

## Simple estimation of absolute free energies for biomolecules

F. Marty Ytreberg<sup>a)</sup> and Daniel M. Zuckerman<sup>b)</sup>*Department of Computational Biology, School of Medicine, University of Pittsburgh, Pittsburgh, Pennsylvania 15213 and Department of Environmental and Occupational Health, Graduate School of Public Health, University of Pittsburgh, Pittsburgh, Pennsylvania 15213*

(Received 15 August 2005; accepted 19 January 2006; published online 10 March 2006)

One reason that free energy difference calculations are notoriously difficult in molecular systems is due to insufficient conformational overlap, or similarity, between the two states or systems of interest. The degree of overlap is irrelevant, however, if the absolute free energy of each state can be computed. We present a method for calculating the absolute free energy that employs a simple construction of an exactly computable reference system which possesses high overlap with the state of interest. The approach requires only a physical ensemble of conformations generated via simulation and an auxiliary calculation of approximately equal central-processing-unit cost. Moreover, the calculations can converge to the correct free energy value even when the physical ensemble is incomplete or improperly distributed. As a “proof of principle,” we use the approach to correctly predict free energies for test systems where the absolute values can be calculated exactly and also to predict the conformational equilibrium for leucine dipeptide in implicit solvent. © 2006 American Institute of Physics. [DOI: [10.1063/1.2174008](https://doi.org/10.1063/1.2174008)]

### I. INTRODUCTION

Knowledge of the free energy for two different states or systems of interest allows the calculation of solubilities,<sup>1,2</sup> determines binding affinities of ligands to proteins,<sup>3,4</sup> and determines conformational equilibria (e.g., Ref. 5). Free energy differences ( $\Delta F$ ) therefore have potential application in structure-based drug design where current methods rely on *ad hoc* protocols to estimate binding affinities.<sup>6,7</sup>

Poor “overlap,” the lack of configurational similarity between the two states or systems of interest, is a key cause of computational expense and error in  $\Delta F$  calculations. The most common approach to improve overlap in free energy calculations (used in thermodynamic integration and free energy perturbation) is to simulate the system at multiple hybrid, or intermediate stages (e.g., Refs. 8–12). However, the simulation of intermediate stages greatly increases the computational cost of the  $\Delta F$  calculation.

Here, we address the overlap problem by calculating the absolute free energy for each of the end states, thus avoiding the need for any configurational overlap. Our method relies on the calculation of the free energy difference between a reference system (where the exact free energy can be calculated, either analytically or numerically) and the system of interest.

Such use of a reference system with a computable free energy has been used successfully in solids, where the reference system is generally a harmonic or Einstein solid,<sup>13,14</sup> and liquid systems, where the reference system is usually an ideal gas.<sup>15,16</sup> The scheme has also been applied to molecular systems by Stoessel and Nowak, using a harmonic solid in Cartesian coordinates as a reference system.<sup>17</sup>

Other approaches to calculate the absolute free energies of molecules have been developed. Meirovitch and co-workers calculated absolute free energies for peptides in vacuum, for liquid argon, and water using the hypothetical scanning method.<sup>18,19</sup> Computational cost has thus far limited the approach to peptides with 60 degrees of freedom.<sup>20</sup> The “mining minima” approach, developed by Gilson and co-workers, estimates the absolute free energy of complex molecules by attempting to enumerate the low-energy conformations and estimating the contribution to the configurational integral for each.<sup>21,22</sup> Anharmonic effects can be included.<sup>23</sup> The mining minima method can, in principle, include potential correlations between the torsions and bond angles or lengths, and uses an approximate method to compute local partition functions. Other investigators have estimated absolute free energies for molecules using harmonic or quasiharmonic approximations;<sup>23–25</sup> however, as discussed in Refs. 23 and 24, local minima can deviate substantially from a parabolic shape.

We introduce, apparently for the first time, a reference system which is constructed to have high overlap with fairly general molecular systems. The approach can make use of either *internal* or *Cartesian* coordinates. For biomolecules, using internal coordinates greatly enhances the accuracy of the method since internal coordinates are tailored to the description of conformations. Furthermore, *all degrees of freedom and their correlations* are explicitly included in the method.

Our method differs in several ways from the important study of Stoessel and Nowak:<sup>17</sup> (i) we use internal coordinates for molecules which are key for optimizing the overlap between the reference system and the system of interest; (ii) we may use a nearly arbitrary reference potential because only a numerical reference free energy value is needed, not an analytic value; and (iii) there is no need, in cases we have

<sup>a)</sup>Electronic mail: fmy1@pitt.edu<sup>b)</sup>Electronic mail: dmz@cbb.pitt.edu

studied, to use multistage methodology to find the desired free energy due to the overlap built into the reference system.

We consider this report a “proof of principle” for our reference system method. After introducing the method, it is tested on single- and double-well two-dimensional systems and on a methane molecule where absolute free energy estimates can be compared to exact values. The method is then used to compute the absolute free energy of the alpha and beta conformations for leucine dipeptide (ACE-(leu)<sub>2</sub>-NME) in implicit solvent, using *all* 115 degrees of freedom, correctly calculating the free energy difference  $\Delta F_{\text{alpha}\rightarrow\text{beta}}$ . Extensions of the method to larger systems are then discussed.

## II. REFERENCE SYSTEM METHOD

### A. The fundamental relations

The absolute free energy of the system of interest (“phys” for physical) is defined using the partition function  $Z_{\text{phys}}$ ,

$$F_{\text{phys}} = -k_B T \ln Z_{\text{phys}} \\ = -k_B T \ln \left[ \int d\mathbf{x} e^{-\beta(U_{\text{phys}}(\mathbf{x}) + K_{\text{phys}}(\mathbf{x}))} \right], \quad (1)$$

where  $T$  is the system temperature,  $\beta = 1/k_B T$ ,  $U_{\text{phys}}$  and  $K_{\text{phys}}$  are, respectively, the physical potential energy (i.e., simulation force field) and the kinetic energy, and  $\mathbf{x}$  represents the full set of configurational coordinates (internal or Cartesian).<sup>26</sup> The kinetic energy term can be integrated exactly to obtain<sup>26</sup>

$$Z_{\text{phys}} = \left[ \frac{1}{h^{3N}} \frac{8\pi^2}{\sigma C^0} \prod_{i=1}^N (2\pi k_B T m_i)^{3/2} \right] \int d\mathbf{x} e^{-\beta U_{\text{phys}}(\mathbf{x})}, \quad (2)$$

where  $m_i$  is the mass of atom  $i$ ,  $h$  is Planck’s constant,  $C^0$  is the standard concentration,  $\sigma$  is the symmetry number,<sup>22</sup>  $N$  is the number of particles in the system, and the integral is defined to be the configurational partition function. For the method used in this study the absolute free energy of the system of interest is calculated using a reference system (“ref”), and the following relationships are used:

$$Z_{\text{phys}} = Z_{\text{ref}} \frac{Z_{\text{phys}}}{Z_{\text{ref}}}, \quad (3)$$

$$F_{\text{phys}} = F_{\text{ref}} + \Delta F_{\text{ref}\rightarrow\text{phys}},$$

where  $F_{\text{ref}}$  is the trivially computable free energy of the reference system and  $\Delta F_{\text{ref}\rightarrow\text{phys}}$  is the free energy difference between the reference and physical system which can be calculated using standard techniques.

For this report, we include estimates of the configurational integral only, i.e., the leading constant factor in square brackets in Eq. (2) is not included in our results. Ignoring the constant is not a limitation since, for the conformational free energies studied here, the term cancels for free energy differences.

### B. The reference energy and its normalization

The trivial identities of Eq. (3) suggest that arbitrary reference systems can be used in our approach. To be con-

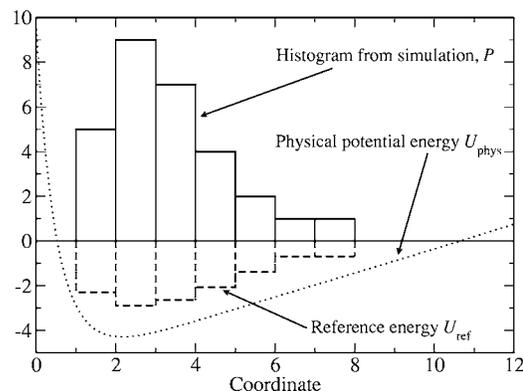


FIG. 1. Depiction of how the reference potential energy  $U_{\text{ref}}$  is calculated for a one-coordinate system. First the coordinate is binned, creating a histogram  $P$  (solid bars) populated according to a simulation. Then Eq. (4) is used to calculate reference energies for each coordinate bin (dashed bars). A hypothetical physical potential is shown as a dotted curve for comparison to  $U_{\text{ref}}$ . For a multicoordinate system  $U_{\text{ref}}$  would be the sum of the single-coordinate reference potential energies.

crete and anticipate the procedure used, our discussion below will assume that a finite-length simulation of the system of interest has been performed—from which histograms of the coordinates have been generated. For the molecular systems studied in this report, ordinary Langevin dynamics simulations are performed using standard force fields. The reference potential energy can be constructed from a wide variety of histograms, as discussed below. Denoting the computed histograms over all coordinates as  $P(\mathbf{x})$ , we define

$$U_{\text{ref}}(\mathbf{x}) \equiv -k_B T \ln P(\mathbf{x}), \quad (4)$$

where  $P(\mathbf{x})$  is the normalized probability of a particular configuration corresponding to a set of histogram bins (see Fig. 1). For example, if all coordinates are binned as independent, then

$$P(\mathbf{x}) = \prod_{i=1}^{N_{\text{coords}}} P_i(x_i), \quad (5)$$

where  $P_i(x_i)$  is the binned probability distribution (histogram) for the  $i$ th coordinate and there are  $N_{\text{coords}}$  degrees of freedom in the system. If all coordinates are binned as pairwise correlated, then

$$P(\mathbf{x}) = \prod_{\{i,j\}} P_{ij}(x_i, x_j), \quad (6)$$

where  $\{i,j\}$  is a set of pairs in which each coordinate occurs exactly once and  $P_{ij}(x_i, x_j)$  is the probability for two particular coordinate values from the two-dimensional histogram for these coordinates. It is also possible to use an arbitrary combination of independent and correlated coordinates—so long as each coordinate occurs in only one  $P$  factor.

We emphasize that the final computed free energy values include all correlations embodied in the true potential  $U_{\text{phys}}$ . This is true regardless of whether or how coordinates are correlated in the reference potential.

A schematic of how  $U_{\text{ref}}$  is computed for a one-coordinate system is shown in Fig. 1. The coordinate histogram is first determined (solid bars) using a simulation trajectory, then Eq. (4) is used to calculate  $U_{\text{ref}}$  (dashed bars). A

possible physical potential is also included (dotted line) for comparison to  $U_{\text{ref}}$ . For a system containing many degrees of freedom, the process is carried out for all coordinates, based on Eqs. (5) and (6) or another correlation scheme.  $U_{\text{ref}}$  is the sum of all the appropriate terms, consistent with Eq. (4) and the binning choice.

The free energy of the reference system can now be calculated via the reference partition function

$$Z_{\text{ref}} = \int d\mathbf{x} e^{-\beta U_{\text{ref}}(\mathbf{x})} = \int d\mathbf{x} P(\mathbf{x}). \quad (7)$$

In practice, we normalize the histogram for each coordinate to 1 independently by summing over all histogram bins. So, for a particular bond length  $r_1$ , that is binned as independent, we account for the Jacobian factor [see Eq. (11)] by defining  $\xi = r_1^3/3$ , and then

$$Z_{\xi} = \int d\xi P(\xi) = \sum_{N_{\text{bin}}} \Delta\xi P(\xi) = 1, \quad (8)$$

where  $\Delta\xi$  is the histogram bin size and  $N_{\text{bin}}$  is the number of bins in the  $r_1$  histogram. (Binning choices are discussed below.) Similar relationships are used for all coordinates. Thus the reference free energy  $F_{\text{ref}}=0$  and Eq. (3) becomes

$$F_{\text{phys}} = \Delta F_{\text{ref} \rightarrow \text{phys}} \quad (F_{\text{ref}} \equiv 0). \quad (9)$$

### C. Using the physical and reference ensembles

With the reference potential energy  $U_{\text{ref}}$  defined in Eq. (4) and the physical potential energy  $U_{\text{phys}}$  given by the force field, which may include implicit solvation energies, Boltzmann-distributed snapshots from both the reference and physical systems can be utilized to calculate  $F_{\text{phys}} = \Delta F_{\text{ref} \rightarrow \text{phys}}$ . Here, we simply use free energy perturbation<sup>8</sup> from the reference to the physical systems,

$$F_{\text{phys}} = -k_B T \ln \langle e^{-\beta(U_{\text{phys}} - U_{\text{ref}})} \rangle_{\text{ref}} \\ \doteq -k_B T \ln \left( \frac{1}{N_{\text{ref}}} \sum_{i=1}^{N_{\text{ref}}} e^{-\beta(U_{\text{phys}}(\mathbf{x}_i) - U_{\text{ref}}(\mathbf{x}_i))} \right), \quad (10)$$

where  $N_{\text{ref}}$  is the number of structures  $\mathbf{x}_i$  in the reference ensemble, the “ $\doteq$ ” symbol denotes a computational estimate, and  $\langle \dots \rangle_{\text{ref}}$  represents a canonical average using structures from the reference ensemble only. It is important to note that, while other choices for computing  $F_{\text{phys}}$  are possible, such as Bennett’s method,<sup>5,27–31</sup> Eq. (10) is the only choice which relies solely on configurations drawn from the reference ensemble which are, by construction, sampled canonically and without dynamical trapping. We also note that “unidirectional” estimates like that of Eq. (10) have been analyzed extensively (e.g., Refs. 32 and 33) and may be amenable to error-reduction techniques;<sup>34,35</sup> however, we have applied the perturbation approach here to keep our initial analysis as straightforward as possible. Staged free energy methods such as thermodynamic integration<sup>36</sup> and adaptive integration<sup>37</sup> may also be used.

### D. The physical ensemble and construction of the reference system

The method used in this report relies on simple histograms for all degrees of freedom (in principle, with internal or Cartesian coordinates) based on a “physical ensemble” of conformations generated via molecular dynamics, Monte Carlo, or other sampling scheme. The histograms define a reference system with a free energy that is trivially computable, as described in Sec. II B. We emphasize that an analytical solution need not be available; a precise numerical evaluation is more than adequate. A well-sampled ensemble of reference system configurations is then readily generated and used to compute the free energy difference via Eq. (10).

The first step in our approach to constructing the reference system is to generate a physical ensemble (i.e., a trajectory) by simulating the system of interest using standard molecular dynamics, Monte Carlo, or other sampling techniques. The trajectory produced by the simulation is used to generate histograms for all coordinates as described below. In creating histograms, note that constrained coordinates, such as bond lengths involving hydrogens constrained by RATTLE,<sup>38</sup> need not be binned since these coordinates do not change between configurations. Such coordinate constraints are not required in the method, however.

If internal coordinates are used (such as for the molecules in this study), care must be taken to account for the Jacobian factors. Using internal coordinates with bond lengths  $r$ , bond angles  $\theta$ , and dihedrals  $\omega$ , the volume element in the configurational integral of Eq. (2) is given by<sup>23</sup>

$$d\mathbf{x} = \prod_{i=1}^{N-1} r_i^2 dr_i \prod_{i=1}^{N-2} \sin \theta_i d\theta_i \prod_{i=1}^{N-3} d\omega_i \\ = \prod_{i=1}^{N-1} d(r_i^3/3) \prod_{i=1}^{N-2} d(-\cos \theta_i) \prod_{i=1}^{N-3} d\omega_i, \quad (11)$$

where  $N$  is the number of atoms in the system. Thus, when using internal coordinates, the simplest strategy to account for the Jacobian is to bin according to a set of rules: bond lengths are binned according to  $r^3/3$ , bond angles are binned according to  $\cos \theta$ , and dihedrals are binned according to  $\omega$  (i.e., the same as Cartesian coordinates).

### E. Generation of the reference ensemble

Once the histograms are constructed and populated using the physical ensemble, the reference ensemble is generated. To generate a single reference structure, for each coordinate, choose a histogram bin according to the probability associated with that bin. Then a coordinate value is chosen at random uniformly within the bin according the Jacobian factor in Eq. (11)—e.g., for a bond length  $r$ , one chooses uniformly in the variable ( $r^3/3$ ). The process is repeated for every degree of freedom in the system. By repeating the entire procedure, one can generate as many reference structures as desired (i.e., the reference ensemble).

TABLE I. Absolute free energy estimates obtained using our reference system approach for cases where the absolute free energy can be determined exactly. In all cases, the estimate is in excellent agreement with the exact free energy. The uncertainty, shown in parentheses [e.g., 3.14 (0.05) = 3.14 ± 0.05] is the standard deviation from five independent simulations. The results for the two-dimensional systems are in  $k_B T$  units and methane results have units of kcal/mole. The table shows estimates of the configurational integral in Eq. (2), i.e., the constant term is not included in the estimate.

System	Exact	Estimate
Two-dimensional single-well <sup>a</sup>	-1.1443	-1.1449(0.0003)
Two-dimensional double-well <sup>a</sup>	5.4043	5.4058(0.0003)
Methane molecule	10.932	10.934(0.002)

<sup>a</sup>Reference 39.

## F. Summary of the reference system method

In summary, the method is implemented by first constructing properly normalized histograms for all internal (or Cartesian) coordinates based on a physical ensemble of structures. An ensemble of reference structures is then chosen at random from the histograms. The reference energy [ $U_{\text{ref}}$  of Eq. (4)] and physical energy ( $U_{\text{phys}}$  from the force field) must be calculated for each structure in the reference ensemble. Finally, Eq. (10) is used to calculate the desired absolute free energy of the system of interest.

The central-processing-unit (CPU) cost of the method, above that of the initial “physical” trajectory, is one physical energy evaluation for each of the  $N_{\text{ref}}$  reference structures, plus the less expensive cost of generating reference structures.

## III. RESULTS

To test the effectiveness of the reference system method we first estimated the absolute free energy for three test systems where the free energy is known exactly. We chose the two-dimensional potentials from Ref. 39 and a methane molecule in vacuum. Finally, we used the method to estimate the absolute free energies of the alpha and beta conformations of the 50-atom leucine dipeptide (ACE-(leu)<sub>2</sub>-NME) and compared the free energy difference obtained via our method with an independent estimate. In all cases, the free energy estimate computed by our approach is in excellent agreement with independent results.

### A. Simple test systems

We first studied the two-dimensional single- and double-well potentials from Ref. 39,

$$U_{\text{phys}}^{\text{single}}(x, y) = (x + 2)^2 + y^2, \quad (12)$$

$$U_{\text{phys}}^{\text{double}}(x, y) = \frac{1}{10} \{ ((x - 1)^2 - y^2)^2 + 10(x^2 - 5)^2 + (x + y)^4 + (x - y)^4 \}.$$

Table I shows the excellent agreement between the reference system estimates and the exact free energies (obtained analytically) for the two-dimensional potentials used in this study [Eq. (12)]. The physical simulations used Metropolis

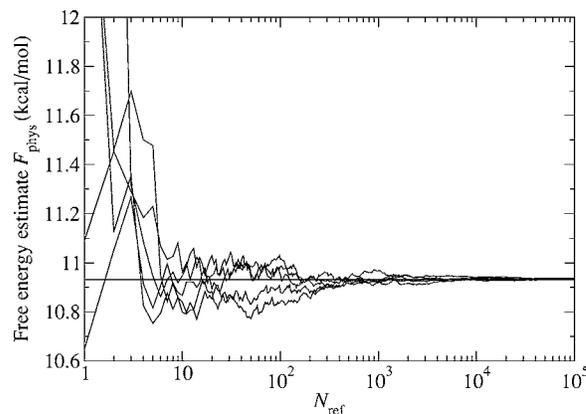


FIG. 2. Absolute free energy for methane estimated by the reference system method as a function of the number of reference structures  $N_{\text{ref}}$  used in the estimate. The solid horizontal line is the exact free energy obtained by numerical integration. Five independent simulations are shown on a log scale to clearly show the convergence of the free energy estimate. The results shown were obtained using Eq. (10) with 100 bins for each degree of freedom, i.e., the estimates for the absolute free energy of methane in Table I are the values shown here for  $N_{\text{ref}}=1\,000\,000$ .

Monte Carlo with  $k_B T=1.0$  and  $1 \times 10^6$  snapshots in the physical and reference ensembles. For all two-dimensional simulations, both coordinates were treated with full correlations—i.e., two-dimensional histograms were used—and the bin sizes were chosen such that the number of bins ranged from 100 to 1000. The error shown in Table I in parentheses is the standard deviation from five independent estimates using five separate physical ensembles—and thus five different reference systems. Good estimates were also obtained using fewer snapshots—e.g., we obtained  $F=-1.142$  (0.003) for the single-well potential and  $F=5.408$  (0.007) for the double-well potential using 10 000 snapshots in both the physical and reference ensembles.

Table I also shows the excellent agreement between the reference system estimates and the exact value of the free energy for methane in vacuum. Methane trajectories were generated using TINKER 4.2 (Ref. 40) with the all-atom optimized potentials for liquid simulations (OPLS-AA) force field.<sup>41</sup> The temperature was maintained at 300.0 K using Langevin dynamics with a friction coefficient of 91.0 ps<sup>-1</sup> and a time step of 0.5 fs. The physical ensemble was created by generating five 10.0 ns trajectories with snapshots saved every 0.1 ps. Using the 100 000 methane structures in the physical ensemble, the reference system was generated by binning internal coordinates into histograms. The absolute free energy was then estimated by generating 100 000 structures for the reference ensemble and using Eq. (10). All coordinates were binned as independent using 100 bins per coordinate, thus only one-dimensional histograms were required. The uncertainty shown in parentheses in Table I is the standard deviation from five independent estimates using the five separate methane trajectories—and thus five different reference systems.

Figure 2 shows the convergence behavior of the reference system method for methane. Five independent absolute free energy estimates are shown as a function of the number of reference structures used in the estimate. Each of the five

simulations uses the same protocol as described above, i.e., the absolute free energy estimates in Table I are the values shown in Fig. 2 for  $N_{\text{ref}}=100\,000$ .

Methane was chosen as a test system because intramolecular interactions are due only to bond lengths and angles. In the OPLS-AA force field no nonbonded terms are present in the potential energy  $U_{\text{phys}}$  for methane, and thus the exact absolute free energy can be computed numerically without great difficulty. For methane, a configuration is determined by (i) four bond lengths, which are independent of each other and all of the other coordinates in the force field, and (ii) five bond angles which are correlated to one another but not to the bond lengths. Thus the exact partition function  $Z_{\text{meth}}$  is a product of four bond length partition functions  $Z_r$  and one angular partition function  $Z_\theta$ ,

$$Z_{\text{meth}} = Z_r Z_\theta,$$

$$Z_r = \int_0^\infty dr e^{-\beta U_{\text{phys}}(r)}, \quad (13)$$

$$Z_\theta = \int_0^\pi d\theta_1 d\theta_2 d\theta_3 d\theta_4 d\theta_5 e^{-\beta U_{\text{phys}}(\theta_1, \theta_2, \theta_3, \theta_4, \theta_5)}.$$

$U_{\text{phys}}(r)$  is harmonic, and thus  $Z_r$  was computed analytically using parameters from the force field. For  $U_{\text{phys}}(\theta_1, \theta_2, \theta_3, \theta_4, \theta_5)$  the correlations between angles must be taken into account; thus  $Z_\theta$  was estimated numerically using TINKER to evaluate  $U_{\text{phys}}$  in the five-dimensional integral. We found that  $F_{\text{meth}} = -k_B T \ln Z_{\text{meth}} = 10.932$  kcal/mol as shown in Table I.

Methane was also used to show that the method correctly computes the free energy even when the physical ensemble is incorrect or incomplete. In our studies we found that the correct free energy is obtained using our method even when the histogram for each coordinate was assumed to be flat, i.e., without the use of a physical ensemble (data not shown).

Choosing the size of the histogram bins is an important consideration. Figure 3 shows the large “sweet spot” where bins are large enough to be well populated, and yet small enough to reveal histogram features. The figure shows results for the absolute free energy for a methane molecule using 10 000 structures in both the physical and reference ensembles [ $N_{\text{phys}}=N_{\text{ref}}=10\,000$ , (dashed curve)] and 100 000 structures in both ensembles [ $N_{\text{phys}}=N_{\text{ref}}=100\,000$  (solid curve)]. The small vertical scale of two kcal/mol and the logarithmic horizontal scale emphasize that there is a wide range of bin sizes that produce excellent results for the reference system approach. Error bars are the standard deviation of five independent simulations. The solid horizontal line shows the exact free energy and the curves are free energy estimates, using Eq. (10) as a function of the number of bins used for the histograms for all degrees of freedom. From this plot it is clear that one should choose at least 50 bins, and that the maximum number of bins that should be used depends on the number of snapshots in the physical ensemble—more snapshots in the physical ensemble means one can use more bins for the reference system.

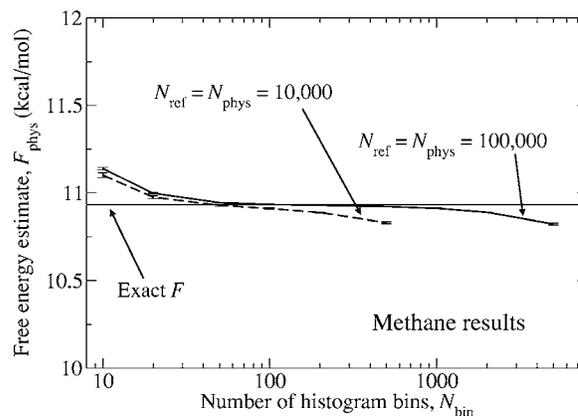


FIG. 3. Absolute free energy for methane estimated by the reference system method as a function of the number of histogram bins used for each degree of freedom. The plot shows the “sweet spot” where histogram bins are small enough to reveal histogram features, yet large enough to give sufficient population in each bin. The results are shown with a vertical scale of 2 kcal/mol and on a log scale to emphasize the wide range of bin sizes that produce excellent results for the reference system approach. The results shown were obtained using Eq. (10) for a methane molecule using  $N_{\text{phys}}=N_{\text{ref}}=10\,000$  (dashed curve) and  $N_{\text{phys}}=N_{\text{ref}}=100\,000$  (solid curve). The solid horizontal line shows the exact free energy and the error bars are the standard deviations of five independent trials. The plot demonstrates that at least 50 bins should be used for each independent coordinate and that the maximum number of bins depends on the number of snapshots in the physical ensemble.

## B. Leucine dipeptide

Table II shows the agreement for leucine dipeptide (ACE-(leu)<sub>2</sub>-NME) between the free energy difference  $\Delta F_{\text{alpha}\rightarrow\text{beta}}$  as predicted by the reference system method and that as predicted via long simulation. The leucine dipeptide physical ensembles were generated using TINKER 4.2 (Ref. 40) with the OPLS-AA force field.<sup>41</sup> The temperature was maintained at 500.0 K (to enable an independent  $\Delta F$  estimate via repeated crossing of the free energy barrier between alpha and beta configurations), using Langevin dynamics with a friction coefficient of 5.0 ps<sup>-1</sup>. Generalized Born Surface area<sup>42</sup> (GBSA) implicit solvation was used, and RATTLE was utilized to maintain all bonds involving hydrogens at their ideal lengths<sup>38</sup> allowing the use of a 2.0 fs time step.

We calculated reference systems and computed absolute

TABLE II. Absolute free energy estimates of the alpha ( $F_{\text{alpha}}$ ) and beta ( $F_{\text{beta}}$ ) conformations obtained using the reference system method for leucine dipeptide with GBSA solvation, in units of kcal/mol. The independent measurement for the free energy difference was obtained via a 1.0  $\mu\text{s}$  unconstrained simulation. The uncertainty for the absolute free energies, shown in parentheses, is the standard deviation from five independent 10.0 ns leucine dipeptide simulations using  $1 \times 10^6$  reference structures in the reference ensemble. The uncertainty for the free energy differences is obtained by using every possible combination of  $F_{\text{alpha}}$  and  $F_{\text{beta}}$ , i.e., 25 independent estimates. The standard error associated with the  $\Delta F_{\text{alpha}\rightarrow\text{beta}}$  reference system estimate is 0.18 kcal/mol, reflecting the 25 independent estimates. The table shows estimates of the configurational integral in Eq. (2), i.e., the constant term is not included in the estimate.

System	Estimate (kcal/mol)	Independent estimate
$F_{\text{alpha}}$	87.3(0.7)	...
$F_{\text{beta}}$	86.3(0.7)	...
$\Delta F_{\text{alpha}\rightarrow\text{beta}}$	-1.0(0.9)	-0.85(0.05)

free energies of the alpha and beta conformations based on five 10.0 ns trajectories. For all simulations, backbone torsions were constrained using a flat-bottomed harmonic restraint (zero force if the torsion angles were within the allowed range, and harmonic otherwise), namely, for alpha,  $-105 < \phi < -45$  and  $-70 < \psi < -10$ , and for beta,  $-125 < \phi < -65$  and  $120 < \psi < 180$ . The reference system was generated using 100 000 snapshots from the physical ensemble, then free energy estimates were obtained by generating 1 000 000 structures for the reference ensemble for each estimate. All 115 internal coordinates [excludes bond lengths constrained by RATTLE (Ref. 38)] were binned as independent with 50 bins for each coordinate. The uncertainty shown in parentheses is the standard deviation from the five independent estimates using the five separate trajectories, i.e., five different physical ensembles and five different reference systems.

Since independent estimates of the absolute free energies of the alpha and beta conformations of leucine dipeptide are not available, we calculated the free energy difference  $\Delta F_{\text{alpha} \rightarrow \text{beta}} = -0.85$  (0.05) kcal/mol via a 1.0  $\mu\text{s}$  unconstrained simulation. The uncertainty of the independent estimate was obtained using block averages. The temperature was chosen to be 500.0 K which allowed around 1500 crossings of the free energy barrier between the alpha and beta conformations, providing an accurate independent estimate. As can be seen in Table II, our estimated free energy difference is in good agreement with the independent value obtained via long simulation.

We emphasize that the nearly kcal/mol fluctuations observed in our leucine dipeptide estimates are completely independent of the magnitude of the free energy difference of the same order. That is, for a similar sized system and similar CPU investment, one would expect similar uncertainty, even for a very large free energy difference. This, indeed, is the motivation for performing absolute free energy calculations. We believe, moreover, that efficiency improvements will be achieved beyond the data in this initial report.

Figure 4 shows the convergence behavior of the reference system method for leucine dipeptide. Five free energy estimates are shown as a function of the number of reference structures used in the estimate for (a) the alpha configuration and (b) the beta configuration. Each of the five simulations uses the same protocol as described above.

The leucine dipeptide calculations also demonstrate two important aspects of the particular reference system defined in this study: (i) the reference system has good overlap with the physical system and (ii) the reference system is broader than the physical system. Figure 5 shows a scatter plot of the  $\chi_2$  torsions of each residue for both the physical and reference ensembles. Each ensemble contains 100 000 structures. The figure clearly shows the excellent overlap between the reference and physical ensembles, as can be seen by the similarity between the two plots. In addition, the reference ensemble scatter plot has data in the region  $(-60, -60)$  which does not exist in the physical ensemble, showing that the reference system is “broader” than the physical system.

Figure 6 shows a histogram of the distance between the  $C_\delta$  atom of residue one and the  $C_\alpha$  of residue 2 for the same

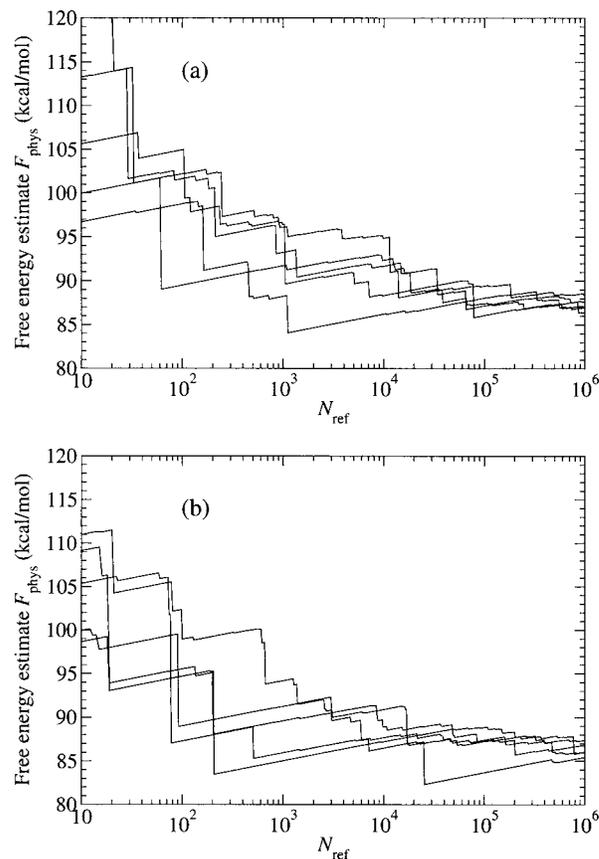


FIG. 4. Free energy for leucine dipeptide estimated by the reference system method as a function of the number of reference structures  $N_{\text{ref}}$  used in the estimate. Five independent simulations are shown on a log scale to demonstrate the convergence behavior of the free energy estimate for (a) the alpha configuration and (b) the beta configuration. The results shown were obtained using Eq. (10) with 50 bins for each degree of freedom.

ensembles as Fig. 5. The figure again shows how the reference system has both excellent overlap with the physical system and is also broader than the physical system.

#### IV. DISCUSSION

The present results raise a number of questions regarding the reference system approach to computing absolute free energies—in particular, regarding the use of correlations, the importance of the physical ensemble, and the potential for application to larger systems.

##### A. Correlation of coordinates

How can correlations among coordinates be used to increase the method’s effectiveness? One may choose to bin coordinates as independent (i.e., one-dimensional histograms) or with correlations (i.e., multidimensional histograms). For example, in peptides, one may choose to bin all sets of backbone  $\phi, \psi$  torsions as correlated, and all other coordinates (bond lengths, bond angles, other torsions) as independent. It might always seem advantageous to bin some coordinates (at least backbone torsions) as correlated, since reference structures drawn randomly from the histograms will be less likely to have steric clashes. On the other hand, including correlations with small bin sizes is impractical. As

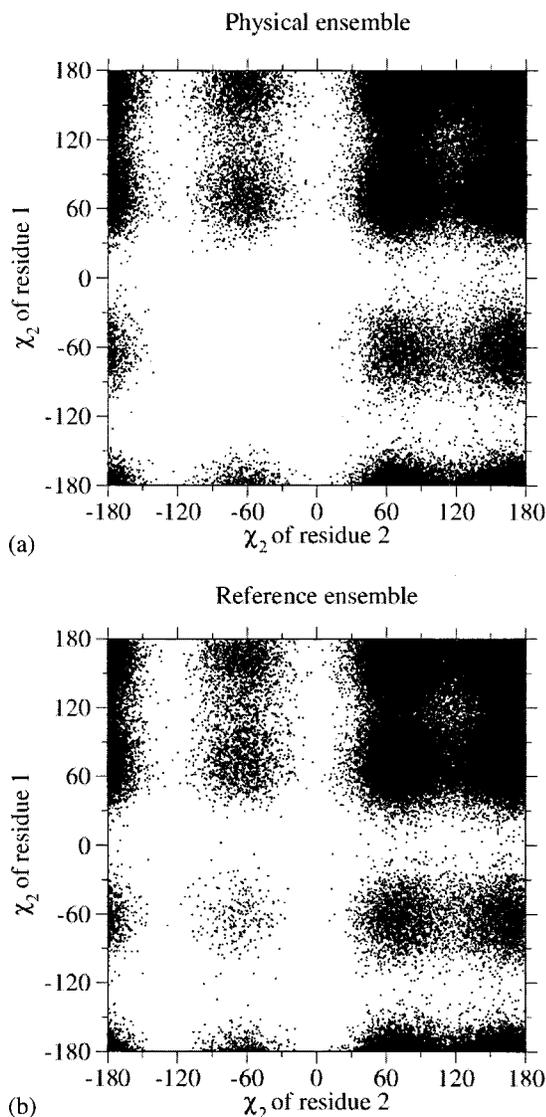


FIG. 5. Scatter plots of the two  $\chi_2$  torsions of each residue for leucine dipeptide. The results are shown for both physical and reference ensembles containing 100 000 structures each. The figure shows that (i) the reference system has good overlap with the physical system, as can be seen by the similarity between the two plots, and (ii) the reference system is more broadly distributed than the physical system, as evidenced by the data at  $(-60, -60)$  for the reference system that is not present for the physical system.

an example, imagine that for the leucine dipeptide molecule used in this study, one binned the four  $\phi, \psi$  backbone torsions as correlated. If 50 bins for each torsion were used (as should be done according to the discussion below), then there would be  $50^4 = 6\,250\,000$  multidimensional bins to populate, which is simply not feasible.

There does appear to be an important advantage to eliminating at least some correlations from the original physical ensemble: namely, a larger portion of conformational space becomes available to the reference ensemble (see Figs. 5 and 6). Since coordinates for the reference structures are drawn randomly and independently, it is possible to generate reference structures that are in entirely different energy basins than those in the physical ensemble. It is thus *possible* to overcome the inadequacies of the physical ensemble by bin-

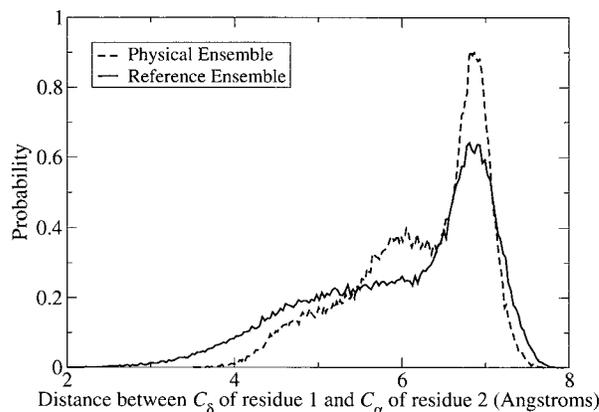


FIG. 6. Histogram of the distance between the  $C_\delta$  of residue 1 and the  $C_\alpha$  of residue 2 for leucine dipeptide. The results are shown for both reference and physical ensembles containing 100 000 structures each. The figure shows that (i) the reference system has good overlap with the physical system and (ii) the reference system is broader than the physical system.

ning internal coordinates *independently*. The optimal (presumably) limited use of correlations will be considered in future work.

Regardless of the degree of correlations included in  $U_{\text{ref}}$ , we emphasize that final results fully include correlations in the physical potential  $U_{\text{phys}}$ .

## B. Quality of the physical ensemble

Since the reference ensemble is generated by drawing at random from histograms which, in turn, were generated from the physical ensemble, a natural question to ask is how complete the physical ensemble needs to be. The surprising answer is that, for our reference system method, the physical ensemble does not need to be complete, or even correct (properly distributed). Since Eqs. (3) and (9) are valid for arbitrary reference systems, the convergence of the free energy estimate to the correct value is guaranteed, in the limit of infinite sampling ( $N_{\text{ref}} \rightarrow \infty$ ), regardless of the quality of the physical ensemble. The “trick” is that the ensemble for the reference system must be converged, which can be achieved with much less expense since there is no dynamical trapping. Unlike the typical case for molecular mechanics simulation, we sample the reference ensemble “perfectly”—there is no possibility of being trapped in a local basin. By construction, since all coordinate values are generated exactly according to the reference distributions, the reference ensemble can only suffer from statistical (but not systematic) error. For example, it was possible to obtain the correct free energy for methane based on 10 000 reference structures even when the histogram for each coordinate was assumed to be flat, i.e., without the use of a physical ensemble (data not shown).

It is important to note that, while convergence to the correct free energy is guaranteed for any choice of reference system, the efficiency of the method could be dramatically reduced if the reference system does not overlap well with the physical system.

Given the fact that the physical ensemble need not be correct, it is easy to imagine a modified method that does not

require simulation, but instead populates the histogram bins using the “bare” potential for each internal coordinate (e.g., Gaussian histograms for bond lengths and angles). Of course, the conformational state must be defined explicitly, with upper and lower limits for coordinates. Allowed ranges for the torsions (especially  $\phi$ ,  $\psi$ ) are naturally obtainable via, e.g., Ramachandran propensities (e.g., Ref. 43), and reasonable ranges for bond lengths and angles could be chosen to be, e.g., several standard deviations from the mean.

### C. Extension to larger systems

While the initial results of our reference system method are promising, a naive implementation of the method will find difficulty with large systems (as do all absolute and relative free energy methods). For our method, the difficulty with including a very large number of degrees of freedom is due to the fact that if one does not treat all correlations in the backbone, then steric clashes will occur frequently when generating the reference ensemble.

However, it is possible to extend the method to larger peptides, still include all degrees of freedom, and bin all coordinates independently (important for broadening configurational space, as discussed above), by using a “segmentation” technique motivated by earlier works.<sup>44,45</sup> Consider generating reference structures for a ten-residue peptide in the alpha helix conformation. Due to the large number of backbone torsions, most of the reference structures chosen at random will not be energetically favorable. However, if one breaks the peptide into two pieces, then one can generate many structures for each segment, and only “keep” energetically likely segment structures. The selected structures may be joined to form full structures which are reasonably likely to have low energy. For example, if one generates  $10^5$  structures for each of the two segments and keeps only  $10^3$  of those, then one only need evaluate  $10^3 \times 10^3 = 10^6$  full structures out of a possible  $10^5 \times 10^5 = 10^{10}$ . A statistically correct segmentation strategy is currently being investigated by the authors for use in large peptides.

Another strategy which may prove useful for larger systems is to use the reference system method with multistage simulation. Multistage simulation requires the introduction of a hybrid potential energy parameterized by  $\lambda$ , e.g.,

$$U_\lambda = \lambda U_{\text{phys}} + (1 - \lambda) U_{\text{ref}}. \quad (14)$$

Thus,  $U_0 = U_{\text{ref}}$  and  $U_1 = U_{\text{phys}}$ . Simulations are performed using the hybrid potential energy  $U_\lambda$  (and thus a hybrid force field, if using molecular dynamics) at intermediate  $\lambda$  values between 0 and 1. Conventional free energy methods such as thermodynamic integration or free energy perturbation can then be used to obtain  $F_{\text{phys}}$ .

We also believe that including correlations, such as suggested by Eq. (6) and possibly other ways, may be useful. The inclusion of correlations should improve the overlap between the reference and physical ensembles—thereby reducing the amount of sampling required in the reference system, hence improving efficiency. This also will be explored in

future work. (We also remind the reader that the final free energy value includes the full correlations in  $U_{\text{phys}}$ , regardless of  $U_{\text{ref}}$ .)

The method could prove useful in future protein-ligand binding studies. In the simplest approach, one could freeze all degrees of freedom except for the ligand and side-chain degrees of freedom in the binding site. While the absolute free energy would be unphysical, the approach could permit comparison of ligands or protein mutations with little or no conformational similarity.

In principle, it is possible to extend the reference system method to include explicitly solvated biomolecules. However, as with all absolute free energy methods, the addition of the solvent degrees of freedom causes the free energy estimate to converge much more slowly than without explicit solvent. Thus, we feel that the method described in this study will find use primarily in implicitly solvated biomolecules.

### V. CONCLUSIONS

In conclusion, we have introduced and tested a simple method for calculating absolute free energies in molecular systems. The approach relies on the construction of an ensemble of reference structures (i.e., the reference system) that is designed to have high overlap with the physical system of interest. The method was first shown to reproduce exactly computable absolute free energies for simple systems and then used to correctly predict the stability of leucine dipeptide conformations using all 115 degrees of freedom.

Some strengths of the approach are that (i) the reference system is built to have good overlap with the system of interest by using internal coordinates and by using a single equilibrium ensemble from Monte Carlo or molecular dynamics; (ii) the absolute free energy estimate is guaranteed to converge to the correct value, whether or not the physical ensemble is complete and, in fact, it is possible to estimate the absolute free energy without the use of a physical ensemble; (iii) the method explicitly includes all degrees of freedom employed in the simulation; (iv) the reference system need only be numerically computable, i.e., the exact analytic result is not needed; and (v) the method can be trivially extended to include the use of multistage simulation. The CPU cost of the approach, beyond that for short trajectories of the physical system of interest, is one energy call for each reference structure, plus the less expensive cost of generating the reference ensemble.

In the present “proof of principle” report, our method was used to study conformational equilibria; however, we feel that the simplicity and flexibility of the method may find broad use in computational biophysics and biochemistry for a wide variety of free energy problems. We have also described a segmentation strategy, currently being pursued, to use the approach in much larger systems.

### ACKNOWLEDGMENTS

The authors would like to thank Edward Lyman, Ronald White, Srinath Chelvarajah, and Hagai Meirovitch for many fruitful discussions. The authors thank the Departments of Computational Biology and Environmental and Occupational

Health at the University of Pittsburgh, and the National Institutes of Health (F32 GM073517) for support.

- <sup>1</sup>A. Grossfield, P. Ren, and J. W. Ponder, *J. Am. Chem. Soc.* **125**, 15671 (2003).
- <sup>2</sup>J. W. Pitera and W. F. van Gunsteren, *J. Phys. Chem. B* **105**, 11264 (2001).
- <sup>3</sup>S. B. Singh, D. E. Wemmer, and P. A. Kollman, *Proc. Natl. Acad. Sci. U.S.A.* **91**, 7673 (1994).
- <sup>4</sup>B. C. Oostenbrink, J. W. Pitera, M. M. van Lipzig, J. H. N. Meerman, and W. F. van Gunsteren, *J. Med. Chem.* **43**, 4594 (2000).
- <sup>5</sup>F. M. Ytreberg and D. M. Zuckerman, *J. Phys. Chem. B* **109**, 9096 (2005).
- <sup>6</sup>B. K. Shoichet, *Nature (London)* **432**, 862 (2004).
- <sup>7</sup>J. Y. Trosset and H. A. Scheraga, *J. Comput. Chem.* **20**, 412 (1999).
- <sup>8</sup>R. W. Zwanzig, *J. Chem. Phys.* **22**, 1420 (1954).
- <sup>9</sup>D. Beveridge and F. DiCapua, *Annu. Rev. Biophys. Chem.* **18**, 431 (1989).
- <sup>10</sup>W. L. Jorgensen and C. Ravimohan, *J. Chem. Phys.* **83**, 3050 (1985).
- <sup>11</sup>W. Yang, R. Bitetti-Putzer, and M. Karplus, *J. Chem. Phys.* **120**, 2618 (2004).
- <sup>12</sup>J. A. McCammon, *Curr. Opin. Struct. Biol.* **2**, 96 (1991).
- <sup>13</sup>W. G. Hoover, S. G. Gray, and K. W. Johnson, *J. Chem. Phys.* **55**, 1128 (1971).
- <sup>14</sup>D. Frenkel and A. J. C. Ladd, *J. Chem. Phys.* **81**, 3188 (1984).
- <sup>15</sup>W. G. Hoover and F. H. Ree, *J. Chem. Phys.* **47**, 4873 (1967).
- <sup>16</sup>L. M. Amon and W. P. Reinhardt, *J. Chem. Phys.* **113**, 3573 (2000).
- <sup>17</sup>J. P. Stoessel and P. Nowak, *Macromolecules* **23**, 1961 (1990).
- <sup>18</sup>S. Chelvaraja and H. Meirovitch, *Proc. Natl. Acad. Sci. U.S.A.* **101**, 9241 (2004).
- <sup>19</sup>R. P. White and H. Meirovitch, *Proc. Natl. Acad. Sci. U.S.A.* **101**, 9235 (2004).
- <sup>20</sup>S. Chelvaraja and H. Meirovitch, *J. Chem. Phys.* **1022**, 054903 (2004).
- <sup>21</sup>M. S. Head, J. A. Given, and M. K. Gilson, *J. Phys. Chem. A* **101**, 1609 (1997).
- <sup>22</sup>M. K. Gilson, J. A. Given, B. L. Bush, and J. A. McCammon, *Biophys. J.* **72**, 1047 (1997).
- <sup>23</sup>C. E. Chang and M. K. Gilson, *J. Am. Chem. Soc.* **126**, 13156 (2004).
- <sup>24</sup>M. Karplus and J. N. Kushick, *Macromolecules* **14**, 325 (1981).
- <sup>25</sup>J. Carlsson and J. Aqvist, *J. Phys. Chem. B* **109**, 6448 (2005).
- <sup>26</sup>C. E. Chang, M. J. Potter, and M. K. Gilson, *J. Phys. Chem. B* **107**, 1048 (2003).
- <sup>27</sup>C. H. Bennett, *J. Comput. Phys.* **22**, 245 (1976).
- <sup>28</sup>M. R. Shirts and V. S. Pande, *J. Chem. Phys.* **122**, 144107 (2005).
- <sup>29</sup>M. R. Shirts, E. Bair, G. Hooker, and V. S. Pande, *Phys. Rev. Lett.* **91**, 140601 (2003).
- <sup>30</sup>G. E. Crooks, *Phys. Rev. E* **61**, 2361 (2000).
- <sup>31</sup>N. Lu, D. A. Kofke, and T. B. Woolf, *J. Comput. Chem.* **25**, 28 (2004).
- <sup>32</sup>D. M. Zuckerman and T. B. Woolf, *Phys. Rev. Lett.* **89**, 180602 (2002).
- <sup>33</sup>D. M. Zuckerman and T. B. Woolf, *J. Stat. Phys.* **114**, 1303 (2004).
- <sup>34</sup>D. M. Zuckerman and T. B. Woolf, *Chem. Phys. Lett.* **351**, 445 (2002).
- <sup>35</sup>F. M. Ytreberg and D. M. Zuckerman, *J. Comput. Chem.* **25**, 1749 (2004).
- <sup>36</sup>T. P. Straatsma and J. A. McCammon, *J. Chem. Phys.* **95**, 1175 (1991).
- <sup>37</sup>M. Fasnacht, R. H. Swendsen, and J. M. Rosenberg, *Phys. Rev. E* **69**, 056704 (2004).
- <sup>38</sup>H. C. Andersen, *J. Comput. Phys.* **52**, 24 (1983).
- <sup>39</sup>F. M. Ytreberg and D. M. Zuckerman, *J. Chem. Phys.* **120**, 10876 (2004).
- <sup>40</sup>J. W. Ponder and F. M. Richard, *J. Comput. Chem.* **8**, 1016 (1987); <http://dasher.wustl.edu/tinker>.
- <sup>41</sup>W. L. Jorgensen, D. S. Maxwell, and J. Tirado-Rives, *J. Am. Chem. Soc.* **117**, 11225 (1996).
- <sup>42</sup>W. C. Still, A. Tempczyk, and R. C. Hawley, *J. Am. Chem. Soc.* **112**, 6127 (1990).
- <sup>43</sup>S. C. Lovell, I. W. Davis, W. B. Arendall III, P. I. W. de Bakker, J. M. Word, M. G. Prisant, J. S. Richardson, and D. C. Richardson, *Proteins* **50**, 437 (2003).
- <sup>44</sup>K. D. Gibson and H. A. Scheraga, *J. Comput. Chem.* **8**, 826 (1987).
- <sup>45</sup>A. R. Leach, K. Prout, and D. P. Dolata, *J. Comput.-Aided Mol. Des.* **2**, 107 (1988).