

Calculation of absolute protein-ligand binding free energy using distributed replica sampling

Tomas Rodinger, P. Lynne Howell, and Régis Pomès

Citation: *J. Chem. Phys.* **129**, 155102 (2008); doi: 10.1063/1.2989800

View online: <http://dx.doi.org/10.1063/1.2989800>

View Table of Contents: <http://jcp.aip.org/resource/1/JCPSA6/v129/i15>

Published by the [American Institute of Physics](#).

Related Articles

Mode of bindings of zinc oxide nanoparticles to myoglobin and horseradish peroxidase: A spectroscopic investigations

J. Appl. Phys. **110**, 024701 (2011)

A water-swap reaction coordinate for the calculation of absolute protein–ligand binding free energies

JCP: BioChem. Phys. **5**, 02B611 (2011)

A water-swap reaction coordinate for the calculation of absolute protein–ligand binding free energies

J. Chem. Phys. **134**, 054114 (2011)

The axial methionine ligand may control the redox reorganizations in the active site of blue copper proteins

JCP: BioChem. Phys. **4**, 11B601 (2010)

The axial methionine ligand may control the redox reorganizations in the active site of blue copper proteins

J. Chem. Phys. **133**, 175101 (2010)

Additional information on *J. Chem. Phys.*

Journal Homepage: <http://jcp.aip.org/>

Journal Information: http://jcp.aip.org/about/about_the_journal

Top downloads: http://jcp.aip.org/features/most_downloaded

Information for Authors: <http://jcp.aip.org/authors>

ADVERTISEMENT

**AIP**Advances

Submit Now

Explore AIP's new
open-access journal

- Article-level metrics now available
- Join the conversation! Rate & comment on articles

Calculation of absolute protein-ligand binding free energy using distributed replica sampling

Tomas Rodinger,^{1,2,3,a)} P. Lynne Howell,^{1,2} and Régis Pomès^{1,2,b)}

¹*Molecular Structure and Function, The Hospital for Sick Children, 555 University Ave., Toronto, Ontario M5G 1X8, Canada*

²*Department of Biochemistry, University of Toronto, 1 King's College Circle, Toronto, Ontario M5S 1A8, Canada*

³*Institute of Biomaterials and Biomedical Engineering, University of Toronto, 164 College Street, Toronto, Ontario M5S 3G9, Canada*

(Received 24 May 2007; accepted 3 September 2008; published online 16 October 2008)

Distributed replica sampling [T. Rodinger *et al.*, *J. Chem. Theory Comput.* **2**, 725 (2006)] is a simple and general scheme for Boltzmann sampling of conformational space by computer simulation in which multiple replicas of the system undergo a random walk in reaction coordinate or temperature space. Individual replicas are linked through a generalized Hamiltonian containing an extra potential energy term or bias which depends on the distribution of all replicas, thus enforcing the desired sampling distribution along the coordinate or parameter of interest regardless of free energy barriers. In contrast to replica exchange methods, efficient implementation of the algorithm does not require synchronicity of the individual simulations. The algorithm is inherently suited for large-scale simulations using shared or heterogeneous computing platforms such as a distributed network. In this work, we build on our original algorithm by introducing Boltzmann-weighted jumping, which allows moves of a larger magnitude and thus enhances sampling efficiency along the reaction coordinate. The approach is demonstrated using a realistic and biologically relevant application; we calculate the standard binding free energy of benzene to the L99A mutant of T4 lysozyme. Distributed replica sampling is used in conjunction with thermodynamic integration to compute the potential of mean force for extracting the ligand from protein and solvent along a nonphysical spatial coordinate. Dynamic treatment of the reaction coordinate leads to faster statistical convergence of the potential of mean force than a conventional static coordinate, which suffers from slow transitions on a rugged potential energy surface. © 2008 American Institute of Physics. [DOI: [10.1063/1.2989800](https://doi.org/10.1063/1.2989800)]

INTRODUCTION

Although computational power continues to grow, performing simulations of large and complex biomolecular systems remains a challenge. Typical simulations of biopolymers in explicit solvent involve on the order of 10^4 atoms and are limited to timescales on the order of nanoseconds on a single CPU. The need to attain converged thermodynamic averages for systems of increasing size and complexity has fueled the continual development of methods to spread the computational work over multiple CPUs and has led to new algorithms designed to improve sampling efficiency of rugged energy surfaces.¹⁻⁵

The speed at which a single system can be simulated can be dramatically improved by parallelizing the workload over multiple CPUs. In a typical parallelized molecular simulation, each CPU concentrates on calculations relating to a region of space.⁶ This approach, which scales well with large systems, has enabled simulations of the ribosome (2.64×10^6 atoms for a few nanoseconds)⁷ and of an entire virus

(a million atoms for 50 ns).⁸ Such simulations require expensive supercomputers or dedicated computing clusters. An alternative to large-scale simulations is to run independent single-CPU simulations of a molecular system on as many CPUs as possible, a strategy embodied by distributed computing. In such trivially parallelizable calculations, simultaneous simulations may be employed in an effort to sample phase space as exhaustively as possible. Such extensive sampling is useful if the reaction coordinate for the process of interest is unknown (such as in protein folding). Distributed computing has recently been applied to the study of a polyalanine helix, using 20 000 distributed CPUs for a combined simulation time of 800 μ s.⁹

Concomitant to the use of multiple CPUs, new algorithms have been devised in the past decade to further improve sampling efficiency by helping overcome potential energy barriers. This approach is the basis for several generalized-ensemble algorithms that allow random walks in temperature.¹⁰⁻¹² One such method, replica exchange (RE), also known as multiple Markov chain or parallel tempering, is explicitly designed to run on multiple processors.¹³⁻¹⁵ Multiple noninteracting replicas of a system, each governed by the same potential energy function but differing in temperature, are simulated at once on separate CPUs. Periodi-

^{a)}Present address: Zymeworks, Inc., 540-1385 West 8th Ave., Vancouver, British Columbia V6H 3V9, Canada

^{b)}Author to whom correspondence should be addressed. Tel.: 416-813-5686. FAX: 416-813-5022. Electronic mail: pomes@sickkids.ca.

cally, the simulations are halted and replicas i and j with neighboring temperatures T_i and T_j are swapped. The RE method has enabled a significant increase in the size and complexity of the systems studied.^{14,16–25} The method was used to fold a 46 amino acid protein domain²⁶ and was extended to two-dimensional random walks in pressure and temperature.²⁴ The principal drawback of RE is that all the replicas must run synchronously so that an efficient implementation of the algorithm requires a dedicated and homogeneous cluster. In addition, because the number of replicas required scales as the square root of the number of degrees of freedom in the system,¹³ simulations of complex systems typically require a large number of CPUs.

An important application of large-scale biomolecular simulations, the calculation of protein-ligand binding free energies involves sampling a potential energy landscape which is too large and too rugged for exhaustive sampling. In such cases, a reaction coordinate must be chosen to restrict sampling to meaningful regions of conformational space. To this end, constraints are used to progressively transform the system from an initial state to a final state along the chosen coordinate. Techniques for such calculations include free energy perturbation,²⁷ thermodynamic integration (TI),²⁸ and umbrella sampling (US),²⁹ all of which can be applied to physical or nonphysical (alchemical) transformations.^{1–5} These simulations are typically performed using windowing to restrict sampling to small intervals of the reaction coordinate. The free energy change for the whole reaction is then obtained by combining the results obtained for each window. Since each window can be simulated independently, such approaches are naturally suited for multiple CPUs and distributed computing, a strategy recently used to calculate the hydration free energies of the naturally occurring amino acids.³⁰

However, imposing artificial restraints on the reaction coordinate increases the ruggedness of the potential energy. While this is not a problem in situations where no significant barriers exist in degrees of freedom orthogonal to the transformation path (for example, in the calculation of hydration free energies of simple solutes³¹), in general, ruggedness in the degrees of freedom perpendicular to the reaction coordinate can lead to systematic sampling errors by trapping the system in a high-energy state.^{32,33} In protein-ligand simulations, this situation may occur if the orientation or the conformation of the bound ligand is incorrect. In general, the formation of a molecular complex results in a decrease in the conformational freedom of protein and ligand. In addition, both bound and unbound states may consist of several conformations that are populated significantly. Systematic approaches to avoid the need to cross rotational energy barriers have recently been proposed to help alleviate this problem.^{34,35}

Alternatively, treating the transformation coordinate as a dynamic variable decreases ruggedness and increases the probability of crossing energy barriers. This is because the higher dimensionality of phase space results in additional routes by which barriers can be avoided. Barriers that exist at one coordinate position (e.g., in the bound state) may vanish at other coordinate positions (the unbound state), enabling

transitions that would be nearly impossible without dynamic coordinate movement. Thus, a dynamic coordinate allows a ligand to move out of a protein binding site, reorient, and move back in. Two techniques designed to induce a random walk in the transformation parameter are adaptive umbrella sampling^{36–40} and RE along the parameter (also known as Hamiltonian exchange).^{20,41} Both methods aim to achieve uniform sampling along the reaction coordinate. While RE does not require adaptation to achieve perfect sampling uniformity, as noted above the algorithm requires a large dedicated cluster.

To circumvent this limitation and make large-scale simulations more practical, we recently introduced distributed replica (DR) sampling, a simple and general scheme for efficient Boltzmann sampling of conformational space.³² As in RE, multiple replicas of the system covering a preassigned range in temperature or reaction coordinate are simulated independently. However, instead of pairwise exchanges, stochastic moves of individual replicas are considered one at a time. The coupling between replicas is attained by a generalized Hamiltonian containing an extra potential energy bias that depends on the distribution of all replicas and acts to enforce a target distribution. Like RE, DR does not require adaptation. The algorithm leads to a random walk with an efficiency comparable to RE.³² However, by avoiding the need for all replicas to run synchronously, DR is inherently suited for a shared or inhomogeneous computing cluster, or even a large-scale distributed network. Thus, DR combines an efficient (barrier-crossing) sampling algorithm with accessible large-scale computing. In particular, DR can scale up to any number of replicas running on heterogeneous platforms in a trivial way. Another advantage of DR over RE is that the magnitude of stochastic moves need not be restricted to preassigned values (as required in an exchange process). Instead, moves of arbitrary magnitude can be considered, which makes it possible to further improve the efficiency of the random walk.

In our previous work, we used a simple model system containing a potential energy barrier orthogonal to the reaction coordinate.³² The use of replica simulation techniques, such as DR or RE, to achieve a random walk along the reaction coordinate, was shown to lead to the correct Boltzmann distribution much more quickly than independent simulations.^{32,33} Here we utilize the DR algorithm to its full advantage as follows: (1) We build on the original algorithm by introducing Boltzmann-weighted jumping, which enables stochastic moves of arbitrary amplitudes and (2) we apply the approach to a challenging large-scale biomolecular simulation by computing the absolute binding affinity of a molecular ligand to a protein. The approach is demonstrated in an engineered T4 lysozyme protein in which the single-point mutation L99A results in a buried hydrophobic cavity able to bind benzene and other similarly sized hydrophobic molecules.^{42,43} This system is a good benchmark for computational studies of protein-ligand binding free energies because high-resolution structures of the enzyme have been determined by x-ray crystallography both in its apo and complexed forms.⁴³ In addition, the binding strengths of a num-

ber of these ligands have been measured experimentally,⁴² providing values to which the computed binding free energies can be compared.

In the remainder of this article, we describe new methodological developments and their application to the calculation of the absolute binding free energy of benzene to T4 lysozyme. DR is combined with TI to compute the potential of mean force (PMF) for the decoupling of protein-ligand interactions in four-dimensional (4D) space.^{31,44} This combination of approaches enables us to (i) calculate the ligand insertion/extraction free energy in a single series of simulations, (ii) maximize the use of available computing resources, and (iii) facilitate the sampling of conformational degrees of freedom whose transition probabilities vary with the extent of coupling between protein and ligand via a random walk in the coupling parameter. The improvement of sampling efficiency over simulations in which the transformation coordinate is static is examined. Since we aim to design methodology and execution protocols to be used on shared clusters or large-scale distributed computing platforms, where CPU resources change unpredictably, we also examine the effect of CPU availability on the calculated free energy result.

THEORY AND METHODS

Distributed replica sampling and Boltzmann-weighted jumping

Consider N noninteracting copies (or “replicas”) of a system governed by an identical potential energy function, $E(\mathbf{q}_i, \lambda_i)$, where \mathbf{q}_i represents atomic coordinates of the atoms in replica i and λ_i is the coupling parameter for the reaction coordinate of interest (the reaction in question may be either an alchemical or a spatial transformation). The DR method makes use of an additional potential energy term $D(\lambda_1, \lambda_2, \dots, \lambda_N)$, henceforth referred to as the distributed replica potential energy (DRPE), which enforces the distribution of replicas across the range of the transformation coordinate (i.e., an energy penalty is associated with a nonideal distribution). The generalized Hamiltonian for all replicas together with the DRPE is given by

$$H_\lambda = \sum_{i=1}^N E(\mathbf{q}_i, \lambda_i) + D(\lambda_1, \lambda_2, \dots, \lambda_N). \quad (1)$$

The weight factor for a state $\mathbf{X} = \{\mathbf{q}_1, \lambda_1, \mathbf{q}_2, \lambda_2, \dots, \mathbf{q}_N, \lambda_N\}$ is given by

$$W(X) = \exp(-\beta H_\lambda). \quad (2)$$

We consider one λ move at a time. Suppose that the λ value of replica m is to be changed from λ_m to $\lambda_m + \delta\lambda_m$, thus taking state X to state X'

$$\begin{aligned} X &= \{\mathbf{q}_1, \lambda_1, \dots, \mathbf{q}_m, \lambda_m, \dots, \mathbf{q}_N, \lambda_N\} \\ \rightarrow X' &= \{\mathbf{q}_1, \lambda_1, \dots, \mathbf{q}_m, \lambda_m + \delta\lambda_m, \dots, \mathbf{q}_N, \lambda_N\}. \end{aligned}$$

In order for the exchange process to converge toward the equilibrium distribution, it is sufficient to impose the detailed balance condition on the transition probability $p(X \rightarrow X')$

$$W(X)p(X \rightarrow X') = W(X')p(X' \rightarrow X). \quad (3)$$

From Eqs. (1)–(3), we have

$$\frac{p(X \rightarrow X')}{p(X' \rightarrow X)} = \exp(-\beta\Delta), \quad (4)$$

where

$$\begin{aligned} \Delta &= E(\mathbf{q}_m, \lambda_m + \delta\lambda_m) - E(\mathbf{q}_m, \lambda_m) \\ &\quad + D(\lambda_1, \lambda_2, \dots, \lambda_m + \delta\lambda_m, \dots, \lambda_N) \\ &\quad - D(\lambda_1, \lambda_2, \dots, \lambda_m, \dots, \lambda_N). \end{aligned} \quad (5)$$

This can be satisfied using the Metropolis Monte Carlo criterion

$$p_{\text{accept}} = \min[1, \exp(-\beta\Delta)]. \quad (6)$$

Another way to satisfy Eq. (4) is to consider at once the energy change associated with a jump of replica i from its current position λ_i to all possible states in a discretized λ space. The N possible states have λ values of Λ_j , where j is the index for the possible states and E_j is the potential energy change for jumping to state j . A separate Δ_j is calculated for each state

$$\Delta_j = E(\mathbf{q}_i, \Lambda_j) + D(\lambda_1, \lambda_2, \dots, \Lambda_j, \dots, \lambda_N). \quad (7)$$

The normalized probability p_j of each of the possible states j is then given by

$$p_j = \frac{\exp(-\beta\Delta_j)}{\sum_{k=1}^N \exp(-\beta\Delta_k)}. \quad (8)$$

The new value of the transformation parameter is appropriately chosen out of the N possible states based on their probabilities; λ_i is updated. This process will be referred to as a Boltzmann-weighted jump; see Table I for an example calculation. Theoretically, a Boltzmann-weighted jump can take the transformation coordinate of a replica from one place to any other in a single move.

The DRPE function D is calculated using the following three algorithmic steps. First, the λ (or β) values for all replicas are sorted in ascending order. The following holds true for the new order, $\lambda_{i,\text{sorted}}$:

$$\lambda_{i,\text{sorted}} > \lambda_{i-1,\text{sorted}} \quad \text{for } i = 2 \text{ to } N. \quad (9)$$

Second, the spacing system is transformed to a uniform unit spacing arrangement to give $\lambda_{i,\text{unit}}$

$$\lambda_{i,\text{unit}} = f^{-1}(\lambda_{i,\text{sorted}}), \quad (10)$$

where f^{-1} is the inverse of a function f which maps the replica index to the nominal λ value of that replica (i.e., the λ position where the replica started). Note that since replica indices are integers, f is constructed by linearly interpolating between adjacent points. Finally, the function D is calculated as

TABLE I. Example Boltzmann-weighted jump calculation for replica 3.

Index j	1	2	3	4	5	6
Nominal λ positions of replicas (Λ_j)	0.0	0.1	0.2	0.4	0.6	1.0
λ position before jump (λ_j)	0.2	0.0	0.1	1.0	0.1	0.6
Sorted λ before jump ($\lambda_{j,\text{sorted}}$)	0.0	0.1	0.1	0.2	0.6	1.0
Linearized λ before jump ($\lambda_{j,\text{unit}}$)	1	2	2	3	5	6
DRPE scaling constants (c_1, c_2)			0.1, 0.1			
DRPE (D) value before jump			First term=1.6; Second term=0.4, Total=2.0			
DRPE (D) value if λ_3 moves to Λ_j	2.7	2.0	1.1	0.0	1.1	2.0
System energy of replica 3 if λ_3 moves to Λ_j	100.3	100.0	100.2	100.6	101.0	102.5
Total energy change for jump (Δ_j)	103.0	102.0	101.3	100.6	102.1	104.5
$\exp(-\beta\Delta_j)$	0.018	0.097	0.311	1.000	0.082	0.002
$\sum_{k=1}^N \exp(-\beta\Delta_k)$			1.510			
Probability (p_j) to jump $\lambda_3 \rightarrow \Lambda_j$	0.012	0.064	0.206	0.662	0.054	0.001
Condition to accept jump	$0.000 \leq R < 0.012$	$0.012 \leq R < 0.076$	$0.076 \leq R < 0.283$	$0.283 \leq R < 0.945$	$0.945 \leq R < 0.999$	$0.999 \leq R < 1.000$
Random number ($0 \leq R < 1$)			0.78			
New position of λ_3			0.4			
DRPE (D) value after jump			First term=0.0; Second term=0.0; Total=0.0			

$$D = c_1 \sum_{i=1}^N \sum_{j=1}^N [(\lambda_{i,\text{unit}} - \lambda_{j,\text{unit}}) - (i - j)]^2 + c_2 \left[\sum_{i=1}^N \lambda_{i,\text{unit}} - \sum_{i=1}^N i \right]^2, \quad (11)$$

where c_1 and c_2 are parameters that scale the DRPE function for adjusting the severity of the penalty associated with non-ideal distribution of replicas. The first term in Eq. (11) enforces the spacing between replicas and has no effect on their absolute positions. The second term prevents a concerted drift of all λ_i values away from the region of interest.

A DR simulation is realized as follows. Initially, each replica i is created at a different position, λ_i , spanning the transformation coordinate, and is optionally equilibrated. The spacing between adjacent λ_i values is chosen based on the application at hand and may be uniform or nonuniform. The following two steps are then iterated. First, each replica is run as an independent molecular dynamics (MD) or Monte Carlo simulation at a fixed λ_i value for a set number of steps or period of time. Second, periodically, one replica is considered for a λ move which is accomplished using either the Metropolis Monte Carlo [Eq. (6)] or the Boltzmann-weighted jumping [Eq. (8)] formalism.

Calculation of absolute binding free energies

To compute the absolute binding free energy between ligand (benzene) and protein receptor (T4 lysozyme), we use a method consisting of three separate computational steps: The ligand is extracted from the binding site, brought to a standard concentration (1M), and inserted into bulk water. The thermodynamic pathway (depicted in Fig. 1) takes advantage of the efficiency gains afforded by alchemical trans-

formations in 4D space. Furthermore, the pathway provides a correction to a well-defined standard state so as to make comparisons to experimental data possible.

The main focus of this paper is on the first step, which is by far the most computationally demanding and challenging, and can therefore benefit the most from DR sampling. In the first step, TI in four spatial dimensions is used to decouple a

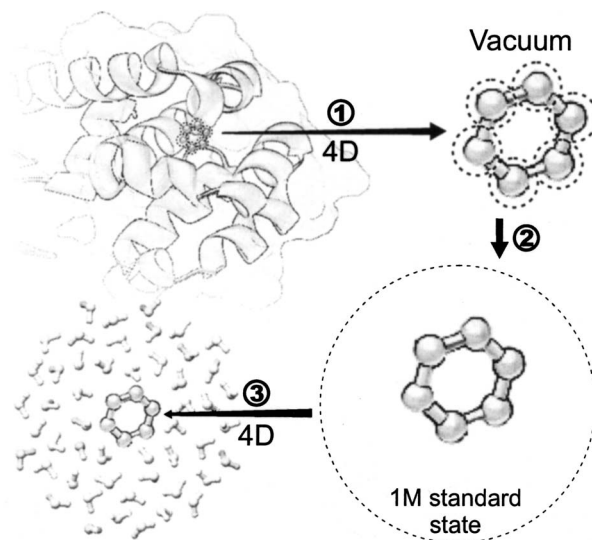


FIG. 1. Thermodynamic pathway used to calculate the absolute binding free energy of a ligand (benzene) to an enzyme (T4 lysozyme). (1) The ligand is extracted from the binding site into vacuum along a fourth spatial dimension; spatial restraints on the heavy atoms of the ligand limit the ligand's mobility once it leaves the binding site. (2) The free energy change associated with the removal of the restraints is accounted for, leaving the ligand inside a single boundary that mimics a 1M standard state. (3) The hydration free energy of the ligand is computed by inserting the ligand into a water droplet along the fourth dimension; note that this step does not depend on solute concentration; this calculation is performed with a harmonic potential that keeps the ligand near the center of the water droplet.

ligand molecule bound to a receptor protein and bring it into a noninteracting state (vacuum). Using a fourth dimension as the coupling parameter provides an efficient means of modulating nonbonded interactions similar to alchemical soft-core type methods (where the parameter, λ , drives the system between initial and final states).⁴³ Here, the transformation parameter is w , a spatial coordinate that describes the separation distance between the ligand and receptor along the 4D axis. This unphysical spatial coordinate is used to decouple protein-ligand interactions, which are expressed as a function of pairwise atomic separations computed in four spatial dimensions (x , y , z , and w). The system is fully coupled (ligand inserted in the binding site) at $w=0$ and fully decoupled (ligand in an unbound state) at $w=\infty$. The 4D method was designed to avoid steep steric barriers along the reaction pathway as well as to decrease the number of simulation windows needed to carry out the required sampling of the full coordinate. The formalism of using 4D space for free energy calculations by US and TI has been described in detail elsewhere.^{31,44} The method was shown to provide a simple and general reaction pathway for decoupling all nonbonded interactions in a single calculation, thus providing an alternative to more conventional methods in which the decoupling of Coulombic and short- and long-range Lennard-Jones interactions is achieved in separate steps.⁴⁵ Furthermore, TI leads to smaller statistical sampling errors than US in this type of application.^{31,46}

Throughout the extraction step, a spherical restraint placed on each heavy atom of the ligand and centered at the atom's mean bound position exerts a force F as follows:

$$F = \begin{cases} k(r - r_{\text{off}})^2, & r > r_{\text{off}} \\ 0, & r \leq r_{\text{off}} \end{cases}, \quad (12)$$

where k is a force constant (or stiffness of the restraint), r is the distance of the heavy atom from the sphere's center, and r_{off} is the restraint radius. Note that for these restraints, r is calculated based on the atoms' three-dimensional coordinates only. Thus, the restraint continues to act even when the ligand is extracted into the fourth dimension. r_{off} values are made large enough such that when the ligand is fully bound (and thus maximally restrained by the receptor), the restraints exert a negligible influence. However, as the ligand travels into the fourth dimension, and protein-ligand interactions vanish, the restraints prevent the ligand from leaving the vicinity of the binding site (in three dimensions). Without these restraints, the ligand would not sample a well-defined volume and the calculation of a standard binding free energy would not be possible. Furthermore, a ligand undergoing a random walk in coordinate space, as in a DR simulation, would be unlikely to return to the original binding pocket once it leaves, resulting in an irreversible transformation.

The restraints imposed on the heavy atoms of the ligand are removed in step two of the overall binding free energy calculation (see Fig. 1), and at the same time, the effective concentration of the ligand is brought to a 1M standard state. The importance of a well-defined standard state and its theoretical basis have been described elsewhere.⁴⁷⁻⁴⁹ Briefly, the entropy of the dissociated state depends on the volume sampled by the ligand (i.e., its concentration) after it is ex-

tracted from the protein. Therefore, the volume sampled by the extracted ligand directly affects the outcome of the binding free energy calculation and must be accounted for to make the comparison with experimental binding constants possible. The free energies for removing the imposed restraints are calculated as follows. The restrained ligand is placed inside a spherical boundary with radius 7.3 Å (volume of 1660 Å³ corresponding to a 1M concentration). Initially, the individual restraints acting on the ligand's heavy atoms prevent it from interacting with this boundary. One by one, the atomic restraints are released by increasing their r_{off} parameter until they are no longer interacting with the ligand (the average force F and pressure P exerted by each atom on the restraint boundary go to zero). Once all the atomic restraints have been released, the ligand is allowed to sample the 1M standard state boundary. The free energy associated with releasing the restraints is calculated by integrating the average pressure exerted by the ligand atom on the restraint as the latter changes in volume

$$\Delta G = \int \langle P \rangle dV = \int \langle F \rangle dr_{\text{off}}, \quad (13)$$

where $\langle P \rangle$ is the average pressure exerted on the restraint by the heavy atom, V is the restraint's volume calculated from r_{off} , and $\langle F \rangle$ is the average magnitude of the outward force [see Eq. (12)] exerted on the boundary by the heavy atom over the course of a long MD simulation.

The third step of the binding free energy calculation involves determining the hydration free energy of the ligand, which can be done routinely in a similar way as in step one, by extracting the ligand from a droplet of water along the fourth dimension. Note that the free energy change associated with step three (an insertion) is the negative of the calculated result (an extraction). DR sampling was not used for this step as the potential energy landscape is not rugged and a random walk is not beneficial. Instead, we applied TI along a static 4D coordinate, with each discrete value of the coordinate simulated independently. The methodology for this step was described in detail previously.³¹

In general, an additional step may be required to account for the free energy change for filling the cavity left behind after extraction of the ligand with water. In the case of the L99A mutant of T4 lysozyme, this step is not necessary because in the apo form, the cavity exists in a fully dehydrated state.⁵⁰

Finally, a symmetry factor correction of $k_B T \ln 2 = 0.42$ kcal/mol is required since the restraints placed on the individual carbon atoms of the benzene ligand restrict its orientation to one out of two possible binding modes. Benzene actually has 12 possible binding modes due to its symmetry; however, the restraints imposed on its carbon atoms are loose enough to allow free rotation of the molecule around the axis normal to its plane; the restraints only prevent flipping of the benzene plane.

Simulation protocol

For all MD simulations, version c28b1 of the CHARMM program,⁵¹ with modifications to allow TI in 4D space,³¹ was

used and all parameters were taken from the CHARMM force field (version 27 for the protein and version 22 for the benzene ligand).⁵² A time step of 2 fs was employed. Bonds containing H atoms were subjected to holonomic constraints using the SHAKE algorithm. The leap-frog integrator with the Langevin algorithm, with a friction constant of 5 ps⁻¹ acting on all nonhydrogen atoms, was used to integrate the equations of motion.

The L99A single-point mutant of T4 lysozyme (162 amino acids; 1290 atoms) with bound benzene ligand was taken from the Protein Databank (accession code 181L) and inserted into a TIP3P water cylinder large enough to completely encompass the protein with a water margin of at least 5 Å around the protein. 40 ps of equilibration of the water itself was performed (with protein and ligand frozen) at a temperature of 3000 K. This was followed by 40 ps of equilibration at room temperature (298 K). A spherical region of interest with a radius of 18 Å centered at the center of mass of the benzene molecule was then defined. All protein atoms outside of this region were frozen. All water molecules outside of a slightly bigger concentric sphere with a radius of 21 Å were deleted and a spherical quartic potential boundary was placed to constrain the remaining water molecules inside this region. A stiff spherical harmonic boundary with a radius of 19 Å, concentric with the other boundaries, was imposed to prevent any protein atom from approaching the water-vacuum interface. Furthermore, a loose harmonic potential with a force constant of 0.2 kcal/mol Å was placed on the C_δ of Arg80 because otherwise the high flexibility of this side chain leads to local unraveling of the protein fold. An additional 60 ps of equilibration was performed. Finally, a 200 ps simulation was run in order to determine the size of the spatial restraints acting on the ligand during the extraction stage (Fig. 1). We first calculated the mean position and spatial distribution of each carbon atom of the benzene ligand. The r_{off} parameter for each restraint was chosen such that the probability of an interaction occurring between the heavy atom and the restraint is less than 0.1%. The force constant k for all atomic restraints imposed on the ligand was 1000 kcal/mol Å.

From the initial equilibrated system, 61 replicas of the system differing only in the w coordinate of the benzene molecule were created. The nominal w positions for the replicas were as follows:

$$0.1, 0.2, \dots, 3.7, 3.8, 4.0, 4.2, 4.5, 4.8, 5.2, 5.6, 6.1, 6.6, 7.2, \\ 7.8, 8.5, 9.2, 10, 11, \dots, 19, 20 \text{ \AA}$$

(a simulation at $w=0$ Å was deemed unnecessary to the PMF calculation since the projected force along the w -axis would be zero at all times). Four different variations of the simulation protocol were performed: (1) Independent replicas with fixed w parameters (static treatment), (2) DR sampling using Monte Carlo moves on a busy cluster, (3) DR sampling using Boltzmann-weighted jumps on a busy cluster, and (4) DR sampling using Boltzmann-weighted jumps on a free cluster. We define a “busy” cluster as one having less available CPUs than there are replicas; otherwise the cluster is available. For each replica in protocols (2)–(4), Monte Carlo moves or

Boltzmann-weighted jumps along w were attempted every 200 MD steps (that is, every 0.4 ps). The DRPE constants were set to $c_1=0.008$ kcal/mol and $c_2=2.0$ kcal/mol. The simulations were run on a shared cluster where the number of available CPUs fluctuated between as little as 10 to over 61. At times when the number of available CPUs was less than the number of replicas, replicas took turns running. A replica running on a particular CPU executed for a maximum of 6 h (or about 24 ps of simulated time) before forfeiting the CPU to a nonrunning replica. Simulation protocol number (4) was run on the same busy cluster. However, provisions were made to emulate a free cluster. This was accomplished by forcing a given running replica to forfeit the CPU to a nonrunning replica immediately after completing its requisite 200 MD steps. This kept all replicas “up to date” as would be the case on a free cluster, where all replicas run simultaneously though not necessarily synchronously.

The hydration free energy calculation was performed using a TIP3P water sphere of 20 Å in radius which contained 1034 water molecules. A spherical quartic potential boundary was used to contain the system. The benzene solute was placed at the center and a harmonic restraint was imposed on the solute to keep it at the center of the water droplet (force constant of 10 kcal/mol Å). No cutoffs for nonbonded interactions were imposed. Independent simulations were performed at the following ligand w -coordinates: 0.1, 0.2, ..., 3.9, 4.0, 4.5, 5.0, ..., 9.5, 10, 11, ..., 19, 20 Å for a total combined sampling time of 0.89 ns. The average force acting on the solute in the fourth dimension was integrated with respect to w to yield the PMF for the extraction of benzene from water. The free energy change for removing the solute completely from $w=0$ to ∞ was estimated by extrapolating the PMF to large w (see Ref. 31). The calculated hydration free energy corresponds to the negative of this result.

RESULTS AND DISCUSSION

Extraction of benzene from T4 lysozyme

The 4D force acting on the benzene ligand as a function of the progression of the simulation, computed using DR sampling and Boltzmann-weighted jumping on a busy cluster (simulation protocol 3), is shown for each discrete w position in Fig. 2. Since the starting structures for all replicas were generated by copying a single system that was equilibrated with benzene in the binding site (see the Simulation Protocol section), it is expected to take some time before the replicas reach equilibrium. Many of the windows exhibit an initial drift in the force over time as the simulation progresses toward equilibrium.

Statistical noise in the force measurements over time is evident in most windows, especially those in regions of w between 0.3 and 4.8 Å. This is typical of 4D PMF calculations due to ruggedness in that region.³¹ Ruggedness is mostly caused by the repulsive part of the Lennard-Jones potential (steric interactions), which is short range and becomes negligible beyond 3.3–3.5 Å (the van der Waals radius of the largest atoms of the solute). Therefore, any random walk (e.g., DR) simulation extending significantly

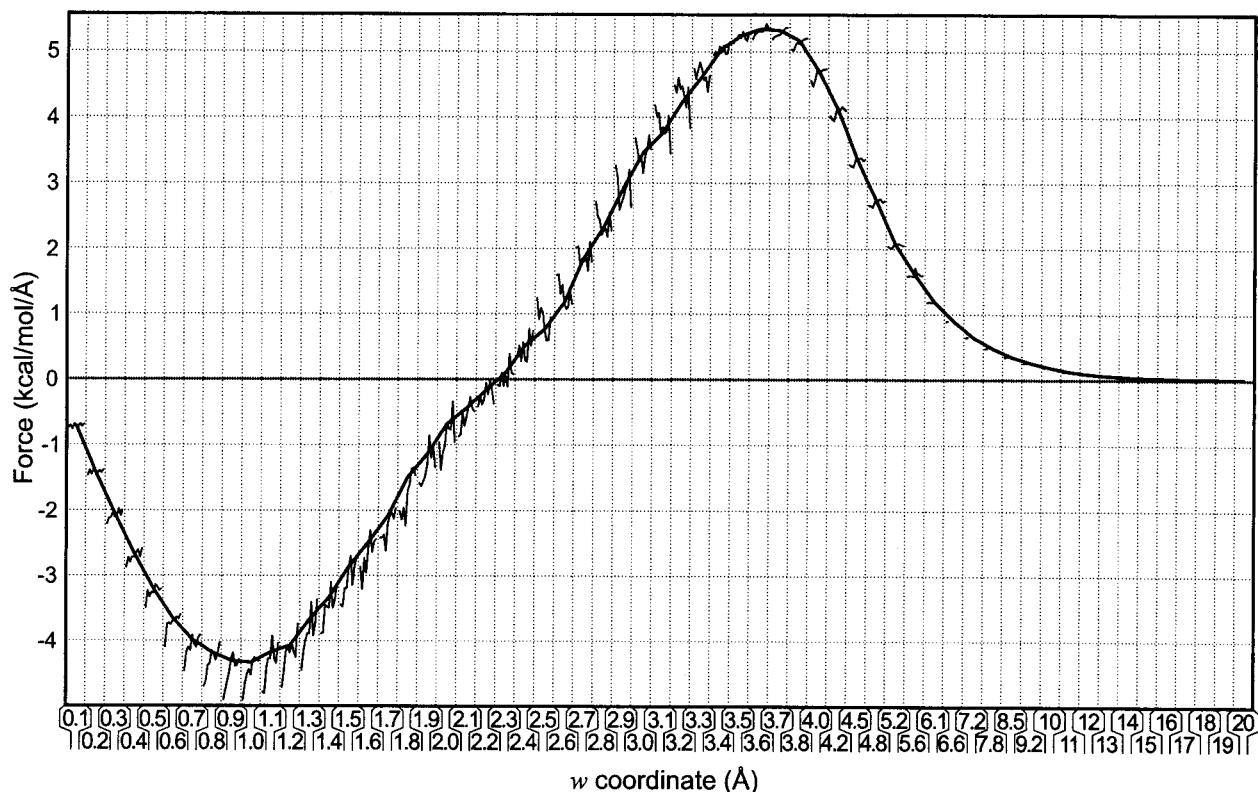


FIG. 2. The 4D force acting on the benzene ligand computed using DR sampling and Boltzmann-weighted jumping on a busy cluster (simulation protocol 3; see the Simulation Protocol section above). Each column shows, for a particular w , the force acting on the ligand as a function of the progression of the simulation. Note that since each replica moves in a random-walk fashion along w , the data for each particular discrete w position represent a composite of all replicas that visited that position. Data are shown for all simulated w positions (note the nonlinear scale). The thick line represents the average sampled force taken from the second half of the data at each w and serves to guide the eye.

beyond that range will help overcome ruggedness. Detailed analysis of the structural transitions occurring in the binding site shows that several amino acid side chains undergo infrequent rotameric transitions that can give rise to the statistical noise observed. In particular, Ile78, Met102, and Leu118, which are located in the binding pocket, closely affect the forces felt by the ligand in the binding site as they change conformation. Satisfactory convergence of the average force in regions beyond 4.8 Å is attained quickly while significantly more sampling is required at smaller w values. Some of the replicas that were sampling regions of w beyond 6.1 Å were stopped early to make CPU resources available to the other replicas. The remainder of the w range was treated dynamically for the entire duration of the simulation to take advantage of the improved sampling efficiency on rugged energy surfaces that DR offers. Regions of w between 0.1 and 4.5 Å were each sampled for approximately 1.6 ns. Regions of w beyond this range were sampled for much less time (see Fig. 3). Together, a total of 73 ns of sampling was performed. The desired and attained sampling profiles, as shown in Fig. 3, match well, with a root mean square deviation (RMSD) between them of 0.069.

Sampling efficiency

Figure 4 shows the mean force computed over the course of the second half of the simulation. Data from the first half of the simulation were not used as it is considered far from equilibrium. Simulation protocols (1)–(4), as outlined in the

Theory and Methods section, are compared. The three protocols using DR sampling give very similar results in terms of quantitative agreement and smoothness of the curve. In contrast, the plot derived from independent replicas appears to be plagued with statistical noise (especially near $w=1$ Å). The noise arises from the fact that some of the replicas get trapped in local energy minima and cannot escape on the time scale of the simulation, resulting in systematic sampling

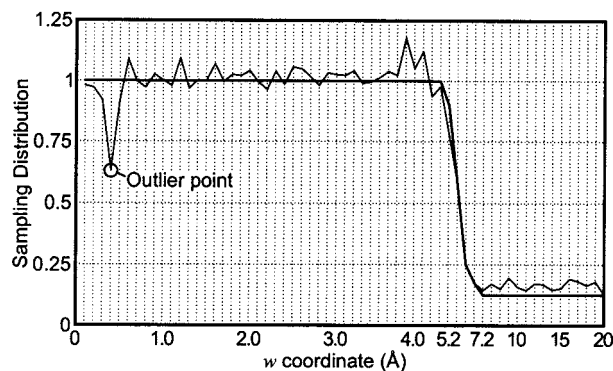


FIG. 3. Sampling distribution over w (note nonlinear scale) resulting from all replicas. The number of counts sampled was normalized such that the desired number of counts is 1 for replica 1. The thick line represents the desired sampling profile. The thin line is the profile attained using DR sampling and Boltzmann-weighted jumping on a busy cluster (simulation protocol 3; see the Simulation Protocol section above). Note that some replicas at large w values were stopped much earlier than others. The outlier point (at $w=0.4$ Å) resulted from data corruption that occurred as a result of a full disk and should not be considered an artifact of DR sampling.

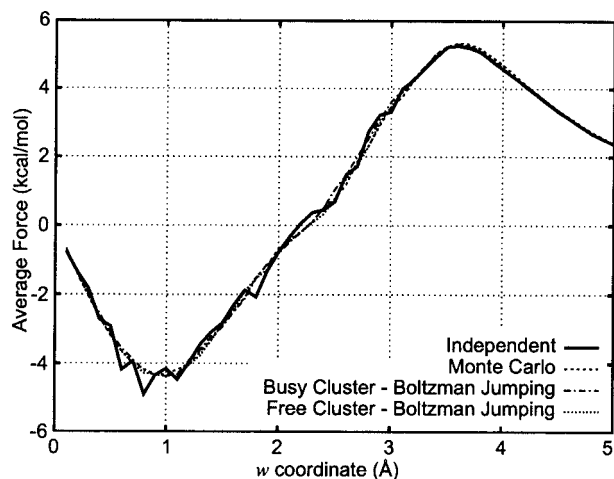


FIG. 4. Average force along the w -axis as a function of w . Results are shown for the four simulation protocols. Only the range from $w=0$ to 5 \AA is shown. Beyond 5 \AA , the four curves are nearly identical and smoothly taper to zero.

errors and lack of convergence of the mean force. This phenomenon was clearly demonstrated by simple test cases in earlier publications.^{32,33} The DR method provides better efficiency in circumventing the barriers that trap the replicas, thus leading to a smoother and statistically converged mean force. A similar qualitative improvement would be obtained with adaptive approaches, inasmuch as they achieve a random walk along the reaction coordinate. The sampling efficiency gain from using DR, or other random-walk approaches, is expected to be even greater in systems with more complex cavities and ligands, where the corresponding energy surface is more rugged.

The random-walk movement of two representative replicas is illustrated in Fig. 5. We measure the mobility in the coupling coordinate w by calculating the average change in the linearized version of the coordinate [i.e., $f^{-1}(w)$; see Eq. (10)] per move/jump attempt. We call this measure the productivity ratio. Note that the productivity ratio is equivalent to the move acceptance probability in the case of Monte Carlo moves. The productivity ratios for DR with Monte

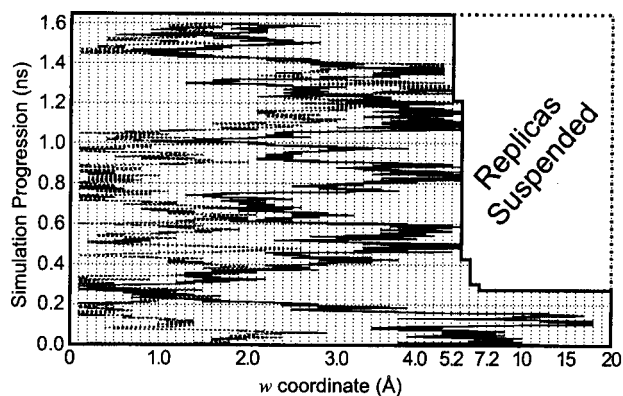


FIG. 5. w coordinates (note the nonlinear scale) of two representative replicas as a function of the progression of the simulation. Random-walk behavior is demonstrated. Note that some replicas at large w values were stopped much earlier than others as this region requires much less sampling. The protocol prevents replicas from entering the region of suspended replicas as shown.

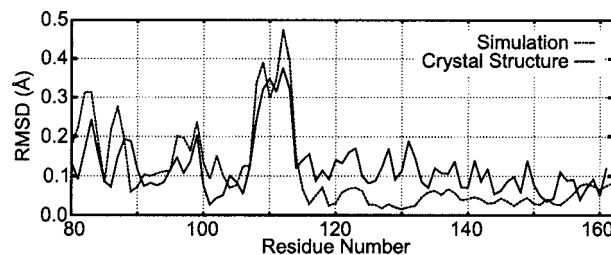


FIG. 6. Main-chain positional shifts in the benzene-bound complex relative to the apoprotein. Alignment was performed on the basis of main-chain atoms. The value plotted for each residue is the RMSD of the shifts in the three backbone atoms. The N -terminal domains (residues 1–79) were not used in the alignment and are not included in the figure. Plots derived from the crystal structures (Ref. 43) (solid) as well as the simulation results (dashed) are shown. Alignment and distance calculations were performed using the PROFIT program (Ref. 59).

Carlo moves and DR with Boltzmann-weighted jumping were 0.63 and 0.85, respectively, showing the improvement of w mobility in the latter technique. These 4D simulations benefit from Boltzmann-weighted jumping because the space (in the fourth dimension) that the ligand moves through is mostly empty and even a large jump has a good chance of being accepted. The productivity ratio cannot be greater than 1 for Monte Carlo moves, but Boltzmann-weighted jumps have no such limit and can provide a substantial benefit to some applications. When applied to the simple test case described in the original DR sampling publication,³² Monte Carlo moves and Boltzmann-weighted jumping yield productivity ratios of 0.65 and 3.35, respectively.

Parameter mobility can also be modulated by the frequency of attempted moves. There are no restrictions on the intervals between replica move attempts, although some optimal interval will exist for a given application. Frequent move attempts allow greater mobility of λ_i values, but at the cost of increased overhead (network communication and calculation of energies). Very frequent moves are unlikely to lead to improved sampling as they do not give the orthogonal degrees of freedom a chance to relax. In the case of Monte Carlo moves, the distance by which λ_i changes at one time is not restricted and can be optimized for the application at hand. In Boltzmann-weighted jumping, λ space must be discretized *a priori* and only those discrete values of λ can be visited. There is no limit on how fine the discretization can be (there can be many more discrete λ values than there are replicas), although overhead is associated with the calculation of Δ_j [Eq. (5)] for each discrete state j during each λ jump attempt.

Although in our simulations the extraction of the ligand is performed through an unphysical fourth spatial dimension, the endpoints, which represent the bound and unbound states, are physically meaningful. We computed averages of atomic coordinates from snapshots taken from the simulation. The positional shifts of main-chain atoms that occur when the protein transitions from the bound state (approximated by $w=0.1 \text{ \AA}$) to the apo form (approximated by $w=20 \text{ \AA}$) were computed to yield a plot that can be compared directly to the equivalent analysis of the crystal structures (Fig. 6). The largest shifts occur in residues of the binding pocket (indices in the area of 110; see Ref. 43 for more

details). The efficiency of a random-walk approach helps to achieve thorough sampling of conformational space. In T4 lysozyme, simulations on the order of 1 ns per window, coupled with an efficient sampling scheme such as DR, are capable of reproducing the largest conformational changes that occur in the protein upon ligand binding.

Estimating statistical convergence

One of the difficult questions to answer when calculating thermodynamic averages from molecular simulations is as follows: When has statistical convergence been reached? Or, has the system reached equilibrium yet and can accumulation now start? Free energy simulations are often divided into two parts: The equilibration part and the production (or sampling) part. A rigorous method for calculating where this division should be has been proposed.⁵³ However, this approach is artificial, as there is no point in time where we can say that equilibration has suddenly ended. Equilibrium is approached as the simulation progresses and may not be reached in practical simulation times. Furthermore, data generated by the equilibration part of the simulation are often completely discarded.⁵³ This is unfortunate as complex systems typically require long equilibration times. Consider the current test system where a ligand is extracted from the binding site of a protein receptor. All 61 replicas spanning the full range from bound to unbound states were generated from the benzene-bound crystal structure of T4 lysozyme. The conformational changes that the protein undergoes as the ligand moves from the bound to the unbound state are not reflected in these initial replicas (the protein conformations all resemble the bound structure). Equilibration must occur before the individual replicas can relax to reflect the extent of ligand binding. In the following analysis, we make no presumptions about when equilibrium has been reached. Instead, we assume that equilibrium is not reached on the simulated time scale and we predict the binding free energy via extrapolation to an infinitely long simulation. In this process, we make use of all generated data.

Each force sample taken from any replica is stamped with both a time and a w -coordinate. A time t_0 is chosen as the division point between the equilibration and production phase. The equilibration data are temporarily discarded while the production data are consolidated into chronological order. A window of data containing all points with a time stamp ranging from t_0 to $t_0 + 144$ ps is then extracted. Here, the window size (144 ps) was arbitrarily chosen as 1/10th of the total simulation time per window. From these data, the force samples are sorted by their w -coordinate into 61 separate groups. Force samples within each group are averaged, resulting in 61 data points that can be used to compute the mean force versus w . Integrating the mean force with respect to w yields the work or free energy change for extracting the ligand from its binding site. In Fig. 7, we plot the ligand extraction free energy versus t_0 . In effect, the independent variable t_0 controls a sliding window of production data from which the free energy change is computed. The free energy of extraction computed from a running average in this way, eventually at a large t_0 , should fluctuate around a stable value

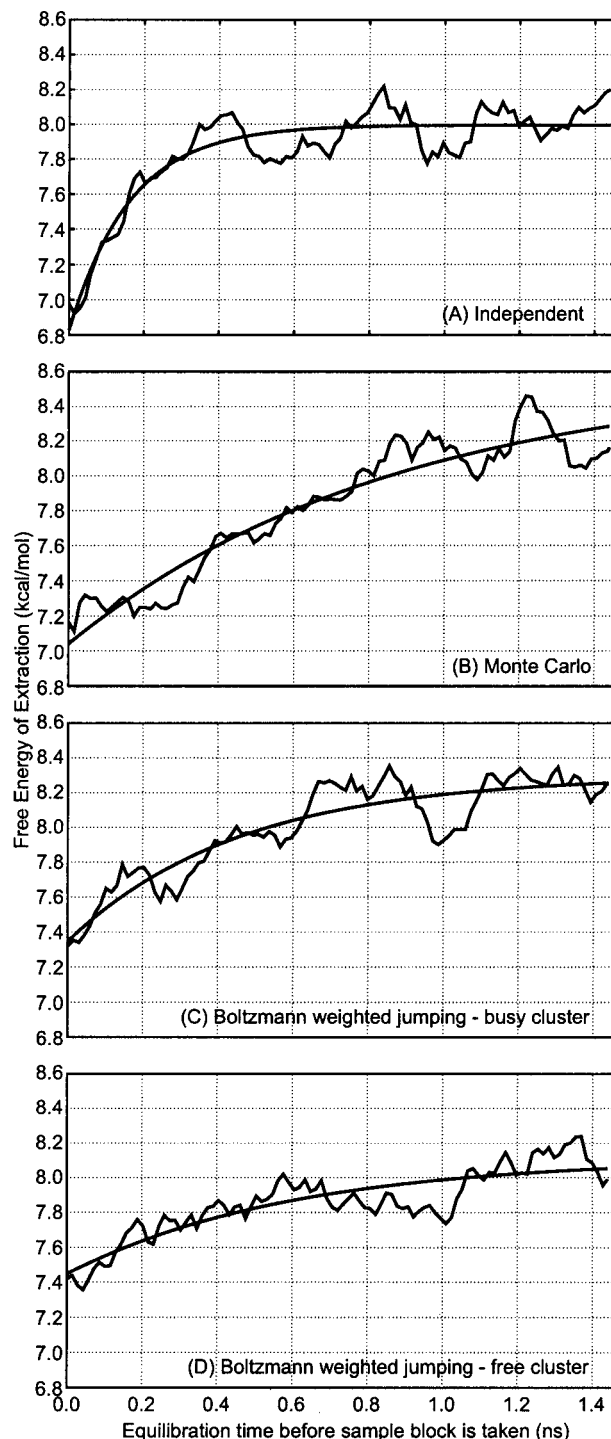


FIG. 7. Free energy of extraction of benzene from T4 lysozyme calculated from a block of sample data (0.144 ns in duration for each replica) vs the amount of equilibration time that had elapsed before that block was taken. Results are shown for the four simulation protocols: (a) Independent simulations, DR sampling with (b) Monte Carlo moves, (c) Boltzmann-weighted jumping on a busy cluster, and (d) Boltzmann-weighted jumping on a free cluster. Best fit curves of the form $\Delta G = a \exp(-bt_0) + c$ were applied (see Table II for parameters a , b , and c).

(or asymptote). The corresponding plot is expected to follow a decay profile as t_0 is increased. For simplicity, here we assume that the process by which each replica tends from the initial configuration (which resembles the bound structure) toward a configuration more reflective of the replica's w coordinate is governed by a single dominating barrier and

TABLE II. Summary of the free energy data for the calculation of absolute binding free energy of benzene to T4 lysozyme using simulation protocols (1)–(4)—see the Simulation Protocol section. Units are in kcal/mol unless otherwise indicated.

Protocol used	(1) Independent	(2) Monte Carlo	(3) Boltzmann (busy)	(4) Boltzmann (available)
Fit parameters ^a <i>a</i>	-2.84 (11%)	-1.82 (4.6%)	-1.30 (5.9%)	-0.846 (6.0%)
<i>b</i> (per picosecond)	6.15×10^{-3} (8.8%)	1.13×10^{-3} (16%)	2.18×10^{-3} (13%)	1.65×10^{-3} (21%)
<i>c</i>	8.00 (0.18%)	8.59 (1.4%)	8.30 (0.45%)	8.12 (0.65%)
Time to 95% equilibrium (ns)	0.487	2.65	1.37	1.82
Extraction free energy (step 1)	8.00	8.59	8.30	8.12
Standard state correction (step 2)			-0.98	
Hydration free energy (step 3)			1.09	
Symmetry factor correction→ $RT \ln(2)$			0.42	
Total free energy of unbinding	8.53	9.12	8.83	8.65
Calculated binding free energy	-8.53	-9.12	-8.83	-8.65
Experimental hydration free energy ^b			-0.89	
Experimental binding free energy ^c			-5.19	

^aPercentage errors are given.

^bReference 58.

^cReference 42.

therefore follows two-state kinetics. In this simplified view, any property measured during the equilibration process is expected to follow an exponential decay process. In practice, we find that an exponential decay function of the form $\Delta G = a \exp(-bt_0) + c$ fits the data very well. A Levenberg–Marquardt solver was used to calculate the constants *a*, *b*, and *c*. $\Delta G = c$ gives the extrapolated free energy at infinite time. The fit parameters, together with the time required to reach 95% of the way to equilibrium, are listed in Table II.

In Fig. 7, we compare the plots from the four simulation protocols. Consistent with a lower productivity ratio (0.63), Monte Carlo moves lead to the slowest convergence toward the steady state (it takes 2.65 ns to get 95% of the way there). The two Boltzmann-weighted jumping protocols achieve higher productivity ratios (0.84), although the jump magnitude is limited by the fact that the ligand exists in a dense medium at small *w* and because large excursions of a replica in one step are prohibited by the DRPE. With higher productivity ratios, the Boltzmann-weighted jumping schemes achieve faster progression toward the steady state (1.37 and 1.82 ns to complete 95% of the decay process) as expected. The rate at which convergence is reached is consistent with another study of T4 lysozyme,³⁴ in which the authors found that with a static coordinate, about 5 ns of simulation time per window was required to achieve adequate convergence. However, with a clever scheme to decompose space and thus eliminate the need to cross difficult barriers (a similar aim to that of DR sampling), convergence was achieved in about 1 ns per window. The simulation with independent replicas appears to reach a steady state in the shortest amount of time (0.487 ns to achieve 95% of the way to steady state). However, this is unlikely to be the true equilibrium but more probably represents a trapped state where several replicas are caught in local energy minima (as discussed above; see Fig. 4).

Local conformational transitions occurring on longer

time scales than those accessed in the reported simulations are possible. Evidence from a run with 19 ns per window (data not shown) suggests that rotameric isomerizations of several amino acid side chains (including those of Ile 78, Met 102, and Leu 118, as noted above) occur infrequently, over time scales in the tens of nanoseconds. Such time scales might therefore be required to improve statistical convergence of the free energy result. Generally, the length of a simulation required to sample all relevant configurations cannot currently be predicted.

Influence of cluster availability

Each time a different replica is allocated to a particular CPU, overhead is associated in transferring restart data and simulation parameters through the network. In the ideal case when there are as many available CPUs as there are replicas, all replicas can run at once and no reallocation is necessary. However, when fewer CPUs are available, replicas must take turns running. In the limit where a replica change is allowed each time a replica move is considered, the simulation is energetically equivalent to one with enough CPUs for each replica, albeit with more network overhead. This is the procedure that was used to emulate a free cluster (simulation protocol 4). When replicas are allowed to retain the CPU for a longer time, other nonrunning replicas will not be given a chance to move. In extreme cases, this will hamper the mobility of the transformation parameter and lead to slow convergence. Furthermore, significant deviations from the desired sampling profile may be induced. This becomes more problematic the busier the cluster is. Despite being run on a busy cluster, the simulation using protocol 3 achieves a mean force profile showing a very good agreement with that of simulation protocol 4 (RMSD=0.10; see Fig. 4) and ap-

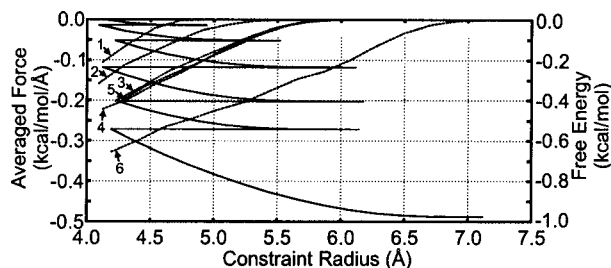


FIG. 8. Force acting inward by the spherical restraints as a function of restraint radius (r_{off}) and the free energy change associated with the expansion process. The individual free energies for expanding the restraints are stacked on top of each other. The free energy sum is -0.98 kcal/mol.

proaches statistical convergence of the steady state at about the same rate (Table II). This shows that the method is fairly robust and insensitive to CPU availability.

Calculation of the absolute binding free energy

Figure 8 depicts the free energy associated with releasing the six restraints originally placed on the carbon atoms of the benzene ligand. The free energy for the removal of the restraints and thus the correction to reach the standard state representing a $1M$ ligand concentration was determined to be -0.98 kcal/mol. The hydration free energy of benzene was calculated to be 1.09 kcal/mol. Finally, a symmetry factor correction for benzene of 0.42 kcal/mol was applied. The overall free energy results for the absolute binding free energy calculations are summarized in Table II.

Although a detailed investigation of the sources of systematic errors is beyond the scope of this paper, we include the experimental result for the hydration free energy of benzene and the binding free energy of benzene to T4 lysozyme for comparison. The computed hydration free energy of benzene (1.09 kcal/mol) is an overestimate of the experimental value (-0.89 kcal/mol) by 1.97 kcal/mol. Since hydration free energy calculations can be performed routinely using current methods to very high statistical precision,^{30,54,55} most of the error (in the hydration and thus presumably also in the extraction free energy) is attributed to systematic biases introduced by the approximations in the force field and by limitations in the setup of the system (finite size effects, frozen atoms, and long-range nonbonded cutoffs). The computed binding free energies of benzene (ranging from -8.53 to -9.12 kcal/mol) indicate that the ligand binds tighter to the protein than measured by experiment (-5.19 kcal/mol) by 3.3 – 3.9 kcal/mol. Results from similar calculations performed in our laboratory involving other ligands of T4 lysozyme (data not shown) demonstrate very similar systematic errors. The results suggest that better agreement with experiment may be achieved via an improvement in the simulation setup, including better treatment of long-range interactions, dynamic treatment of the entire protein, solvation in a larger droplet or periodic box, and optimized force field parameters. Other computational results for the binding free energy of benzene to T4 lysozyme have recently been reported in the literature, for example, -5.96 kcal/mol (Ref. 56) and -5.14 kcal/mol.⁵⁷ Nevertheless, the results attained in this work demonstrate that meth-

odologies for efficient Boltzmann sampling of conformational space improve the rate at which statistically precise free energy results can be attained in protein-ligand systems.

CONCLUSIONS

Through the calculation of the absolute binding free energy of benzene to T4 lysozyme, we have demonstrated the usefulness of DR sampling combined with a 4D coordinate in simulations of biological scale and complexity.

The 4D separation distance between ligand and protein is a suitable transformation coordinate that minimizes energy barriers and allows rapid removal and insertion of the ligand into the binding site. In addition, an improvement in efficiency is also achieved as all nonbonded interactions are decoupled at once, without resorting to the separation of Lennard-Jones and Coulombic interactions as is commonly done.⁴⁵ Given a set amount of CPU cycles, this approach minimizes the statistical error in the result because the degrees of freedom that do not affect the 4D force are not sampled redundantly.

DR sampling allows conformational space to be searched more effectively (as compared to standard TI), improving the rate of convergence of free energy calculations while making the most of available computational resources. Dynamic treatment provides new paths for energy barriers (such as rotation of a ligand in a binding site) to be circumvented. In addition, when a replica sampling one region of the transformation coordinate discovers a new conformation, the effect of this new conformation is propagated to other locations along the coordinate as the replica undergoes a random walk. We have built on the original DR sampling algorithm by introducing Boltzmann-weighted jumping. Boltzmann-weighted jumping improves the mobility of the transformation coordinate by 33% in the present system and leads to faster convergence of the PMF. Boltzmann-weighted jumps cannot be readily applied in RE algorithms.

DR sampling is a novel approach designed for shared clusters or large-scale distributed computing platforms. Full utilization of available CPU resources is realized automatically. In contrast, efficient implementation of RE type algorithms on such platforms is very difficult to realize as CPUs would sit idle waiting for other replicas to finish a simulation segment so that swap events can occur. The simulation protocol presented here achieves convergence of the binding free energy calculation in practical wall clock time and on computer hardware accessible to most research groups. Finally, the method is shown to be insensitive to fluctuation in CPU availability. This is especially important when DR is implemented on a distributed or shared computing system where CPU availability changes unpredictably.

ACKNOWLEDGMENTS

We gratefully acknowledge the Canadian Institutes of Health Research (Grant No. MOP43998) for support. T.R. was funded in part by a NSERC Canada Graduate Scholarship. P.L.H. and R.P. are CRCP chairholders.

- ¹T. Rodinger and R. Pomès, *Curr. Opin. Struct. Biol.* **15**, 164 (2005).
- ²I. Andricioaei, D. Asthagiri, T. L. Beck, C. Chipot, E. Darve, C. Dellago, G. Hummer, N. Lu, A. E. Mark, A. Z. Panagiotopoulos, V. S. Pande, A. Pohorille, L. R. Pratt, M. S. Shell, T. Simonson, and T. B. Woolf, in *Free Energy Calculations: Theory and Applications in Chemistry and Biology*, edited by C. Chipot and A. Pohorille (Springer-Verlag, Berlin, 2007).
- ³C. Chipot and D. A. Pearlman, *Mol. Simul.* **28**, 1 (2002).
- ⁴W. F. van Gunsteren, X. Daura, and A. E. Mark, *Helv. Chim. Acta* **85**, 3113 (2002).
- ⁵T. Simonson, G. Archontis, and M. Karplus, *Acc. Chem. Res.* **35**, 430 (2002).
- ⁶L. Kale, R. Skeel, M. Bhandarkar, R. Brunner, A. Gursoy, N. Krawetz, J. Phillips, A. Shinozaki, K. Varadarajan, and K. Schulten, *J. Comput. Phys.* **151**, 283 (1999).
- ⁷K. Y. Sanbonmatsu, S. Joseph, and C. S. Tung, *Proc. Natl. Acad. Sci. U.S.A.* **102**, 15854 (2005).
- ⁸P. L. Freddolino, A. S. Arkhipov, S. B. Larson, A. McPherson, and K. Schulten, *Structure (London)* **14**, 437 (2006).
- ⁹E. J. Sorin, Y. M. Rhee, M. R. Shirts, and V. S. Pande, *J. Mol. Biol.* **356**, 248 (2006).
- ¹⁰B. A. Berg and T. Neuhaus, *Phys. Lett. B* **267**, 249 (1991).
- ¹¹A. P. Lyubartsev, A. A. Martsinovski, S. V. Shevkunov, and P. N. Vorontsovvellyaminov, *J. Chem. Phys.* **96**, 1776 (1992).
- ¹²E. Marinari and G. Parisi, *Europhys. Lett.* **19**, 451 (1992).
- ¹³K. Hukushima and K. Nemoto, *J. Phys. Soc. Jpn.* **65**, 1604 (1996).
- ¹⁴Y. Sugita and Y. Okamoto, *Chem. Phys. Lett.* **314**, 141 (1999).
- ¹⁵M. C. Tesi, E. J. J. van Rensburg, E. Orlandini, and S. G. Whittington, *J. Stat. Phys.* **82**, 155 (1996).
- ¹⁶A. Mitsutake and Y. Okamoto, *J. Chem. Phys.* **121**, 2491 (2004).
- ¹⁷A. Mitsutake and Y. Okamoto, *Chem. Phys. Lett.* **332**, 131 (2000).
- ¹⁸A. Mitsutake, Y. Sugita, and Y. Okamoto, *J. Chem. Phys.* **118**, 6676 (2003).
- ¹⁹A. Mitsutake, Y. Sugita, and Y. Okamoto, *J. Chem. Phys.* **118**, 6664 (2003).
- ²⁰Y. Sugita, A. Kitao, and Y. Okamoto, *J. Chem. Phys.* **113**, 6042 (2000).
- ²¹Y. Sugita and Y. Okamoto, *Chem. Phys. Lett.* **329**, 261 (2000).
- ²²A. E. García and J. N. Onuchic, *Proc. Natl. Acad. Sci. U.S.A.* **100**, 13898 (2003).
- ²³H. Nymeyer and A. E. García, *Proc. Natl. Acad. Sci. U.S.A.* **100**, 13934 (2003).
- ²⁴D. Paschek and A. E. García, *Phys. Rev. Lett.* **93**, 238105 (2004).
- ²⁵A. E. García and K. Y. Sanbonmatsu, *Proteins* **42**, 345 (2001).
- ²⁶S. Gnanakaran, H. Nymeyer, J. Portman, K.-Y. Sanbonmatsu, and A. E. García, *Curr. Opin. Struct. Biol.* **13**, 168 (2003).
- ²⁷R. W. Zwanzig, *J. Chem. Phys.* **22**, 1420 (1954).
- ²⁸J. G. Kirkwood, *J. Chem. Phys.* **3**, 300 (1935).
- ²⁹G. M. Torrie and J. P. Valleau, *Chem. Phys. Lett.* **28**, 578 (1974).
- ³⁰M. R. Shirts, J. W. Pitera, W. C. Swope, and V. S. Pande, *J. Chem. Phys.* **119**, 5740 (2003).
- ³¹T. Rodinger, P. L. Howell, and R. Pomès, *J. Chem. Phys.* **123**, 034104 (2005).
- ³²T. Rodinger, P. L. Howell, and R. Pomès, *J. Chem. Theory Comput.* **2**, 725 (2006).
- ³³C. Neale, T. Rodinger, and R. Pomès, *Chem. Phys. Lett.* **460**, 375 (2008).
- ³⁴D. L. Mobley, J. D. Chodera, and K. A. Dill, *J. Chem. Phys.* **125**, 084902 (2006).
- ³⁵G. Jayachandran, M. R. Shirts, S. Park, and V. S. Pande, *J. Chem. Phys.* **125**, 084901 (2006).
- ³⁶E. Darve and A. Pohorille, *J. Chem. Phys.* **115**, 9169 (2001).
- ³⁷E. Darve, M. A. Wilson, and A. Pohorille, *Mol. Simul.* **28**, 113 (2002).
- ³⁸J. Hénin and C. Chipot, *J. Chem. Phys.* **121**, 2904 (2004).
- ³⁹M. Mezei, *J. Comput. Phys.* **68**, 237 (1987).
- ⁴⁰C. Bartels and M. Karplus, *J. Comput. Chem.* **18**, 1450 (1997).
- ⁴¹C. J. Woods, J. W. Essex, and M. A. King, *J. Phys. Chem. B* **107**, 13703 (2003).
- ⁴²A. Morton, W. A. Baase, and B. W. Matthews, *Biochemistry* **34**, 8564 (1995).
- ⁴³A. Morton and B. W. Matthews, *Biochemistry* **34**, 8576 (1995).
- ⁴⁴R. Pomès, E. Eisenmesser, C. B. Post, and B. Roux, *J. Chem. Phys.* **111**, 3387 (1999).
- ⁴⁵B. O. Brandsdal and A. O. Smalås, *Protein Eng.* **13**, 239 (2000).
- ⁴⁶D. Trzesniak, A. P. E. Kunz, and W. F. van Gunsteren, *ChemPhysChem* **8**, 162 (2007).
- ⁴⁷M. K. Gilson, J. A. Given, B. L. Bush, and J. A. McCammon, *Biophys. J.* **72**, 1047 (1997).
- ⁴⁸J. A. McCammon, *Curr. Opin. Struct. Biol.* **8**, 245 (1998).
- ⁴⁹S. Boresch, F. Tettinger, M. Leitgeb, and M. Karplus, *J. Phys. Chem. B* **107**, 9535 (2003).
- ⁵⁰M. D. Collins, G. Hummer, M. L. Quillin, B. W. Matthews, and S. M. Gruner, *Proc. Natl. Acad. Sci. U.S.A.* **102**, 16668 (2005).
- ⁵¹B. R. Brooks, R. E. Bruccoleri, B. D. Olafson, D. J. States, S. Swaminathan, and M. Karplus, *J. Comput. Chem.* **4**, 187 (1983).
- ⁵²A. D. MacKerell, Jr., D. Bashford, M. Bellott, R. L. Dunbrack, Jr., J. D. Evanseck, M. J. Field, S. Fischer, J. Gao, H. Guo, S. Ha, D. Joseph-McCarthy, L. Kuchnir, K. Kuczera, F. T. K. Lau, C. Mattos, S. Michnick, T. Ngo, D. T. Nguyen, B. Prodhom, W. E. Reiher III, B. Roux, M. Schlenkrich, J. C. Smith, R. Stote, J. Straub, M. Watanabe, J. Wiórkiewicz-Kuczera, D. Yin, and M. Karplus, *J. Phys. Chem. B* **102**, 3586 (1998).
- ⁵³W. Yang, R. Bitetti-Putzer, and M. Karplus, *J. Chem. Phys.* **120**, 2618 (2004).
- ⁵⁴Y. Q. Deng and B. Roux, *J. Phys. Chem. B* **108**, 16567 (2004).
- ⁵⁵J. L. Maccallum and D. P. Tieleman, *J. Comput. Chem.* **24**, 1930 (2003).
- ⁵⁶Y. Q. Deng and B. Roux, *J. Chem. Theory Comput.* **2**, 1255 (2006).
- ⁵⁷J. Hermans and L. Wang, *J. Am. Chem. Soc.* **119**, 2707 (1997).
- ⁵⁸J. Hine and P. K. Mookerjee, *J. Org. Chem.* **40**, 292 (1975).
- ⁵⁹G. D. Smith, PROFIT, a locally written program for orienting one protein molecule onto another by a least-squares procedure, Hauptman-Woodward Medical Research Institute, Buffalo, NY, 1993.