- (13) Phongphanphanee, S.; Yoshida, N.; Hirata, F. *J. Am. Chem. Soc. Commun.* **2008**, *130*, 1540–1541.
- (14) Phongphanphanee, S.; Yoshida, N.; Hirata, F. *J. Phys. Chem. B* **2010**, *114*, 7967–7973.
- (15) Phongphanphanee, S.; Rungrotmongkol, T.; Yoshida, N.; Hannongbua, S.; Hirata, F. *J. Am. Chem. Soc.* **2010**, *132*, 9782–9788.
- (16) Kiyota, Y.; Hiraoka, R.; Yoshida, N.; Maruyama, Y.; Imai, T.; Hirata, F. *J. Am. Chem. Soc. Commun.* **2009**, *131*, 3852–3853.
- (17) Kiyota, Y.; Yoshida, N.; Hirata, F. *J. Mol. Liq.* **2011**, *159*, 93–98.
- (18) Imai, T.; Oda, K.; Kovalenko, A.; Hirata, F.; Kidera, A. *J. Am. Chem. Soc.* **2009**, *131*, 12430–12440.
- (19) Roux, B.; Yu, H. A.; Karplus, M. *J. Phys. Chem.* **1990**, *94*, 4683–4688.
- (20) Hirata, F. *Molecular Theory of Solvation*; Kluwer: Dordrecht, The Netherlands, 2003.
- (21) Kovalenko, A.; Hirata, F. *J. Phys. Chem. B* **1999**, *103*, 7942–7957.
- (22) Yoshida, N.; Imai, T.; Phongphanphanee, S.; Kovalenko, A.; Hirata, F. *J. Phys. Chem. B* **2009**, *113*, 873–886.
- (23) Kovalenko, A.; Hirata, F. *J. Chem. Phys.* **2000**, *112*, 10391–10402.
- (24) Minezawa, N.; Kato, S. *J. Chem. Phys.* **2007**, *126*, 054511 (15 pages).
- (25) Kovalenko, A.; Ten-No, S.; Hirata, F. *J. Comput. Chem.* **1999**, *20*, 928–936.
- (26) Kinoshita, M.; Okamoto, Y.; Hirata, F. *J. Chem. Phys.* **1997**, *107*, 1586–1599.
- (27) Kruse, S. W.; Zhao, R.; Smith, D. P.; Jones, D. N. M. *Nat. Struct. Biol.* **2003**, *10*, 694–700.
- (28) Dennis, E. A. *J. Biol. Chem.* **1994**, *269*, 13057–13060.
- (29) Nicolas, J. P.; Lin, Y.; Lambeau, G.; Ghomashchi, F.; Lazdunski, M.; Gelb, M. H. *J. Biol. Chem.* **1997**, *272*, 7173–7181.
- (30) Argiolas, A.; Pisano, J. J. *J. Biol. Chem.* **1983**, *258*, 13697–13702.
- (31) Chang, D. T.; Oyang, Y.; Lin, J. *Nuc. Acids Res.* **2005**, *33*, W233–W238.
- (32) Hansen, J. P.; McDonald, I. R. *Theory of Simple Liquids*, 3rd ed.; Academic: London, 2006.
- (33) Chandler, D.; Andersen, H. C. *J. Chem. Phys.* **1972**, *57*, 1930–1937.
- (34) Harano, Y.; Imai, T.; Kovalenko, A.; Kinoshita, M.; Hirata, F. *J. Chem. Phys.* **2001**, *114*, 9506–9511.
- (35) Chandler, D.; McCoy, D. J.; Singer, S. J. *J. Chem. Phys.* **1986**, *85*, 5971–5977.
- (36) Kovalenko, A.; Hirata, F. *Chem. Phys. Lett.* **1998**, *290*, 237–244.
- (37) Kovalenko, A.; Hirata, F. *J. Chem. Phys.* **1999**, *110*, 10095–10112.
- (38) Percus, J. K. *Phys. Rev. Lett.* **1962**, *8*, 462–463.
- (39) Kovalenko, A.; Hirata, F. *Chem. Phys. Lett.* **2001**, *346*, 496–502.
- (40) Wang, J. M.; Cieplak, P.; Kollman, P. A. *J. Comput. Chem.* **2000**, *21*, 1049–1074.
- (41) Wang, J.; Wolf, R. M.; Caldwell, J. W.; Kollman, P. A.; Case, D. A. *J. Comput. Chem.* **2004**, *25*, 1157–1174.
- (42) Jorgensen, W. L.; Chandrasekhar, J.; Madura, J. D.; Impey, R. W.; Klein, M. L. *J. Chem. Phys.* **1983**, *79*, 926–935.
- (43) Pettersen, E. F.; Goddard, T. D.; Huang, C. C.; Couch, G. S.; Greenblatt, D. M.; Meng, E. C.; Ferrin, T. E. *J. Comput. Chem.* **2004**, *25*, 1605–1612.
- (44) Laaksonen, L. *J. Mol. Graphics* **1992**, *10*, 33–34.
- (45) Bergman, D. L.; Laaksonen, L.; Laaksonen, A. *J. Mol. Graphics Modell.* **1997**, *15*, 301.
- (46) Yoshida, N.; Hirata, F. *J. Comput. Chem.* **2006**, *27*, 453–462.
- (47) Yoshida, N.; Kiyota, Y.; Hirata, F. *J. Mol. Liq.* **2011**, *159*, 83–92.

# Unraveling the Molecular Mechanism of Enthalpy Driven Peptide Folding by Polyol Osmolytes

Regina Gilman-Politi and Daniel Harries\*

Institute of Chemistry and The Fritz Haber Center, The Hebrew University, Jerusalem 91904, Israel

**S** Supporting Information

**ABSTRACT:** Many polyols and carbohydrates serve in different organisms as protective osmolytes that help to stabilize proteins in their native, functional state, even under a variety of environmental stresses. However, despite their important role, much of the molecular mechanism by which these osmolytes exert their action remains elusive. We have recently shown experimentally that, although polyols and carbohydrates are excluded from protein and peptide interfaces, as also expected for the known entropic “crowding” mechanism, the osmolyte folding action can in fact primarily be enthalpic in nature. To follow this newly resolved enthalpically driven stabilization mechanism, we report here on molecular dynamics simulations of a model peptide that can fold in solution into a  $\beta$ -hairpin. In agreement with experiments, our simulations indicate that sorbitol, a representative polyol, promotes peptide folding by preferential exclusion. At the molecular level, simulations further show that peptide stabilization can be explained by sorbitol’s perturbation of the solution hydrogen bonding network in the peptide first hydration shells. Consequently, fewer hydrogen bonds between peptide and solvating water are lost upon folding, and additional internal peptide hydrogen bonds are formed in the presence of sorbitol, while internal peptide and water-associated hydrogen bonds are strengthened, resulting in stabilization of the peptide folded state. We further find that changes in water orientational entropy are reduced upon folding in sorbitol solution, reflecting the struggle of water molecules to maintain optimal hydrogen bonding in the presence of competing polyols. By providing first molecular underpinnings for enthalpically driven osmolyte stabilization of peptides and proteins, this mechanism should allow a better understanding of the variety of physical forces by which protective osmolytes act in biologically realistic solutions.

## INTRODUCTION

Protein stability and activity sensitively depend on myriad modulators of environmental solvent conditions, including hydration levels, ion concentrations, and pH. In efforts to maintain protein function and integrity, one of the important ways that living organisms combat such environmental stresses involves the accumulation of molecularly small cosolutes termed osmolytes.<sup>1–3</sup> Naturally occurring osmolytes are cosolutes that can be typically grouped into three major classes: polyols, amino acids, and combinations of methylamines with urea.<sup>1</sup> Of these, the addition of “protective” osmolytes to protein solutions shifts the thermodynamic equilibrium of folding toward more compact, native states. While protein folding by osmolytes has been a subject of numerous investigations,<sup>4–13</sup> much of the underlying molecular mechanism remains unknown.

Protective osmolytes are generally excluded from protein–water interfaces, and it is this preferential exclusion from macromolecular surfaces that necessarily confers thermodynamic stability to proteins.<sup>10,14–20</sup> Molecular crowding due to excluded volume interactions has been widely invoked to explain how osmolytes can shift the folding equilibrium toward the more folded state.<sup>21–25</sup> According to this mechanism, the restriction of protein conformations to allow larger free volume for the added osmolytes destabilizes the unfolded state with respect to the native conformation.<sup>23,26</sup> Crowding has been useful in explaining the protein stabilizing effect of macromolecular solutes, such as polymers and other proteins, based on steric interactions that are entropic in nature.<sup>27–30</sup>

In contrast to entropically driven steric “crowding”, recent evidence indicates that, when molecularly small solutes are involved, the stabilizing mechanism may be enthalpically dominated.<sup>31,32</sup> For example, our recent experiments indicate that sugars and polyols can drive peptide folding primarily through diminishing the enthalpic loss involved in the folding process. Moreover, the added entropic contribution to folding wrought by osmolytes is negative, and disfavors folding. Interestingly, these effects are dependent on the molecular size of the osmolytes, as also expected for a crowding mechanism. These findings require new molecular mechanisms that can explain how osmolytes can confer a favorable enthalpic contribution to folding, while concurrently remaining preferentially excluded from the peptide–solution interface.

Here, we employ molecular dynamics (MD) simulations to explore the molecular origins of the stabilizing mechanism of protective osmolytes. We follow a model 16-residue peptide that has been shown experimentally to fold into a  $\beta$ -hairpin from a disordered state.<sup>32,33</sup> At room temperature and pH 7, about half of the peptide population is found in the folded state ( $\Delta G_{\text{fold}} \approx 0$ ). In the presence of polyols and carbohydrates, however, the equilibrium is shifted, and the peptide primarily adopts the folded state. Our strategy was to separately simulate the folded and unfolded states of the peptide in two different solvating solutions: pure water and aqueous solutions of osmolyte. We focus on

Received: June 30, 2011

Published: September 22, 2011

sorbitol as an important representative of polyol osmolytes; sorbitol is one of the largest polyols, is highly soluble, and has been experimentally shown to have one of the strongest effects on peptide folding in this group. Because the sorbitol and water translational and orientational relaxation around the peptide (on the order of hundreds of picoseconds) is much faster than the peptide folding–unfolding times (estimated at several hundreds of nanoseconds), we have been able to exploit this difference in time scales to separately average the properties of solution microstates around the folded and unfolded populations.

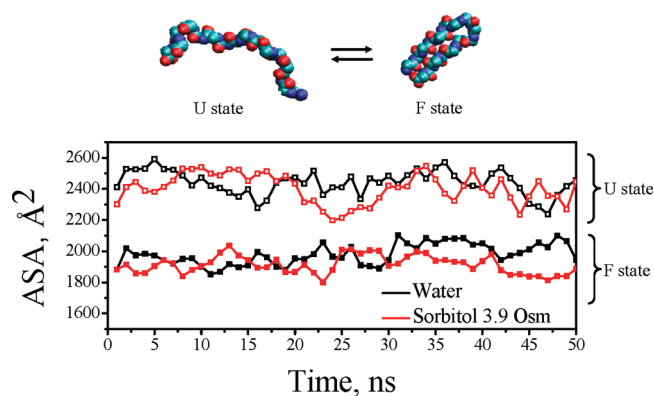
By analyzing the changes in the solvating environment of the folded and unfolded states in the presence of osmolytes, and comparing these to the differences in pure water found around both states, we characterized a previously unknown enthalpy driven mechanism for polyol stabilization of peptides. Our results reveal that the key driving forces for peptide stabilization by osmolytes involve the reduced loss upon folding of the number of hydrogen bonds located in the first and second solvation layers around the peptide, the increased strength of the hydrogen bonds that remain, and the larger number of internal peptide hydrogen bonds created in the presence of the polyol.

## RESULTS

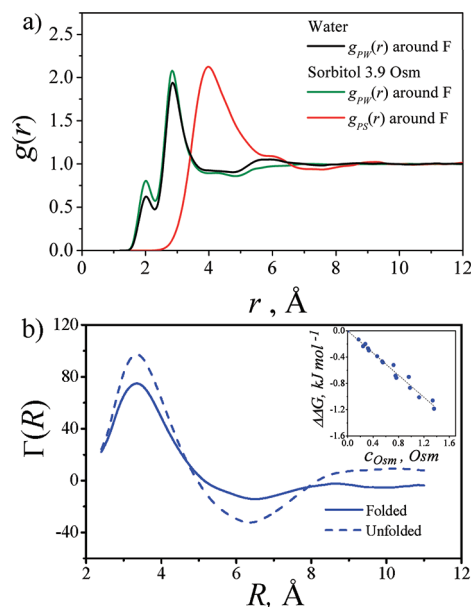
Changes in the model peptide's secondary structure in both water and sorbitol solution (at 3.9 Osm) can be followed by tracking the solvent accessible surface area (ASA) for the folded (F) and unfolded (U) states over the MD trajectory, as shown in Figure 1 and detailed in the Methods. The U state undergoes larger fluctuations in ASA than the F state (with a standard deviation of 78 vs 68 Å<sup>2</sup> in pure water, and 95 vs 59 Å<sup>2</sup> in sorbitol solution, respectively); however, on the basis of average ASA alone, there is no discernible difference between the U states in water and in the presence of sorbitol. Similarly, the F state also showed no significant difference in ASA for the first 30 ns simulated, but at longer times we find more compact structures in the presence of sorbitol than in pure water (smaller ASA by, on average, 100 Å<sup>2</sup> over the last 20 ns). While these observations already qualitatively match the known peptide stabilizing capacity of sorbitol, additional thermodynamic measures can provide a direct link between solvation thermodynamics and sorbitol-induced peptide stability, as we discuss in the following sections.

**Sorbitol is Preferentially Excluded from Peptide Interfaces.** In general, stabilizing cosolutes are found to be “excluded” from the macromolecular peptide interfaces and thereby increase the chemical potential of the macromolecule; in the limit of infinite peptide dilution, this necessarily also implies that the peptide is “preferentially hydrated”.<sup>19</sup> To assess osmolyte exclusion as well as the distance and size of different solvation layers around the peptide in the F and U states, we follow the structure of solution using the radial distribution function,  $g_{Px}(r)$ , representing the local densities with respect to the bulk of species  $x$  at a distance  $r$  from the peptide (P), where  $x$  represents the chemical species in solution: peptide (P), solute (S), or water (W) (see also Methods). Figure 2a shows  $g_{PW}(r)$  and  $g_{PS}(r)$  for the local densities of water's oxygen and sorbitol's center of mass, respectively, where  $r$  represents the shortest distance from any atom of the folded peptide states.

In contrast to the results often reported for proteins using MD simulations and neutron diffraction experiments,<sup>34,35</sup> the peptide–water distribution function  $g_{PW}(r)$  shows not one, but two prominent hydration peaks, Figure 2a. This apparent discrepancy



**Figure 1.** Solvent accessible surface area (ASA) calculated over 50 ns of MD simulation for the peptide folded (F) and unfolded (U) states in pure water (black line) and in the presence of sorbitol aqueous solutions at a concentration of 3.9 Osm (red line). Each data point represents an average ASA taken over the course of 1 ns. The upper panel shows a schematic of the F and U states.



**Figure 2.** Solution structure around the folded state and preferential interaction coefficients for the folded and unfolded states. (a) Radial distribution function,  $g(r)$ , of water oxygen's and sorbitol's centers of mass around the peptide. (b) Preferential interaction coefficient,  $\Gamma$ , of water near the peptide. The inset in b shows experimental data for the folding free energy versus sorbitol's osmolyte concentration at  $T = 298$  K, reproduced from ref 32. Both  $g(r)$  and  $\Gamma(R)$  are plotted versus the distance from any peptide atom.

results from the way we calculate  $g(r)$ : here, distances are measured with respect to each peptide atom, rather than from a single point (such as the peptide center of mass) that tends to smear out these two hydration peaks. The first peak in  $g_{PW}(r)$ , representing the first hydration layer, appears at  $r \approx 2$  Å and is mainly associated with water molecules localized close to the charged amino acid Lys, as well as Asn, Ser, and Thr, which show the largest average number of neighboring water molecules within this first hydration layer, see Supporting Information Figure S2. These water molecules primarily orient with oxygen toward the

peptide, Supporting Information Figure 3S. The second peak in  $g_{PW}(r)$  at  $r \approx 3 \text{ \AA}$  corresponds to the second hydration layer and shows a more random distribution of water molecules, with some preference to waters pointing with hydrogens toward the peptide, particularly in the U state, Supporting Information Figure 3S.

Interestingly, the peptide–sorbitol pair correlation function  $g_{PS}(r)$  shows a first peak at  $r \approx 4 \text{ \AA}$ , further than the first two hydration layers, representing sorbitol exclusion at these short distances from the peptide interface, Figure 2a. This exclusion is in agreement with other studies that have found that polyols are preferentially excluded from protein interfaces.<sup>8,10,20</sup> There are only minor differences in  $g_{PW}(r)$  for the U and F states (see Supporting Information, Figure 4S), but for both states in the presence of sorbitol, the local densities of water in the first hydration layers relative to the bulk, seen as the height of the  $g_{PW}(r)$  peaks, are somewhat higher than in pure water. This difference translates into  $\sim 11$  water molecules found in the vicinity of the peptide folded state ( $r \leq 3.6 \text{ \AA}$ ) in the presence of sorbitol compared with  $\sim 10$  in pure water.

The net exclusion or accumulation of water or osmolytes near the solvated peptide (at infinite dilution) can be quantified either by the preferential hydration coefficient  $\Gamma$  or the preferential interaction coefficient for osmolyte  $\Gamma_S$  because these two are necessarily related.<sup>8,16,36</sup> We focus here on the preferential hydration  $\Gamma$ , which we have found in experiments to remain constant over a wide range of concentrations.<sup>32,37</sup> To find  $\Gamma$  in simulations, we use the operational definition<sup>19,38</sup>

$$\Gamma(R) = N_W \left( 1 - \frac{N_S/N_W}{n_S/n_W} \right) \quad (1)$$

Here, we have defined a vicinal volume surrounding the peptide satisfying  $r < R$ , and a bulk domain for  $r \geq R$  within the simulation box.<sup>9</sup> The values  $N_S$  and  $N_W$  represent the number of sorbitol and water molecules, respectively, within the vicinal volume, whereas  $n_S$  and  $n_W$  are the number of sorbitol and water molecules in the bulk domain. The preferential interaction coefficient  $\Gamma$  emerges as the converged value of  $\Gamma(R)$  when  $R$  is large enough. A positive value for  $\Gamma(R)$  represents preferential accumulation of water (or equivalently, exclusion of sorbitol) from the peptide interface, whereas a negative value indicates the depletion of water or accumulation of sorbitol.

Determining  $\Gamma$  from simulations is notoriously difficult.<sup>39</sup> Sufficient statistical sampling requires trajectories of 10 ns or longer, and convergence of values at a large enough  $R$  should be confirmed. Even then, preferential interaction coefficients are highly sensitive to the particular force fields used.<sup>40</sup> Indeed, we found that due to the overall sensitivity of  $\Gamma(R)$ , the profiles do not usually reach convergence for any single nanosecond segment along the MD trajectory. However, when we average over 15 ns out of the final 20 ns in the trajectory, values of  $\Gamma(R)$  were well converged, as shown in Figure 2b. To ensure that only conformations representative of the U and F states are included in the averaging, the analysis of  $\Gamma(R)$  was performed only on frames with a peptide ASA value that was lower than the mean value plus one standard deviation for the F state and higher than the mean value less one standard deviation for the U state. For the analyzed trajectory,  $\Gamma(R)$  values for the U state converge at a large  $R$  to a positive value,  $\Gamma_U \cong 9$ , indicating preferential hydration, while for the F state  $\Gamma$  is slightly negative,  $\Gamma_F \cong -5$ , Figure 2b. Convergence of  $\Gamma$  values for the U state to a more positive value than for the F state indicates stronger sorbitol exclusion from U versus F conformations.

Thermodynamically, preferential interaction coefficients are directly related to the changes in peptide stability imposed by osmolytes. Our previous experiments<sup>32</sup> have shown that the peptide folding free energy  $\Delta G_{UF}$  changes linearly with solute Osmolal concentration  $c_{Osm}$ , see inset in Figure 2b. Analogous changes in protein folding free energy due to cosolute addition are commonly described using the relation  $\Delta G_{UF}(c_{Osm}) = \Delta G_{UF}(c_{Osm} = 0) + mc_{Osm}$ , where the defined  $m$  value describes the constant slope in  $\Delta G_{UF}(c_{Osm})$ .<sup>19,41</sup> Importantly, this linearity in  $\Delta G_{UF}(c_{Osm})$  translates into a constant change in the number of solute-excluding water molecules,  $\Delta\Gamma_{UF}$ , upon folding.<sup>42</sup> Thus, the  $\Gamma$  values evaluated in simulations for the different peptide states, F and U, can be used to calculate the change in peptide folding free energy upon the addition of solute,  $\Delta\Delta G_{UF}(c_{Osm})$ , as follows:

$$\begin{aligned} \Delta\Delta G_{UF}(c_{Osm}) &= mc_{Osm} = \frac{RT}{55.6}(\Gamma_F - \Gamma_U)c_{Osm} \\ &= \frac{RT}{55.6}\Delta\Gamma_{UF}c_{Osm} \end{aligned} \quad (2)$$

where 55.6 is the number of moles of water in 1 kg and  $RT$  is the thermal energy per mole. Experimentally, we have used the variation in  $\Delta G_{UF}(c_{Osm})$  with sorbitol concentration, Figure 2b inset, to determine  $\Delta\Gamma_{UF} = -19$ . This indicates that, in the folding process, the preferential hydration of the peptide drops by 19, or alternatively, that 19 water molecules on average are “released” upon peptide folding, independent of sorbitol concentration.<sup>32</sup> This number corresponds closely with  $\Delta\Gamma_{UF} = \Gamma_F - \Gamma_U = -14$  that we find in simulations, as described above. We note that our simulations use a somewhat higher concentration of sorbitol than in experiments, allowing us to observe larger changes in solution structure and to gain statistically significant and convergent results. Because experiments show a highly constant preferential hydration over a very large concentration regime, we expect the same trends to be valid also for the concentrations used in simulations.

An additional validation of the peptide simulations can be made by using the “transfer model” formalism.<sup>13,36,43</sup> It has been shown experimentally that changes in peptide free energy upon transfer from water to solutions containing cosolutes can be dissected into a sum of contributions from different amino acid side chains and backbone. These changes due to solvation can be translated into changes in folding free energy due to the presence of cosolutes, as long as values of accessible surface areas are available for both U and F states. Using the experimental values reported by Bolen and Auton,<sup>36</sup> we estimated the changes in solvation energies for the simulated folded and unfolded states. Specifically, using the exposed surface areas of side chains and the backbone in the different conformations of the U and F states of the peptide in water, we determine the average ensemble difference  $\Delta\Delta G_{UF}(c_{Osm})$ , see Figure 5S (Supporting Information). This procedure allowed us to derive the average  $m$  value using eq 2. Using the transfer model, we find  $m = -646.14 \text{ J mol}^{-1}$ , which, is close to the experimental value of  $-841.36 \text{ J mol}^{-1}$ .

Taken together, the correspondence of experimental and predicted  $m$  values, as well as the close match that we find in the experimental and simulated  $\Delta\Gamma_{UF}$ , further support our MD simulation as good models for the peptide in its solvating environment, allowing us to explore additional properties that are hardly accessible in experiments. Specifically, while  $\Delta\Gamma_{UF}$  lets

us quantify the changes in preferential hydration upon folding in sorbitol solution, further analysis is required to reveal the underlying molecular origins for the exclusion. In the next sections, we proceed to the principal findings from our work to follow the possible molecular mechanism of peptide stabilization.

**Hydrogen Bonding Plays a Major Role in Peptide Stabilization by Osmolyte.** The experimentally determined enthalpic contribution to peptide folding in the presence of sorbitol<sup>32</sup> ( $\Delta\Delta H_{UF} = -3.3 \text{ kJ mol}^{-1}$  for the folding in 1 Osm sorbitol solutions) could originate from a variety of forces acting in solution, such as altered van der Waals (vdW) interactions, electrostatic forces within media of an altered dielectric constant, and variations in hydrogen bonding. To begin to unravel the important contributions to the system enthalpy in simulations, we first calculated the total force-field energies, arising from Coulomb ( $U_{El}$ ) and Lennard-Jones ( $U_{vdW}$ ) contributions, Table 1. Examining the differences in the force field energies between the F and U states, we find a positive value for  $\Delta U_{vdW}$  representing a Lennard-Jones contribution to the potential energy that disfavors folding, Table 1. In contrast to  $\Delta U_{vdW}$ ,  $\Delta U_{El}$  shows a much larger, negative contribution and hence favors folding. The same trend is observed in the presence of sorbitol; however, the contributions of  $\Delta U_{vdW}$ , as well as  $\Delta U_{El}$  to the potential energy are smaller than in pure water. We conclude, therefore, that the dominating forces in the simulations leading to favorable folding are electrostatic. Importantly, hydrogen bonding is represented within the empirical force field primarily by electrostatic interactions, making it a potential source that drives folding.

**Table 1. Differences between F and U States in Force Field Energies Arising from van der Waals or Coulomb Interactions, As Well As the PV Work Involved**

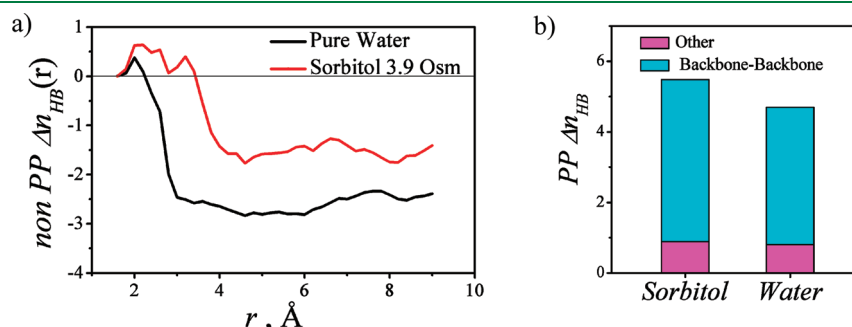
	$\Delta U_{vdW}^a$ (kJ/mol)	$\Delta U_{El}^b$ (kJ/mol)	$\Delta U_{Tot}^c$ (kJ/mol)	$P\Delta V^d$ (kJ/mol)
water	212.86	-1562.37	-1349.51	$-1.56 \times 10^{-3}$
sorbitol 3.9 Osm	171.34	-765.35	-594.01	$163.00 \times 10^{-3}$

<sup>a</sup>Total force-field energies arising from Lennard-Jones contributions. <sup>b</sup>Total force-field energies arising from Coulomb contributions. <sup>c</sup>Total force-field energies arising from Lennard-Jones and Coulomb contributions. <sup>d</sup>PV work calculated from the simulation for the transition from the U to the F state at a constant pressure of 1 atm and a volume corresponding to the same amount of water molecules around F and U states.

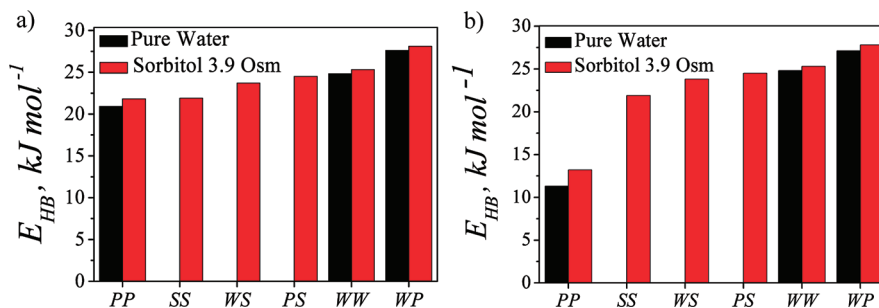
We also report in Table 1 values of  $\Delta U_{Tot}$  representing the sum of  $\Delta U_{vdW}$  and  $\Delta U_{El}$  to the potential energy, as well as  $P\Delta V$ , where  $P$  is the pressure (1 atm in the simulation) and  $\Delta V$  is the difference in the volume that includes the same amount of water molecules around F and U states. Under the simulation isobaric conditions, this  $P\Delta V$  term represents the difference between the enthalpy and the energy of the system. We find that this difference is very small relative to  $\Delta U_{Tot}$ , Table 1, and is also very small compared to the experimental value of  $\Delta\Delta H_{UF}$ . This allows us in the following discussion to equally speak about differences in the energy and the enthalpy.

To isolate the hydrogen bonding contribution to the changes in energy, we further analyzed changes to the hydrogen bonding network upon peptide folding in pure water and in the presence of sorbitol. Specifically, we enumerate the total number of hydrogen (H) bonds in the peptide solution and the number of internal peptide (PP) hydrogen bonds that are formed or lost upon folding. We use a geometric H bond definition, so that a hydrogen bond exists between two molecules if the oxygen–oxygen distance is less than  $d = 3.5 \text{ \AA}$ , and the  $O \cdots O-H$  angle is smaller than  $\theta = 30^\circ$ , as previously suggested.<sup>44–46</sup> This definition, as originally described for pure water, conveniently delineates contacts that form a prominent peak in the probability density of water–water nearest-neighbors as mapped in the  $d-\theta$  plane. We have further verified that this definition is applicable to these contacts in the presence of polyols and sugars, as well as for water–sorbitol H bond interactions,<sup>46</sup> and that it also applies for the ternary system that includes the peptide, see the Supporting Information Figure 6S. The H bonds accounted for in the peptide environment include contributions from the different pairs of chemical species: water–water (WW), water–sorbitol (WS), water–peptide (WP), sorbitol–sorbitol (SS), and sorbitol–peptide (SP). For PP contacts, which also include nonhydroxyl types of H bonds, we use a slightly different definition described and evaluated by Thornton and McDonald.<sup>47</sup>

To evaluate changes in the number of H bonds in the different peptide states and solution conditions, we compared differences in numbers of bonds with respect to the average number of the same type of H bond in the bulk. By only counting the excess or deficit of bonds with respect to the bulk,  $n_{HB}$ , we are able to compare simulated systems that are slightly different in size or density. For this purpose, the bulk was defined as the volume satisfying  $7.4 > r > 8.8 \text{ \AA}$  from the peptide, where we find convergence in the difference in excess or deficit of bonds, for both the F and U states. Figure 3a indicates the difference



**Figure 3.** Difference in the number of hydrogen bonds in the folded and unfolded states in pure water and in the presence of sorbitol. (a) Hydrogen bonds in the peptide environment calculated with respect to the bulk versus the distance from any peptide atom. Hydrogen bonds in the peptide environment include water–water, water–sorbitol, water–peptide, sorbitol–sorbitol, and sorbitol–peptide hydrogen bonds. (b) Number of internal peptide hydrogen bonds.



**Figure 4.** Average energies of each hydrogen bond in pure water and in the presence of sorbitol in the folded (a) and unfolded (b) states. All hydrogen bonds considered are within  $r \leq 4 \text{ \AA}$ , where  $r$  is the shortest distance from any peptide atom.

between F and U states in non-PP H bonds (non-PP  $\Delta n_{\text{HB}}(r)$ ) in the peptide environment, plotted versus the distance  $r$  from any peptide atom. Negative values represent a loss of H bonds, while positive values represent gains in the number of hydrogen bonds upon folding.

Surprisingly, our analysis shows that by folding, the solvating environment loses hydrogen bonds in both pure water and in the presence of sorbitol. However, the loss is smaller in the presence of sorbitol, Figure 3a. Concomitantly, the number of internal peptide H bonds added in the U→F transition (PP  $\Delta n_{\text{HB}}$ ) is also larger than in pure water, Figure 3b. We further dissected the number of internal peptide H bonds into backbone–backbone hydrogen bonds and all other PP H bonds. While the number of backbone–backbone H bonds is significantly larger in the presence of sorbitol, there is only a small difference in the corresponding numbers of nonbackbone–backbone H bonds. The trend for the changes in non-PP  $\Delta n_{\text{HB}}(r)$  upon sorbitol addition is consistent with our experimental results that show a diminished unfavorable enthalpic contribution to folding in the presence of polyols.<sup>32</sup> This finding suggests an important contribution of H bonds in solution to the stabilizing effect of sorbitol.

To fully appreciate the effect of sorbitol requires that we not only follow the change in the numbers of H bonds but also account for variations in the strength of hydrogen bonds in the presence and absence of the osmolyte. Energies of each hydrogen bond ( $E_{\text{HB}}$ ) were estimated using the correlation between H bond length and its strength, parametrized by Espinosa et al.,<sup>48</sup> and expressed as  $E_{\text{HB}} = 2.5 \times 10^4 \exp(-3.6 \times d(\text{H} \cdots \text{O}))$ , where  $d$  denotes the distance between the hydrogen atom and acceptor atom and  $E_{\text{HB}}$  is bond energy given in kilojoules per mole. We note that phenomenological energies derived using this parametrization are used here, as in other studies, only to provide a semiquantitative measure for the strength of H bonds.<sup>49–51</sup> Our analysis, therefore, relies only on the relative strength of the hydrogen bonds and not on their absolute values.

With the exception of PP hydrogen bonds, we find only small differences in H bond energies between the folded and unfolded states. Figure 4 shows the average energies for each class of H bond that are in close proximity to the peptide ( $r \leq 4 \text{ \AA}$ ), both in pure water and in the presence of sorbitol. The small number of internal peptide H bonds that form in the U state, as well as the fact that these hydrogen bonds are relatively weak (compare 11.3 kJ/mol for U with 20.9 kJ/mol for F state in water) indicate an advantage to folding both in pure water and in the presence of sorbitol. In both states (U and F), we find that the WP contacts form the strongest H bonds (27.6 kJ/mol for F and 27.1 kJ/mol for U state in water). The relative strength of these H bonds

grows in the presence of sorbitol (28.1 kJ/mol for F and 27.8 kJ/mol for U state in sorbitol). These WP H bonds are much stronger than PS hydrogen bonds (24.5 kJ/mol for both states). These rather weak PS bonds may, at least partly, explain sorbitol's exclusion from the peptide surface, and its preference to remain hydrated in the bulk solution, where the water–sorbitol hydrogen bond energy is  $\sim 23.8$  kJ/mol.

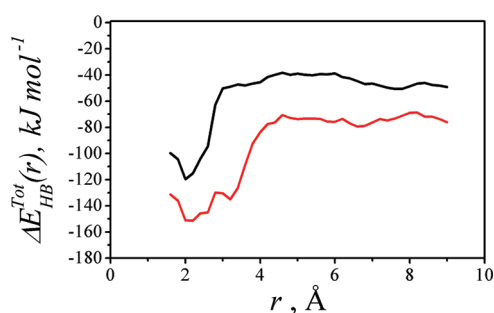
Interestingly, WW and PP hydrogen bonds are stronger in the presence of sorbitol than in pure water. This implies that water molecules released upon folding in the presence of sorbitol create hydrogen bonds with other water molecules in the bulk that are stronger than those that they form in pure water. In fact, every WW and WP hydrogen bond created in the presence of sorbitol gains an additional 0.5–0.8 kJ/mol over that formed in pure water.

Generally, hydrogen bond energies in simulations that use empirical force fields are not uniquely defined, as they arise in these models from electrostatic forces that originate from collective particle-charge interactions of many water molecules. We have, however, verified that another possible measure of hydrogen bond strength confirms our conclusions, as derived from the phenomenological estimates of Espinosa et al.<sup>48</sup> Specifically, we evaluated the force-field energy distribution for WW and WS contact pairs defined as H bonded, measured with respect to their energy at infinite separation. Using this alternate definition, we find that, in agreement with our conclusions from Figure 4, WW contacts are strengthened in the presence of sorbitol and that WW contacts are overall stronger than WS contacts, see Table 1S (Supporting Information).

To properly compound changes in H bond strength and changes in their numbers, we weigh the number of hydrogen bonds by the estimated energy of each hydrogen bond at each distance from the peptide. The resulting sum of hydrogen bond energies that includes all perturbations in H bonds around the peptide up to a distance  $r$  are compiled in Figure 5 for differences between the U and F states ( $\Delta E_{\text{HB}}^{\text{Tot}}$ ) in the presence and absence of sorbitol. We find that despite the loss in H bonds upon folding in the peptide environment, the peptide is overall enthalpically stabilized by changes in H bonding due to folding. In addition, we find larger peptide stabilization as a result of folding in the presence of the osmolyte, so that the H bond energy contributions converge to a value of  $\sim -72$  kJ/mol in sorbitol versus  $\sim -48$  kJ/mol in pure water.

While the choice of water and osmolyte model and force field may affect our results, particularly the hydrogen bond energies and solvation entropies, we do not expect other empirical force fields to yield qualitatively different results. Specifically, several studies using a variety of force fields for aqueous carbohydrate





**Figure 5.** Cumulative hydrogen bond energy differences between the folded and unfolded states in pure water (black line) and in the presence of sorbitol (red line) calculated as the product of the number of hydrogen bonds and the estimated energy of each type of hydrogen bond. Summation is performed as a function of distance  $r$  from the peptide.

binary solutions, including our own tests of additional force fields, have found similar effects of polyols on water structuring.<sup>5</sup> Moreover, another MD study of polyol interaction with large folded proteins that used a different force field (GROMOS96) has found several trends that are in common with our analysis, including a polyols-induced increase in the number of PP H bonds with a concomitant decrease in the number of PW bonds, increased protein hydration numbers, and stronger PW hydrogen bonds.<sup>20</sup>

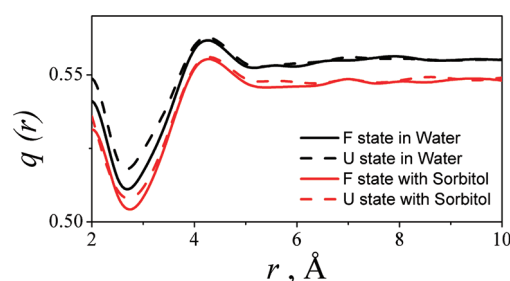
#### Entropic Contributions and Water Orientational Order.

An important and perhaps counterintuitive experimental result is that, in the presence of sorbitol, folding becomes *less* entropically favored ( $\Delta\Delta S = -9 \text{ J/mol K}$  for folding at 1 Osm sorbitol).<sup>32</sup> To follow the possible sources of this decrease in the entropic contribution, we first studied the water tetrahedral structural order parameter  $q(r)$ , in pure water and in the presence of sorbitol, at distance  $r$  around folded and unfolded peptide conformations. This order parameter is a metric that has been used to quantify the tendency of a water molecule and its four nearest neighbors to adopt a tetrahedral arrangement<sup>46,52</sup> and is defined for the  $i$ th water molecule as

$$q_i = 1 - \frac{3}{8} \sum_{j>k} \left( \cos \psi_{ijk} + \frac{1}{3} \right)^2 \quad (3)$$

where  $\psi_{ijk}$  is the angle formed between the central oxygen atom  $i$  and two neighboring atoms  $j$  and  $k$  (belonging to either water, polyol, or peptide hydroxyl oxygen or nitrogen), and the sum extends over four nearest neighbors. If oxygens (or nitrogens) are arranged in a perfect tetrahedral arrangement, then  $q = 1$ , while for an uncorrelated distribution of oxygens/nitrogens,  $q = 0$ .

We find a lower average value of  $q$  for sorbitol solutions relative to pure water, see Figure 6. We have also found a similar trend for the ordering of water around other polyol osmolytes in binary solutions in the absence of a peptide.<sup>46</sup> The decrease in  $q(r)$  values for  $r < 4 \text{ \AA}$  indicates that the peptide further imposes a destructuring effect on hydrating waters found in its close proximity. A similar conclusion was reported previously by Czapiewski and Zielkiewicz, who investigated the structural and dynamic properties of water in the solvation shell formed around different conformations of a polypeptide chain in pure water.<sup>53</sup> By following the two-particle entropy around the peptide core, they concluded that water around the peptide is (locally) less structured than in the bulk.



**Figure 6.** Tetrahedral structural order parameter,  $q$ , plotted as a function of distance  $r$  from the peptide. Order parameter  $q(r)$  is shown in pure water (black) and in the presence of sorbitol (red) for folded (full lines) and unfolded (dashed lines) states.

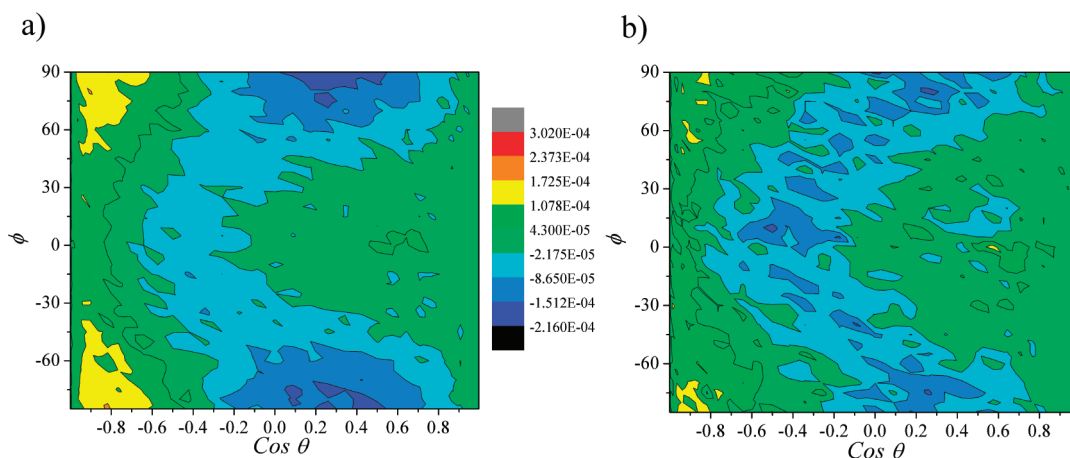
The order parameter further indicates that water is less tetrahedrally structured around F than U conformations, Figure 6. Moreover, the difference between  $q$  values in the peptide's vicinity for F and U states in pure water is larger than in the presence of sorbitol. These findings suggest that the presence of sorbitol imposes a disordering effect on water, so that both the peptide F and U states no longer alter water structuring to the same extent as in pure water.

To further investigate the structural properties of water within the first solvation layers, we determine the water angular probability distribution in term of  $\cos \theta$  and  $\phi$ , corresponding to the two angles that describe the orientation of water molecules with respect to the peptide within 4  $\text{\AA}$  from any peptide atom, as described in the Methods. Figure 7 shows the *difference* in the water angular probability distribution between folded and unfolded states, in pure water and in the presence of sorbitol.

The difference in water's angular distribution is very close to zero in the presence of sorbitol, indicating a similar distribution of water molecules' orientations in close vicinity to F and U states. However, in pure water, we find a considerably larger difference, for example, in the population of water orientations with  $-0.6 < \cos \theta < -0.9$  and  $45 < |\phi| < 90$ . Semiquantitatively, this larger probability difference in water orientations translates to a larger orientational entropy decrease for folding in pure water through the known expression for information entropy,<sup>54</sup>  $S = -k_B \sum_i P_i \ln P_i$  that relates the probability of accessible states  $i$  to the entropy  $S$ . These results also correlate well with the difference in  $q$  values around F and U states, as discussed above.

There are, in fact, many potential sources of entropy in this three-component system. As we have shown, an important entropic contribution relates to the changes in water orientation due to the addition of a peptide to solution. Therefore, to accurately evaluate the part of the solvation entropy that is due to the orientational degree of freedom of water molecules with respect to the peptide ( $s_{PW}^2$ ), we calculate  $s_{PW,o}^2$  and  $s_{PW,r}^2$  that describe the orientational and radial parts of the solvation entropy, respectively (see Methods for details). Summarized in Table 2, we show  $s_{PW,o}^2$  and  $s_{PW,r}^2$  calculated for the F and U states and compare the differences between folded and unfolded states,  $\Delta s_{PW,o}^2$  and  $\Delta s_{PW,r}^2$ , to the experimentally determined entropy of folding.

In water, the values of the orientational entropy  $s_{PW,o}^2$  for the folded ( $-807.50 \text{ J/mol K}$ ) and unfolded ( $-1152.87 \text{ J/mol K}$ ) states are much lower than in the presence of sorbitol ( $-89.06 \text{ J/mol K}$  and  $-131.99 \text{ J/mol K}$  for F and U states, respectively). The larger orientational freedom in water may be explained by our previous findings, indicating that, even in the absence of



**Figure 7.** Difference in water angular probability distribution between folded and unfolded states in pure water (a) and in the presence of sorbitol (b).

**Table 2.** Different Contributions to the Solvation Entropy<sup>a</sup>

		$s_{PW,o}^2$	$s_{PW,r}^2$	$\Delta s_{PW,o}^2$	$\Delta s_{PW,r}^2$	$\Delta S_{exp}$
pure water	folded	-807.50	-3.78	345.37	0.45	44
	unfolded	-1152.87	-4.23			
sorbitol 3.9 Osm	folded	-89.06	1.49	42.93	2.07	35
	unfolded	-131.99	-0.58			

<sup>a</sup>All contributions are given in units of J/mol K.  $s_{PW,o}^2$  denotes the orientational parts and  $s_{PW,r}^2$  the radial parts.

peptide, the binary sorbitol solution is more disordered than pure water.<sup>46</sup> Thus, the addition of the peptide to the binary solution cannot further increase water's orientational disorder to the same extent as it does in pure water.

It has been established that the solvation entropy  $s_{PW}^2$  is typically strongly dominated by the orientational part.<sup>55</sup> Accordingly, we find that the orientational entropy gain for folding in water ( $\Delta s_{PW,o}^2 = 345.37$  J/mol K) is larger than in the presence of sorbitol ( $\Delta s_{PW,o}^2 = 42.93$  J/mol K), while changes in the radial part are significantly smaller, Table 2. This result is also in good qualitative agreement with our experimental results that show diminished favorable entropy for folding once osmolytes (including sorbitol) are added.<sup>32</sup> These findings point to water orientational degree of freedom as an important contribution that should be considered, together with additional possible entropic contributions to the folding free energy of the peptide.

## DISCUSSION

Recently, the mechanism by which polyols and sugar osmolytes impact peptide stability has been shown to be enthalpically driven.<sup>32</sup> The current study aims to resolve the molecular origins of this enthalpic mechanism, by simulating the folded ( $\beta$ -hairpin, F state) and unfolded (U) states of a model peptide in pure water and in the presence of sorbitol.<sup>33,56,57</sup> Sorbitol was shown experimentally to significantly shift the thermodynamic equilibrium of the peptide, making it a convenient model for stabilizing polyol osmolytes.

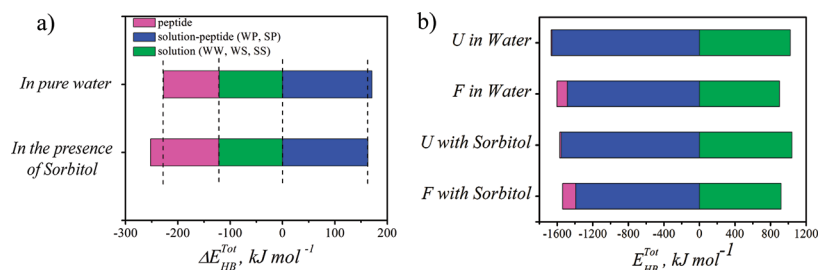
Using molecular dynamics simulations, we find that solution structure is modified upon peptide addition, but these changes are somewhat different in the presence and absence of sorbitol. In the following sections, we describe these differences and show how these modulations can explain the additional enthalpic peptide stabilization found in the presence of sorbitol.

**Sorbitol Alters the Peptide Hydration Layer.** Whereas recent works have argued that the presence of osmolytes alters protein stability indirectly by affecting water structure,<sup>5-7,20</sup> a case has also been made that it is the direct interaction between osmolyte and protein, or steric (excluded volume) interactions, that mediate the osmolyte effect.<sup>9,12,58,59</sup> The interplay of direct interactions and hydration has been particularly important in explaining the action of various solutes such as urea<sup>51,60</sup> or TMAO.<sup>61</sup> Our simulations indicate that while polyols remain excluded from the first solvation layer around peptide, they drive peptide stabilization primarily by impacting that first hydration layer. Specifically, local changes in concentrations and in hydrogen bonding relative to the bulk upon sorbitol addition are all shown to be limited to the peptide's first hydration layers, see Figures 2, 3a, 5, and 6.

The higher values of peptide–water  $g(r)$  hydration peaks in the presence of sorbitol compared with their value in pure water suggest a larger accumulation of water molecules around the peptide when sorbitol is present, see Figure 2. It may be tempting to speculate that this accumulation will result in an increase in the number of available hydrogen bonds for peptide–water interactions in the presence of sorbitol; however, in simulations we find an opposite trend, as further discussed in the following.

The overall change in the number of sorbitol-excluding water molecules (preferential hydration) upon peptide folding from the U to F states in our simulations is close to the values derived experimentally ( $\Delta\Gamma_{UF} \approx -14$  and  $-19$  in simulation and experiments, respectively).<sup>32</sup> This release of sorbitol-excluding waters should inevitably incur an enthalpic contribution, analogous to the heat of dilution associated with adding pure water to a binary aqueous sorbitol solution,  $\Delta H_m^{dil}$ .<sup>62-64</sup> Interestingly,  $\Delta H_m^{dil}$  for the corresponding number of waters released into 1 *m* sorbitol solution is  $\sim -0.225$  kJ/mol.<sup>64</sup> This value can account for only  $\sim 10\%$  of the peptide folding enthalpy change found experimentally in the presence of sorbitol. This result highlights the stark difference between the released osmolyte-excluding waters and pure bulk water. Specifically, the peptide interfacial waters can be expected to be much different than bulk waters, both structurally and in their interactions with peptide and solution.

**Solution–Peptide and Internal Peptide Hydrogen Bonds Drive Peptide Stabilization.** It is instructive to dissect the converged cumulative hydrogen bond energies, shown in Figure 5, into solution–peptide H bonds (including WP and SP), solution H bonds (including WW, WS, and SS), and internal

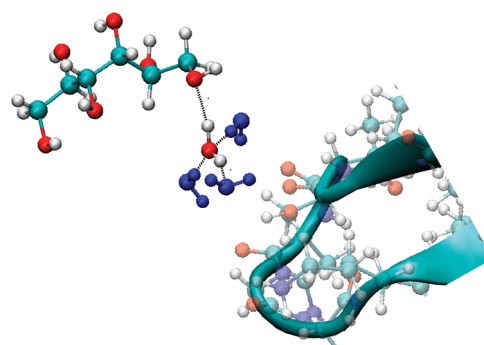


**Figure 8.** Cumulative hydrogen bond energies dissected into internal peptide, solution–peptide, and solution hydrogen bonds. (a) Differences between F and U states in pure water and in 3.9 Osm sorbitol. (b) F and U states in pure water and in the presence of sorbitol. Solution–peptide hydrogen bonds include WP and SP hydrogen bonds. Solution hydrogen bonds include WW, WS, and SS hydrogen bonds.

peptide H bonds. Figure 8a shows the energy differences between folded and unfolded states  $\Delta E_{HB}^{Tot}$  in each of these categories. Solution H bond energy differences are almost the same in pure water and aqueous sorbitol solution. However, the energies of solution–peptide and internal peptide hydrogen bonds significantly change in the presence of sorbitol. These changes are manifested in reduced solution–peptide hydrogen bond energy losses, as well as an increased internal peptide hydrogen bond energy gain upon folding in the presence of sorbitol, Figure 8.

**Changes in Cumulative Hydrogen Bond Energies Track the Changes in Number of Hydrogen Bonds.** Figure 8b shows the dissected cumulative hydrogen bond energies for the U and F peptide states in pure water and in the presence of sorbitol. While we found an increase in the energy of all H bond types in the presence of sorbitol (Figure 4), we concurrently found lower total solution–peptide H bond energies in the presence of sorbitol for both F and U states, Figure 8b. This finding reflects the smaller overall number of solution–peptide hydrogen bonds that form in the presence of sorbitol. This smaller number of bonds is also responsible for the decrease in the number of hydrogen bonds lost upon folding in sorbitol, see Figure 3a. Some of these H bonds lost are peptide–water hydrogen bonds that are the strongest H bonds formed in this system. To compensate for the loss of these H bonds as a result of folding, the system tends to increase the number of internal peptide hydrogen bonds, significantly strengthens these (seen as shorter bond length  $d$ ), and in addition creates stronger hydrogen bonds between peptide and water (27.8 and 28.1 kJ/mol for U and F states, respectively). We further find that the increase in internal peptide hydrogen bond enthalpy upon folding in the presence of sorbitol, Figure 8a,b, is mostly due to the higher number of backbone–backbone H bonds, see Figure 3b. This finding is also consistent with the work of Bolen and co-workers,<sup>4</sup> who have recently shown an enhancement of protein backbone–backbone hydrogen bonding interactions upon dilution from a good solvent (urea solution) to poorer (osmolytes containing) solvents, such as sarcosine or TMAO aqueous solutions.

**Reduced Numbers of Available Hydrogen Bonds in the Presence of Sorbitol Arise from Polyol’s Ability to Form H Bonds.** Polyols are known to compete with water’s own tendency to create optimal hydrogen bonds, thereby leading to a smaller number of solution–peptide hydrogen bonds that can form in their presence.<sup>46,53</sup> The smaller number of potential hydrogen bonds in the peptide vicinity in the presence of sorbitol also results in a smaller number of these bonds that can be lost as a result of folding. Overall, we have previously shown that polyols tend to participate in forming weaker, more distorted hydrogen



**Figure 9.** Typical snapshot from MD simulation of peptide in sorbitol solution, showing a water molecule that forms a weak H bond contact with sorbitol but three strong, optimized H bonds with waters around it.

bonds, in this way disrupting the water hydrogen bonding network and allowing less tetrahedral arrangements.<sup>20,32,52</sup> In concert, water molecules tend to optimize remaining water–water H bonds by creating more linear and shorter contacts leading to stronger interactions.<sup>20,46</sup> Indeed, we find here that sorbitol strengthens water–water as well as water–peptide hydrogen bonds, Figure 4. This impact of osmolytes on solution structure is somewhat similar to that shown in previous simulation studies of another known protective osmolyte, TMAO, showing a modest enhancement of the water structure near TMAO as well as an increase in water–water hydrogen bonding.<sup>7</sup>

The emerging picture of solution structure can be summarized in a typical simulation snapshot, Figure 9. The image shows a water molecule in the peptide’s first hydration layer that forms a weak hydrogen bond contact with sorbitol that is located further from the peptide but makes three much stronger hydrogen bonds with water molecules within that solvation layer. These H bonds with neighboring waters are even closer to optimal than water in the bulk, see Figure 7S. We find that this type of configuration is statistically favored and that sorbitol tends to reduce the number of hydrogen bonds with which the peptide can potentially interact. Sorbitol addition may, therefore, be viewed as introducing a competition between peptide and osmolyte for hydrogen bonding with water. Due to this competition, some potential H bonds to the peptide are lost, forcing the peptide to optimize remaining H bonds with water or within itself. This view also explains how water accumulation in the peptide solvation shell, discussed in Results, coincides with fewer yet stronger H bonds between water and peptide in the presence of sorbitol. Interestingly, Collins and Washbaugh<sup>65</sup> suggested an analogous mechanism for the alteration of water structuring at interfaces due to solutes to describe the action of different salt solutes.

Our results indicate that peptide-sorbitol H bonds are weaker and less favorable than other H bonds that can form in the ternary mixtures. This weaker interaction leads to sorbitol's exclusion and to preferential hydration of the peptide. In terms of polymer theory, polyol solutions are poorer solvents to the peptide than water, and therefore promote its collapse.<sup>66,67</sup> This forces the peptide to optimize internal (primarily backbone–backbone) hydrogen bonds. These findings gain support from the experimental studies of Bolen and coauthors,<sup>4,12,68</sup> showing that the unfavorable interactions between osmolytes and the peptide backbone raise the free energy of the U state, thereby shifting the thermodynamic equilibrium toward the native state. Non-hydrogen bonded “hydrophobic” interactions between nonpolar parts of the peptide are important for the collapse, but seem to be less altered by the presence of sorbitol than the hydrogen bonding network.

**Hydrogen Bond Optimization Is Consistent with a Decrease in Peptide Solvation Entropy in the Presence of Sorbitol.** To follow the contribution of water structural properties within the peptide hydration shell to the solvation entropy, we used the two-particle approximation to configurational entropy. We focused here on two terms of the water solvation entropy, the orientational and radial entropy contributions, and followed changes in peptide solvation entropy as a result of folding in sorbitol solution. This analysis revealed restriction in water orientations in the presence of sorbitol and a dominant contribution of orientational entropy ( $\Delta\Delta s_{PW,o}^2 = -302.44 \text{ J/mol K}$  compared with  $\Delta\Delta s_{PW,r}^2 = 1.62 \text{ J/mol K}$ ). The decrease in solvation entropy results from the optimization of hydrogen bonds wrought by the presence of sorbitol that concomitantly restricts the orientational freedom of water. This entropic contribution agrees with the trends found experimentally, showing a decrease in the total favorable entropic contribution to folding as a result of sorbitol addition ( $\Delta\Delta S = -9 \text{ J/mol K}$  for folding at 1 Osm sorbitol).<sup>32</sup>

Other terms in the total entropy that are more hardly accessible computationally could be important, but previous studies have shown that their contribution is typically limited. For example, estimates of the term depending on water–water interactions,  $s_{WW}^2$ , were calculated by Zielkiewicz and Czapiewski for water within a peptide solvation layer;<sup>53</sup> the study concluded that the local structure of the solvating water changes only slightly compared to that of bulk water. We suggest, therefore, that this term is expected to have a smaller influence also on the solvation entropy in the presence of sorbitol.

**The Possible Contribution of Crowding.** In addition to these entropic contributions, there are other sorbitol-related entropic terms that have not been explicitly determined here. For example, a calculation of depletion entropy can be made on the basis of the number of water molecules released as a result of peptide folding in sorbitol solution<sup>32</sup> according to Asakura–Oosawa theory.<sup>69</sup> A simple estimate results in  $\Delta\Delta S_{\text{dep}} = \Pi\Delta V = \Pi\Delta\Gamma_{UF}\nu N_{Av} \approx 10 \text{ J/mol K}$  for a solution osmotic pressure of  $\Pi = 3.9 \text{ Osm}$  at  $T = 298 \text{ K}$ , where  $\Delta V$  is the change in the osmolyte's free volume due to folding, as determined in ref 32, and  $\nu$  is the volume of a water molecule. This value can account only for a small part of the total change in entropy and would suggest that the folded state is entropically more favored in the presence of sorbitol, contrary to what we found experimentally. This suggests that mechanisms that rely purely on steric interactions do not play the dominant role in sorbitol's action.

Experiments further showed that the enthalpically driven mechanism together with the entropic penalty associated with the native state stabilization are strongly osmolyte size-dependent.<sup>32</sup> Thus, the entropic penalty grows as cosolute is varied from smaller polyols to the larger carbohydrates in the order glycerol < sorbitol < trehalose. This, again, is in contrast to a pure volume exclusion mechanism that would dictate a stronger but more favorable entropic contribution to folding when larger molecular crowders are present at the same mole concentration. These findings are, however, in agreement with recent experiments on the stabilization (or destabilization) of DNA in the presence of polyethylene glycols of various molecular weights.<sup>70</sup> These experiments showed a size-dependent exclusion that could be dissected into the extent of monomer–macromolecule “chemical interaction” effect and a steric effect. Interestingly, the steric effect becomes stronger with the volume of the crowder molecule, while the chemical interaction contribution is proportional to macromolecule and cosolute interacting surface areas. In agreement, we find that for the polyols we have tested,<sup>32</sup> the entropic penalty grows with cosolute size but that the slope of this change becomes less steep as cosolute size becomes larger, possibly suggesting a larger entropically favorable steric contribution for the larger cosolutes.

## CONCLUSION

We have used MD simulations to analyze the mechanism of peptide stabilization by sorbitol. Contrary to common wisdom, the peptide stabilization imposed by polyols was found experimentally to be enthalpically and not entropically driven. The emerging molecular mechanism shows that sorbitol stabilizes the native folded state of peptides by changing their immediate solvation layer. These changes lead to a decrease in the number of hydrogen bonds lost as a result of folding, optimization of existing hydrogen bonds, and an increase in the number of internal peptide hydrogen bonds. While the effects of sorbitol described in this study for a short 16-amino-acid-long peptide are relatively small, the significant accumulation of small contributions from solute–peptide interactions over extended macromolecular interfaces should have a profound effect on stabilization of the larger proteins typically found in cellular environments. It will be interesting to find out if the described mechanism is common to all peptides with different sequences and secondary structures, as well as to proteins, in solutions that contain polyols or other osmolytes.

## METHODS

**All-Atom Molecular Dynamics Simulations.** To compare with our previous experimental work,<sup>32</sup> we simulate here a model 16-residue peptide (sequence: Ac-KKYTVSINGKKITVSI) that can fold to a  $\beta$ -hairpin structure. Unfolded and folded peptide states were simulated in pure water and in the presence of sorbitol. To select initial configurations for subsequent MD simulations, we first employed the CHARMM empirical force field with implicitly included water to distinguish between the two primary states: folded and unfolded. An all-atom representation of the peptide was used in these simulations together with the SASA implicit water parameters, at two temperatures: 300 K and 400 K. Statistics of peptide dynamics showed two major peptide conformational populations when dissected by their accessible surface area (ASA), see Supporting Information Figure 1S. These results indicate conformations and changes in peptide folding that are similar to those found using NMR measurements,<sup>33</sup> indicating

Table 3. Parameters for Simulation Runs

peptide state	sorbitol (Osm)	number of sorbitol molecules	number of water molecules
folded (F)	0	0	3769
	3.9	150	2174
unfolded (U)	0	0	3747
	3.9	150	2169

that for different effective temperatures, the peptide occupies a different set of accessible conformations, changing from a more compact and folded  $\beta$ -hairpin structure at low temperatures to more extended “unfolded” structures with a larger ASA at high temperatures. One of the most probable structures from each state was used for further all-atom MD simulations with explicit water performed using NAMD.<sup>71</sup>

Table 3 lists the entire set of the MD simulations performed. Conformations of the folded and unfolded states were placed in a cubic box of TIP3P water molecules and three chloride counterions (representing pure water solvent solutions). In addition, each state was also immersed in osmolyte solutions by inserting sorbitol molecules in a cubic box of TIP3P water molecules with a single peptide molecule, corresponding to concentrations of approximately 3.9 Osm (3.86 *m*), as detailed in Table 3. All interactions were subject to the CHARMM27 force field<sup>72,73</sup> and used without further modifications. Bonds were kept at a constant length for solutes and solvent molecules using the SHAKE algorithm. All simulations were performed within the NPT ensemble, at  $T = 298$  K (using Langevin dynamics algorithm as implemented in NAMD) and  $P = 1$  bar (maintained using the Nose–Hoover Langevin piston method), within a cubic box with fluctuating length of ca.  $L = 48$  Å and periodic boundary conditions. After initial energy minimization of 1000 steps and 100 ps of MD equilibration, 50 ns MD simulation trajectories were collected. Of these, the last 13–15 ns were used for further analysis, with collection steps every 0.5 ps, resulting in well converged averages for all calculated distributions. The time step in all simulations was 2 fs. Electrostatic calculations were performed using the Ewald particle-mesh summation with 1 Å grid spacing. The van der Waals interactions were truncated smoothly with a cutoff of 12 Å and a switching distance of 10 Å. MD trajectory analysis was performed using VMD.<sup>74</sup>

**Radial Distribution Functions.** Radial distribution functions,  $g_{xy}(r)$ , assess local densities of atom type  $y$  at a distance  $r$  from an atom of type  $x$ . This  $g_{xy}(r)$  is calculated as

$$g_{xy}(r) = \frac{y(r',r)}{\rho_{y,\text{bulk}} \delta V(r',r)} \quad (4)$$

where  $r$  is the radius of the solvation shell,  $y(r',r)$  is the number of  $y$  molecules found between  $r'$  and  $r$ ,  $\delta V(r',r)$  is the volume of the shell ranging from  $r'$  to  $r$ , and  $\rho_{y,\text{bulk}}$  is the bulk density of  $y$ . The volume  $\delta V(r',r)$  was calculated using the Monte Carlo method for determining volume by randomly placing 1000 points in the simulation box and computing the ratio of hits within a shell to the total number of points. Unless otherwise stated, in all of our reported calculations, we have evaluated local densities of water oxygen's and sorbitol's center of mass at a distance  $r$ , corresponding to the

shortest distance from *any atom* of folded or unfolded peptide states.

**Osmolyte Force-Field Validation.** We have previously validated the sorbitol parameter set used here<sup>46</sup> by comparing densities from binary mixture simulations of sorbitol at 2.4 M to the densities extrapolated from experimental data<sup>75</sup> (1.14 gr/cm<sup>3</sup>). The experimental value differs by 3% from the value found in our simulations (1.11 gr/cm<sup>3</sup>). In addition, we have calculated another experimentally available thermodynamic property, the Kirwood–Buff integral for sorbitol<sup>46</sup>  $G_{SS} = \int (g_{SS} - 1) dv$ , where  $S$  stands for solute, and  $g_{SS}$  is the solute–solute pair correlation function, measured as a function of the distance between a central polyol hydroxyl oxygen and an osmolyte-representing atom. We found that in simulations  $G_{SS} = -0.2$ , very close to the previously published experimental value of  $-0.23$ .<sup>76</sup>

**Calculation of the Two-Body Peptide–Water Contribution to the Solvation Entropy.** The position and orientation of each water molecule with respect to the peptide were determined from the MD trajectories and used within the two-particle approximation to find the contribution of the local structure of water to the solvation entropy, as has been previously developed and described in detail (see refs 53, 55, 77). Within this approximation, the two-particle contribution to entropy evaluated relative to a completely random orientation of water molecules for a two-component system has the form:<sup>55,77,78</sup>

$$s^2 = s_{PP}^2 + s_{PW}^2 + s_{WW}^2 \quad (5)$$

where  $s_{PP}^2$  describes peptide–peptide interactions and therefore is irrelevant in our one-peptide simulations. The  $s_{WW}^2$  term is harder to access and is expected to be less important for the differences between solutions with and without peptides, as further detailed in the Discussion. Finally,  $s_{PW}^2$  describes the protein–water interactions given by

$$s_{PW}^2 = -k_B N_W \rho \int g_{PW}^2 \ln g_{PW}^2 d\bar{r} + k_B N_W \rho \int (g_{PW}^2 - 1) d\bar{r} \quad (6)$$

where  $\rho$  denotes the number density of the peptide, taking into account the volume within 4 Å from each peptide atom,  $g_{PW}^2$  is the two-body peptide–water distribution function,  $N_W$  is the number of water molecules within the same volume, and  $k_B$  is Boltzmann's constant.

The two-body contribution to the solvation entropy can be further separated into an orientational  $s_{PW,o}^2$  and a radial (or “non-orientational”)  $s_{PW,r}^2$  part by writing  $g_{PW}^2$  as a product,  $g_{PW}^2 = g(r) a(r,\theta,\phi)$ , of the radial distribution function  $g(r)$  and a function  $a(r,\theta,\phi)$  that describes the orientation of water molecules (the angular distribution function), normalized to  $4\pi$ . The angles  $\theta$  and  $\phi$  determine the position of the surrounding water molecules relative to the closest peptide atom.<sup>78</sup> The angle  $\theta$  describes the angle between the dipole moment of the water molecule,  $d$ , and the vector originating at the water molecule's oxygen atom and ending at the closest peptide atom,  $r_{OP}$ . The angle  $\phi$  is formed between the plane spanned by the water molecule and the plane spanned by the  $d$  and  $r_{OP}$  vectors.<sup>78</sup> Using the product that gives  $g_{PW}^2$  in the expression for  $s_{PW}^2$ , we thus obtain  $s_{PW}^2 = s_{PW,o}^2 + s_{PW,r}^2$

where

$$s_{PW,r}^2 = -4\pi k N_W \rho \int_0^\infty [g(r) \ln g(r) - g(r) + 1] r^2 dr \quad (7)$$

$$s_{PW,o}^2 = -k N_W \rho \int_0^\infty r^2 g(r) dr \times \int_0^\pi \int_0^{2\pi} a(r, \theta, \phi) \ln a(r, \theta, \phi) \sin \theta d\theta d\phi \quad (8)$$

In our evaluations of these quantities, we used the following integration steps:  $\delta r = 0.2 \text{ \AA}$  and  $\delta\theta = \delta\phi = \pi/36$ .

## ■ ASSOCIATED CONTENT

**S Supporting Information.** Figures for the probability distribution of accessible surface area of the peptide under different conditions, the average number of water molecules around the different peptide amino acids, radial distribution functions for water oxygens in the peptide vicinity, properties of hydrogen bonds and their probability distributions, and the average energy of hydrogen bonds in the bulk. This material is available free of charge via the Internet at <http://pubs.acs.org>.

## ■ AUTHOR INFORMATION

### Corresponding Author

\*Tel.: 972-2-6585484. Fax: 972-2-6513742. E-mail: [daniel@fh.huji.ac.il](mailto:daniel@fh.huji.ac.il).

## ■ ACKNOWLEDGMENT

We thank Hirsh Nanda for his help with setting up the implicit solvent simulations and Liel Sapir for his help with the transfer free energy calculations. The financial support from the Israel science foundation (ISF grant Nos. 1011/07 and 1012/07) is gratefully acknowledged. The Fritz Haber Research Center is supported by the Minerva Foundation, Munich, Germany.

## ■ REFERENCES

- Yancey, P. H.; Clark, M. E.; Hand, S. C.; Bowlus, R. D.; Somero, G. N. Living with Water-Stress - Evolution of Osmolyte Systems. *Science* **1982**, *217* (4566), 1214–1222.
- Willmer, P. Biochemical adaptation - Mechanism and process in physiological evolution. *Science* **2002**, *296* (5567), 473–473.
- Wood, J. M. Osmosensing by bacteria: Signals and membrane-based sensors. *Microbiol. Mol. Biol. Rev.* **1999**, *63* (1), 230.
- Holthauzen, L. M. F.; Rosgen, J.; Bolen, D. W. Hydrogen Bonding Progressively Strengthens upon Transfer of the Protein Urea-Denatured State to Water and Protecting Osmolytes. *Biochemistry* **2010**, *49* (6), 1310–1318.
- Sharp, K. A.; Madan, B.; Manas, E.; Vanderkooi, J. M. Water structure changes induced by hydrophobic and polar solutes revealed by simulations and infrared spectroscopy. *J. Chem. Phys.* **2001**, *114* (4), 1791–1796.
- Freda, M.; Onori, G.; Santucci, A. Hydrophobic hydration and hydrophobic interaction in aqueous solutions of tert-butyl alcohol and trimethylamine-N-oxide: a correlation with the effect of these two solutes on the micellization process. *Phys. Chem. Chem. Phys.* **2002**, *4* (20), 4979–4984.
- Zou, Q.; Bennion, B. J.; Daggett, V.; Murphy, K. P. The molecular mechanism of stabilization of proteins by TMAO and its ability to counteract the effects of urea. *J. Am. Chem. Soc.* **2002**, *124* (7), 1192–1202.

(8) Rosgen, J.; Pettitt, B. M.; Bolen, D. W. An analysis of the molecular origin of osmolyte-dependent protein stability. *Protein Sci.* **2007**, *16* (4), 733–743.

(9) Athawale, M. V.; Dordick, J. S.; Garde, S. Osmolyte trimethylamine-N-oxide does not affect the strength of hydrophobic interactions: Origin of osmolyte compatibility. *Biophys. J.* **2005**, *89* (2), 858–866.

(10) Rosgen, J.; Pettitt, B. M.; Bolen, D. W. Protein folding, stability, and solvation structure in osmolyte solutions. *Biophys. J.* **2005**, *89* (5), 2988–2997.

(11) Auton, M.; Bolen, D. W. Application of the transfer model to understand how naturally occurring osmolytes affect protein stability. *Osmosensing Osmosignaling* **2007**, *428*, 397–418.

(12) Bolen, D. W.; Baskakov, I. V. The osmolyte effect: Natural selection of a thermodynamic force in protein folding. *J. Mol. Biol.* **2001**, *310* (5), 955–963.

(13) O'Brien, E. P.; Ziv, G.; Haran, G.; Brooks, B. R.; Thirumalai, D. Effects of denaturants and osmolytes on proteins are accurately predicted by the molecular transfer model. *Proc. Natl. Acad. Sci. U. S. A.* **2008**, *105* (36), 13403–13408.

(14) Kornblatt, J. A.; Kornblatt, M. J. The effects of osmotic and hydrostatic pressures on macromolecular systems. *Biochim. Biophys. Acta—Protein Struct. Mol. Enzymol.* **2002**, *1595* (1–2), 30–47.

(15) Cayley, S.; Record, M. T. Roles of cytoplasmic osmolytes, water, and crowding in the response of *Escherichia coli* to osmotic stress: Biophysical basis of osmoprotection by glycine betaine. *Biochemistry* **2003**, *42* (43), 12596–12609.

(16) Timasheff, S. N. In disperse solution, “osmotic stress” is a restricted case of preferential interactions. *Proc. Natl. Acad. Sci. U. S. A.* **1998**, *95* (13), 7363–7367.

(17) Bolen, D. W. Effects of naturally occurring osmolytes on protein stability and solubility: issues important in protein crystallization. *Methods* **2004**, *34* (3), 312–322.

(18) Gibbs, J. W. On the equilibrium of heterogeneous substances. *Trans. Connecticut Acad.* **1876/78**, *3* (108–248), 343–542.

(19) Parsegian, V. A. Protein-water interactions. *Int. Rev. Cytol. Survey Cell Biol.* **2002**, *215*, 1–31.

(20) Liu, F. F.; Ji, L.; Zhang, L.; Dong, X. Y.; Sun, Y. Molecular basis for polyol-induced protein stability revealed by molecular dynamics simulations. *J. Chem. Phys.* **2010**, *132*, 22.

(21) Parsegian, V. A.; Rand, R. P.; Rau, D. C. Osmotic stress, crowding, preferential hydration, and binding: A comparison of perspectives. *Proc. Natl. Acad. Sci. U. S. A.* **2000**, *97* (8), 3987–3992.

(22) Linhananta, A.; Hadizadeh, S.; Plotkin, S. S. An Effective Solvent Theory Connecting the Underlying Mechanisms of Osmolytes and Denaturants for Protein Stability. *Biophys. J.* **2011**, *100* (2), 459–468.

(23) Saunders, A. J.; Davis-Searles, P. R.; Allen, D. L.; Pielak, G. J.; Erie, D. A. Osmolyte-induced changes in protein conformation equilibria. *Biopolymers* **2000**, *53* (4), 293–307.

(24) Patel, C. N.; Noble, S. M.; Weatherly, G. T.; Tripathy, A.; Winzor, D. J.; Pielak, G. J. Effects of molecular crowding by saccharides on alpha-chymotrypsin dimerization. *Protein Sci.* **2002**, *11* (5), 997–1003.

(25) O'Connor, T. F.; Debenedetti, P. G.; Carbeck, J. D. Simultaneous determination of structural and thermodynamic effects of carbohydrate solutes on the thermal stability of ribonuclease A. *J. Am. Chem. Soc.* **2004**, *126* (38), 11794–11795.

(26) Pincus, D. L.; Hyeon, C.; Thirumalai, D. Effects of trimethylamine N-oxide (TMAO) and crowding agents on the stability of RNA hairpins. *J. Am. Chem. Soc.* **2008**, *130* (23), 7364–7372.

(27) McPhie, P.; Ni, Y. S.; Minton, A. P. Macromolecular crowding stabilizes the molten globule form of apomyoglobin with respect to bath cold and heat unfolding. *J. Mol. Biol.* **2006**, *361* (1), 7–10.

(28) Cheung, M. S.; Klimov, D.; Thirumalai, D. Molecular crowding enhances native state stability and refolding rates of globular proteins. *Proc. Natl. Acad. Sci. U. S. A.* **2005**, *102* (13), 4753–4758.

(29) Zhou, H. X.; Rivas, G. N.; Minton, A. P. Macromolecular crowding and confinement: Biochemical, biophysical, and potential physiological consequences. *Ann. Rev. Biophys.* **2008**, *37*, 375–397.

- (30) Stanley, C. B.; Strey, H. H. Osmotically induced helix-coil transition in poly(glutamic acid). *Biophys. J.* **2008**, *94* (11), 4427–4434.
- (31) Jiao, M.; Li, H. T.; Chen, J.; Minton, A. P.; Liang, Y. Attractive Protein-Polymer Interactions Markedly Alter the Effect of Macromolecular Crowding on Protein Association Equilibria. *Biophys. J.* **2011**, *99* (3), 914–923.
- (32) Politi, R.; Harries, D. Enthalpically driven peptide stabilization by protective osmolytes. *Chem. Commun.* **2010**, *46* (35), 6449–6451.
- (33) Maynard, A. J.; Sharman, G. J.; Searle, M. S. Origin of beta-hairpin stability in solution: Structural and thermodynamic analysis of the folding of model peptide supports hydrophobic stabilization in water. *J. Am. Chem. Soc.* **1998**, *120* (9), 1996–2007.
- (34) Soper, A. K.; Bruni, F.; Ricci, M. A. Site-site pair correlation functions of water from 25 to 400 degrees C: Revised analysis of new and old diffraction data. *J. Chem. Phys.* **1997**, *106* (1), 247–254.
- (35) Mark, P.; Nilsson, L. Structure and dynamics of the TIP3P, SPC, and SPC/E water models at 298 K. *J. Phys. Chem. A* **2001**, *105* (43), 9954–9960.
- (36) Auton, M.; Bolen, D. W. Predicting the energetics of osmolyte-induced protein folding/unfolding. *Proc. Natl. Acad. Sci. U. S. A.* **2005**, *102* (42), 15065–15068.
- (37) Courtenay, E. S.; Capp, M. W.; Anderson, C. F.; Record, M. T. Vapor pressure osmometry studies of osmolyte-protein interactions: Implications for the action of osmoprotectants in vivo and for the interpretation of “osmotic stress” experiments in vitro. *Biochemistry* **2000**, *39* (15), 4455–4471.
- (38) Ghosh, T.; Kalra, A.; Garde, S. On the salt-induced stabilization of pair and many-body hydrophobic interactions. *J. Phys. Chem. B* **2005**, *109* (1), 642–651.
- (39) Shukla, D.; Shinde, C.; Trout, B. L. Molecular Computations of Preferential Interaction Coefficients of Proteins. *J. Phys. Chem. B* **2009**, *113* (37), 12546–12554.
- (40) Ploetz, E. A.; Benteñis, N.; Smith, P. E. Developing force fields from the microscopic structure of solutions. *Fluid Phase Equilib.* **2010**, *290* (1–2), 43–47.
- (41) Fersht, A. R. *Structure and Mechanism in Protein Science*; 3rd ed.; Palgrave Macmillan U. K.; W. H. Freeman: New York, 1999.
- (42) Harries, D.; Rosgen, J. A practical guide on how osmolytes modulate macromolecular properties. In *Biophysical Tools for Biologists: Vol 1 in Vitro Techniques*; Correia, J. J., Detrich, H. W., Eds.; Elsevier: New York, 2008; Vol. 84, pp 679.
- (43) Auton, M.; Bolen, D. W. Application of the transfer model to understand how naturally occurring osmolytes affect protein stability. *Osmosensing Osmosignaling* **2007**, *428*, 397–418.
- (44) Luzar, A.; Chandler, D. Effect of environment on hydrogen bond dynamics in liquid water. *Phys. Rev. Lett.* **1996**, *76* (6), 928–931.
- (45) Kumar, R.; Schmidt, J. R.; Skinner, J. L. Hydrogen bonding definitions and dynamics in liquid water. *J. Chem. Phys.* **2007**, *126*, 20.
- (46) Politi, R.; Sapir, L.; Harries, D. The Impact of Polyols on Water Structure in Solution: A Computational Study. *J. Phys. Chem. A* **2009**, *113* (26), 7548–7555.
- (47) McDonald, I. K.; Thornton, J. M. Satisfying Hydrogen-Bonding Potential in Proteins. *J. Mol. Biol.* **1994**, *238* (5), 777–793.
- (48) Espinosa, E.; Molins, E.; Lecomte, C. Hydrogen bond strengths revealed by topological analyses of experimentally observed electron densities. *Chem. Phys. Lett.* **1998**, *285* (3–4), 170–173.
- (49) Arnold, W. D.; Sanders, L. K.; McMahon, M. T.; Volkov, R. V.; Wu, G.; Coppens, P.; Wilson, S. R.; Godbout, N.; Oldfield, E. Experimental, Hartree-Fock, and density functional theory investigations of the charge density, dipole moment, electrostatic potential, and electric field gradients in L-asparagine monohydrate. *J. Am. Chem. Soc.* **2000**, *122* (19), 4708–4717.
- (50) Galvez, O.; Gomez, P. C.; Pacios, L. F. Variation with the intermolecular distance of properties dependent on the electron density in hydrogen bond dimers. *J. Chem. Phys.* **2001**, *115* (24), 11166–11184.
- (51) Stumpe, M. C.; Grubmüller, H. Aqueous urea solutions: Structure, energetics, and urea aggregation. *J. Phys. Chem. B* **2007**, *111* (22), 6220–6228.
- (52) Lee, S. L.; DeBenedetti, P. G.; Errington, J. R. A computational study of hydration, solution structure, and dynamics in dilute carbohydrate solutions. *J. Chem. Phys.* **2005**, *122*, 20.
- (53) Czapiewski, D.; Zielkiewicz, J. Structural Properties of Hydration Shell Around Various Conformations of Simple Polypeptides. *J. Phys. Chem. B* **2010**, *114* (13), 4536–4550.
- (54) Jaynes, E. T. Information Theory and Statistical Mechanics. *Phys. Rev.* **1957**, *106* (4), 620–630.
- (55) Zielkiewicz, J. Two-particle entropy and structural ordering in liquid water. *J. Phys. Chem. B* **2008**, *112* (26), 7810–7815.
- (56) Griffiths-Jones, S. R.; Maynard, A. J.; Searle, M. S. Dissecting the stability of a beta-hairpin peptide that folds in water: NMR and molecular dynamics analysis of the beta-turn and beta-strand contributions to folding. *J. Mol. Biol.* **1999**, *292* (5), 1051–1069.
- (57) Searle, M. S.; Griffiths-Jones, S. R.; Skinner-Smith, H. Energetics of weak interactions in a beta-hairpin peptide: Electrostatic and hydrophobic contributions to stability from lysine salt bridges. *J. Am. Chem. Soc.* **1999**, *121* (50), 11615–11620.
- (58) Paul, S.; Patey, G. N. Structure and interaction in aqueous urea-trimethylamine-N-oxide solutions. *J. Am. Chem. Soc.* **2007**, *129* (14), 4476–4482.
- (59) Zhang, Y. J.; Cremer, P. S. Chemistry of Hofmeister anions and osmolytes. In *Annu. Rev. Phys. Chem.*; Leone, S. R., Cremer, P. S., Groves, J. T., Johnson, M. A., Richmond, G., Eds.; Annual Reviews: Palo Alto, CA, 2010; Vol. 61.
- (60) Canchi, D. R.; Garcia, A. E. Backbone and Side-Chain Contributions in Protein Denaturation by Urea. *Biophys. J.* **2011**, *100* (6), 1526–1533.
- (61) Hu, C. Y.; Lynch, G. C.; Kokubo, H.; Pettitt, B. M. Trimethylamine N-oxide influence on the backbone of proteins: An oligoglycine model. *Proteins: Struct. Funct. Bioinf.* **2010**, *78* (3), 695–704.
- (62) Blackburn, G. M.; Lilley, T. H.; Walmsley, E. Aqueous-Solutions Containing Amino-Acids and Peptides. 13. Enthalpy of Dilution and Osmotic Coefficients of Some N-Acetyl Amino-Acid Amides and Some N-Acetyl Peptide Amides at 298.15 K. *J. Chem. Soc., Faraday Trans. I* **1982**, *78*, 1641–1665.
- (63) Gaffney, S. H.; Haslam, E.; Lilley, T. H.; Ward, T. R. Homotactic and Heterotactic Interactions in Aqueous-Solutions Containing Some Saccharides - Experimental Results and an Empirical Relationship between Saccharide Solvation and Solute Solute Interactions. *J. Chem. Soc., Faraday Trans. I* **1988**, *84*, 2545–2552.
- (64) Li, L.; Zhu, L. Y.; Qiu, X. M.; Sun, D. Z.; Di, Y. Y. Concentration effect of sodium chloride on enthalpic interaction coefficients of D-mannitol and D-sorbitol in aqueous solution. *J. Therm. Anal. Calorim.* **2007**, *89* (1), 295–301.
- (65) Collins, K. D.; Washabaugh, M. W. The Hofmeister effect and the behaviour of water at interfaces. *Q. Rev. Biophys.* **1985**, *18* (4), 323–422.
- (66) de Gennes, P.-G. *Scaling Concepts in Polymer Physics*; Cornell University Press: Ithaca, NY, 1979; p 315.
- (67) Rubinstein, M.; Colby, R. H. *Polymer Physics*; Oxford University Press: Oxford, U. K., 2003; p 440.
- (68) Ferreon, J. C.; Ferreon, A. C. M.; Bolen, D. W.; Hilser, V. J. Discrepancy between conformational stabilities obtained from native state hydrogen-deuterium exchange and denaturant-induced unfolding. *Biophys. J.* **2003**, *13* (2), 84A–13A.
- (69) Asakura, S.; Oosawa, F. On Interaction between Two Bodies Immersed in a Solution of Macromolecules. *J. Chem. Phys.* **1954**, *22*, 1255–1256.
- (70) Knowles, D. B.; LaCroix, A. S.; Deines, N. F.; Shkel, I.; Record, M. T. Separation of preferential interaction and excluded volume effects on DNA duplex and hairpin stability. *Proc. Natl. Acad. Sci. U. S. A.* **2011**, *108* (31), 12699–12704.
- (71) Phillips, J. C.; Braun, R.; Wang, W.; Gumbart, J.; Tajkhorshid, E.; Villa, E.; Chipot, C.; Skeel, R. D.; Kale, L.; Schulten, K. Scalable molecular dynamics with NAMD. *J. Comput. Chem.* **2005**, *26* (16), 1781–1802.
- (72) Brooks, B. R.; Brucoleri, R. E.; Olafson, B. D.; States, D. J.; Swaminathan, S.; Karplus, M. ChARMm - a Program for Macromolecular

Energy, Minimization, and Dynamics Calculations. *J. Comput. Chem.* **1983**, *4* (2), 187–217.

(73) MacKerell, A. D.; Bashford, D.; Bellott, M.; Dunbrack, R. L.; Evanseck, J. D.; Field, M. J.; Fischer, S.; Gao, J.; Guo, H.; Ha, S.; Joseph-McCarthy, D.; Kuchnir, L.; Kuczera, K.; Lau, F. T. K.; Mattos, C.; Michnick, S.; Ngo, T.; Nguyen, D. T.; Prodhom, B.; Reiher, W. E.; Roux, B.; Schlenkrich, M.; Smith, J. C.; Stote, R.; Straub, J.; Watanabe, M.; Wiorkiewicz-Kuczera, J.; Yin, D.; Karplus, M. All-atom empirical potential for molecular modeling and dynamics studies of proteins. *J. Phys. Chem. B* **1998**, *102* (18), 3586–3616.

(74) Humphrey, W.; Dalke, A.; Schulten, K. VMD: Visual molecular dynamics. *J. Mol. Graphics* **1996**, *14* (1), 33–8.

(75) Blodgett, M. B.; Ziemer, S. P.; Brown, B. R.; Niederhauser, T. L.; Woolley, E. M. Apparent molar volumes and apparent molar heat capacities of aqueous adonitol, dulcitol, glycerol, meso-erythritol, myo-inositol, D-sorbitol, and xylitol at temperatures from (278.15 to 368.15) K and at the pressure 0.35 MPa. *J. Chem. Thermodyn.* **2007**, *39* (4), 627–644.

(76) Rosgen, J., Molecular basis of osmolyte effects on protein and metabolites. In *Osmosensing and Osmosignaling*; Elsevier Academic Press Inc: San Diego, CA, 2007; Vol. 428, pp 459–486.

(77) Lazaridis, T.; Paulaitis, M. E. Simulation Studies of the Hydration Entropy of Simple, Hydrophobic Solutes. *J. Phys. Chem.* **1994**, *98* (2), 635–642.

(78) Bergman, D. L.; Lyubartsev, A. P.; Laaksonen, A. Topological and spatial aspects of the hydration of solutes of extreme solvation entropy. *Phys. Rev. E* **1999**, *60* (4), 4482–4495.



# IDSite: An Accurate Approach to Predict P450-Mediated Drug Metabolism

Jianing Li,<sup>†</sup> Severin T. Schneebeli,<sup>†</sup> Joseph Bylund,<sup>†</sup> Ramy Farid,<sup>‡</sup> and Richard A. Friesner<sup>\*,†</sup>

<sup>†</sup>Department of Chemistry, Columbia University, New York, New York

<sup>‡</sup>Schrödinger, Inc., 120 W. 45th St., New York, New York

 Supporting Information

**ABSTRACT:** Accurate prediction of drug metabolism is crucial for drug design. Since a large majority of drugs' metabolism involves P450 enzymes, we herein describe a computational approach, IDSite, to predict P450-mediated drug metabolism. To model induced-fit effects, IDSite samples the conformational space with flexible docking in *Glide* followed by two refinement stages using the Protein Local Optimization Program (PLOP). Sites of metabolism (SOMs) are predicted according to a physical-based score that evaluates the potential of atoms to react with the catalytic iron center. As a preliminary test, we present in this paper the prediction of hydroxylation and O-dealkylation sites mediated by CYP2D6 using two different models: a physical-based simulation model and a modification of this model in which a small number of parameters are fit to a training set. Without fitting any parameters to experimental data, the physical IDSite scoring recovers 83% of the experimental observations for 56 compounds with a very low false positive rate. With only four fitted parameters, the fitted IDSite was trained with a subset of 36 compounds and successfully applied to the other 20 compounds, recovering 94% of the experimental observations with high sensitivity and specificity for both sets.

## INTRODUCTION

It is crucial to understand how potential drugs are metabolized in the body, because human metabolism has profound impacts on the bioactivity and the safety profiles of drug candidates. On one hand, metabolism can convert these compounds into their active forms, which interact with the therapeutic targets; on the other hand, metabolism eliminates the compounds by converting them into inactive excretable metabolites. Sometimes the metabolic modifications also lead to toxicity, which can cause unexpected failures in the later phases of drug development. Furthermore, the metabolic behavior of drug compounds is also highly related to other critical issues such as food–drug interactions, drug–drug interactions, and personalized medication.<sup>1–3</sup> Given the enormous impact of metabolism on drug bioavailability and toxicity, it is important to determine metabolites in the early stage of the drug discovery process. However, to obtain such information experimentally is often a very lengthy and expensive process. Therefore, it would be extremely useful if one could use computational methods to predict the metabolic decomposition of drug candidates.

Since cytochrome P450 enzymes (CYP) are involved in a large majority of drug metabolism pathways, many computational studies have been published attempting to predict P450-mediated metabolism using a variety of methods and models. For a recent review, see the work of Afzelius et al.<sup>4</sup> These previous studies mainly focused on the important P450 isoforms 2D6, 2C9, and 3A4, aiming to predict the primary metabolites of drug compounds. Several ligand-based methods have been developed during the past decade, making predictions based on hydrogen abstraction energies estimated with semiempirical quantum mechanics<sup>5</sup> or DFT methods.<sup>6</sup> Although such ligand-based

methods are very fast, it is often necessary to consider the interaction between the enzyme and the substrate in order to reach high accuracy (for example, >80% agreement with experiments) in the predictions. It is possible to include a limited amount of enzyme-specific information by making descriptors of ligand-based models dependent on the nature of the enzyme.<sup>7–9</sup> Such approaches have been successfully implemented in software packages such as MetaSite, and some were reported to recover up to 86% of the experimental observations.<sup>10</sup> On the other hand, molecular dynamics (MD) or induced-fit docking simulations in combination with transition state calculations at the QM/MM or semiempirical quantum level were used to predict metabolites for a few ligands.<sup>11,12</sup> Other promising methods based on molecular docking have been implemented as well,<sup>13–17</sup> which determine the predictions using a reactivity model and/or distance cutoffs from the reactive iron center.

Traditional empirical ligand-based approaches to the prediction of P450 SOMs rely primarily on implicit estimation of the intrinsic site reactivity to the compound I oxo species, coupled with a heuristic attempt to take into account the ability of the ligand to bind to the P450 active site. While such methods can yield some discrimination of true positives from false positives when a sufficiently large training set is employed,<sup>5,8,18</sup> the precision of the approach is fundamentally limited, as the treatment of protein–ligand binding is highly approximate. Methods such as MetaSite<sup>10</sup> provide some incorporation of P450 structural information but employ a much smaller training set and fewer empirical parameters; the overall results appear to actually be less

**Received:** July 5, 2011

**Published:** September 02, 2011

accurate than a ligand-based approach employing an extensive data set. The problem is again that the MetaSite algorithm for modeling the reactive protein–ligand complex does not rigorously evaluate the binding energy or perform a thorough conformational search, severely limiting the predictive capability that can be attained.

The method described in the present work (IDSite) represents a qualitatively different approach from those discussed above, as well as from other efforts in the literature.<sup>6,8–10</sup> First, the goal is to actually generate an accurate structure for the protein–ligand complex that enables reactivity at a specified site; this requires construction of a good approximation to a transition state structure for both aliphatic and aromatic sites of reaction. Second, the relative binding affinity, as compared to alternative structures for both the site in question and for other sites, has to be computed with a respectable degree of precision, on the order of a few kilocalories per mole. Finally, the relative intrinsic barrier height of the reaction (combined with the relative binding affinity to produce an overall relative barrier), as compared to other possible reactions of the molecule, must be estimated, to within  $\sim 1$  kcal/mol. These are extraordinarily daunting tasks, given that the P450 isoforms present large, complex active site regions with substantial capability for induced-fit conformational changes, a necessary condition for them to accommodate the wide range of exogenous ligands with which they need to interact to perform their biological function.

The algorithms in IDSite employ a novel model for the total energy of the protein–ligand complex, which has recently been shown to provide remarkably accurate predictions for side chains and loops,<sup>19</sup> and a sophisticated algorithm for generating converged induced-fit structures which combines docking, conformational search, and hybrid Monte Carlo (MC) methods based on MD trajectories. The algorithm enables a hierarchical search which addresses the various length scales of the problem, including the small correlated motions provided by the MD trajectories which we have found are absolutely necessary to produce useful rank ordering of structures, particularly for larger ligands. Constraints are employed in conjunction with these simulation algorithms to enforce appropriate transition state structures. The energy model enables the targets of a few kilocalories per mole accuracy in relative binding affinity to be reached. Finally, a quantum chemically based model is employed to calculate relative intrinsic reactivities and again is shown below to yield outstanding performance. On the basis of such a high level of success, it is documented below in predicting true positive SOMs versus false positives.

To our knowledge, these results represent the first reliable and accurate computation of binding poses and transition states for a wide range of drug-like molecules interacting with an important human P450 isoform. There are a few previous papers in which structures are generated via QM/MM calculations.<sup>12,20,21</sup> However, these typically address a very small number of ligands (usually one); the ligands are typically simpler and smaller than those treated here; and the sampling algorithms are much less extensive. We believe that these structures can be very useful in practical drug design applications, in situations where modification of P450 metabolic properties for candidates in later stages of lead optimization is required. The availability of an atomic level three-dimensional structure, as well as the ability to predict the structural and energetic effects of chemical modification of the molecule, provides a new tool for chemists to rationally engineer desirable metabolic properties into clinical candidates. Extension

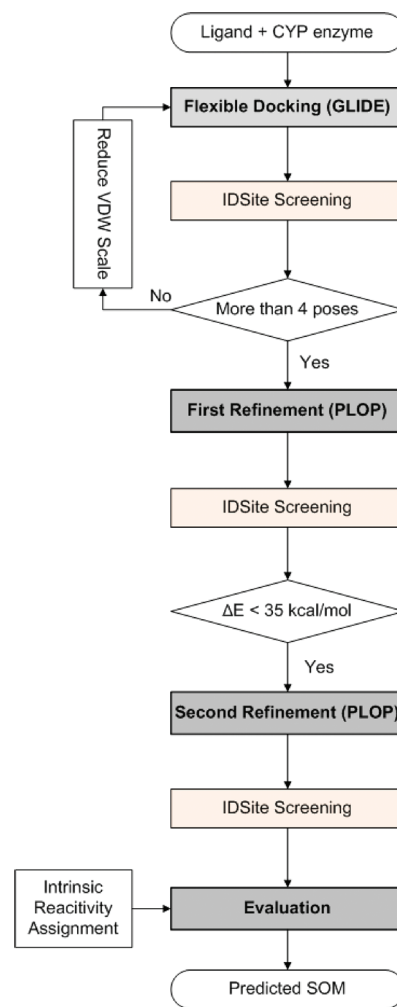


Figure 1. IDSite workflow.

of our methods to other P450 isoforms such as 1A2, 2C9, and 3A4, which is currently in progress, will enhance the utility of our approach for this important application.

## METHODS AND MATERIALS

**IDSite Methodology.** IDSite combines the docking program *Glide*<sup>22</sup> and the protein structure modeling program PLOP (Protein Local Optimization Program, available as the protein refinement module in the protein modeling package *Prime* of Schrödinger, Inc.<sup>23</sup>) to model induced-fit effects and to predict sites of metabolism. IDSite consists of three hierarchical sampling stages and one final scoring stage (Figure 1). It begins with flexible *Glide* docking calculations, which place the ligand into the active site. Following the docking stage, two refinement stages in PLOP are carried out to refine the protein side-chain and ligand orientations. At the end of each sampling stage, the generated/refined poses are screened on the basis of their structures and energies and clustered according to the similarity of the ligand conformation. Finally, the refined lowest energy poses are used to predict the sites of metabolism on the basis of a physical score, which is dependent on the energies of the poses as well as the intrinsic chemical reactivities of the potential sites of metabolism.

Table 1. IDSite Filters in the Screening for CYP2D6

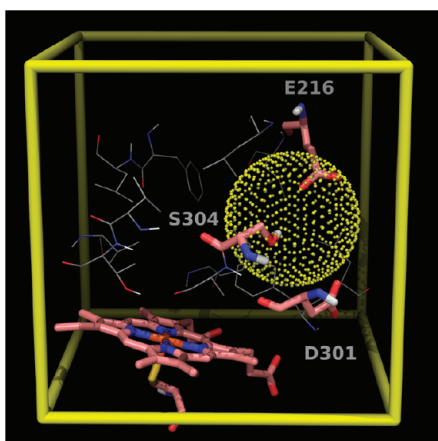
stage	filters applied in the screening at the end of the stage
Glide Docking and PLOP Refinement 1	<p>Poses that fulfill any of the criteria below are removed:</p> <ol style="list-style-type: none"> <li>(1) The distance of the basic nitrogen to the ferryl oxygen is less than 5.0 Å</li> <li>(2) The distance of the basic nitrogen to the negative charged oxygen (in Glu216 or Asp301) is greater than 5.5 Å</li> <li>(3) More than two heavy atoms from the ligands are further than 16.0 Å away from the heme iron</li> <li>(4) More than one heavy atom from the ligand is closer than 1.0 Å to the receptor</li> <li>(5) More than six heavy atoms from the ligand are closer than 1.8 Å to the receptor</li> <li>(6) No heavy atom in the ligand is within 5.0 Å to the heme iron</li> </ol> <p>For PLOP refinement 1: All of the poses are ranked with PLOP energies. Poses with energy higher than 35 kcal/mol compared to the lowest energy pose are removed.</p>
PLOP refinement 2	<p>Poses that fulfill any of the criteria below are removed:</p> <ol style="list-style-type: none"> <li>(1) The distance between the constrained atom and the ferryl oxygen is outside the optimal range, which is from 1.65 to 2.60 Å for sp<sup>3</sup> atoms and from 1.60 to 2.08 Å for sp<sup>2</sup> atoms</li> <li>(2) The distance of the basic nitrogen to the ferryl oxygen is less than 4.8 Å</li> <li>(3) The distance of any polar atom to the ferryl oxygen is less than 3.2 Å</li> <li>(4) The distance of the constrained salt bridge (between the basic nitrogen and the oxygen from Glu216 or Asp301) is greater than 3.6 Å; the angle of the salt bridge (N–H–O) is less than 140°</li> <li>(5) More than two heavy atoms from the ligands are either further than 14.5 Å or closer than 1.6 Å from the heme iron</li> <li>(6) The pose has at least one distorted cyclohexane ring.</li> </ol>

IDSite is able to use knowledge about specific conserved interactions to perform efficient sampling and accelerate the calculations. For example, in the case of CYP2D6, a typical substrate always contains a basic center (e.g., an amine nitrogen) that binds to one of the two acidic residues, Glu216 or Asp301. IDSite constrains such salt bridges to reduce the sampling cost associated with the docking and refinement stages. Filters are applied during the screening at the end of each stage in IDSite to reduce the number of poses passed to further refinement or evaluation (Table 1). The following is a detailed description of each stage of IDSite.

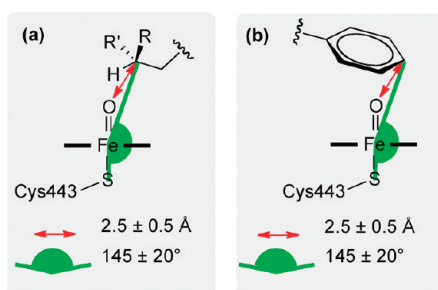
We have constructed our sampling and scoring algorithms with the intention of approximating the correct transition state structure of the protein–ligand complex and associated activation energy, which would lead to reactivity of the target atom of the ligand. There are two components of the problem: finding the transition state in reasonable CPU time (a daunting task for a large, complex ligand when induced fit effects are important) and estimating the free energy of activation associated with the transition state. The VSGB 2.0 energy function, with constraints to enforce a suitable geometry for the reaction to take place (and some other constraints as well to facilitate sampling, as described in the text below), is minimized to generate these structures for the various possible candidate reactant heavy atoms. We use the classical force field and solvation model to produce a “reactant” structure which is optimally positioned for the targeted chemical reactivity. The activation barrier from a precomputed quantum chemical fragment calculation, as described in the following text, is then added to the VSGB 2.0 energy to estimate the relative energy barrier for converting such a structure into products. This is an approximation to a more rigorous approach such as using QM/MM methods to generate the reactant, transition state, and product structures. Note that it is only important that relative free energies of the various potential sites of reaction are calculated

with reasonable accuracy, as the most reactive (lowest activation free energy) site is always used as a reference point (i.e., the energy function for this site is subtracted from the energy function for the candidate site) in our assessment of the metabolic contribution of each site. Finally, in applying the above protocol, the VSGB 2.0 energy must be calculated using a structure with the constraints in place; otherwise the structure would minimize to something that is not a suitable starting point for reaction. The constraints introduce some strain energy into the structure, but this strain energy is an appropriate component of the activation free energy as it does cost energy to create a suitable reactive structure.

*1. Glide Docking.* Starting from the ligand and the protein receptor structures, IDSite carries out flexible ligand docking with *Glide*.<sup>24,25</sup> The flexible ligand docking protocol generates a large number of ligand conformations that are then docked into the rigid receptor. The first step in *Glide* docking is to define the binding box and calculate the receptor grid. As in *Glide*, in IDSite the binding site is defined as a box centered at the center of selected residues or a ligand (if the structure contains a ligand). Because we start from the apo structure of CYP2D6 (PDB ID: 2F9Q; see below for details about the protein preparation), the center of the binding box is selected as the centroid of the residues Glu216, Asp301, Thr309, and Phe483. The box dimension on each side is set to 10 Å for the inner box and 20 Å for the outer box. After the grid generation, IDSite samples the conformations of freely rotatable bonds and rings with *Glide* Standard Precision (SP). In order to increase sampling, IDSite uses reduced van der Waals (VDW) radii and skips the default filtering with a rough score within *Glide* (also referred to as expanded sampling). Similar poses are clustered according to their RMSD (cutoff 2.0 Å). Finally, a postdocking minimization is performed, and the top 60 minimized poses according to the *Glide* SP score are retained. These poses are then screened to



**Figure 2.** Definition of the binding box (yellow cube) and the positional constraint (yellow dotted sphere) in IDSite for CYP2D6.

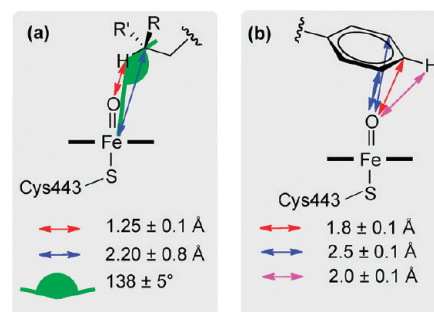


**Figure 3.** Constraints applied to the heme region in the *first* refinement stage. The ferryl oxygen is a “dummy” atom (1.6 Å above the heme iron), only used to define the constraints in the IDSite calculations. (a) Constraints for  $sp^3$  carbons. (b) Constraints for  $sp^2$  carbons.

remove the poses with obvious steric clashes, with too many atoms outside the inner binding box, or without atoms close to the heme iron (Table 1). The remaining poses are then passed to the first refinement stage.

IDSite uses reduced VDW radii for nonpolar atoms both in the protein receptor and the ligand, so that slight steric clashes are tolerated during the docking stage. For the protein receptor, the VDW scaling factor is fixed at 0.40, while for the ligand, the scaling factor starting from 0.80 is adaptively adjusted until at least four valid poses are found. With highly flexible ligands and relatively high scaling factors, *Glide* often finds only a handful of valid poses, and even fewer survive after IDSite screening. However, if the scaling factor is set too low, the docked poses may contain too many serious steric clashes, which can cause problems in the subsequent minimization. If IDSite fails to find enough valid poses, the scaling factor is adjusted, and the number of poses to pass the initial docking phase in *Glide* is increased accordingly to augment sampling.

Since a typical CYP2D6 substrate forms a highly conserved salt bridge with either Glu216 or Asp301,<sup>26</sup> IDSite employs this conserved interaction to reduce the sampling cost of the CYP2D6 docking in the following way: IDSite adds a positional constraint to ensure that the generated poses fulfill at least part of the preferred conserved interactions. The positional constraint defines a spherical region in the receptor that is within 4.0 Å of the center of the Glu216, Asp301, and Ser304 residues (Figure 2). It is required that during docking and postdocking



**Figure 4.** Constraints applied to the heme region in the *second* refinement stage. The ferryl oxygen is a “dummy” atom (1.6 Å above the heme iron), only used to define the constraints in the IDSite calculations. (a) Constraints for  $sp^3$  carbons. (b) Constraints for  $sp^2$  carbons.

minimization each pose should maintain at least one hydrogen-bond donor inside the spherical region. If the ligand contains other hydrogen-bond donors except for the basic nitrogen, the constrained docking is likely to generate poses that form hydrogen bonds instead of the salt bridge to Glu216 or Asp301. However, IDSite is able to distinguish these poses and filter them via an additional salt bridge filter in the pose screening (Table 1), so that only the poses with a stable salt bridge are allowed to pass to the refinement stage.

**2. PLOP Refinements.** The refinement of the docked poses includes multiple parallel Monte Carlo Minimization (MCM) simulations in PLOP. For each pose from the previous stage (the docking or first refinement stage), IDSite finds all of the heavy atoms in the ligand close to the heme iron. For each of these atoms, distance and angular harmonic constraints are applied in order to force sampling of the conformations that potentially lead to metabolism. The optimal distances and angles of the constraints were obtained from hydroxylation transition state geometries with a heme model system at the B3LYP/LACVP\* level using *Jaguar*.<sup>27</sup> The detailed nature of the employed constraints is shown for both  $sp^3$  and  $sp^2$  type carbons in Figures 3 and 4. The constraints are then employed in the minimization step but were not included in the energy used for the acceptance step of the MCM simulations. PLOP uses the overlap factor (the ratio of distance between two atom centers to the sum of their atomic radii) to quickly reject randomized structures with serious steric clashes (defined as the overlap factor being lower than a specific cutoff). PLOP repeats the random attempts until a structure with tolerable clashes is generated, after which a constrained minimization using the truncated Newton method is performed. The acceptance or rejection of the minimized structure is decided by the Metropolis criteria based on the energy calculated in the VSGB 2.0 model. (Performing the minimization step before testing the acceptance criteria violates detail balance, but this is not an issue, as we are interested only in low energy structures and not the population/ensemble distribution.) The simulations run until a certain number of accepted structures are collected.

In order to sample the various degrees of freedom in the conformational space, IDSite employs three types of randomized moves in the MCM simulations: side-chain rotation, rigid body translation/rotation, and hybrid moves.

**Side-Chain Moves.** By varying the dihedral angles of the rotatable bonds, IDSite uses side chain MC moves in PLOP to sample the selected side-chain conformations of the protein and

Table 2. Comparison of Settings in the First and Second Refinement Stages

	PLOP refinement 1	PLOP refinement 2
number of residues to sample (including the ligand)	12	40
number of accepted structures for each job	maximum of 8 times the number of rotatable bonds and 24	maximum of 20 times the number of rotatable bonds and 60
types and probabilities of MCM moves	side chain: 0.50 rigid body: 0.10 hybrid: 0.40	side chain: 0.70 rigid body: 0.10 hybrid: 0.20

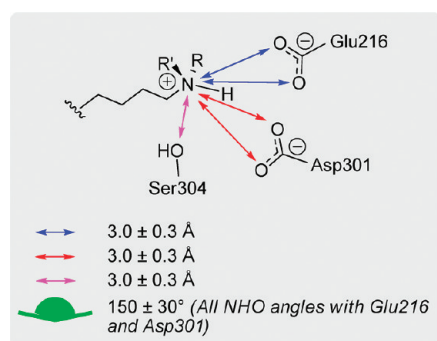


Figure 5. Constraints applied to the salt bridge region of CYP2D6 in the first refinement stage.

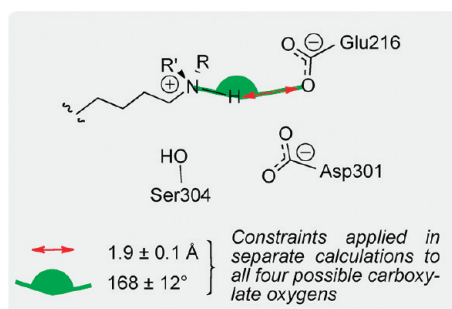


Figure 6. Constraints applied to the salt bridge region of CYP2D6 in the second refinement stage.

of the ligand. Up to three close residues ( $C_\beta$  distance within  $6 \text{ \AA}$ ) are allowed to rotate collectively, but the moves of the protein residues and those of the ligand are separated. In each attempted movement, the conformations of the selected side chains (from the protein/ligand) are either changed by random perturbations or assigned to a randomly selected rotamers from a library. For an attempt with a random perturbation, the displacement of each dihedral angle is the sum of a large rotation ( $N$  times  $60^\circ$  with  $N$  as a random integer between 1 and 5) and a random perturbation from  $0$  to  $30^\circ$ . For a rotamer library attempt, a side-chain conformation is updated with a random rotamer from a high resolution side-chain library for protein residues,<sup>28</sup> and from a homogeneous library at  $10^\circ$  resolution for the ligand. If a structure with tolerable overlaps is generated in an attempt, it is minimized and sent to subsequent stages for judgment of acceptance. Each side-chain move takes less than 15 s and is the fastest among all of the three move types.

**Rigid Body Moves.** Rigid body moves are used to sample the translational and rotational space of the ligand. Multiple attempts

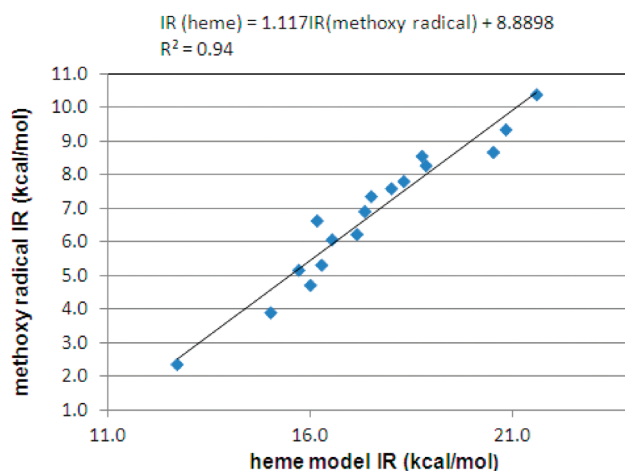


Figure 7. Correlation between the intrinsic reactivities calculated with the methoxy radical model and the heme model (17 sites from 9 selected fragment compounds; details are shown in Supporting Information).

with reduced VDW radii are applied, as it is quite common to fail in searching for a clash-free conformation in a single rigid body moving attempt (especially when the ligand is large and flexible and the binding pocket is relatively small). Each rigid body move includes 1000 attempts, and each attempt performs a translation along a random vector and a rotation around a random axis, with less than  $0.5 \text{ \AA}$  and  $60^\circ$  displacement, respectively. In addition, the VDW radii are reduced (scaling factor 0.8) to soften the Lennard-Jones potential, so that mild steric clashes are allowed, which are likely to be resolved by the subsequent minimization. The rigid body move usually takes 20 to 40 s per move.

**Hybrid Monte Carlo Moves.** The hybrid Monte Carlo (HMC) move<sup>29</sup> in PLOP performs simultaneous sampling for the selected residues in the protein side chains and backbone as well as the ligand. Each HMC move performs a 5 ps, constant energy molecular dynamic simulation (starting at 900 K) on all of the atoms in the selected residues. The molecular dynamics simulation uses a RESPA based integration of short-range forces with a time step of 1 fs and updates long-range forces with a Verlet integration every fifth step.<sup>30</sup> Taking up to 15 min per move, the HMC is the most expensive among all three types of moves in PLOP.

Considering the different costs for the three types of moves, the frequency of deployment of each move type in the various refinement stages is adjustable according to the sampling requirements. Two stages of refinement with different combinations of moves and constraints are carried out in the hierarchical sampling. Using more HMC moves, the first refinement stage applies loose distance constraints between an atom in question

Table 3. Summary of Results for the Training Set

symbol	compound name	physical IDSite			fitted IDSite		
		TP	FP	FN	TP	FP	FN
1	4-methoxyamphetamine	1	0	0	1	0	0
2	amitriptyline	2	2	0	2	0	0
3	aprindine	4	0	1	5	0	0
4	brofaromine	1	0	0	1	0	0
5	bufuralol	0	1	1	1	0	0
6	carvedilol	1	0	2	2	0	1
7	cinnarizine	0	2	1	0	2	1
8	clomipramine	1	0	1	1	0	1
9	codeine	1	0	0	1	0	0
10	desipramine	2	0	0	2	0	0
11	dextromethorphan	1	0	0	1	0	0
12	dihydrocodeine	1	1	0	1	0	0
13	ethylmorphine	1	0	0	1	0	0
14	flunarizine	1	0	0	1	0	0
15	fluperlapine	1	0	0	1	0	0
16	hydrocodone	1	0	0	1	0	0
17	imipramine	2	0	0	2	0	0
18	indoramine	1	0	0	1	0	0
19	MDMA	1	0	0	1	0	0
20	methamphetamine	1	0	0	1	2	0
21	methoxyphenamine	2	0	0	2	0	0
22	metoprolol	1	0	1	2	0	0
23	mexiletine	2	0	1	2	0	1
24	mianserin	1	0	0	1	0	0
25	mirtazapine	0	1	1	1	1	0
26	nortriptyline	1	1	0	1	0	0
27	ondansetron	2	0	0	1	0	1
28	paroxetine	1	0	0	1	0	0
29	perhexiline	2	0	0	2	0	0
30	propafenone	1	1	0	1	1	0
31	propranolol	2	2	0	2	1	0
32	tamoxifen	1	0	0	1	0	0
33	terfenadine	3	0	0	3	0	0
34	tiracizine	1	2	0	1	1	0
35	tropisetron	2	0	1	3	0	0
36	venlafaxine	1	0	0	1	0	0
	total	47	13	10	52	8	5

(from the ligand) to the ferryl oxygen. It is designed to “pull” the close atom (identified from the docking poses) toward the heme iron, to estimate the likelihood that the atom can approach the iron and react with the ferryl oxygen. When an atom in the ligand is forced to be proximate to the ferryl oxygen under the constraints, the rest of the ligand and the surrounding protein residues have to adjust their conformations accordingly. The adjustments for some poses are easy and for some others are difficult, depending upon the specific geometrical issues and energetics of the protein–ligand interactions for the trajectory connecting particular starting and target poses. Resulting poses with steric clashes or distorted structures can be identified by their high energies and discarded in the IDSite energy and structure screening. (Table 1). The low energy poses after screening, mostly with favorable interactions

Table 4. Result Summary for the Test Set

symbol	compound name	physical IDSite			fitted IDSite		
		TP	FP	FN	TP	FP	FN
37	atomoxetine	0	1	1	1	2	0
38	bicifadine	1	2	0	1	0	0
39	bupranolol	1	0	0	1	0	0
40	carteolol	1	1	0	1	0	0
41	chlorpromazine	1	0	0	1	0	0
42	EMAMC	1	0	0	1	0	0
43	encainide	1	1	0	1	1	0
44	harmaline	1	0	0	1	0	0
45	harmine	1	1	0	1	1	0
46	ibogaine	1	0	0	1	0	0
47	MAMC	1	0	0	1	0	0
48	MMAMC	1	0	0	1	0	0
49	MOPPP	1	0	0	1	0	0
50	oxycodone	1	0	0	1	0	0
51	spirosulfonamide	2	0	0	2	0	0
52	timolol	2	0	2	4	0	0
53	tolterodine	0	1	1	1	1	0
54	tramadol	1	1	0	1	1	0
55	tyramine	2	0	0	2	0	0
56	zotepine	1	0	0	1	0	0
	total	21	8	4	25	6	0

between the protein and the ligand, are passed to the second refinement stage. Mainly focusing on side-chain sampling, the second refinement stage applies tight constraints that force the structure to form special conformations similar to that of the transition states obtained from DFT calculations of model systems. The second refinement stage is used to further refine the poses and distinguish the potential of each atom in question to be oxidized. The comparison of the settings for these two refinement stages with PLOP are shown in Table 2, while the constraints are illustrated in Figures 3 and 4. There are approximately 39 protein residues, identified to be important for ligand binding by mutagenesis experiments<sup>31</sup> or are adjacent to these key residues, that are sampled during the refinement stages. At the end of each refinement stage, all of the poses sampled in that stage are screened and clustered for further refinement or evaluation (Table 1).

For CYP2D6, harmonic constraints are also applied to force the basic nitrogen to interact with the acidic residues, Glu216 and Asp301 (Figures 5 and 6), as they are believed to play important roles in substrate binding to CYP2D6 from mutagenesis experiments.<sup>26,32,33</sup>

**3. Evaluation.** Herein, we present two scoring models to evaluate the potential sites of metabolism and to determine the predictions. Our first scoring model (referred to as physical IDSite) is based on the following assumptions: (1) For hydroxylation of an aliphatic chain carbon, the P450-hydrogen abstraction step is rate determining.<sup>34,35</sup> (2) For hydroxylation of aromatic rings, the electrophilic attack of compound I on the aromatic ring is rate determining.<sup>34,35</sup> (3) All reaction intermediates before the rate determining step are in equilibrium.<sup>36</sup> Given these assumptions, the relative rates of product formation depend only on the relative transition state free energies of the rate determining (RD) transition states ( $\Delta G^\ddagger$ ) according to the

Curtin–Hammett principle. These can then simply be written as

$$\Delta G^\ddagger = \Delta G_{\text{bind}} + \Delta G_{\text{RD-step}}^\ddagger \quad (1)$$

where  $\Delta G_{\text{bind}}$  is the binding free energy of the substrate into the reactive conformation in the P450 active site and  $\Delta G_{\text{RD-step}}^\ddagger$  is the activation barrier of the RD step.

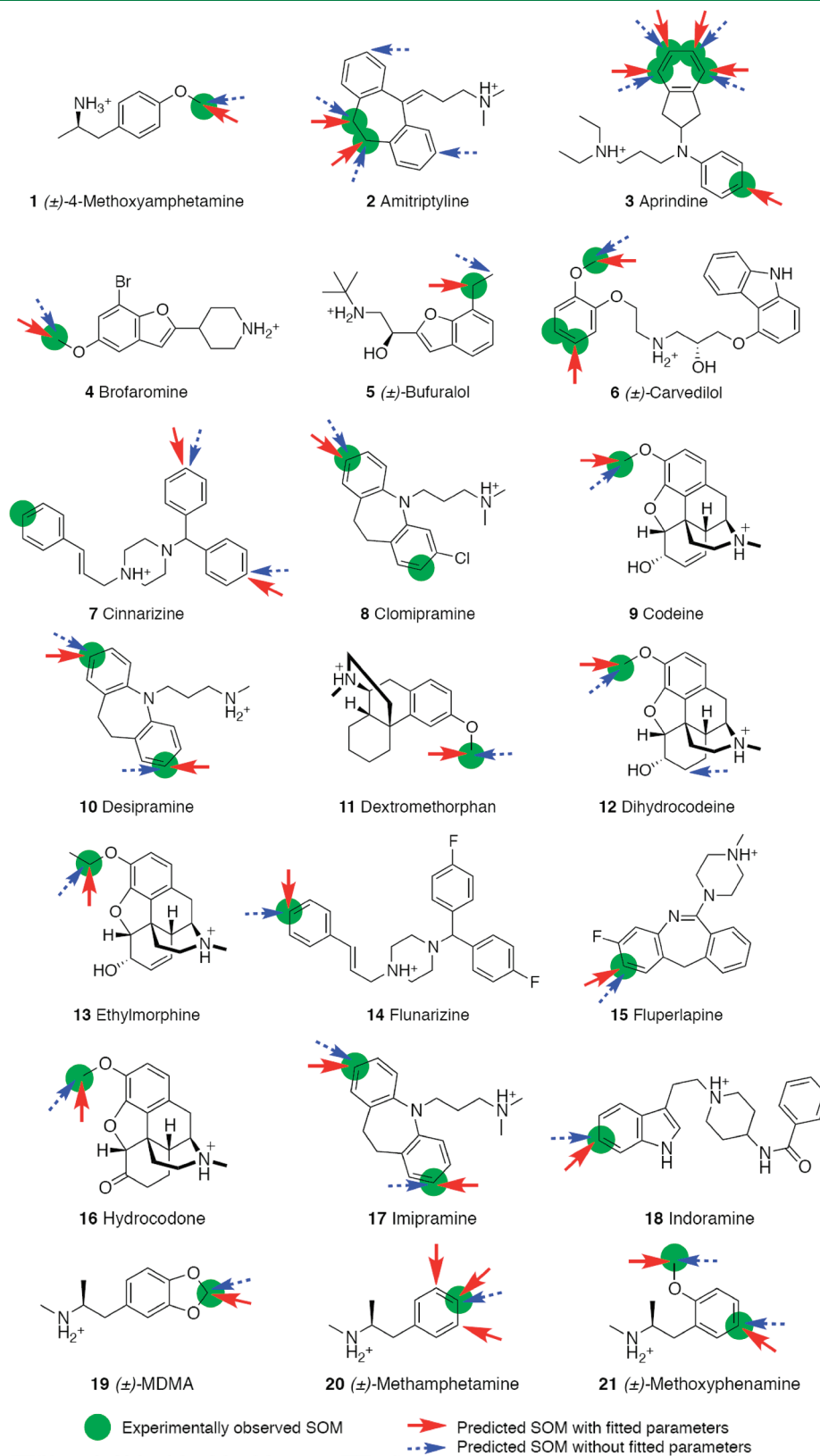


Figure 8. Continued

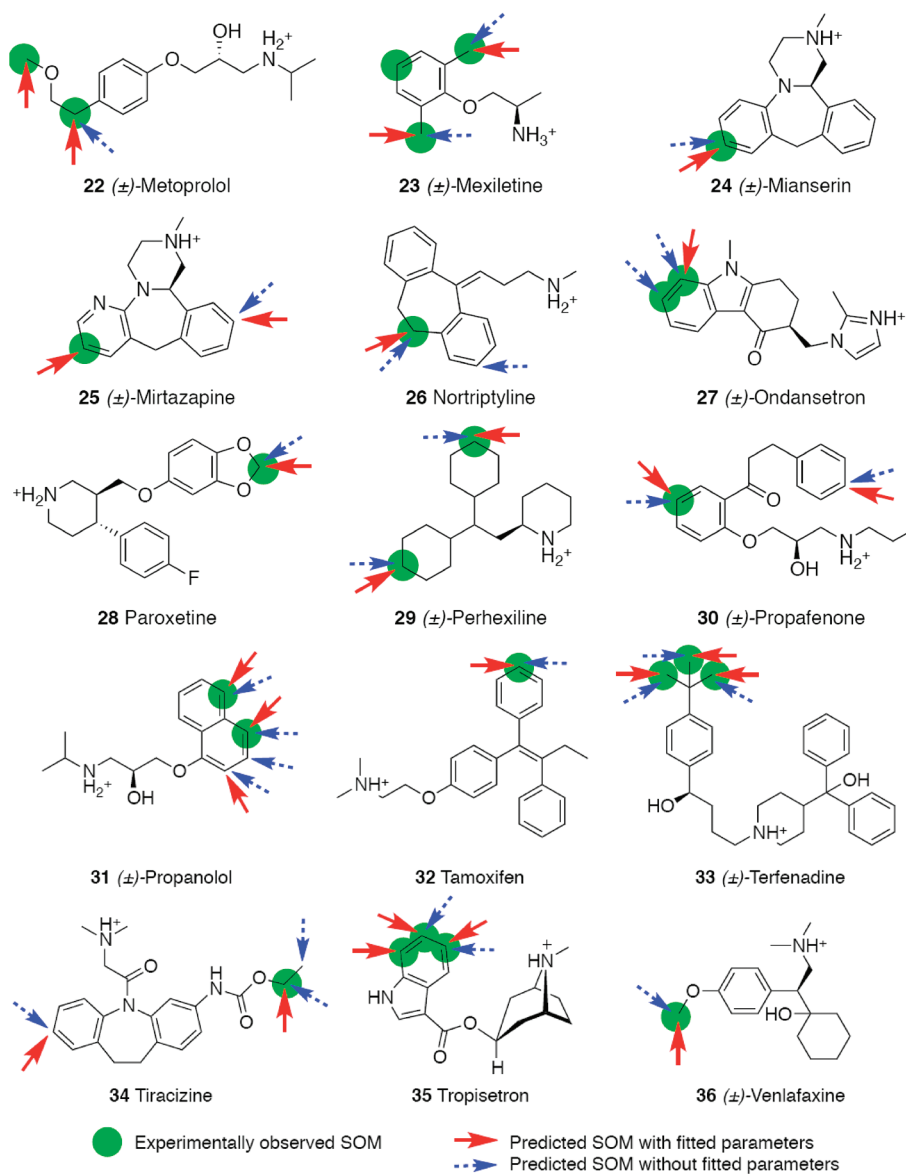


Figure 8. IDSite predicted results for the training set.

In the present application of IDSite, we attempt to calculate only relative, as opposed to absolute, site reactivity for a given ligand. Absolute site reactivity for the ligand can typically be obtained via inexpensive experiments. However, detailed metabolic chemistry is often more difficult to determine, and an accurate three-dimensional structure leading to reactions at each metabolic site is not available given the severe challenge of obtaining a crystal structure of a P450 isozyme with the ligand bound in the reactive conformation. Prediction of the most highly reactive site, followed by the identification of all sites with relative reactivities sufficiently large to be experimentally detected along the dominant metabolic pathway, coupled to structural prediction for each relevant reactive geometry, complements current experimental practice and facilitates compound modification in situations where P450 metabolism needs to be altered to confer improved metabolic properties on a candidate drug molecule.

In the physical IDSite model, the relative binding energies of various docked poses are calculated from the PLOP VSBG 2.0 energies of these poses, while the barriers for the RD steps are estimated from the corresponding activation barriers of model compounds with a methoxy radical (calculated at the DFT level).  $E_{\text{pose}}$  in eq 3, calculated in PLOP, estimates the protein–ligand interactions when a potential site is forced to approach the catalytic center in a certain pose with a transition state-like conformation. On the basis of the linear correlation (Figure 7) between the methoxy radical activation barriers and the corresponding activation barriers with the heme system, we approximated the real activation barrier for each potential site of metabolism from the intrinsic reactivity calculated with the methoxy radical model according to eq 2.

$$\text{IR}(\text{heme}) = 1.117 \times \text{IR}(\text{methoxyradical}) + \text{constant} \quad (2)$$



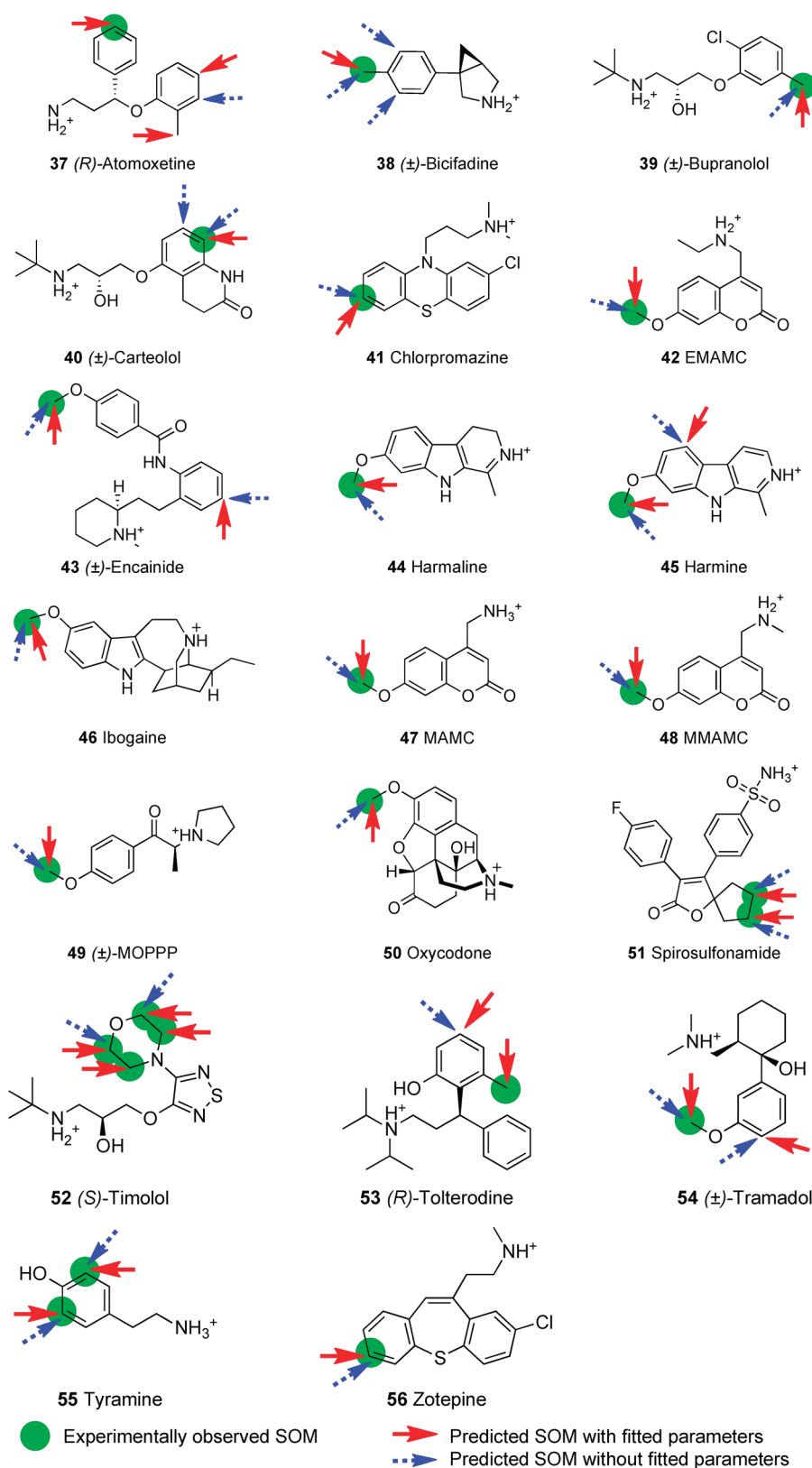


Figure 9. IDSite predicted results for the test set.

With the constant from eq 2 ignored, the relative  $\Delta G^\ddagger$  for each potential site (approximated as the score  $E$ ) is then calculated as the Boltzmann weighted average over the energies of all contributing

poses, where angle brackets represent the Boltzmann averages (eq 3). A term describing the configurational entropy of equivalent hydrogen atoms at 298 K, proportional to the logarithm of the

number of symmetrically equivalent hydrogen atoms, was also included. The  $\Delta G^\ddagger$  values for all symmetrically equivalent sites were set to the lowest  $\Delta G^\ddagger$  of the sites.

$$E = \langle 1.117 \times \text{IR}(\text{methoxyradical}) + E_{\text{pose}} \rangle - kT \ln N_{\text{H}} \quad (3)$$

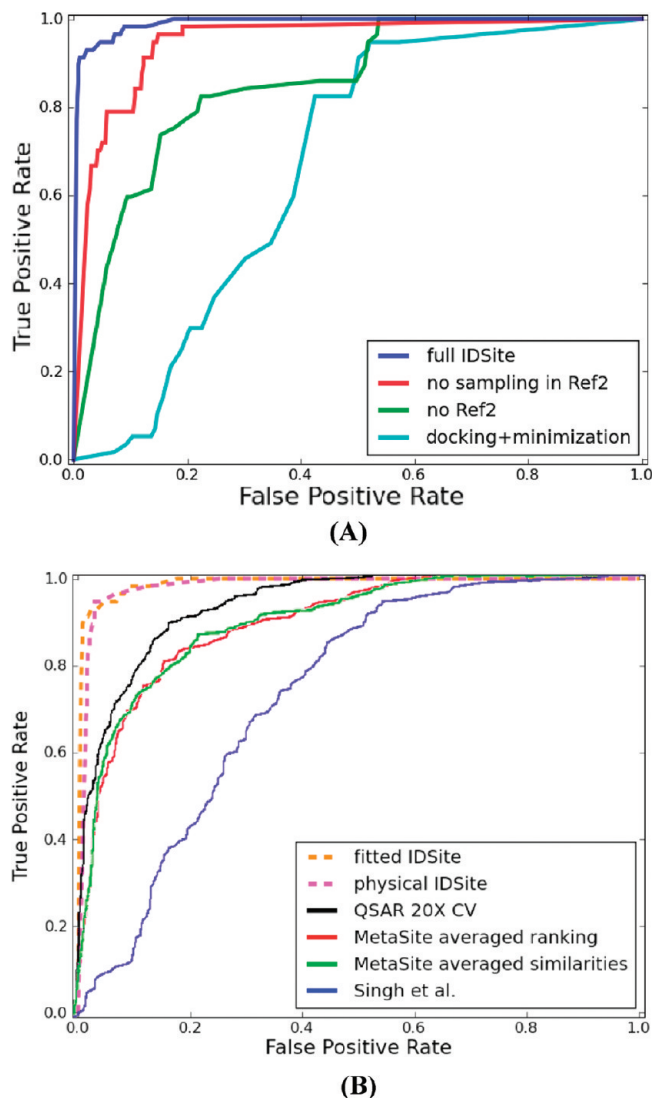
Since (as a rule of thumb) it is difficult to observe a minor metabolite experimentally if it is formed in less than ca. 0.1% yield (which corresponds to ca. 4.75 kcal/mol increase in relative  $\Delta G^\ddagger$  compared to the free energy of the most favored product), we used 4.75 kcal/mol as a cutoff for the prediction; with physical IDSite, any potential sites of metabolism having a relative  $\Delta G^\ddagger$  lower than 4.75 kcal/mol is predicted to be a site of metabolism.

The second scoring model represents an empirically optimized version of the physical model described above with the following changes: (1) The PLOP energy ( $E_{\text{pose}}$ ) is not used directly but rescaled with two parameters as described below, which are fitted to a training set of 36 compounds. (2) Instead of obtaining the scaling coefficient for the methoxy radical intrinsic reactivities from the correlation in Figure 7, we fit it to the training set of 36 compounds. Note that the fitted value for the latter of 1.071 (eq 4) is very similar to the value obtained by correlating the DFT-activation energies (1.117), which further highlights the physical nature of this parameter. (3) The final selection criteria for predictions (score cutoff) were fit to the training set as well. All four fitted parameters were obtained from a fitting algorithm by maximizing the number of true positives over the sum of the numbers of false positives and false negatives.

$$E = \langle 1.071 \times \text{IR}(\text{methoxyradical}) + E_{\text{score}} \rangle - kT \ln N_{\text{H}} \quad (4)$$

As introduced above, instead of directly using the PLOP energy ( $E_{\text{pose}}$ ), eq 4 recalculates the binding contribution ( $E_{\text{score}}$ ) with a linear energy score; the angle brackets again represent Boltzmann averages. If a pose has a PLOP energy ( $E_{\text{pose}}$ ) within 5.26 kcal/mol from the lowest one, the energy score ( $E_{\text{score}}$ ) is zero; otherwise, it is 0.58 times the relative energy. The potential sites that have a relative score within 1.46 kcal/mol of a site predicted to have the highest reactivity are considered to be a site of metabolism.

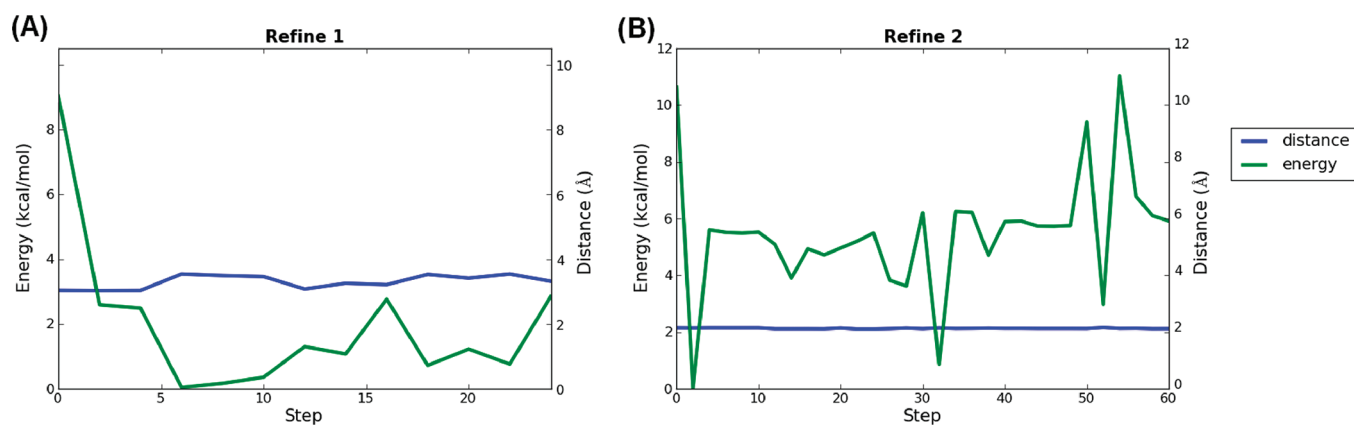
**Reactivity Model.** The sites at which a ligand gets metabolized by a P450 enzyme depends not only on whether the atom in question can approach the heme iron center with the correct geometry but also on the intrinsic chemical reactivity of the site. Assuming that the intrinsic chemical reactivities of the ligand sites are independent of the presence of the enzyme, we estimated the intrinsic reactivities from activation energies of a library of model systems using QM. Since DFT with the B3LYP functional and the 6-31G\* basis set has been shown to give high accuracy for relative energies of transition states,<sup>37</sup> while still allowing for fast calculations, we employed that level of theory for our intrinsic reactivity model. It has been shown that in general, an accurate linear correlation exists between the QM activation energies of hydrogen abstraction reactions with a methoxy radical and the corresponding hydrogen abstraction barriers with an iron–oxo porphyrin species, generally referred to as compound I in the P450 literature.<sup>34</sup> In agreement with previous reports,<sup>38,39</sup> we herein investigated the above-mentioned correlation including aliphatic hydrogen abstraction barriers as well as aromatic ones. As shown in Figure 7, we find good correlation between the methoxy radical and the compound I based activation barriers for both  $sp^2$  and  $sp^3$  hybridized systems ( $R^2 = 0.94$ ),



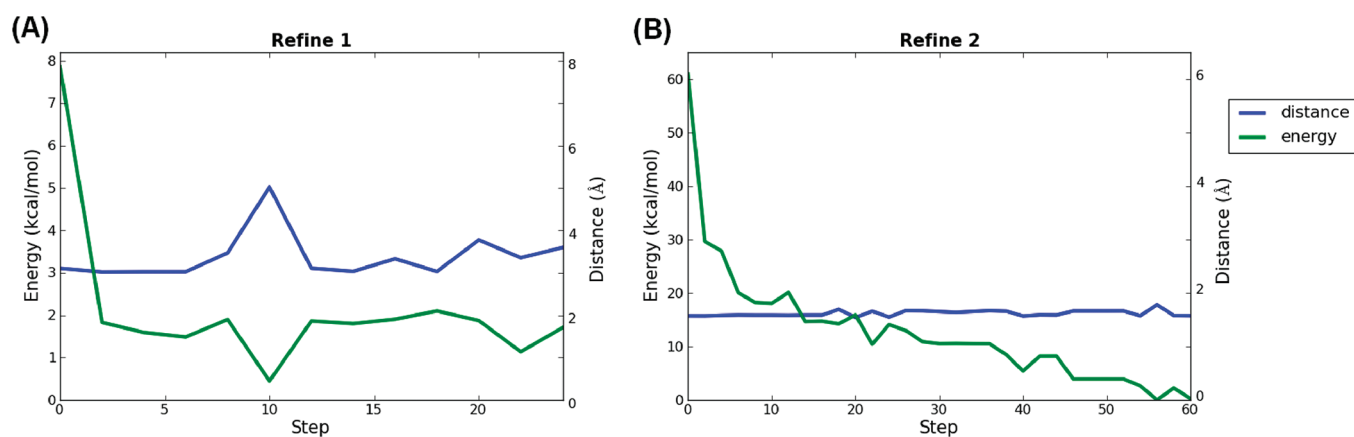
**Figure 10.** (A) ROC curves comparing the full IDSite method to the reduced methods. (B) ROC curves superimposed on the results of Sheridan et al.<sup>8</sup>

which validates the use of the methoxy radical model to estimate the intrinsic reactivities. Therefore, transition states for methoxy radical based hydrogen abstraction reactions were optimized at the B3LYP/LACVP\* level of theory with *Jaguar*<sup>27</sup> for a fragment library consisting of 150 model compounds, 483 distinct hydrogen atoms, and more than 2000 conformations, in order to accurately model all distinct chemical environments. Carbon atom based intrinsic reactivities were then assigned as the Boltzmann weighted activation energies over different transition state conformations. Intrinsic reactivities of the ligand sites were assigned using a simple SMARTS string matching algorithm of the fragment library. Thereby the best matching fragment was determined as the one with (1) the largest number of heavy atoms, (2) the most hydrogen atoms, and (3) the largest sum of atomic numbers.

**Preparation of Protein and Ligands.** The X-ray crystallographic structure of CYP2D6 was obtained from the Protein Data Bank (PDB ID: 2F9Q; 3.0 Å resolution) and contains a well-defined active site above the heme group.<sup>40</sup> We applied the



**Figure 11.** The energy and distance (constrained atom to the ferryl oxygen) changes during the MCM simulation during the first (A) and the second (B) refinement stages for 4-methoxyamphetamine.



**Figure 12.** The energy and distance (constrained atom to the ferryl oxygen) changes during the MCM simulation during the first (A) and the second (B) refinement stages for dextromethorphan.

Protein Preparation Wizard (PPW) of Schrödinger, Inc. to add hydrogen atoms, optimize the hydroxyl orientation, correct the Gln/Asn/His side-chain orientations, and determine the protonation states of titratable residues. PPW also assigned the bond order of the heme group and the iron oxidation state, which defines the iron atom as  $\text{Fe}^{3+}$  covalently bonded to the side chain of Cys443. The positions of all hydrogen atoms were optimized with a constraint of 0.3 Å with the OPLS 2005 force field.

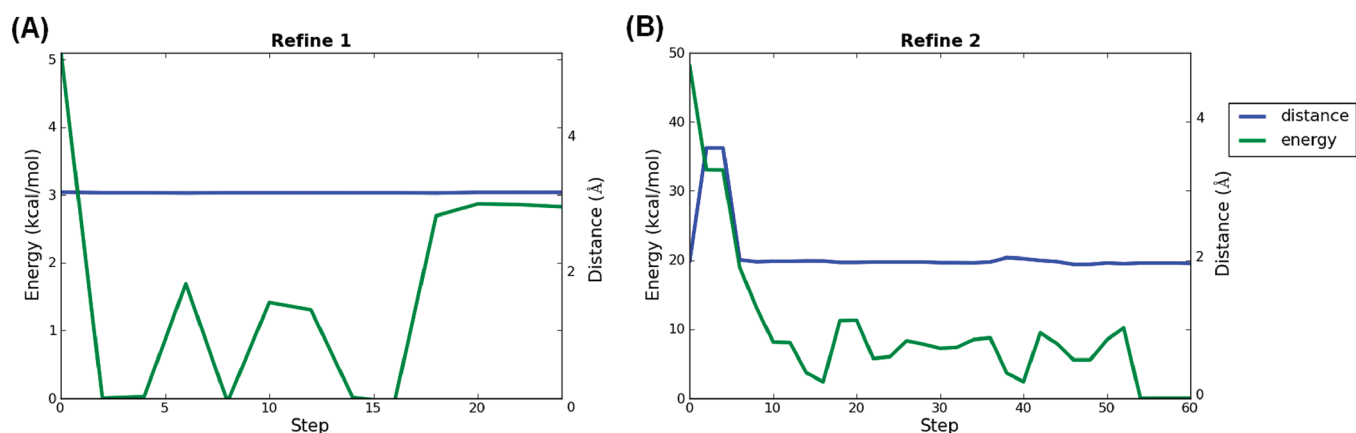
A training set of 36 compounds and a test set of 20 compounds were collected from the experimental literature.<sup>31,41</sup> These compounds mainly undergo O-dealkylation and hydroxylation by CYP2D6. The training and test sets contain 774 and 383 heavy atoms, respectively. Details about the data selection are explained in the Supporting Information. All stereoisomers used in the experiments were enumerated, as were the protonation states at pH = 7.0. All structures were minimized in vacuum using the OPLS 2005 force field, prior to the IDSite calculations.

## RESULTS AND DISCUSSION

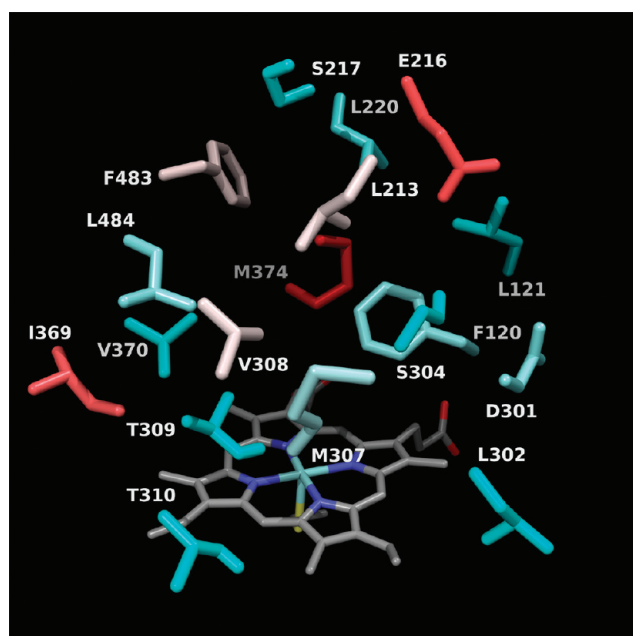
Tables 3 and 4 present the summary of our predicted results with the training set and the test set. The data show that IDSite has high sensitivity and specificity with both IDSite scoring models in predicting the 2D6-mediated metabolism of the 56

compounds: using the physical IDSite scoring, we achieve high sensitivity (0.83) and high specificity (0.98); using the fitted IDSite scoring, we can achieve even higher sensitivity (0.94) and similarly high specificity (0.99). With the fitted IDSite scoring, the results for the training set (sensitivity 0.91 and specificity 0.99) and test set (sensitivity 1.0 and specificity 0.98) are very similar, indicating that for the fitted model, no overfitting to the training set can be detected (see Figures 8 and 9 for IDSite predicted results for the training and test sets).

It is interesting to note that the principal effect of the parameter fitting is to reduce the number of false negatives; the reduction is of similar magnitude in both the training and test sets (there is also some reduction of false positives in the training set, but this is a less prominent result). The principal effect of the parametrization is to take into account the fact that there is some noise in the induced fit calculation energetics, reflected in the 5.26 kcal/mol energy window and scaling factor of 0.58. The noise is a combination of imperfect sampling and residual errors in the continuum solvent free energy model; the parameters suggest that there is a slight overestimation of the relative energetics of poses close in energy. Buffering and scaling the contribution from this term enables a (small) number of secondary sites to be recognized by the model as contributing to the reactivity, without increasing the number of false positives. As



**Figure 13.** The energy and distance (constrained atom to the ferryl oxygen) changes during the MCM simulation during the first (A) and the second (B) refinement stages for fluperlapine.



**Figure 14.** Illustration of the induced-fit effects modeled by IDSite. The cyan–white–red scheme is used to show the side chains from the least changed to the most changed, defined as the maximum mean absolute dihedral angle change for each residue.

noted above, the intrinsic reactivity appears to have less noise associated with it, which is not surprising in view of the fact that it poses a much less demanding sampling challenge.

A second question of interest is whether the various intensive sampling components of the algorithm actually improve the predictive capability. In order to analyze the importance of each sampling stage in IDSite, ROC (Receiver Operating Characteristic) curves were calculated (Figure 10A) to compare three reduced methods using the fitted score to the full method using *physical* and *fitted* scores. As mentioned in the Methods and Materials section, each refinement stage performs a constrained minimization, followed by sampling with MCM simulations. After Glide docking, the prediction can be made after minimization in the first refinement stage (referred to as “docking+minimization”), after the sampling in the first refinement stage

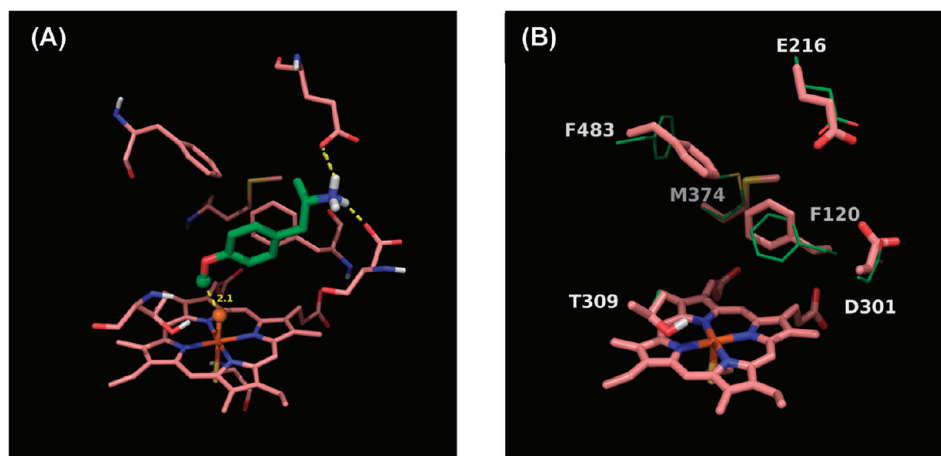
(referred to as “no Ref2”), or after the minimization in the second refinement stage (referred to as “no sampling in Ref2”). A higher energy cutoff (150 kcal/mol, instead of 24 kcal/mol) and distance cutoff (8.0 Å, instead of 2.6 Å for  $sp^3$  and 2.08 Å for  $sp^2$  hybridized atoms) are adjusted for the methods of “docking+minimization” and “no Ref2”. To draw the ROC curves, the scoring cutoff (4.75 and 1.46 kcal/mol are used for the results shown in Table 3 and 4 for the *physical* and *fitted* scores, respectively) is varied at a 0.5 kcal/mol interval from 0.0 to 100 kcal/mol, which represents the true positive rate ( $y$  axis) and the corresponding false positive rate ( $x$  axis) of the methods. The true positive rate and false positive rate are calculated according to eq 5,

$$\text{true positive rate} = \frac{\text{number of true positives}}{\text{number of SOMs observed in experiments}} \quad (5a)$$

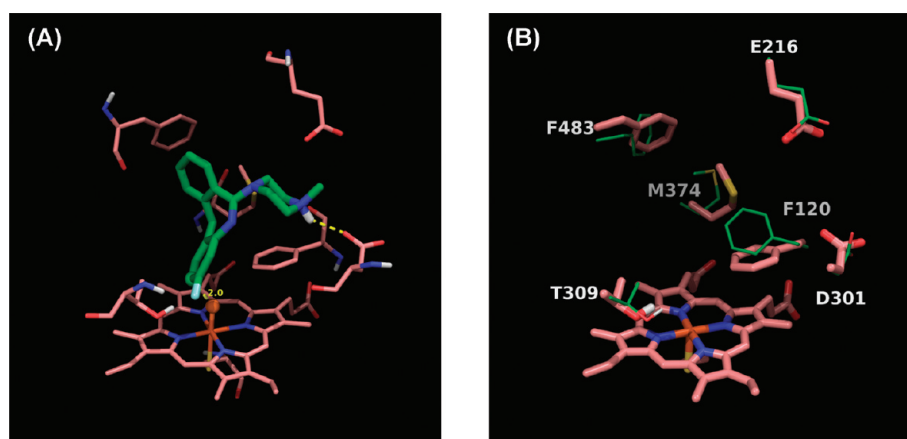
$$\text{false positive rate} = \frac{\text{number of false positives}}{\text{number of non-SOMs observed in experiments}} \quad (5b)$$

where true positives are the SOMs ( $sp^2$  and  $sp^3$  carbon atoms which undergo hydroxylation or O-dealkylation) identified by experiments as well as predicted correctly by IDSite, and the false positives are non-SOMs (nonhydrogen atoms) but mispredicted by IDSite as hydroxylated/dealkylated by CYP2D6. As currently we mainly focus on the typical CYP2D6-mediated hydroxylation and O-dealkylation involving  $sp^2$  and  $sp^3$  carbon atoms with bonded hydrogen atoms, those sites (carbon atoms or heteroatoms) which potentially undergo other metabolic reactions such as N-dealkylation and oxidation are currently considered non-SOMs in our preliminary study.

The ROC curves in Figure 10A indicate that at the same false positive rate (sensitivity), the false negative rate decreases with more sampling, and the full IDSite method always has the lowest false negative rate (the highest specificity) with both scoring models. It is interesting that the *physical* score derived from the basic physical chemistry model is very close to the *fitted* score. For the reduced methods, there is an obvious trend that increasing the sampling efforts yields substantially higher specificity at each stage. This means that using the IDSite scoring models in conjunction with binding requirements, sufficient sampling in



**Figure 15.** (A) The lowest energy pose in the second refinement stage for 4-methoxyamphetamine. Orange sphere = “dummy” ferryl oxygen, green sphere = experimental and predicted SOM. (B) Comparison of side chains important for induced-fit effects. Crystal structure (green, PDB ID: 2F9Q) minimized with the VSGB 2.0 model and superimposed onto the lowest energy pose with 4-methoxyamphetamine (salmon). Large dihedral changes are seen for Asp301 ( $\Delta\chi_2$ ,  $121^\circ$ ), Met374 ( $\Delta\chi_3$ ,  $114^\circ$ ), and Phe483 ( $\Delta\chi_1$ ,  $60^\circ$ ).



**Figure 16.** (A) The lowest energy pose in the second refinement stage for fluperlapine. Orange sphere = “dummy” ferryl oxygen, green sphere = experimental and predicted SOM. (B) Comparison of side chains important for induced fit effects. Crystal structure (green, PDB ID: 2F9Q) minimized with the VSGB 2.0 model and superimposed onto the lowest energy pose with fluperlapine (salmon). Large dihedral changes are seen for Phe120 ( $\Delta\chi_2$ ,  $73^\circ$ ), Glu216 ( $\Delta\chi_1$ ,  $60^\circ$ ), Asp301 ( $\Delta\chi_2$ ,  $64^\circ$ ), Met374 ( $\Delta\chi_3$ ,  $105^\circ$ ), and Phe483 ( $\Delta\chi_2$ ,  $94^\circ$ ).

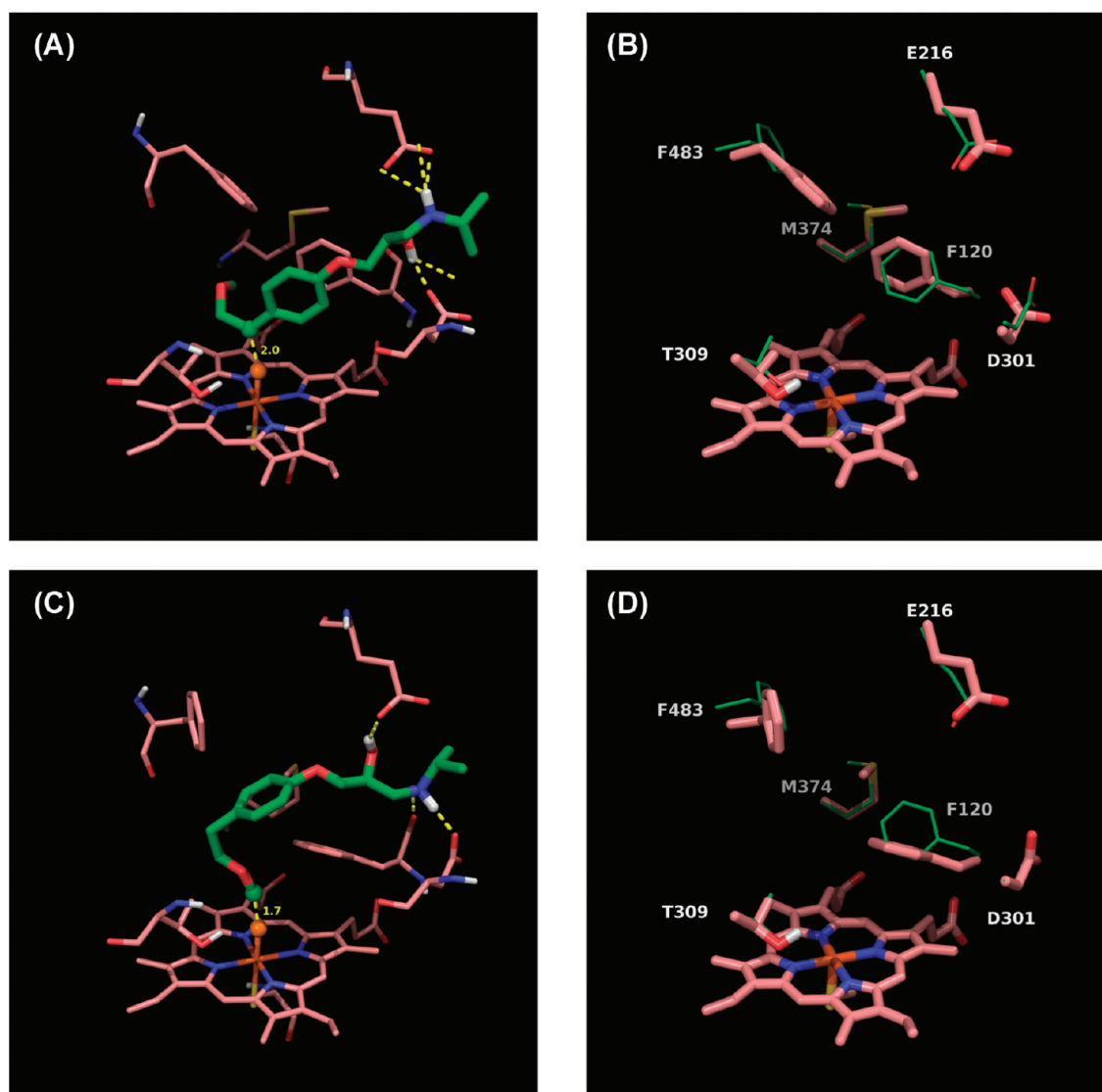
IDSite can specifically identify the sites metabolism observed by the experiments.

In Figure 10B, we compare the physical IDSite and fitted IDSite results to results from Sheridan et al.,<sup>8</sup> who evaluated true positive and false positive rates, using the same ROC metric that we employ, for their test set of CYP2D6 ligands. The test set employed in ref 8 is different in detail from the one we use here, but the types of ligands in both test sets are similar on the basis of examples of test set molecules given in ref 8. Hence, while the comparison is not completely rigorous, it is a reasonable way to estimate relative performance. It can be seen that given the caveat above, both physical IDSite and fitted IDSite substantially outperform both MetaSite and the in-house Merck QSAR-based approached plotted in Figure 10B. To recover 90% of true positives, the QSAR method included roughly 20% false positives, whereas MetaSite included 40% false positives. In contrast, IDSite incorporated only  $\sim 1\%$  false positives. This is a qualitative transformation of performance that has significant implications for use in drug discovery applications, as does the availability of a predicted three-dimensional structure that is likely to be quite accurate.

So far, only the apo enzyme structure of CYP2D6 has been determined by X-ray crystallography. In order to investigate the capability of IDSite in modeling the induced-fit effects and to understand the effects of the hierarchical sampling, several compounds of various sizes and flexibility were selected to analyze the structural and energetic changes at each stage.

It is very common that the poses from docking that have the SOM close to the ferryl oxygen are not among the top poses considered by *Glide* SP scoring. For example, the pose with the shortest distance (1.8 Å) is ranked sixth in the case of 4-methoxyamphetamine; the pose (1.4 Å) that leads to the prediction of O-demethylation is ranked 20th for the case of metoprolol. Further, it is also possible for some cases (e.g., fluperlapine) that none of the poses have the SOM close enough to the ferryl oxygen. Therefore, it is very difficult to make specific predictions with only a small distance cutoff and a few top poses from docking. In order to improve the sensitivity as well as the specificity of the predictions, it appears to be necessary to employ the refinement stages.

Focusing on the distance between the site(s) of metabolism observed experimentally, we investigated the Boltzmann averaged



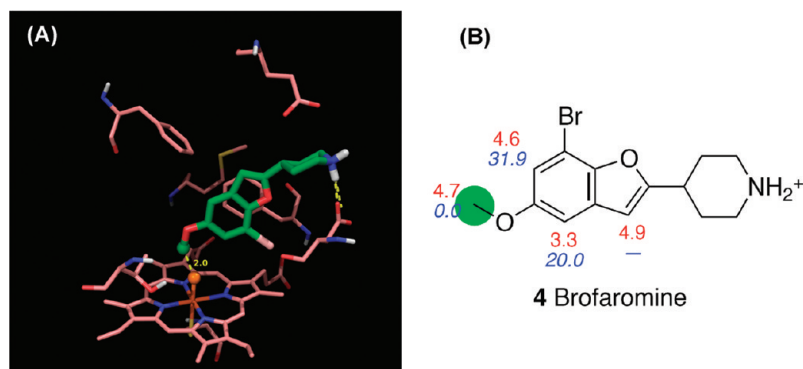
**Figure 17.** (A) The lowest energy poses in the second refinement stage for metoprolol benzylic hydroxylation. (B) Comparison of side chains important for induced fit effects for metoprolol benzylic hydroxylation. (C) The lowest energy poses in the second refinement stage for metoprolol O-dealkylation. (D) Comparison of side chains important for induced fit effects for metoprolol O-dealkylation. For A and C, orange spheres = “dummy” ferryl oxygen, green spheres = experimental and predicted SOMs. For B and D, crystal structure (green, PDB ID: 2F9Q) minimized with the VSGB 2.0 model and superimposed onto the lowest energy poses with metoprolol (salmon). For benzylic hydroxylation, large dihedral changes are seen for Glu216 ( $\Delta\chi_1$ ,  $60^\circ$ ), Asp301 ( $\Delta\chi_2$ ,  $66^\circ$ ), Met374 ( $\Delta\chi_3$ ,  $112^\circ$ ), and Phe483 ( $\Delta\chi_1$ ,  $40^\circ$ ); for O-dealkylation, large dihedral changes are seen for Phe120 ( $\Delta\chi_2$ ,  $67^\circ$ ), Glu216 ( $\Delta\chi_2$ ,  $50^\circ$ ), and Phe483 ( $\Delta\chi_2$ ,  $194^\circ$ ).

energy and distance from the site(s) to the ferryl oxygen over all of the poses sampled at any even numbered step (see Figures 11–13). Given the strong harmonic constraints applied in the refinement stages, the distance change is generally relatively small, as expected. The energy change in the first refinement is usually small ranging from 4 to 25 kcal/mol. However, the energy change during the second refinement stage is quite different for small ligands as compared to large ligands. For ligands as small as 4-methoxyamphetamine, the energy of the poses fluctuated within the range of 12 kcal/mol, and the lowest energy structure was obtained at the early steps. In contrast, the energy can decrease by more than 60 kcal/mol during the sampling of the second refinement stage for flexible or bulky ligands such as fluperlapine. For such cases, it is often not until the end of the simulation that the low energy structure is sampled.

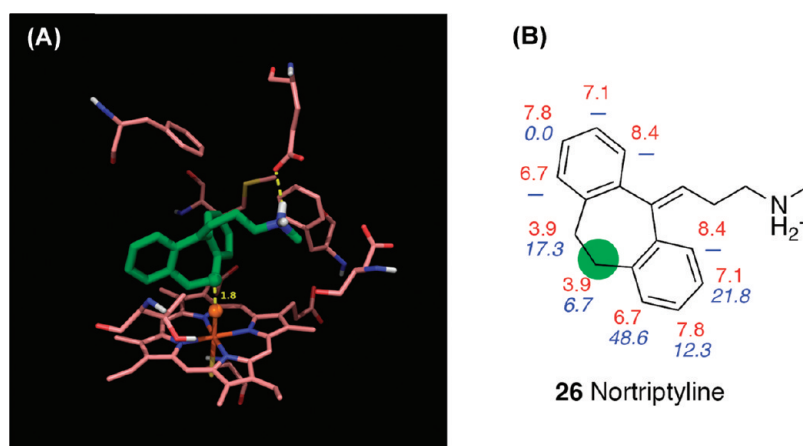
This implies that the second refinement plays an important role in optimizing the structure for bulky or flexible compounds.

Skipping the second refinement, about 40% of the compounds (24/56) in the training set have the same results as obtained from the full protocol, and most of them are small compounds like 4-methoxyamphetamine, MDMA, MAMC, etc. This observation is consistent with our discussion above that links the need for extended refinement to the presence of large, bulky ligands where protein induced-fit effects are significant, and where optimization of the free energy of the reactive binding complex can pose great difficulties due to various types of energy barriers and additional degrees of freedom to explore in the ligand.

**Analysis of Induced-Fit Effects.** P450 enzymes are believed to have high flexibility in adjusting their active site to accommodate a large variety of substrates.<sup>42</sup> In order to model such



**Figure 18.** (A) The lowest energy pose in the second refinement stage for brofaromine. Orange sphere = “dummy” ferryl oxygen, green sphere = experimental and predicted SOM. (B) Intrinsic reactivities (red) for each site and the relative energy (blue) of the poses with the corresponding site constrained to the ferryl oxygen. The SOM observed experimentally is marked with a green circle.



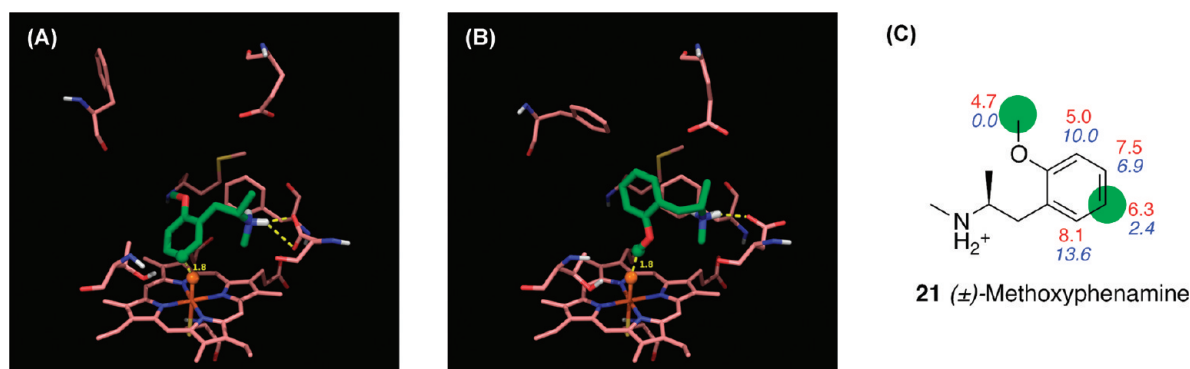
**Figure 19.** (A) The lowest energy pose in the second refinement stage for nortriptyline. Orange sphere = “dummy” ferryl oxygen, green sphere = experimental and predicted SOM. (B) Intrinsic reactivities (red) for each site and the relative energy (blue) of the poses with the corresponding site constrained to the ferryl oxygen. The SOM observed experimentally is marked with a green circle.

induced fit effects, sufficient sampling provided by the two refinement stages of IDSite is critical, as demonstrated in the previous section. In order to further investigate the capability of IDSite in modeling induced-fit effects, we calculated the average absolute change for each dihedral angle of the protein side chains in the binding box in comparison to the minimized crystal structure of CYP2D6. The largest change of all the  $\chi$  angles for each residue is used to represent the change for that residue. Figure 14 illustrates the induced-fit effects by showing the largest change for each residue. Ten of the 18 residues in the binding box have changes greater than  $30^\circ$ . This shows that IDSite is able to model induced-fit effects required to correctly identify the “bio-active” conformation of the ligands by changing the side-chain orientations in the active site. Phe120 and Phe483 with bulky side chains have changes as large as  $40^\circ$  and  $60^\circ$ , respectively. However, the magnitude of their induced fit effects depends highly on the ligand size. Between these two Phe residues in the binding box, Met374 has the most significant change ( $108^\circ$ ) because a small rotation in the Phe side chains can cause a big adjustment in Met374. Compared to the large change of Glu216 ( $88^\circ$ ), the change of Asp301 ( $38^\circ$ ) is relatively smaller due to the shorter side chain.

The above-mentioned trends are illustrated in Figures 15–17, which compare the docked structures leading to the SOM of 4-methoxyamphetamine (PMA), fluperlapine, and metoprolol to

the crystal structure of the apoenzyme minimized with the VSGB 2.0 energy model. Analogous figures can be found for all of our predictions in the Supporting Information. One striking example of induced-fit effects involves Phe120. For small ligands such as PMA, the benzene ring conformation of Phe120 changes only slightly (Figure 15), while it has to move out of the way for larger ligands such as fluperlapine (Figure 16) or metoprolol (Figure 17), therefore rotating by almost  $90^\circ$ . Interestingly, for compounds with multiple sites of metabolism, such as metoprolol (Figure 17), different binding modes leading to different SOMs have very different conformations of the Phe120 side chain as well. Our IDSite docked structures clearly highlight the importance of induced fit effects for CYP2D6 metabolism and therefore explain why it is difficult to accurately predict SOMs with a rigid receptor model.

**The Importance of Structural Effects in Determining SOMs.** The two main competing factors in determining the SOMs with P450 enzymes are the intrinsic reactivities of the ligand sites and the geometric fit of the ligand in the active site. As mentioned in the Methods and Materials section, IDSite considers both of these effects in determining the SOMs, which enables it to select the correct SOM even for difficult cases, where the intrinsic reactivity favors the nonsite of metabolism. For these cases, the structural fit of the ligand with the receptor, i.e., how easily the ligand site can reach the ferryl oxygen, mainly determines



**Figure 20.** The lowest energy pose in the second refinement stage for methoxyphenamine. Orange sphere = “dummy” ferryl oxygen, green sphere = experimental and predicted SOM. (A) Aromatic hydroxylation. (B) O-demethylation. (C) Intrinsic reactivities (red) for each site and the relative energy (blue) of the poses with the corresponding site constrained to the ferryl oxygen. The SOM observed experimentally is marked with a green circle.

the SOM. Therefore, the structures and energies of the poses, with consideration of the receptor, have to be utilized. Three cases are used here to demonstrate the role of a receptor (CYP2D6) in determining the sites of metabolism.

The first case we discuss is brofaromine, for which experiments show that the major metabolic pathway is O-demethylation mediated by CYP2D6.<sup>43</sup> The intrinsic reactivity of the site of metabolism (4.7 kcal/mol) is very close to those of sites on the aromatic rings (nonsites of metabolism, 3.3–4.9 kcal/mol; Figure 18). Due to the receptor geometry, it is impossible for the atoms on the furan ring to get close to the ferryl oxygen while still attaining the salt bridge with either Glu216 or Asp301. Therefore, no qualified poses were found leading to a reaction on the furan ring. Although we found qualified poses for all of the sites on the benzene ring, those poses are all strongly disfavored energetically by more than 20 kcal/mol. This indicates that taking the interactions between the ligand and the receptor into account, IDSite is able to make the prediction of the SOM for brofaromine in good agreement with the experimental observation.

A second interesting case is nortriptyline (Figure 19), since the two sites on the seven-membered aliphatic ring are difficult to distinguish only with their intrinsic reactivity, as they are almost equally reactive. However, experiments show that only the (*E*)-10 site of nortriptyline is metabolized.<sup>44</sup> In IDSite, the poses with the (*Z*)-10 site close to the ferryl oxygen are all at least 10 kcal/mol higher in energy compared to the poses with the (*E*)-10 atom close to the ferryl oxygen. Such an energy gap is large enough for IDSite to correctly determine the (*E*)-isomer as the only metabolite. While structural effects are therefore clearly very important to determine nortriptylene’s SOM, the intrinsic reactivities also play a key role. This is again nicely illustrated with the example of nortriptyline, where a simply structure based method (without considering intrinsic reactivities) would predict the SOM as being an aromatic hydroxylation due to the favorable energy of the corresponding poses. Therefore, IDSite is able to correctly balance the subtle effects stemming from intrinsic reactivity and structural fit.

Methoxyphenamine is another case where the joint effects of intrinsic reactivity and the structural fit lead to the correct predictions. Methoxyphenamine is metabolized through O-demethylation and aromatic hydroxylation mediated by CYP2D6.<sup>45</sup> These two sites not only have very close intrinsic reactivities (5.7 and 6.3 kcal/mol, Figure 20), but their lowest energy poses

also have very similar energies. The non-SOMs are not selected by IDSite, either because of unfavorable intrinsic reactivity or because of high pose energies.

**Computational Times.** On a single 2.2 GHz AMD Opteron Processor 6174, the average CPU time required for a typical IDSite calculation (e.g., with a compound with three rotatable bonds) is about 448 h, of which about 11% of the time is spent on the first refinement stage and 89% on the second refinement. On 20 such processors, the calculation takes 22 h. The initial Glide docking step on a single processor takes about 10 min. The computational cost of PLOP refinement is proportional to the number of rotatable bonds in the compound.

## CONCLUSION

We have developed a novel approach for the prediction of experimentally observable cytochrome P450 sites of metabolism, IDSite, and applied it to a data set for the 2D6 P450 isoform. We obtain remarkably high sensitivity and specificity using a structure-based model, representing a major advance as compared to alternatives in the literature, including various types of ligand-based models. The method delivers not only accurate SOM predictions but also three-dimensional structures of the protein–ligand complex, including induced fit effects (which are quite significant), for every SOM identified by the algorithm.

We selected 2D6 as our initial target because the binding of a ligand positive nitrogen to an acidic group in the protein created an additional constraint that was useful in limiting sampling and achieving reliable poses in the induced fit docking effort. Other important P450 isoforms, such as 1A2, 2C9, and 3A4, may be more difficult to model in this fashion, as they lack such a salt bridge constraint; nevertheless, even if additional sampling effort is required, it should be possible to obtain successful results given the performance of the conformational energy and reactivity models that we have seen in the present work. The development of models for additional isoforms, and additional ligand test sets, is ongoing in our laboratory. Ultimately, predictive use in an active drug discovery project will be required for validation; we look forward to engaging in such tests in the near future.

## ASSOCIATED CONTENT

**S Supporting Information.** Details about our data set selection, figures of all docked structures leading to our predictions, tables with



the activation barriers used to draw Figure 7, and tables with detailed dihedral angle changes due to induced-fit effects. This information is available free of charge via the Internet at <http://pubs.acs.org>.

## AUTHOR INFORMATION

### Corresponding Author

\*E-mail: [rich@chem.columbia.edu](mailto:rich@chem.columbia.edu).

## ACKNOWLEDGMENT

This work was supported by the NIH grant GM-40526 to R.A.F. S.T.S. also thanks the Guthikonda Family for an Arun Guthikonda Memorial Graduate Fellowship. We thank Professor Ronald Breslow, Dr. Tyler Day, Dr. Robert Abel, Dr. Zhiyong Zhou, Michelle L. Hall, and Jing Zhang for the helpful discussions. R.A.F. has a significant financial stake in Schrödinger, Inc., is a consultant to Schrödinger, Inc., and is on the Scientific Advisory Board of Schrödinger, Inc.

## REFERENCES

- (1) Bailey, D. G.; Malcolm, J.; Arnold, O.; Spence, J. D. *Br. J. Clin. Pharmacol.* **1998**, *46*, 101.
- (2) Preskorn, S. H. *Clin. Pharmacokinet.* **1997**, *32*, 1.
- (3) Dresser, G. K.; Spence, J. D.; Bailey, D. G. *Clin. Pharmacokinet.* **2000**, *38*, 41.
- (4) Afzelius, L.; Arnby, C. H.; Broo, A.; Carlsson, L.; Isaksson, C.; Jurva, U.; Kjellander, B.; Kolmodin, K.; Nilsson, K.; Raubacher, F.; Weidolf, L. *Drug Metab. Rev.* **2007**, *39*, 61.
- (5) Singh, S. B.; Shen, L. Q.; Walker, M. J.; Sheridan, R. P. *J. Med. Chem.* **2003**, *46*, 1330.
- (6) Rydberg, P.; Gloriam, D. E.; Zaretski, J.; Breneman, C.; Olsen, L. *ACS Med. Chem. Lett.* **2010**, *1*, 96.
- (7) de Groot, M. J.; Alex, A. A.; Jones, B. C. *J. Med. Chem.* **2002**, *45*, 1983.
- (8) Sheridan, R. P.; Korzekwa, K. R.; Torres, R. A.; Walker, M. J. *J. Med. Chem.* **2007**, *50*, 3173.
- (9) Zaretski, J.; Bergeron, C.; Rydberg, P.; Huang, T.-w.; Bennett, K. P.; Breneman, C. M. *J. Chem. Inf. Model.* **2011**, *51*, 1667.
- (10) Cruciani, G.; Carosati, E.; De Boeck, B.; Ethirajulu, K.; Mackie, C.; Howe, T.; Vianello, R. *J. Med. Chem.* **2005**, *48*, 6970.
- (11) Jones, J. P.; Korzekwa, K. R. In *Methods Enzymol.*; Eric, F. J., Michael, R. W., Eds.; Academic Press: New York, NY, 1996; Vol. 272, p 326.
- (12) Oláh, J.; Mulholland, A. J.; Harvey, J. N. *Proc. Natl. Acad. Sci. U. S. A.* **2011**, *108*, 6050.
- (13) Kirton, S. B.; Kemp, C. A.; Tomkinson, N. P.; St.-Gallay, S.; Sutcliffe, M. J. *Proteins: Struct., Funct., Bioinf.* **2002**, *49*, 216.
- (14) de Graaf, C.; Oostenbrink, C.; Keizers, P. H. J.; van der Wijst, T.; Jongejan, A.; Vermeulen, N. P. E. *J. Med. Chem.* **2006**, *49*, 2417.
- (15) Unwalla, R.; Cross, J.; Salaniwal, S.; Shilling, A.; Leung, L.; Kao, J.; Humblet, C. *J. Comput.-Aided Mol. Des.* **2010**, *24*, 237.
- (16) Vasanathanathan, P.; Hritz, J.; Taboureau, O.; Olsen, L.; Jorgensen, F. S.; Vermeulen, N. P. E.; Oostenbrink, C. *J. Chem. Inf. Model.* **2009**, *49*, 43.
- (17) Rydberg, P.; Hansen, S. M.; Kongsted, J.; Norrby, P. O.; Olsen, L.; Ryde, U. *J. Chem. Theory Comput.* **2008**, *4*, 673.
- (18) Gleeson, M. P.; Davis, A. M.; Chohan, K. K.; Paine, S. W.; Boyer, S.; Gavaghan, C. L.; Arnby, C. H.; Kankkonen, C.; Albertson, N. *J. Comput.-Aided Mol. Des.* **2007**, *21*, 559.
- (19) Li, J.; Abel, R.; Zhu, K.; Cao, Y.; Friesner, R. *Proteins: Struct., Funct., Bioinf.* **2011**, *79*, 2794.
- (20) Bathelt, C. M.; Mulholland, A. J.; Harvey, J. N. *J. Phys. Chem. A* **2008**, *112*, 13149.
- (21) Tian, L.; Friesner, R. A. *J. Chem. Theory Comput.* **2009**, *5*, 1421.
- (22) *Glide*, version 5.6; Schrödinger, Inc.: New York, NY, 2010.
- (23) *Prime*, version 3.0; Schrödinger, Inc.: New York, NY, 2011.

- (24) Halgren, T. A.; Murphy, R. B.; Friesner, R. A.; Beard, H. S.; Frye, L. L.; Pollard, W. T.; Banks, J. L. *J. Med. Chem.* **2004**, *47*, 1750.
- (25) Friesner, R. A.; Banks, J. L.; Murphy, R. B.; Halgren, T. A.; Klicic, J. J.; Mainz, D. T.; Repasky, M. P.; Knoll, E. H.; Shelley, M.; Perry, J. K.; Shaw, D. E.; Francis, P.; Shenkin, P. S. *J. Med. Chem.* **2004**, *47*, 1739.
- (26) Paine, M. J. I.; McLaughlin, L. A.; Flanagan, J. U.; Kemp, C. A.; Sutcliffe, M. J.; Roberts, G. C. K.; Wolf, C. R. *J. Biol. Chem.* **2003**, *278*, 4021.
- (27) *Jaguar*, version 7.6; Schrödinger, Inc.: New York, NY, 2010.
- (28) Xiang, Z. X.; Honig, B. *J. Mol. Biol.* **2001**, *311*, 421.
- (29) Duane, S.; Kennedy, A. D.; Pendleton, B. J.; Roweth, D. *Phys. Lett. B* **1987**, *195*, 216.
- (30) Tuckerman, M.; Berne, B. J.; Martyna, G. J. *J. Chem. Phys.* **1992**, *97*, 1990.
- (31) Wang, B.; Yang, L.-P.; Zhang, X.-Z.; Huang, S.-Q.; Bartlam, M.; Zhou, S.-F. *Drug Metab. Rev.* **2009**, *41*, 573.
- (32) Guengerich, F. P.; Miller, G. P.; Hanna, I. H.; Martin, M. V.; Leger, S.; Black, C.; Chauret, N.; Silva, J. M.; Trimble, L. A.; Yergey, J. A.; Nicoll-Griffith, D. A. *Biochemistry* **2002**, *41*, 11025.
- (33) Guengerich, F. P.; Hanna, I. H.; Martin, M. V.; Gillam, E. M. J. *Biochemistry* **2003**, *42*, 1245.
- (34) Guengerich, F. P. *Chem. Res. Toxicol.* **2001**, *14*, 611.
- (35) Shaik, S.; Kumar, D.; de Visser, S. P.; Altun, A.; Thiel, W. *Chem. Rev.* **2005**, *105*, 2279.
- (36) Wang, Y. H.; Li, Y.; Wang, B. *J. Phys. Chem. B* **2007**, *111*, 4251.
- (37) Schneebeli, S. T.; Hall, M. L.; Breslow, R.; Friesner, R. *J. Am. Chem. Soc.* **2009**, *131*, 3965.
- (38) Olsen, L.; Rydberg, P.; Rod, T. H.; Ryde, U. *J. Med. Chem.* **2006**, *49*, 6489.
- (39) Rydberg, P.; Ryde, U.; Olsen, L. *J. Phys. Chem. A* **2008**, *112*, 13058.
- (40) Rowland, P.; Blaney, F. E.; Smyth, M. G.; Jones, J. J.; Leydon, V. R.; Oxbrow, A. K.; Lewis, C. J.; Tennant, M. G.; Modi, S.; Eggleston, D. S.; Chenery, R. J.; Bridges, A. M. *J. Biol. Chem.* **2006**, *281*, 7614.
- (41) de Groot, M. J.; Ackland, M. J.; Horne, V. A.; Alex, A. A.; Jones, B. C. *J. Med. Chem.* **1999**, *42*, 1515.
- (42) Lill, M. A.; Dobler, M.; Vedani, A. *ChemMedChem* **2006**, *1*, 73.
- (43) Feifel, N.; Kucher, K.; Fuchs, L.; Jedrychowski, M.; Schmidt, E.; Antonin, K. H.; Bieck, P. R.; Gleiter, C. H. *Eur. J. Clin. Pharmacol.* **1993**, *45*, 265.
- (44) Olesen, O. V.; Linnet, K. *Drug Metab. Dispos.* **1997**, *25*, 740.
- (45) Geertsen, S.; Foster, B. C.; Wilson, D. L.; Cyr, T. D.; Casley, W. *Xenobiotica* **1995**, *25*, 895.

## NOTE ADDED AFTER ASAP PUBLICATION

This paper was published on the Web on September 29, 2011, with the second half of Figure 8 missing. The corrected version was reposted on October 5, 2011.

# A Simple Mechanism Underlying the Effect of Protecting Osmolytes on Protein Folding

G. Saladino,<sup>†</sup> M. Marenchino,<sup>‡</sup> S. Pieraccini,<sup>†,§</sup> R. Campos-Olivas,<sup>‡</sup> M. Sironi,<sup>\*,†,§,||</sup> and F. L. Gervasio<sup>\*,‡</sup>

<sup>†</sup>Dipartimento di Chimica Fisica ed Elettrochimica, Università degli Studi di Milano, Via Golgi 19, 20133 Milano, Italy

<sup>‡</sup>Structural Biology and Biocomputing Programme, Spanish National Cancer Research Centre (CNIO), c/Melchor Fernandez Almagro 3, 28029, Madrid, Spain

<sup>§</sup>INSTM Research Unit, Via Golgi 19, 20133 Milano, Italy

<sup>||</sup>Institute of Molecular Science and Technology, Via Golgi 19, 20133 Milano

**S** Supporting Information

**ABSTRACT:** Osmolytes are small organic compounds that confer to the cell an enhanced adaptability to external conditions. Many osmolytes not only protect the cell from osmotic stress but also stabilize the native structure of proteins. While simplified models able to predict changes to protein stability are available, a general physicochemical explanation of the underlying microscopic mechanism is still missing. Here, we address this issue by performing very long all-atom MD simulations, free energy calculations, and experiments on a well-characterized mini-protein, the villin headpiece. Comparisons between the folding free energy landscapes in pure water and osmolyte solutions, together with experimental validation by means of circular dichroism, unfolding experiments, and NMR, led us to formulate a simple hypothesis for the protecting mechanism. Taken together, our results support a novel mechanistic explanation according to which the main driving force behind native state protection is a change in the solvent rotational diffusion.

## INTRODUCTION

Severe environmental conditions, such as extreme temperatures, high osmotic pressure, or high concentrations of urea tend to cause cellular water stress. Many organisms have evolved to respond to these conditions regulating the level of small organic compounds, called osmolytes.<sup>1</sup> Osmolytes have been observed in a wide range of organisms<sup>2</sup> and have been found to accumulate in some species able to survive under harsh conditions,<sup>3,4</sup> such as the so-called “resurrection plants”, able to survive under severe drought.<sup>5</sup> In addition to their ability to control cell water loss or gain,<sup>6–8</sup> some osmolytes are also able to stabilize the native fold of proteins.<sup>9</sup> Bolen and co-workers carefully characterized osmolyte-induced thermotolerance,<sup>10,11</sup> due to the alteration of folded–unfolded equilibria. They also demonstrated that trimethylamine N-oxide (TMAO)<sup>11</sup> can fold natively unfolded proteins. Despite the wide variety of proteins in living organisms, only a few, generally interchangeable,<sup>12</sup> osmolyte molecules exist,<sup>1,13</sup> suggesting a universal underlying mechanism. Contrasting theories have been proposed involving either direct<sup>14,15</sup> or indirect interactions with proteins,<sup>16,17</sup> with the latter one recently prevailing due to the observed exclusion of osmolytes, with the exception of the denaturing urea,<sup>18</sup> from the protein surface, a phenomenon coined the “osmophobic effect”.<sup>19</sup> Recently, we used simulations and free energy methods to study the effect of the osmolyte glycine betaine (GB) on a small  $\beta$ -hairpin peptide, observing the expected increased stability of the native fold.<sup>20,21</sup> Nevertheless, a simple yet universal explanation of the microscopic mechanism of osmoprotection has not been found. In the search for such a general explanation, here, we combine simulations and experiments to study the effect of osmolytes on a more realistic and

well-characterized mini-protein, the human villin headpiece C-terminal helical subdomain (HP35).<sup>22</sup> HP35 has a well-defined secondary and tertiary structure and is one of the smallest peptides that folds cooperatively.<sup>23</sup> It has been the subject of several computational<sup>24,25</sup> and experimental<sup>26–29</sup> studies. In the following, the effect of two different osmoprotectants and urea on the folding of HP35 were investigated using 1.5- $\mu$ s-long unbiased all-atoms MD simulations and massive bias exchange molecular dynamics simulations (BEMD),<sup>30</sup> as well as calorimetry, circular dichroism (CD), and NMR experiments. The experimentally validated free energies, together with a careful structural analysis, allowed us to outline a clear and simple picture of the osmolyte protecting mechanism.

## MATERIALS AND METHODS

The HP35 structure was retrieved from the Protein Data Bank (PDB code: 1UNC).<sup>22</sup> The protein was solvated with TIP3P water molecules<sup>31</sup> in a 50 Å cubic box and neutralized with Cl<sup>−</sup> ions. To obtain the mixed-solvent systems, an appropriate number of water molecules was replaced with GB or TMAO molecules to obtain a 1 M solution. Simulations were run using the GROMACS<sup>32</sup> package combined with the PLUMED<sup>33</sup> plugin, which implements BEMD. As in other collective variables (CV)-based techniques, biasing the evolution of the system along a few variables approximating the reaction coordinate, the convergence of metadynamics can be severely affected by neglecting slow CVs. The BEMD method complements the

**Received:** July 8, 2011

**Published:** September 20, 2011

metadynamics technique introducing a replica exchange algorithm, compensating for this eventual neglect and allowing for a larger number of CVs with respect to standard metadynamics. Albeit BEMD, at difference with the more computationally expensive PTmetaD,<sup>21</sup> might have convergence problems in complex systems, it has already been shown to converge well in the case of HP35 folding.<sup>34</sup> What is more, we have carefully checked the convergence of the free energy profiles reconstructed from the blank replica as a function of simulation time. The Amber99SB\*-ILDN<sup>35,36</sup> force field was used, including backbone corrections.<sup>37</sup> Particle-mesh Ewald was used with a cutoff of 0.8 nm. All bond lengths were constrained to equilibrium distances using the LINCS<sup>38</sup> algorithm. After minimization, the systems were relaxed with 1 ns NPT dynamics, at 320 K and 1 atm, using the V-Rescale<sup>39</sup> algorithm for temperature coupling and a Berendsen barostat.<sup>40</sup> The BEMD runs were performed with the same collective variables (CVs) used in ref 34, i.e., the number of backbone hydrogen bonds, salt bridges, and hydrophobic contacts; the correlation of the backbone dihedral angles; and the fraction of secondary structure, and a neutral replica, on which no bias was applied. Each BEMD simulation required considerably longer simulations than those used in ref 34 to converge (>300 ns). This is most probably due to the different version of the Amber force-field used. Analyses were performed on the neutral replica, whose free energy profiles along the CV were reconstructed from the unbiased probability distribution of the states. The free energy of unfolding was calculated by integrating the density of the folded (F) and unfolded (U) states according to

$$\Delta G_{\text{unfold}} = k_B T \log \left( \frac{\int_{\text{F}} ds \exp\left(-\frac{G(s)}{k_B T}\right)}{\int_{\text{U}} ds \exp\left(-\frac{G(s)}{k_B T}\right)} \right) \quad (1)$$

Preferential coefficients ( $\Gamma_{\text{XP}}$ ) were calculated using the approach developed by Baynes and Trout.<sup>41</sup> According to this approach,  $\Gamma_{\text{XP}}$  can be evaluated defining two domains, a bulk domain (I) and a protein domain (II), and calculating

$$\Gamma_{\text{XP}} = \left\langle n_{\text{X}}^{\text{II}} - n_{\text{W}}^{\text{II}} \left( \frac{n_{\text{X}}^{\text{I}}}{n_{\text{W}}^{\text{I}}} \right) \right\rangle \quad (2)$$

where  $n_{\text{X,W}}^{\text{II}}$  is the number of water (W) or osmolyte (X) molecules in the I and II domains. The solvent density function (SDF) that describes how the molecules of osmolyte are distributed around the protein is, in principle, equivalent to the radial distribution function but takes into account the shape and volume of the protein. The SDF for a generic molecule X is computed as

$$\rho_{\text{X}}(r) = \frac{X(r, r')}{V(r, r')} \quad (3)$$

where  $r$  is the radius of the solvation shell,  $X(r, r')$  is the number of X molecules found from  $rr$  to  $r'$ , and  $V(r, r')$  is the volume of the shell from  $r$  to  $r'$ . The number of molecules  $X(r, r')$  was obtained calculating  $n_{\text{X}}^{\text{II}}$  for different  $r$  values. The volume  $V(r, r')$  was calculated on the basis of the grid-based solvent-accessible methodology of ref 42. Bulk dielectric constants were calculated, according to Neumann's formulation,<sup>43</sup> from the fluctuations of the total dipole moment  $\langle M^2 \rangle$  following the approach reported in ref 44. Three different systems, comprising only the mixed-solvent,

were simulated by standard MD with the same parameters as described before, for a total production phase of 55 ns. The rotational correlation function was calculated using the same systems, following the derivation of Lipari and Szabo for NMR relaxation times<sup>45,46</sup> with a first-order Legendre polynomial. A detailed explanation of the procedure is reported in ref 44 and the references therein.

Human villin headpiece subdomain HP35, LSIED FTQAF GMTPA AFSAL PKWKQ QNLKK EKGLF, was synthesized by Proteogenix (France) with a purity >95%. All CD measurements were performed on a JASCO-810 dichrograph equipped with a Peltier thermoelectric temperature controller. CD spectra of HP35 in water at a concentration of 70  $\mu\text{M}$  were recorded between 190 and 260 nm, with a 0.1-cm-path-length quartz cuvette (Hellma), a 50 nm/min scanning speed, an averaging time of 4 s, and a bandwidth of 1 nm. The spectra shown are the averages of three scans. Thermal denaturation experiments were performed at constant heating rates of 1  $^{\circ}\text{C}/\text{min}$  by following the ellipticity at 222 nm from 5 to 90  $^{\circ}\text{C}$  with a total sample concentration of 50  $\mu\text{M}$ . The analysis of the thermal unfolding curve was performed by nonlinear least-squares fitting according to a two-state model.<sup>47</sup> Equilibrium urea denaturation was monitored by CD in the wavelength range of 210–260 nm and at seven different temperatures between 10 and 40  $^{\circ}\text{C}$ . HP35 solutions at a 50  $\mu\text{M}$  concentration were mixed with varying amounts of stock solution containing 8 M urea. Unfolding was monitored in the range of 0–7 M urea. The urea unfolding profile of HP35 is described by the change of the dichroic signal at 222 nm as a function of denaturant concentration. Chemical denaturation data were analyzed by nonlinear least-squares fitting of the observed CD signal  $[\theta]_t$  to a two-state model of a single unfolding transition between folded (F) and unfolded (U) states:<sup>48</sup>

$$[\theta]_t = \alpha_i([\theta]_{\text{U}} - [\theta]_{\text{F}}) + [\theta]_{\text{F}} \quad (4)$$

where  $[\theta]_{\text{F}}$  is the ellipticity at which the molecule is fully folded and  $[\theta]_{\text{U}}$  is the ellipticity of the fully unfolded molecule. The fractional population of the unfolded form ( $\alpha_i$ ) is determined from the equilibrium constant for unfolding:

$$K_{\text{U}_i} = \exp\left(-\frac{\Delta G_i}{RT}\right) \quad (5)$$

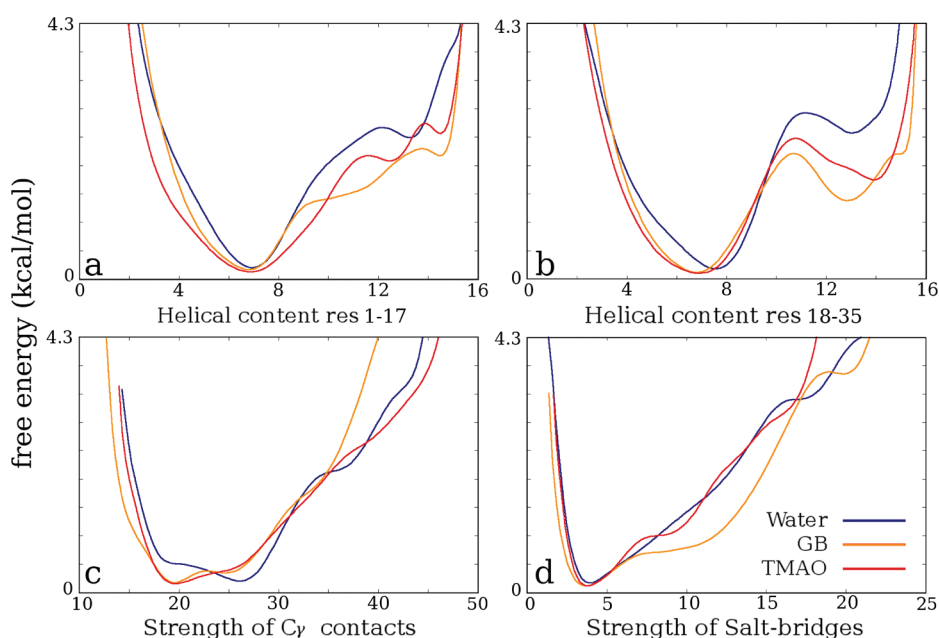
where  $R$  is the gas constant, which equals 1.98 cal/mol, and  $T$  is the absolute temperature.  $\Delta G_i$  is calculated using the linear extrapolation model (LEM):<sup>49</sup>

$$\Delta G_i = \Delta G_0 - m_{\text{urea}}[\text{urea}] \quad (6)$$

where  $\Delta G_0$  is the standard free energy of unfolding in the absence of denaturant and  $m_{\text{urea}}$  is the slope, which characterizes the change in  $\Delta G_i$  with  $[\text{urea}]$ . The denaturant concentration midpoint of the transition,  $[\text{urea}]_{(1/2)}$ , is equal to  $\Delta G_0/m$ . The combined effect of urea and the osmolyte TMAO (or GB) on unfolding free energies was modeled as being linear in both cosolvents:<sup>50</sup>

$$\Delta G_i = \Delta G_0 - m_{\text{urea}}[\text{urea}] - m_{\text{osmolyte}}[\text{osmolyte}] \quad (7)$$

Equation 7 was globally fitted to unfolding transitions in mixtures of urea and osmolyte to yield the free energy of unfolding in the absence of both cosolvents. In order to correctly determine the  $\Delta G$  in the presence of osmolyte, the unfolded fraction was calculated by using the  $[\theta]_{\text{U}}$  value derived



**Figure 1.** Free energy profile as a function of the helical content of residues 1–17 (a) and residues 18–35 (b) of the protein. Free energy profile as a function of the strength of hydrophobic contacts (c) and the strength of salt bridges (d). When the osmolyte is added to the solution, structures with a higher helical content become more populated as conformations with a higher number of salt bridges and less hydrophobic contacts. See ref 34 for the exact definition of the CVs. The typical error due to the convergence of the free energy profiles is reported in Supporting Information Figure S2.

in the absence of osmolytes with urea.<sup>51</sup> Heat capacity change ( $\Delta C_p$ ) for HP35 unfolding was measured by globally fitting the thermal and chemical denaturation data to the Gibbs–Helmholtz equation:<sup>52</sup>

$$\Delta G(T) = \Delta H_m \left( 1 - \frac{T}{T_m} \right) - \Delta C_p \left[ (T_m - T) + T \ln \left( \frac{T}{T_m} \right) \right] \quad (8)$$

where  $\Delta G(T)$  is  $\Delta G$  at temperature  $T$ ,  $T_m$  is the midpoint of the thermal unfolding curve, and  $\Delta H_m$  is the enthalpy change for unfolding measured at  $T_m$ .

## RESULTS AND DISCUSSION

The availability of high-resolution experiments on HP35 folding enables a careful validation of the computational results.

Here, we use previously reported simulations of HP35 in pure water<sup>53</sup> in good agreement with experiments as a reference for the simulations of the osmolyte solutions: a 1.5- $\mu$ s-long fully atomistic MD simulation at 298 K starting from the lowest energy NMR structure (PDB code: 1UNC)<sup>22</sup> and massive BEMD simulations at 298 K and 320 K, close to the experimental melting temperature (see the Supporting Information), were used to reconstruct a fully converged free-energy landscape of HP35 folding. We used the recently described Amber99SB\*-ILDN<sup>35</sup> force field, including several improvements.<sup>36,37</sup>

We repeated the BEMD simulations in the presence of 1 M GB, 1 M TMAO, and, for comparison, 1 M urea, a denaturant. The folded minimum in water (Figure 1) is, in every case, narrow and centered around the values typical of the native structure. The minima in osmolyte solutions are generally broader, and it can be seen (Figure 1c,d) that the osmolytes weaken the hydrophobic core and strengthen the salt bridges. The weakening of the hydrophobic core corresponds to a slight increase of

the exposed surface in the folded state and a more sizable increase in the unfolded ensemble, leading to a  $\Delta$ SASA in good agreement with the observed increase of the heat capacity  $\Delta C_p$  (see Supporting Information Table S1). The protein in osmolyte solutions adopts more helical conformations, as can be seen from the free energy profiles (Figure 1a,b) showing lower minima at higher helical values ( $\sim 14$ ). Since the typical value for the native state is  $\sim 8$ , the higher helical content is found mainly in the unfolded ensemble. From 2D FES (Figure S1, Supporting Information), it is clear that the N and N' free energy basins observed in the pure water simulations<sup>53</sup> are merged in the presence of the osmolytes and that HP35 is more flexible. Another alternative explanation is that the N' state becomes the most stable native structure over the more rigid N state.

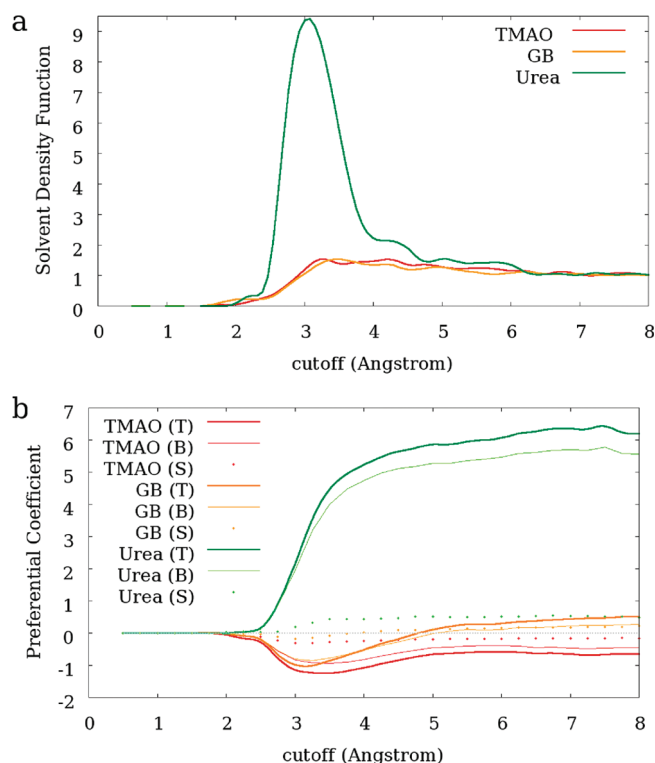
As expected, the free energy profiles of HP35 in urea are significantly different. The FE profiles (Figure S3, Supporting Information) show a narrower minimum corresponding to the folded state, while cluster analysis of the main structures indicates that in urea a partial disruption of the hydrophobic core takes place, with a strong destabilization of helix 3, in contrast to the effect observed for the stabilizing osmolytes. Using 3D FES (Figure S4 and S5, Supporting Information), the unfolding free energy ( $\Delta G_{\text{unfold}}$ ) was calculated by integrating the densities of the folded and unfolded states. The resulting values are  $-0.04$  kcal/mol for the simulation in water, 0.61 and 0.69 kcal/mol, respectively, for 1 M GB and TMAO (Table 1), and  $-0.5$  kcal/mol for urea. As expected, the  $\Delta G_{\text{unfold}}$  is lower for the urea solution, since unfolding is favored with respect to water, and positive for the two protecting osmolytes.

In order to assess whether or not GB and TMAO engage in direct interactions with the protein backbone, we analyzed the distribution of osmolyte molecules around the protein. Calculating the solvent density function (SDF) for GB and TMAO molecules, no relevant peak was observed, suggesting the

Table 1. Thermodynamic Parameters for HP 35 in Pure Water and in NaCl or Osmolite Solutions<sup>a</sup>

	$\Delta H$	$\Delta C_p$	$T_m$	$\Delta G_{47^\circ\text{C}}^b$	$\Delta\Delta G_{47^\circ\text{C}}$	$\Delta G_{47^\circ\text{C}}^{\text{calcd}^c}$	$\Delta\Delta G_{47^\circ\text{C}}^{\text{calcd}}$
H <sub>2</sub> O	24.8 ± 0.9	0.37 ± 0.06	44 ± 0.1	−0.24		−0.04	
NaCl 0.66 M	26.7 ± 0.8	0.44 ± 0.05	49 ± 0.1	0.16	0.40	0.69	0.73
GB 1 M	30.7 ± 3.7	0.72 ± 0.22	49 ± 0.1	0.19	0.43	0.61	0.65
TMAO 1 M	33.4 ± 1.5	0.96 ± 0.09	51 ± 0.1	0.39	0.63	0.69	0.73

<sup>a</sup> Values are in kcal mol<sup>−1</sup> for  $\Delta H$ , kcal mol<sup>−1</sup> K<sup>−1</sup> for  $\Delta C_p$ , °C for  $T_m$ , and kcal mol<sup>−1</sup> for  $\Delta G$ . <sup>b</sup> Obtained from the experiments employing the Gibbs–Helmholtz equation. <sup>c</sup> Obtained from the calculated free energy surfaces.



**Figure 2.** (a) Solvent density function for GB, TMAO, and urea. Only a very slight increase with respect to the bulk limit is observable for the osmolytes with a preferred distance of 3.8 Å. The absence of a well-defined prominent peak suggests that GB and TMAO are excluded from the protein surface, as demonstrated experimentally. The typical 2.8 Å peak is clearly recognizable for urea, confirming its proximity to the protein surface. (b) Preferential coefficient for HP35 in 1 M solutions of GB, TMAO, and urea: total (T), side chain contribution (S), and backbone contribution (B). The negative values in the region 3–5 Å clearly show a preference of the two osmolytes for the bulk domain.

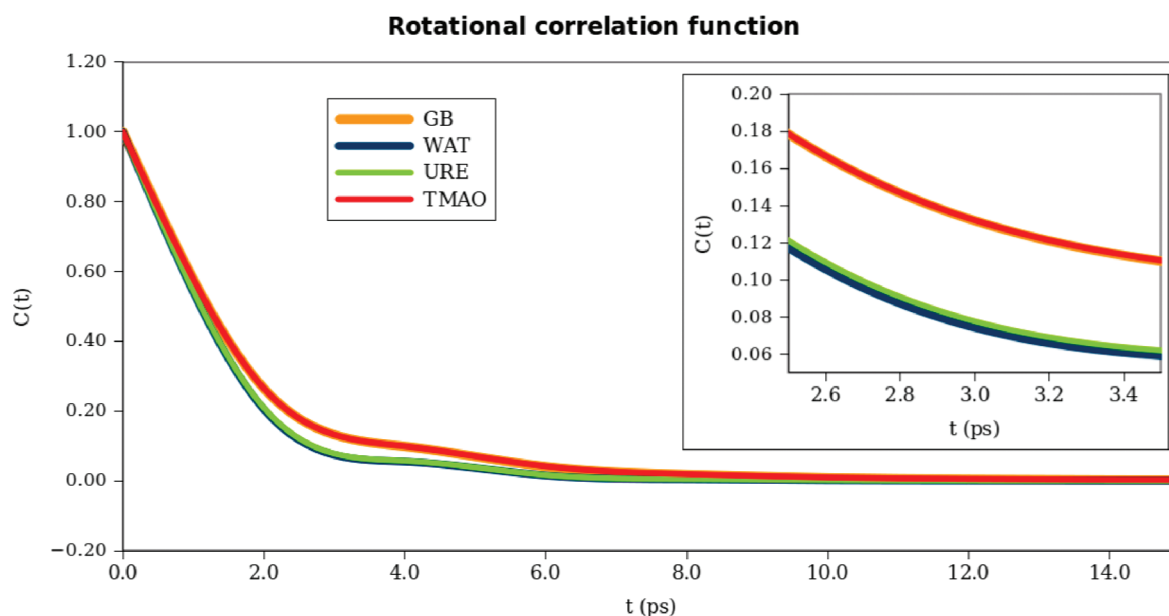
absence of direct contacts between the osmolytes and the protein (Figure 2a), in agreement with refs 16 and 20. The preferential coefficient  $\Gamma_{\text{XP}}$  is calculated to confirm the proposed osmophobic effect<sup>16,19</sup> (see Figure 2b). Choosing a cutoff of 4 Å for the boundary between protein and bulk domains, we obtained a value of −0.52 for GB and −1.08 for TMAO, in agreement with the suggested osmophobic effect. As a comparison, the corresponding value for HP35 in 1 M urea solution is 5.22, confirming urea contacts with the protein.<sup>54</sup> For all molecules, the most relevant contribution to  $\Gamma_{\text{XP}}$  came from the backbone, in agreement with previous results.<sup>55</sup>

The observed differences in the FE profiles due to the osmolytes and the lack of direct interactions with the protein are in agreement with the “indirect” hypothesis. This, together

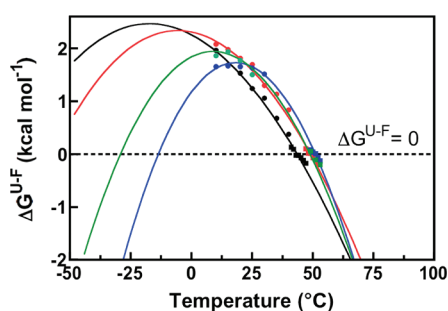
with the proposed changes in the water structure<sup>16,56</sup> due to the osmolytes, led us to investigate whether or not a shift of the dielectric constant would explain the protecting effects. An increase of the  $\epsilon$  value was reported for several osmolytes, including TMAO, GB, taurine, and sarcosine.<sup>57,58</sup> We calculated the static dielectric constant  $\epsilon$  according to Neumann’s formulation,<sup>43,44</sup> obtaining a value of 98.0 (±0.2) for TIP3P water and significantly higher values for the 1 M osmolyte solutions: 106.4 (±0.2) for GB and 103.0 (±0.2) for TMAO, in agreement with the experiments. For comparison, the  $\epsilon$  for the 1 M urea solution was 95.2 (±0.2), similar to that of pure water. These results, suggesting an increased polarity of the solution, are not consistent with the observed increase of salt bridges and hydrophobic core relaxation. However, when we examined more in detail the properties of the solution in the previously defined protein domain, we found a possible explanation to the discrepancy. Osmolytes also affect the rotational dynamics of water molecules both in the bulk and in the protein domain. Indeed, the rotational diffusion is significantly reduced, as shown by the calculated rotational correlation function, showing higher correlation times for water molecules in both 1 M GB and TMAO (Figure 3).

Thus, the high dipole moment of the osmolyte molecules has two different effects on the solution: on the one hand, it causes a considerable increase in the overall dielectric constant; on the other hand, it tends to align the water dipoles. Since, as we have seen the osmolytes are excluded from the protein surface, the lower rotational diffusion of water in the protein first solvation shells has the effect of reducing the local  $\epsilon$  of the solution. The calculated  $\epsilon$  around the protein is significantly smaller, even compared to that of pure water, 83.0 for GB and 90.3 for TMAO.

These results are in agreement with both NMR<sup>59</sup> and IR<sup>60</sup> observations. <sup>1</sup>H NMR data<sup>59</sup> demonstrate a decrease of  $T_1$  relaxation time (i.e., an increase of the rotational correlation time  $\tau_C$ <sup>61</sup>) for several osmolytes, including TMAO, GB, sarcosine, sorbitol, and trehalose. Lower  $T_1$  relaxation times have been ascribed to a more “ice-like” behavior of water (it is to be noted that the  $\epsilon_0$  of ice is 3.19<sup>62</sup>), confirmed by the shifts in NIR spectra.<sup>59</sup> Very recently a similar observation was obtained by 2D infrared spectroscopy on TMAO solutions.<sup>60</sup> Bakulin and co-workers demonstrate the slower rotational reorientation of water molecules around TMAO molecules, supporting the results of our calculations. Much slighter variations, occasionally in the same direction of those observed for protecting osmolytes, are registered for urea, suggesting that the source of the different effect of urea resides mainly in its interactions with the protein. Most of the reported features of osmolyte behavior (e.g., folded state protection, osmophobic effect and backbone repulsion) can be solidly explained in the context of an “ice-like” shift in the aqueous solvent dynamical behavior, due to GB or TMAO addition. The slowing down of water rotational diffusion is highly



**Figure 3.** Rotational correlation function for the simulated systems: pure water and 1 M solutions of TMAO, GB, and urea. An increase of the rotational correlation time can be observed for the two protecting osmolytes.



**Figure 4.** Protein stability curves for HP35 in water (black), 0.66 M NaCl (red), 1 M GB (green), or 1 M TMAO (blue). Squares represent unfolding free energies measured directly from the transition zones of the thermal denaturation curves shown in Figure 2 at an HP 35 concentration of 50  $\mu\text{M}$ . Circles represent  $\Delta G_{\text{U}}^{\text{H}_2\text{O}}$  values determined from an analysis of urea denaturation curves determined at various temperatures. Solid lines show the best fit to the Gibbs–Helmholtz equation.

consistent with previous hypotheses describing the osmolytes effect as a “water-structuring” effect.<sup>63,64</sup> The water molecules’ rotational diffusion slowdown is perceived as an average effect by the protein itself, and the change in water rotational properties affects the thermodynamic and electrostatic response properties of the solvent.

To confirm this new formulation of the “indirect hypothesis”, we performed a BEMD simulation of HP35 in modified TIP3P water molecules (W79), whose charges were scaled down to reproduce the decreased dielectric constant of water in the 1 M GB solution. The FE profiles were strikingly similar to those obtained for the osmolyte solutions (Figure S6 and S7, Supporting Information). Cluster analysis revealed a high similarity of the most populated conformers in 1 M GB and W79, with a RMSD within 2.1 Å; the salt bridge previously observed in 1 M GB was also observed in W79 (Figure S8, Supporting Information). Hence, the W79 simulation provided further evidence that a

dielectric constant shift (i.e., a rotational diffusion slowdown) in the protein domain can explain most of the features of the osmolyte solution.<sup>16</sup>

To validate computational results, we exploited thermal and chemical denaturation to gain an in-depth thermodynamic description of the effects due to the osmolytes. Equilibrium thermal unfolding measurements were performed on HP35 in water and in 1 M TMAO or GB solutions. The stability was also investigated in a 0.66 M NaCl solution. This salt concentration reduces the dielectric constant to 67.2,<sup>65</sup> similar to the shift observed for water in a 1 M GB solution. HP35 showed a cooperative, sigmoidal transition (Supporting Information Figure S9), and the data fit a two-state model. HP35 in water shows a transition temperature ( $T_m$ ) of 44 °C. TMAO or GB increases the  $T_m$  to 51 and 49 °C, respectively. Similarly, HP35 in 0.66 M NaCl unfolds with a  $T_m$  of 49 °C. At 25 °C, the midpoint of the urea-induced chemical denaturation is 2.9 M in water, 3.6 M in NaCl, and 3.5 M in GB or TMAO. These results indicate a clear stabilization of the native state and are in excellent agreement with the predictions of the simulations. The unfolding reaction of HP35 showed similar  $m$  values in all solutions, suggesting similar cooperativity (Supporting Information Table S2). One explanation is that osmolytes are not directly in contact with the protein backbone, in agreement with the results of the simulations. Consistent with this, proton NMR cross-relaxation (ROESY) experiments were unable to detect any TMAO-HP35 contact (see Supporting Information Figures S12 and S13), indicating the absence of direct and persistent (ms time scale) interactions between osmolyte and protein. From the combination of thermal and chemical denaturation, we obtained the stability plot of HP35 and calculated the unfolding  $\Delta G$ . It is evident (Figure 4 and Table 1) that the osmolytes determine an increase of stability with respect to pure water, in agreement with the calculated values.

## CONCLUSIONS

The pursuit of a universal explanation for the osmoprotectant effect has drawn considerable attention in recent decades, due to

its significant importance for both fundamental and applied science. In recent years, the studies of Bolen and co-workers<sup>11,19,55</sup> have succeeded in the defining a simplified model, based on transfer free energies, with considerable predictive power. However, despite multiple efforts and ever increasing interest, a simple yet general microscopic explanation of the mechanism underlying osmolyte-mediated protein protection mainly remains an open issue. State-of-the-art in silico simulations and experiments allowed us to make a significant step forward toward this goal. Our results support a new flavor of the previously reported “indirect hypothesis” and put forward a very simple explanation: the main driving force behind native state protection is a slowdown of the solvent rotational dynamics. The “slower” solvent behaves around the protein, where the osmolytes are excluded, as a colder or lower dielectric aqueous solvent. This local reduction is consistent with, and explanatory of, all reported theoretical and experimental results.<sup>11,16,19,55</sup> Indeed, the alteration of the solvent is translated into a decreased denaturing power of the water molecules that, acting as a less polar media with the dynamical behavior of a lower temperature solvent, is less effective in interfering with the intraprotein interactions sustaining the native fold. This, in turn, explains not only the osmoprotecting effect and the increase of the melting temperature of proteins but also the significant role of backbone interactions, whose importance was systematically predicted by transfer models.

## ■ ASSOCIATED CONTENT

📄 **Supporting Information.** Additional free energy profiles of HP35 in water and in 1 M solutions, 2D and 3D free energy maps, CD and NMR spectra, and thermal denaturation curves. This information is available free of charge via the Internet at <http://pubs.acs.org/>.

## ■ AUTHOR INFORMATION

### Corresponding Author

\*E-mail: [maurizio.sironi@unimi.it](mailto:maurizio.sironi@unimi.it); [flgervasio@cnio.es](mailto:flgervasio@cnio.es).

## ■ ACKNOWLEDGMENT

We acknowledge support by the Spanish Science and Innovation (MICINN) grant (BIO2010-20166, “AlteredDynamics”). M. Morando is acknowledged for helpful discussions. G.S. acknowledges the European Commission Capacities Area—Research Infrastructures Initiative HPC-EUROPA2 (project number: 228398) for partial support. The Barcelona Supercomputing Center is acknowledged for a generous allocation of computer resources.

## ■ REFERENCES

- (1) Yancey, P. H.; Clark, M. E.; Hand, S. C.; Bowlus, R. D.; Somero, G. N. *Science* **1982**, *217*, 1214–1222.
- (2) Yancey, P. H.; Burg, M. B. *Am. J. Physiol.* **1989**, *257*, 602–607.
- (3) Yancey, P. H. *J. Exp. Biol.* **2005**, *208*, 2819–30.
- (4) Weber, D. J. *Salinity and Water Stress*; Springer Netherlands: Dordrecht, The Netherlands, 2009; Vol. 44, p 236.
- (5) Furini, A.; Koncz, C.; Salamini, F.; Bartels, D. *EMBO J.* **1997**, *16*, 3599–3608.
- (6) Baldwin, W. W.; Myer, R.; Kung, T.; Anderson, E.; Koch, A. L. *J. Bacteriol.* **1995**, *177*, 235–237.

- (7) Cayley, S.; Lewis, B. A.; Guttman, H. J.; Record, M. T.; et al. *J. Mol. Biol.* **1991**, *222*, 281–300.
- (8) Csonka, L. N. *Microbiol. Rev.* **1989**, *53*, 121–147.
- (9) Hochachka, P. W.; Somero, G. N. *Biochemical adaptation: mechanism and process in physiological evolution*; Oxford University Press: New York, 2002; p 466.
- (10) Qu, Y.; Bolen, C. L.; Bolen, D. W. *Proc. Natl. Acad. Sci. U.S.A.* **1998**, *95*, 9268–73.
- (11) Baskakov, I. V.; Bolen, D. W. *J. Biol. Chem.* **1998**, *273*, 4831–4834.
- (12) Moriyama, T.; Garcia-Perez, A.; Olson, A. D.; Burg, M. B. *Am. J. Physiol.* **1991**, *260*, 494–497.
- (13) Burg, M. B.; Ferraris, J. D. *J. Biol. Chem.* **2008**, *283*, 7309–13.
- (14) Xie, G.; Timasheff, S. N. *Biophys. Chem.* **1997**, *64*, 25–43.
- (15) Street, T. O.; Bolen, D. W.; Rose, G. D. *Proc. Natl. Acad. Sci. U.S.A.* **2006**, *103*, 13997–14002.
- (16) Zou, Q.; Bennis, B.; Daggett, V.; Murphy, K. J. *Am. Chem. Soc.* **2002**, *124*, 1192–1202.
- (17) Wang, A.; Bolen, D. W. *Biophys. J.* **1996**, *71*, 2117–2122.
- (18) Canchi, D.; Paschek, D.; Garcia, A. *J. Am. Chem. Soc.* **2010**, *132*, 2338–2344.
- (19) Bolen, D. W.; Baskakov, I. V. *J. Mol. Biol.* **2001**, *310*, 955–63.
- (20) Saladino, G.; Pieraccini, S.; Rendine, S.; Recca, T.; Francescato, P.; Speranza, G.; Sironi, M. *J. Am. Chem. Soc.* **2011**, *133*, 2897–2903.
- (21) Bussi, G.; Gervasio, F. L.; Laio, A.; Parrinello, M. *J. Am. Chem. Soc.* **2006**, *128*, 13435.
- (22) Vermeulen, W.; Vanhaesebrouck, P.; Troys, M. V.; Verschuere, M.; Fant, F.; Goethals, M.; Ampe, C.; Martins, J. C.; Borremans, F. A. M. *Protein Sci.* **2004**, *13*, 1276–1287.
- (23) McKnight, C. J.; Doering, D. S.; Matsudaira, P. T.; Kim, P. S. *J. Mol. Biol.* **1996**, *260*, 126–34.
- (24) Duan, Y.; Kollman, P. A. *Science* **1998**, *282*, 740.
- (25) Zagrovic, B.; Snow, C. D.; Shirts, M. R.; Pande, V. S. *J. Mol. Biol.* **2002**, *323*, 927–937.
- (26) Frank, B. S.; Vardar, D.; Buckley, D. A.; McKnight, C. J. *Protein Sci.* **2002**, *11*, 680–7.
- (27) Kubelka, J.; Henry, E. R.; Cellmer, T.; Hofrichter, J.; Eaton, W. A. *Proc. Natl. Acad. Sci. U.S.A.* **2008**, *105*, 18655–62.
- (28) Tang, Y.; Rigotti, D.; Fairman, R.; Raleigh, D. *Biochemistry* **2004**, *43*, 3264–3272.
- (29) Havlin, R. H.; Tycko, R. *Proc. Natl. Acad. Sci. U.S.A.* **2005**, *102*, 3284–9.
- (30) Piana, S.; Laio, A. *J. Phys. Chem. B* **2007**, *111*, 4553–4559.
- (31) Jorgensen, W. L. *J. Am. Chem. Soc.* **1981**, *103*, 335–340.
- (32) Hess, B.; Kutzner, C.; Spoel, D. V. D.; Lindahl, E. *J. Chem. Theory Comput.* **2008**, *4*, 435–447.
- (33) Bonomi, M.; Branduardi, D.; Bussi, G.; Camilloni, C.; Provasi, D.; Raiteri, P.; Donadio, D.; Marinelli, F.; Pietrucci, F.; Broglia, R. A.; et al. *Comput. Phys. Commun.* **2009**, *180*, 1961–1972.
- (34) Piana, S.; Laio, A.; Marinelli, F.; Troys, M. V.; Bourry, D.; Ampe, C.; Martins, J. C. *J. Mol. Biol.* **2008**, *460*, 460–470.
- (35) Hornak, V.; Abel, R.; Okur, A.; Strockbine, B.; Roitberg, A.; Simmerling, C. *Proteins: Struct. Funct. Bioinf.* **2006**, *65*, 712–725.
- (36) Lindorff-Larsen, K.; Piana, S.; Palmo, K.; Maragakis, P.; Klepeis, J. L.; Dror, R. O.; Shaw, D. E. *Proteins: Struct. Funct. Bioinf.* **2010**, *78*, 1950–8.
- (37) Best, R. B.; Buchete, N. V.; Hummer, G. *Biophys. J.* **2008**, *95*, L07–L09.
- (38) Hess, B.; Bekker, H.; Berendsen, H. J.; Fraaije, J. G. *J. Comput. Chem.* **1997**, *18*, 1463–1472.
- (39) Bussi, G.; Donadio, D.; Parrinello, M. *J. Chem. Phys.* **2007**, *126*, 014101.
- (40) Berendsen, H. J. C.; Postma, J. P. M.; Gunsteren, W. F. V.; DiNola, A.; Haak, J. R. *J. Chem. Phys.* **1984**, *81*, 3684.
- (41) Baynes, B. M.; Trout, B. L. *J. Phys. Chem. B* **2003**, *107*, 14058–14067.
- (42) Beck, D. A.; Alonso, D. O.; Daggett, V. *Biophys. Chem.* **2003**, *100*, 221–237.

- (43) Neumann, M. J. *Chem. Phys.* **1985**, *82*, 5663–5672.
- (44) van der Spoel, D.; Maaren, P. J. V.; Berendsen, H. J. J. *Chem. Phys.* **1998**, *108*, 10220–10230.
- (45) Lipari, G.; Szabo, A. J. *Am. Chem. Soc.* **1982**, *104*, 4559–4570.
- (46) Lipari, G.; Szabo, A. J. *Am. Chem. Soc.* **1982**, *104*, 4546–4559.
- (47) Greenfield, N. J. *Nat. Protoc.* **2007**, *1*, 2527–2535.
- (48) Greenfield, N. J. *Nat. Protoc.* **2006**, *1*, 2733–2741.
- (49) Myers, J.; Pace, C.; Scholtz, J. *Protein Sci.* **1995**, *4*, 2138–2148.
- (50) Mello, C.; Barrick, D. *Protein Sci.* **2003**, *12*, 1522–1529.
- (51) Wu, P.; Bolen, D. *Proteins: Struct. Funct. Bioinf.* **2006**, *63*, 290–296.
- (52) Pace, C.; Laurents, D. *Biochemistry* **1989**, *28*, 2520–2525.
- (53) Saladino, G.; Marenchino, M.; Gervasio, F. L. *J. Chem. Theory Comput.* **2011**, *7*, 2675–2680.
- (54) Stumpe, M. C.; Grubmüller, H. *PLoS Comput. Biol.* **2008**, *4*, e1000221.
- (55) Street, T.; Bolen, D.; Rose, G. *Proc. Natl. Acad. Sci. U.S.A.* **2006**, *103*, 13997–14002.
- (56) Niebuhr, M.; Koch, M. H. *Biophys. J.* **2005**, *89*, 1978–1983.
- (57) Shikata, T.; Itatani, S. *J. Sol. Chem.* **2002**, *31*, 823–844.
- (58) Wyman, J. *Chem. Rev.* **1936**, 213–239.
- (59) Lever, M.; Blunt, J. W.; Maclagan, R. G. *Comp. Biochem. Physiol.* **2001**, *130*, 471–86.
- (60) Bakulin, A. a.; Pshenichnikov, M. S.; Bakker, H. J.; Petersen, C. *J. Phys. Chem. A* **2011**, *115*, 1821–9.
- (61) Yoshida, K.; Ibuki, K.; Ueno, M. *J. Chem. Phys.* **1998**, *108*, 1360.
- (62) Mätzler, C.; Wegmüller, U. *J. Phys. D* **1988**, *21*, 1660–1660.
- (63) Bennion, B. J.; Daggett, V. *Proc. Natl. Acad. Sci. U.S.A.* **2003**, *100*, 5142.
- (64) Bennion, B.; Daggett, V. *Proc. Natl. Acad. Sci. U.S.A.* **2004**, *101*, 6433.
- (65) Haggis, G. H.; Hasted, J. B.; Buchanan, T. J. *J. Chem. Phys.* **1952**, *20*, 1452–1465.