

Hybrid Ensembles for Improved Force Matching

Lee-Ping Wang and Troy Van Voorhis

Department of Chemistry, Massachusetts Institute of Technology

77 Massachusetts Ave.

Cambridge, MA 02139 USA

(Dated: 29 October 2010)

Force matching is a method for parameterizing empirical potentials in which the empirical parameters are fitted to a reference potential energy surface (PES). Typically, training data is sampled from a canonical ensemble generated with either the empirical potential or the reference PES. In this Communication, we show that sampling from either ensemble risks excluding critical regions of configuration space, leading to fitted potentials that deviate significantly from the reference PES. We present a hybrid ensemble which combines the Boltzmann probabilities of both potential surfaces into the fitting procedure, and we demonstrate that this technique improves the quality and stability of empirical potentials.

Molecular mechanics (MM) simulation is a powerful method for investigating the dynamical behavior of complex atomistic systems, but its utility depends critically on the accuracy of the empirical potential being used. A common approach is to fit the parameters in the empirical potential to a high-accuracy reference potential energy surface (PES), which can be obtained from quantum mechanical (QM) calculations. This approach, known as force matching¹⁻⁶ or potential fitting⁷⁻¹⁰, aims to find the optimal parameters k that minimize an objective function χ^2 of the difference between a set of properties \mathbf{Q} computed using the reference PES, and the analogous properties \mathbf{M} computed using the empirical potential:

$$\chi^2 \equiv \int_{\mathbb{R}^{3N}} P(\mathbf{r}) |\mathbf{X}(\mathbf{r}, k)|^2 d\mathbf{r}. \quad (1)$$

Here, the norm squared of the difference vector $\mathbf{X}(\mathbf{r}, k) \equiv \mathbf{M}(\mathbf{r}, k) - \mathbf{Q}(\mathbf{r}, k)$ is integrated over the $3N$ -dimensional configuration space. \mathbf{X} may contain energies, atomistic forces, and/or other computable quantities that one wishes to match. The integral is typically performed by quadrature, or sampling techniques such as the Metropolis algorithm; in the latter case, $P(\mathbf{r})$ is a probability density corresponding to some ensemble.

The ensemble being sampled may be generated with the reference PES; we will call this the QM ensemble and denote the corresponding objective function using χ_{QM}^2 . Accurate sampling of the QM ensemble is generally desirable but costly. Alternatively, training configurations may be rapidly sampled from the empirical (MM) ensemble; we will denote the corresponding objective function using χ_{MM}^2 . Ischtwan and Collins,¹¹ and more recently Akin-Ojo and Wang⁶ have proposed a scheme in which the MM ensemble is resampled after parameterization, with the process repeated in *generations* until self-consistency. Both types of force matching utilizing χ_{QM}^2 and χ_{MM}^2 have been applied widely to parameterize empirical potentials for gas-phase molecules¹¹⁻¹³, condensed phase systems^{2,6,14}, and solids.^{1,3,15-17} The force matching method is not restricted to using MM for the empirical potential or QM as the reference PES; in multiscale coarse graining methods¹⁸, the reference PES is

an atomistic MM potential and is used to parameterize an empirical potential for coarse-grained particles.

In this Communication, we first present a simple heuristic example where force matching using either χ_{MM}^2 or χ_{QM}^2 does not optimally reproduce the QM PES. Instead, configurations from both QM and MM ensembles are important for force matching, and using either ensemble alone is insufficient. As a solution, we propose a hybrid-ensemble approach that combines the probability densities of both ensembles and generates a MM potential that optimally reproduces the QM PES in the heuristic example. Finally, we demonstrate the effects of hybrid-ensemble force matching when it is applied to parameterize MM potentials for a helium dimer and a water hexamer.

As an example, take both the QM PES and the MM potential to be hard repulsive walls as illustrated in Fig. 1. The MM potential has one adjustable parameter σ_{MM} which sets the location of the hard wall, and the optimal match is obtained when $\sigma_{\text{MM}} = \sigma_{\text{QM}}$. We make the simplifying assumption that the wall height is $\gg kT$, such that only the configurations where $r \geq \sigma$ are thermally accessible, and the repulsive region is excluded from the integral in Equation 1.

Consider first the case where the QM PES has a hard wall at $\sigma_{\text{QM}} = 50$, and the MM potential underestimates this value at $\sigma_{\text{MM}} = 30$; we will denote this potential using $\text{MM}_<$ (Fig. 1a). The “error region”, with finite $\mathbf{X}(r, \sigma_{\text{MM}})$, is given by $r \in \{30, 50\}$; the QM-accessible region with finite $P_{\text{QM}}(r)$ is given by $r > 50$, and the MM-accessible region with finite $P_{\text{MM}}(r)$ is given by $r > 30$. The error region overlaps with the MM-accessible region and contributes to χ_{MM}^2 ; thus, χ_{MM}^2 is minimized by increasing σ_{MM} until the optimal value, $\sigma_{\text{MM}} = \sigma_{\text{QM}}$, is obtained. However, the error region does *not* overlap with the QM-accessible region and does not contribute to χ_{QM}^2 ; in fact, any $\sigma_{\text{MM}} < 50$ minimizes χ_{QM}^2 ! Thus, minimizing χ_{QM}^2 risks severely underestimating σ_{MM} .

Now, consider the case where the initial guess is an overestimate ($\sigma_{\text{MM}} = 70$), denoted by $\text{MM}_>$ (Fig. 1b). The error region is now $r \in \{50, 70\}$; the QM-accessible region is the same ($r > 50$), and the MM-accessible region

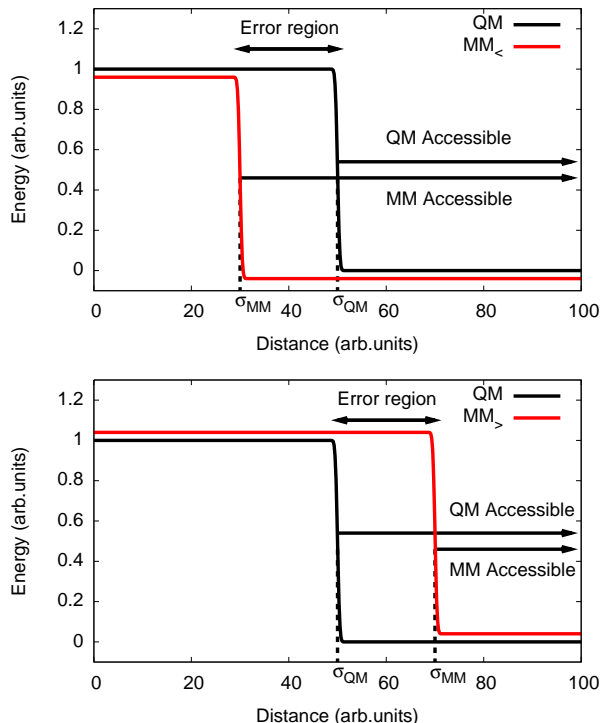


FIG. 1. (a) QM and $MM_{<}$ potential energy curves. (b) QM and $MM_{>}$ potential energy curves. The error region and thermally accessible regions are labeled using arrows.

is $r > 70$. The error region now overlaps with the QM-accessible region but not the MM-accessible region. The fitting situation is reversed; there is a contribution to χ_{QM}^2 but not to χ_{MM}^2 . We now expect minimizing χ_{QM}^2 to effectively optimize σ_{MM} , while any $\sigma_{MM} > 50$ will minimize χ_{MM}^2 . Thus, in a realistic situation where the relationship between the QM PES and the MM potential is unknown, we would expect force matching using χ_{QM}^2 and χ_{MM}^2 to respectively underestimate and overestimate σ ; clearly, both choices are inadequate. The problem is further worsened when approximate functional forms are used for the MM potential.

The key observation is that the optimization fails when the chosen ensemble excludes configurations that are included in the other ensemble. Evidently, it is essential for configurations from both ensembles to contribute to χ^2 in order to produce an accurate MM potential. A simple treatment would be to construct χ^2 using a hybrid ensemble which combines the QM and MM probabilities:

$$\chi_H^2 \equiv \int_{\mathbb{R}^{3N}} (\omega P_{QM}(\mathbf{r}) + (1-\omega)P_{MM}(\mathbf{r})) |\mathbf{X}(\mathbf{r}, k)|^2 d\mathbf{r}. \quad (2)$$

Here, ω is an adjustable mixing parameter which mixes the QM and MM probabilities. In the hard-wall example, the entire error region always has a finite probability density in χ_H^2 , and minimizing χ_H^2 would effectively optimize σ_{MM} starting from any initial value. Note that after

obtaining a perfect fit (i.e. the QM PES and MM potential are exactly the same), the hybrid ensemble reduces to the standard QM or MM ensemble.

Our main results in this Communication are based upon force-matching using χ_H^2 . We now demonstrate the application of hybrid-ensemble force matching to two example systems; the helium dimer and water hexamer. The QM PES, in all cases, is computed at the Hartree-Fock (HF)/3-21G level of theory using Q-Chem.¹⁹ MM simulations were performed using GROMACS.²⁰ Force matching was performed using the ForTune force matching program developed within our group, which minimizes χ^2 using a Newton-Raphson algorithm in similar fashion to previously published software packages.^{21,22} We include both energy and force contributions in χ^2 ; each component is inverse variance-weighted, and the force components are attenuated by $\frac{1}{3^N}$ such that energy and force differences contribute approximately equally. In all optimizations, a penalty proportional to $|\Delta k|^2$ is added to prevent large fluctuations in the parameters when the change in χ^2 is small.

The examples presented in this Communication are designed to highlight the effects of using the hybrid ensemble for force-matching to the QM PES and not intended to reproduce experimental properties. Thus, we will discuss the quality of the parameterization by direct comparison of the QM PES and MM potentials in the former case, and by comparing radial distribution functions from QM and MM molecular dynamics in the latter.

Our first example, the parameterization of an MM potential for a helium dimer, is a direct computational realization of the hard-wall example above. In HF theory, the interaction between helium atoms is well characterized by a van der Waals (vdW) repulsive interaction. Our MM potential uses the Lennard-Jones (LJ) functional form:

$$E_{LJ}(r_{ab}) = 4\epsilon \left(- \left(\frac{\sigma_{ab}}{r_{ab}} \right)^6 + \left(\frac{\sigma_{ab}}{r_{ab}} \right)^{12} \right) \quad (3)$$

In all cases, the LJ well depth ϵ is fixed at 0.01 kJ/mol, and only the X-intercept σ is optimized. Training configurations were obtained from a uniform grid of 800 He-He internuclear separations ranging from 1.5 Å to 4.0 Å. χ_{QM}^2 and χ_{MM}^2 were computed by performing a weighted sum over snapshots, with the appropriate Boltzmann probabilities: $P_{QM,s} = e^{-\beta E_{QM}(\mathbf{r}_s)}$ and $P_{MM,s} = e^{-\beta E_{MM}(\mathbf{r}_s)}$. The hybrid χ_H^2 uses an average probability with ω set to 0.5: $P_{H,s} = \frac{1}{2}(P_{QM,s} + P_{MM,s})$.

Two initial MM potentials, $MM_{<}$ and $MM_{>}$, were chosen such that they respectively underestimate and overestimate the position of the repulsive wall. For each initial potential, we obtained three optimized MM potentials using the three objective functions above; each optimized MM potential was obtained from three generations of the iterative force matching scheme described earlier. Fig. 2 shows the QM PES, initial MM potential, and the optimized MM potentials obtained using each of the three

objective functions. Here, MM_{MM} stands for “optimized MM potential using χ_{MM}^2 ”. The values of σ and χ^2 for the three optimizations are given in Tables S1 and S2.²⁶

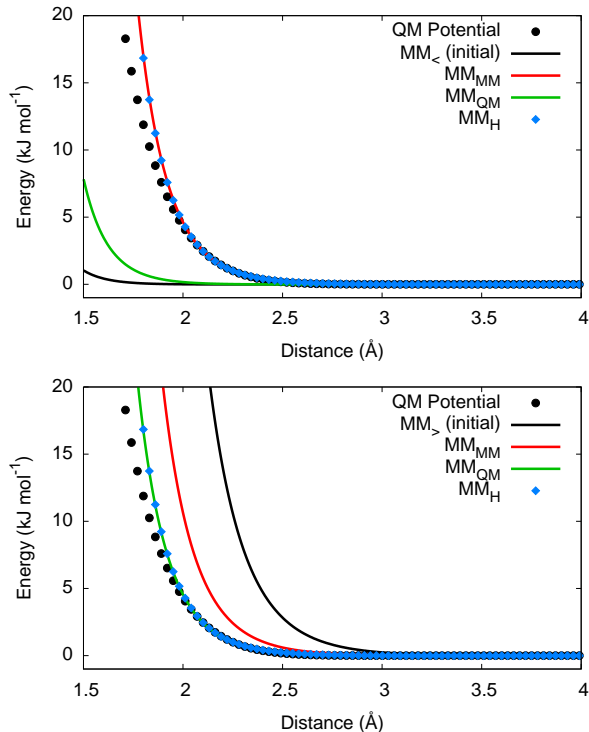


FIG. 2. Optimized MM potentials from force matching, starting from the initial guess (black line) and fitting to the quantum PES (black dots). Initial guesses are $\text{MM}_{<}$ (a) and $\text{MM}_{>}$ (b); optimized MM potentials using χ_{MM}^2 (red), χ_{QM}^2 (green), and χ_{H}^2 (blue) are shown.

From Fig. 2, the role reversal of QM vs. MM weights can be clearly seen. Starting from $\text{MM}_{<}$, using χ_{MM}^2 reproduces the quantum result while using χ_{QM}^2 underestimates the solution; the opposite effect is observed when we start from $\text{MM}_{>}$. In both cases, using χ_{H}^2 produces the correct solution, because the regions that contribute to either χ_{MM}^2 and χ_{QM}^2 are both counted.

It should be noted that in this example, all methods will eventually lead to the correct result; the convergence will just be very slow, taking many generations. This is because the repulsive interaction still contributes a very small amount to χ^2 , and that all regions are perfectly sampled in this ideal case. When more parameters and degrees of freedom are involved, changing the sampling method can lead to different self-consistent MM potentials, as we describe in the following example.

In our second example, our system is a cluster of six water molecules; the intermolecular interactions in small water clusters is highly nontrivial²³ and MM potential development for water remains a highly active field. We chose the flexible SPC/E model²⁴ as our initial MM potential; the MM potential contains 7 empirical parameters. In each generation of force matching, the MM

potential is used to generate 3.0 ns of dynamics in the canonical ensemble ($T = 300\text{K}$), from which 3,000 snapshots are sampled at 1ps time intervals and added to the training data. We utilize WHAM²⁵ to include training data from past generations. A shallow harmonic potential with force constant $0.01 \text{ kJ mol}^{-1} \text{ nm}^{-2}$ was applied to prevent divergent trajectories.

The MM dynamics directly samples the MM ensemble, so in computing χ_{MM}^2 each snapshot has an equal probability. In computing χ_{H}^2 , each snapshot’s probability contains a constant from the MM ensemble, plus a non-Boltzmann factor corresponding to that snapshot’s probability in the QM ensemble:

$$\chi_{\text{H}}^2 = \sum_s \left(\omega \frac{e^{-\beta(E_{\text{QM}}(\mathbf{r}_s) - E_{\text{MM}}(\mathbf{r}_s))}}{S_{\text{QM}}} + (1 - \omega) \frac{1}{S} \right) |\mathbf{X}(\mathbf{r}_s, k)|^2. \quad (4)$$

Here, S is the total number of snapshots, and S_{QM} is the sum of all non-Boltzmann factors.

Three objective functions were used with the following ensembles: χ_{80}^2 ($\omega = 0.8$), χ_{50}^2 ($\omega = 0.5$), and χ_0^2 ($\omega = 0.0$, pure MM). Due to large fluctuations in the non-Boltzmann factors, χ_{100}^2 ($\omega = 1.0$, pure QM) was not usable in force matching and χ_{80}^2 was used as a substitute. Force matching was performed for 19 generations, and the final MM potentials (MM_{80} , MM_{50} , MM_0) were used to obtain radial distribution functions (RDFs) for comparison with an 1.8 ns AIMD reference trajectory, generated using the same temperature and harmonic potential. All three cases produced improved agreement with the AIMD RDFs compared to the initial SPC/E parameters; however, none of the MM potentials reproduced the exact shape of the AIMD RDFs, possibly due to the limitations of the functional form.

Fig. 3 shows the RDFs generated using the three optimized MM potentials compared to the AIMD RDF. MM_{80} underestimates the contact distances and peak positions for all three atom pairs; notably, this occurred starting from an overestimated initial guess. This may have been caused by fitting to components of χ^2 from other intermolecular degrees of freedom (i.e. torques) which couple to this one. MM_0 overestimates the distances, while MM_{50} provides the best estimate with peak positions agreeing with the AIMD result to within 0.1 Å. Thus, the hybrid-ensemble force matching leads to the most accurate MM potential in this more complex case and corresponds well with the earlier examples.

We also performed force matching by sampling from the QM ensemble directly. 50,000 snapshots were sampled at 1ps intervals from the AIMD trajectories, and force matching was used to obtain the optimized potential MM_{100} . We find that MM_{100} underestimates the approach distance, but not as severely as MM_{80} ; this result is surprising as we expect MM_{80} to interpolate between MM_0 and MM_{100} . This may be due to differences in the sampling trajectories; the training data for MM_{100} came from direct AIMD, while the data for MM_{80} came from an MM trajectory and required non-

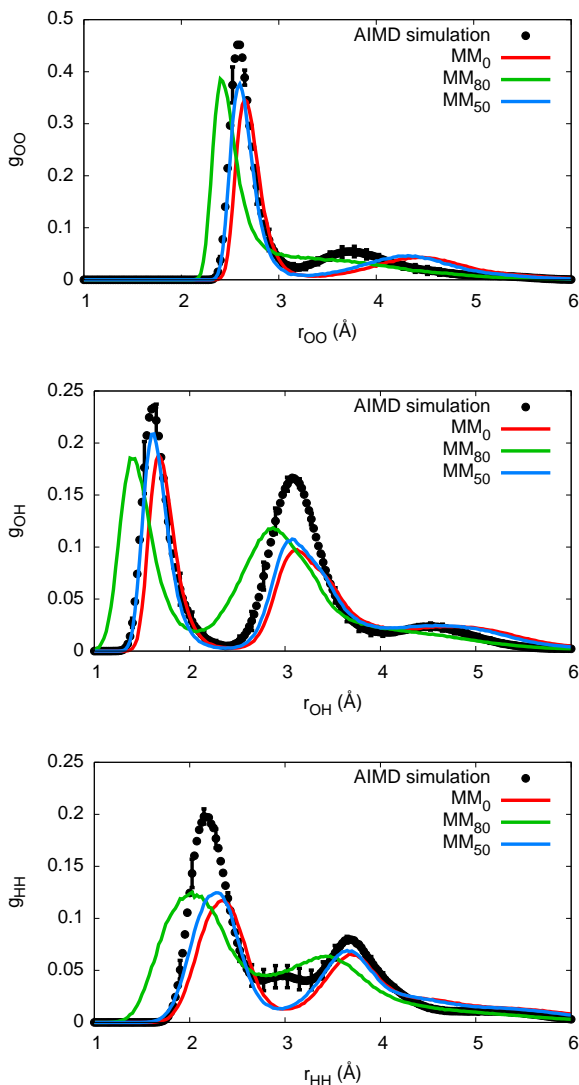


FIG. 3. Radial distribution functions of water for (a) O—O distances, (b) O—H distances, and (c) H—H distances. Thirty AIMD RDFs of length 60 ps were combined to produce the reference curve.

Boltzmann weights. This indicates that more complete sampling may be achieved by combining trajectories from both AIMD and MM dynamics.

In summary, using a single canonical ensemble to obtain training data for force matching can lead to substantial errors in the optimized MM potential. In particular, there is a risk of obtaining a poor fit if χ^2 neglects regions which make large contributions to $|\mathbf{X}(\mathbf{r}, k)|^2$ and have finite probability in *either* the QM or MM ensemble. As a solution, we proposed constructing χ^2 from a hybrid ensemble which combines the probabilities from both ensembles. We demonstrated the advantages of the method by applying force matching to parameterize MM potentials for a helium dimer and a water cluster. In both cases the anomalies associated with single-ensemble objective

functions χ_{QM}^2 and χ_{MM}^2 were shown, and the hybrid χ_{H}^2 provided the correct behavior. Our finding makes intuitive sense: if training data is sampled from the MM ensemble, the MM-accessible configurations will naturally be included in the fit. Meanwhile, QM-accessible configurations that are MM-inaccessible should still be penalized for their absence.

The hybrid ensemble may be applied to improve the parameterization of atomic charges and vdW-type interactions (e.g. LJ, n -6 or Buckingham), which traditionally are the most difficult to parameterize yet are essential for determining the dynamical and bulk properties of any multimolecular system.

This work was funded by ENI S.p.A. as part of the Solar Frontiers Research Program.

- ¹F. ERCOLESSI and J. B. ADAMS, *Europhys. Lett.* **26**, 583 (1994).
- ²T. G. A. YOUNGS, M. G. D. POPOLO, and J. KOHANOFF, *J. Phys. Chem. B* **110**, 5697 (2006).
- ³X. Y. LIU, J. B. ADAMS, F. ERCOLESSI, and J. A. MORIARTY, *Model. Simul. Mater. Sc.* **4**, 293 (1996).
- ⁴S. IZVEKOV, M. PARRINELLO, C. J. BURNHAM, and G. A. VOTH, *J. Chem. Phys.* **120**, 10896 (2004).
- ⁵P. MAURER, A. LAIO, H. W. HUGOSSON, M. C. COLOMBO, and U. ROTH LISBERGER, *J. Chem. Theory Comput.* **3**, 628 (2007).
- ⁶O. AKIN-OJO, Y. SONG, and F. WANG, *J. Chem. Phys.* **129**, 064108 (2008).
- ⁷A. PUKRITTAYAKAMEE, M. MALSHE, M. HAGAN, L. M. RAFF, R. NARULKAR, S. BUKKAPATNUM, and R. KOMANDURI, *J. Chem. Phys.* **130**, 134101 (2009).
- ⁸H. SUN, *J. Phys. Chem. B* **102**, 7338 (1998).
- ⁹P. BROMMER and F. GAHLER, *Philos. Mag.* **86**, 753 (2006).
- ¹⁰G. TOTH, *J. Phys.-Condens. Mat.* **19**, 335222 (2007).
- ¹¹J. ISCHTWAN and M. A. COLLINS, *J. Chem. Phys.* **100**, 8080 (1994).
- ¹²L. M. RAFF, M. MALSHE, M. HAGAN, D. I. DOUGHAN, M. G. ROCKLEY, and R. KOMANDURI, *J. Chem. Phys.* **122**, 084104 (2005).
- ¹³M. MALSHE, R. NARULKAR, L. M. RAFF, M. HAGAN, S. BUKKAPATNAM, and R. KOMANDURI, *J. Chem. Phys.* **129**, 044111 (2008).
- ¹⁴S. IZVEKOV and G. A. VOTH, *J. Phys. Chem. B* **109**, 6573 (2005).
- ¹⁵P. TANGNEY and S. SCANDOLO, *J. Chem. Phys.* **117**, 8898 (2002).
- ¹⁶G. GROCHOLA, S. P. RUSSO, and I. K. SNOOK, *J. Chem. Phys.* **123**, 204719 (2005).
- ¹⁷S. PARAMORE, L. W. CHENG, and B. J. BERNE, *J. Chem. Theory Comput.* **4**, 1698 (2008).
- ¹⁸W. G. NOID, J. W. CHU, G. S. AYTON, V. KRISHNA, S. IZVEKOV, G. A. VOTH, A. DAS, and H. C. ANDERSEN, *J. Chem. Phys.* **128**, 244114 (2008).
- ¹⁹Y. S. *et al.*, *Phys. Chem. Chem. Phys.* **8**, 3172 (2006).
- ²⁰D. V. DER SPOEL, E. LINDAHL, B. HESS, G. GROENHOF, A. E. MARK, and H. J. C. BERENDSEN, *J. Comp. Chem.* **26**, 1701 (2005).
- ²¹P. BROMMER and F. GAHLER, *Model. Simul. Mater. Sc.* **15**, 295 (2007).
- ²²B. WALDHER, J. KUTA, S. CHEN, N. HENSON, and A. E. CLARK, *J. Comput. Chem.* **31**, 2307 (2010).
- ²³C. J. TSAI and K. D. JORDAN, *Chem. Phys. Lett.* **213**, 181 (1993).
- ²⁴H. J. C. BERENDSEN, J. R. GRIGERA, and T. P. STRAATSMAN, *J. Phys. Chem.* **91**, 6269 (1987).
- ²⁵A. M. FERRENBERG and R. H. SWENDSEN, *Phys. Rev. Lett.* **61**, 2635 (1988).
- ²⁶See Supplementary Material Document No. _____ for force field parameters and additional RDF data. For information on Supplementary Material, see <http://www.aip.org/pubservs/epaps.html>. (Supporting Information)