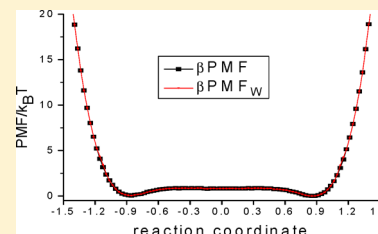


Variance of a Potential of Mean Force Obtained Using the Weighted Histogram Analysis Method

Robert I. Cukier*

Department of Chemistry Michigan State University, East Lansing, Michigan 48824-1322, United States

ABSTRACT: A potential of mean force (PMF) that provides the free energy of a thermally driven system along some chosen reaction coordinate (RC) is a useful descriptor of systems characterized by complex, high dimensional potential energy surfaces. Umbrella sampling window simulations use potential energy restraints to provide more uniform sampling along a RC so that potential energy barriers that would otherwise make equilibrium sampling computationally difficult can be overcome. Combining the results from the different biased window trajectories can be accomplished using the Weighted Histogram Analysis Method (WHAM). Here, we provide an analysis of the variance of a PMF along the reaction coordinate. We assume that the potential restraints used for each window lead to Gaussian distributions for the window reaction coordinate densities and that the data sampling in each window is from an equilibrium ensemble sampled so that successive points are statistically independent. Also, we assume that neighbor window densities overlap, as required in WHAM, and that further-than-neighbor window density overlap is negligible. Then, an analytic expression for the variance of the PMF along the reaction coordinate at a desired level of spatial resolution can be generated. The variance separates into a sum over all windows with two kinds of contributions: One from the variance of the biased window density normalized by the total biased window density and the other from the variance of the local (for each window's coordinate range) PMF. Based on the desired spatial resolution of the PMF, the former variance can be minimized relative to that from the latter. The method is applied to a model system that has features of a complex energy landscape evocative of a protein with two conformational states separated by a free energy barrier along a collective reaction coordinate. The variance can be constructed from data that is already available from the WHAM PMF construction.



1. INTRODUCTION

The potential of mean force (PMF(R)),¹ in contrast to the mechanical potential that underlies it, provides a free energy along some reaction coordinate, R . From the perspective of simulation, it can be obtained from trajectories that are often generated by Molecular dynamics (MD) or Monte Carlo (MC) methods. Typically, MD and MC suffer from the limitation that rough energy landscapes are difficult to sample in their high energy regions, implying the need for what may be unrealistically long simulations. One resolution to this practical difficulty is the use of umbrella-sampling-based window methods^{2,3} that introduce restraints in a reaction coordinate (RC) to improve the sampling in otherwise difficult to sample regions of configuration space. When multiple windows with differing temperatures and/or RC restraint positions are used, the pool of data from all the restrained simulations can be combined in a statistically grounded method known as the Weighted Histogram Analysis Method (WHAM).^{4–7} Variants of umbrella sampling such as Replica Exchange Methods^{8–10} and simulated tempering^{11,12} can also, with suitable modification,¹³ be analyzed using the WHAM and can be applied to reaction coordinate simulations.^{14,15} The generation of a PMF from finite trajectory data does leave the PMF subject to statistical fluctuations in its determination. The purpose of this work is to analyze the sampling error of the WHAM. The approach is designed to obtain insight into the sources of the sampling error as a function of the reaction coordinate.

A number of approaches to the WHAM sampling error have been developed. The generic bootstrap method¹⁶ or estimates based on subsets of data¹⁷ of course can be applied to this sampling data, but no use is made of the specifics of the umbrella-based method and thus no insight relevant to the method is available. For the related but not identical subject of free energy calculations, where the free energy difference ΔF between two thermodynamic states is evaluated, Bennett's method¹⁸ provides an optimal solution for ΔF as well as the associated sampling error under certain assumptions on the statistical properties of the data, as shown by a number of methods.¹⁹ For a PMF, the original focus was on the count fluctuations in the histograms of the window data.^{5,7,13} These histograms are the sampling approximants of the biased (by the applied restraints) window densities $\rho_w^{(b)}(R)$ that would be obtained in infinitely long trajectories. Kastner and Thiel^{20,21} developed an approach based on the mean force ($-\partial \text{PMF}(R)/\partial R$) and obtained error estimates by integrating their error expression for the mean force. They assume that the biased window densities are Gaussian in form and that permits insight into the sources of the sampling error and suggestions for picking window parameters to reduce the error. Berau and Swendson²² analyze the error in free energy differences obtained by WHAM by adding up the contributions from

Received: August 8, 2013

Revised: October 22, 2013

Published: October 31, 2013

errors in the window free energies and from (assumed Gaussian) window density overlap. By minimizing the overall error as a function of the separation relative to the widths of the windows they obtain an optimum separation for the windows. Zhu and Hummer²³ also proceed by working with the mean force and then integrating it to consider the propagation of error through the windows. They emphasize that the errors in the window free energies are strong contributors to the sampling error in addition to those arising from the window count fluctuations. They obtain an expression for the sampling error in the strong force limit, where there is a simple, direct connection between the PMF evaluated at the window restraint positions and the window free energies. A method based on a Bayesian statistics, maximum likelihood approach has been developed that also stresses the sampling error dependence on window free energies.²⁴

In the current approach, the focus will be directly on the potential of mean force. Our basic observation is that for typical umbrella simulations, biased window densities $\rho_w^{(b)}(R)$ that are first neighbors overlap, a requirement of the WHAM, but second and beyond neighbor windows do not overlap to any significant degree. Thus, an approach that is “local” whereby the variance of the overall PMF(R) can be obtained from each window’s $\rho_w^{(b)}(R)$ data in ranges centered on their restrained positions will be useful. As in the work of Kastner and Thiel,^{20,21} we shall assume that the biased window densities are Gaussian in form. Furthermore, we shall assume that the sampling in each window produces a trajectory corresponding to sampling from an independent, identically distributed (i.i.d.) distribution (no correlation of trajectory points). Then, analytic expressions for the variance of the PMF can be found directly.

The PMF variance analysis will be explored based on a simple model consisting of a set of “particles”, each of which experience a double well potential along with a coupling between neighbor particles.^{25,26} This simple model has the flavor of more complex systems, such as proteins, which can be viewed as a set of dihedral coordinates experiencing double well potentials coupled together via excluded volume and electrostatic interactions. The simple model has a collective coordinate, R , that is the average particle position and is suitable as a reaction coordinate. For high temperatures, the PMF has a single well character. As the temperature is lowered, a double well PMF is obtained, as would be the case for a protein with two stable conformational states, parametrized by a collective coordinate, separated by a free energy barrier that arises from the underlying rough potential energy landscape.

The remainder of this paper is organized as follows: Section 2 develops the analytic framework that results in an expression for the reaction coordinate dependence of the variance of a potential of mean force. Section 3 presents the details of the model along with the Langevin equation integrator used to generate trajectory data. Section 4 evaluates the variance expression for the model data and dissects its various components to gain insight into the various sources of the variance. Section 5 summarizes our approximations, and how the model parameters influence the variance. Appendix A examines the generated data to make sure sample points are separated sufficiently to be independent, identically distributed samples.

2. THEORY

2.1. Weighted Histogram Analysis Method. Window simulations bias trajectories, run at temperature $T = 1/(k_B\beta)$,

with a restraint potential $V_w(g(\mathbf{r}))$ to promote sampling of a reaction coordinate $R = g(\mathbf{r})$, with $g(\mathbf{r})$, a function of the configurational coordinates, \mathbf{r} . The WHAM^{5,7} combines the biased window densities $\rho_w^{(b)}(R)$ ($w = 1, 2, \dots, W$) as

$$\rho(R) = \sum_{w=1}^W c_w(R) \rho_w^{(b)}(R) = \sum_{w=1}^W p_w(R) \rho_w^{(u)}(R) \quad (1)$$

with coefficients

$$c_w(R) \equiv 1 / \sum_{w'}^W e^{-\beta(V_w(R) - f_{w'})} \quad (2)$$

that couple the different window simulations. In terms of the unbiased densities, $\rho_w^{(u)}(R)$, with

$$\rho_w^{(u)}(R) = e^{\beta(V_w(R) - f_w)} \rho_w^{(b)}(R) \quad (3)$$

the probabilities $p_w(R)$ are given by

$$p_w(R) = c_w(R) e^{-\beta(V_w(R) - f_w)} \quad (4)$$

The normalization condition for these probabilities is $\sum_{w=1}^W p_w(R) = 1$; they are dimensionless probabilities for fixed values of R — (not R -densities). They can be written as

$$\begin{aligned} p_w(R) &= e^{-\beta(V_w(R) - f_w)} \rho_w^{(u)}(R) / \sum_{w'}^W e^{-\beta(V_{w'}(R) - f_{w'})} \rho_{w'}^{(u)}(R) \\ &= \rho_w^{(b)}(R) / \sum_{w'}^W \rho_{w'}^{(b)}(R) \end{aligned} \quad (5)$$

The WHAM consists of two steps: (1) The determination of the free energies f_w for each window required to unbias the probability densities and (2) construction of PMFs for reaction coordinates based on the unbiased density. The determination of the free energies is obtained from iterative solution of

$$e^{-\beta f_w} = \sum_{w'} \sum_s \left[\frac{e^{-\beta V_{w'}(g(\mathbf{r}_s^{w'}))}}{\sum_k e^{-\beta [V_k(g(\mathbf{r}_s^{w'})) - f_k]}} \right] \quad (6)$$

where $g(\mathbf{r}_s^w) = R_s^w$ denotes the reaction coordinate value for the window w trajectory snapshot s . Using an explicit expression in terms of trajectory snapshots for step (1) is, in principle, more accurate than using a histogram method to determine the f_w . Once the f_w are available, the unbiased WHAM probability density is obtained from

$$\rho(R) = \sum_w \sum_s c_{w,s} \eta_s^w(R) \quad (7)$$

where

$$c_{w,s} = 1 / \sum_{w'} e^{-\beta [V_{w'}(g(\mathbf{r}_s^{w'})) - f_{w'}]} \quad (8)$$

and

$$\begin{aligned} \eta_s^w(R) &= 1 \quad R < R_s^w < R + dR \\ &= 0 \quad \text{otherwise} \end{aligned} \quad (9)$$

bins the reaction coordinate values. The two-step procedure decouples the determination of the free energies that can be done without histograms from evaluation of the PMF that is obtained via histograms.

2.2. WHAM Variance. To construct a variance expression, set

$$\begin{aligned}\beta\Delta\text{PMF}(R) &\equiv \beta\text{PMF}(R) - \beta\text{PMF}(R_0) \\ &= -\ln(\rho(R)/\rho(R_0)) \\ &\equiv -\ln \tilde{\rho}(R)\end{aligned}\quad (10)$$

$$\tilde{\rho}(R) = \sum_{w=1}^W c_w(R) \tilde{\rho}_w^{(b)}(R) = \sum_{w=1}^W p_w(R) \tilde{\rho}_w^{(u)}(R) \quad (11)$$

where R_0 is some reference coordinate value. Then

$$\text{var}[\beta\Delta\text{PMF}(R)] = \text{var}\left[\ln \sum_{w=1}^W p_w(R) \tilde{\rho}_w^{(u)}(R)\right] \quad (12)$$

Using the idea that the $p_w(R)$ probabilities in eq 5 keep the ranges of R more-or-less disjoint, approximate eq 12 as

$$\begin{aligned}\text{var}[\beta\Delta\text{PMF}(R)] &\approx \sum_{w=1}^W \text{var}[\ln p_w(R) \tilde{\rho}_w^{(u)}(R)] \\ &\equiv \sum_{w=1}^W \text{var}[\beta\Delta\text{PMF}_w(R)]\end{aligned}\quad (13)$$

This approximation should be excellent because next-neighbor and further window pairs will have minimal overlap in typical applications, where the desire is to use as small a number of windows as possible. (Its validity is illustrated in section 3.5, see Figure 2.) Assuming that the two variance contributions in eq 13 are statistically independent, that is, using

$$\text{var}[\ln(XY)] = \text{var}[\ln X] + \text{var}[\ln Y] \quad (14)$$

and using propagation of error²⁷ for the log of the first term $\text{var}[\ln(X)] = \text{var}[X]/X^2$ one obtains

$$\text{var}[\ln p_w(R) \tilde{\rho}_w^{(u)}(R)] = \frac{\text{var}[p_w(R)]}{p_w^2(R)} + \text{var}[\beta\Delta\text{PMF}_w(R)] \quad (15)$$

with

$$\beta\Delta\text{PMF}_w(R) \equiv -\ln \tilde{\rho}_w^{(u)}(R) \quad (16)$$

a “local” contribution to the PMF.

Thus

$$\text{var}[\beta\Delta\text{PMF}_w(R)] = \frac{\text{var}[p_w(R)]}{p_w^2(R)} + \text{var}[\beta\Delta\text{PMF}_w(R)] \quad (17)$$

and the overall PMF can be written as

$$\text{var}[\beta\Delta\text{PMF}(R)] = \sum_{w=1}^W \text{var}[\beta\Delta\text{PMF}_w(R)] H_w(R) \quad (18)$$

where

$$\begin{aligned}H_w(R) &= 1 \quad R_w^{\min} < R < R_w^{\max} \\ H_w(R) &= 0 \quad \text{otherwise}\end{aligned}\quad (19)$$

with (R_w^{\min}, R_w^{\max}) essentially the range covered by $p_w^2(R)$ around the window equilibrium position. Thus,

$$\begin{aligned}\text{var}[\beta\Delta\text{PMF}(R)] &= \sum_{w=1}^W \text{var}[\beta\Delta\text{PMF}_w(R)] H_w(R) \\ &\approx \sum_w \text{var}[\beta\Delta\text{PMF}_w(R)] p_w^2(R)\end{aligned}\quad (20)$$

The second form in eq 20 provides the simple expression

$$\text{var}[\beta\Delta\text{PMF}(R)] = \sum_{w=1}^W (\text{var}[p_w(R)] + p_w^2(R) \text{var}[\beta\Delta\text{PMF}_w(R)]) \quad (21)$$

Note that the variance is independent of an origin of energy; $\text{var}[\beta\Delta\text{PMF}(R)] = \text{var}[\beta\text{PMF}(R)]$. In the numerical analysis carried out in section 4, we find that each contribution to the window sum gives similar results using the $p_w^2(R)$ and $H_w(R)$ versions.

2.3. Variance of $\text{PMF}_w(R)$. Consider first the $\text{var}[\beta\Delta\text{PMF}_w(R)]$ term. Use a Gaussian approximation to $\rho_w^{(b)}(R)$ with estimated sample mean $\mu_w = \bar{R}_w$ and standard deviation $\sigma_w^{(b)} = (\bar{R}_w^2 - \bar{R}_w^2)^{1/2}$

$$\rho_w^{(b)}(\mu_w, (\sigma_w^{(b)})^2) = \frac{1}{\sqrt{2\pi(\sigma_w^{(b)})^2}} e^{-((R-\mu_w)^2/2(\sigma_w^{(b)})^2)} \quad (22)$$

The variance of some multivariate function $G(x_1, x_2, \dots, x_n)$ using propagation of error²⁷ is

$$\text{var}[G(x_1, x_2, \dots, x_n)] = \sum_{i=1}^n \left[\left(\frac{\partial G}{\partial x_i} \right) \Big|_{x_i=\langle x_i \rangle} \right]^2 \text{var}[x_i] \quad (23)$$

neglecting covariances. The $\langle x_i \rangle$ denote the infinite trajectory averaged values of these parameters, though of course they can only be estimated as their sample averages, \bar{x}_i . Applied to window densities $\rho_w^{(b)}(R)$, the neglect of correlation between the different windows is quite reasonable in view of the thermal ensemble generation of each window.

From eq 3 and the assumed Gaussian form for $\rho_w^{(b)}(R)$, propagation of error then gives

$$\begin{aligned}\text{var}[\beta\Delta\text{PMF}_w(\mu_w, (\sigma_w^{(b)})^2 | R)] \\ = \text{var}[\beta f_w] + \text{var}[\ln(\tilde{\rho}_w^{(b)}(\mu_w, (\sigma_w^{(b)})^2 | R))]\end{aligned}\quad (24)$$

where the notation emphasizes that the $\text{var}[\beta\Delta\text{PMF}_w]$ will depend parametrically on the R value with R again within a region delimited by (R_w^{\min}, R_w^{\max}) . Working out the second term in eq 24 for a Gaussian yields

$$\begin{aligned}\text{var}[\ln(\tilde{\rho}_w^{(b)}(\mu_w, (\sigma_w^{(b)})^2 | R))] \\ = \frac{1}{(\sigma_w^{(b)})^2} \left[\frac{(R - \mu_w)^2}{(\sigma_w^{(b)})^2} \text{var}[\mu_w] + \frac{1}{(4\sigma_w^{(b)})^2} \right. \\ \left. \left[1 - \frac{(R - \mu_w)^2}{(\sigma_w^{(b)})^2} \right]^2 \text{var}[(\sigma_w^{(b)})^2] \right]\end{aligned}\quad (25)$$

A temperature dependence is implicit in this part of the variance because μ_w and $(\sigma_w^{(b)})^2$ are temperature dependent.

Now suppose the sampling is i.i.d. Then

$$\text{var}[\mu_w] = (\sigma_w^{(b)})^2 / N_w \quad (26)$$

because for an i.i.d. distribution the $\text{var}[\mu_w]$ is actually the variance of the trajectory data itself and this variance is the estimated trajectory variance, $(\sigma_w^{(b)})^2$, divided by the number of sample points, N_w . Also, assuming a normal distribution for the distribution of the true Gaussian width, its sample average variance is^{27,28}

$$\text{var}[(\sigma_w^{(b)})^2] = 2((\sigma_w^{(b)})^2)^2/N_w \quad (27)$$

Using these results in eq 25,

$$\begin{aligned} \text{var}[\ln(\tilde{\rho}_w^{(b)}(\mu_w, (\sigma_w^{(b)})^2|R))] &= \frac{1}{N_w} \left[\frac{(R - \mu_w)^2}{(\sigma_w^{(b)})^2} + \frac{1}{2} \left(1 - \frac{(R - \mu_w)^2}{(\sigma_w^{(b)})^2} \right)^2 \right] \\ &= \frac{1}{2N_w} \left(1 + \frac{(R - \mu_w)^2}{(\sigma_w^{(b)})^2} \right) \end{aligned} \quad (28)$$

For the $\text{var}[f_w]$ we can use a result generated by Zhu and Hummer:²³

$$\begin{aligned} \text{var}[f_{w+1} - f_w] &= [\text{var}(\bar{R}^{w+1}) + \text{var}(\bar{R}^w)] \\ &\quad [k(R_{w+1}^e - R_w^e)/2]^2 \end{aligned} \quad (29)$$

for neighbor window pairs, where \bar{R}^w is the sample average of the window w trajectory and R_w^e is the window w restraint position. Use of propagation of error gives:

$$\text{var}[f_{w+1} - f_w] = \text{var}[f_{w+1}] + \text{var}[f_w] - 2\text{cov}[f_{w+1}, f_w] \quad (30)$$

that is exact because it is a linear relation. Apply this result to neighbor window pairs and neglect the covariance to obtain

$$\begin{aligned} \text{var}[f_{w+1}] + \text{var}[f_w] &= (\text{var}[\bar{R}^{w+1}] + \text{var}[\bar{R}^w]) (k(R_{w+1}^e - R_w^e)/2)^2 \\ (w = 1, 3, \dots) \end{aligned} \quad (31)$$

and partition the variance as

$$\begin{aligned} \text{var}[f_w] &= [\text{var}(\bar{R}^w)] (k(R_{w+1}^e - R_w^e)/2)^2 \\ (w = 1, 2, \dots, W-1) \end{aligned} \quad (32)$$

Of course, as noted in eq 26, $\text{var}[\bar{R}^w] = \text{var}[\mu_w] = (\sigma_w^{(b)})^2/N_w$ so that with the definition $\Delta R_w^e = R_{w+1}^e - R_w^e$

$$\begin{aligned} \text{var}[f_w] &= ((\sigma_w^{(b)})^2/N_w) (k\Delta R_w^e/2)^2 \\ (w = 1, 2, \dots, W-1) \end{aligned} \quad (33)$$

Adding these two variance contributions of eqs 28 and 33 to eq 24 produces

$$\begin{aligned} \text{var}[\beta \Delta \text{PMF}_w(R)] &= \left(\frac{1}{N_w} \right) \left\{ (\beta \sigma_w^{(b)})^2 (k\Delta R_w^e/2)^2 + \frac{1}{2} \left(1 + \frac{(R - \mu_w)^2}{(\sigma_w^{(b)})^2} \right) \right\} \end{aligned} \quad (34)$$

This can be written as

$$\begin{aligned} \text{var}[\beta \Delta \text{PMF}_w(R)] &= \left(\frac{1}{N_w} \right) \left\{ \frac{(\Delta R_w^e/2)^2 (k_w^\beta)^2}{(\sigma_w^{(b)})^2} + \frac{1}{2} \left(1 + \frac{(R - \mu_w)^2}{(\sigma_w^{(b)})^2} \right) \right\} \end{aligned} \quad (35)$$

with $k_w^\beta \equiv 1/\beta(\sigma_w^{(b)})^2$ a “thermal” force constant.

An estimate of the relative contributions of the two terms can be obtained by writing the bracketed term of eq 35 as

$$\begin{aligned} &\left\{ \frac{(\Delta R_w^e/2)^2 (k_w^\beta)^2}{(\sigma_w^{(b)})^2} + \frac{1}{2} \left(1 + \frac{(R - \mu_w)^2}{(\sigma_w^{(b)})^2} \right) \right\} \\ &\approx \left\{ \frac{(\Delta R_w^e/2)^2}{(\sigma_w^{(b)})^2} + \frac{1}{2} \left(1 + \frac{(\Delta R_w^e/2)^2}{(\sigma_w^{(b)})^2} \right) \right\} \end{aligned} \quad (36)$$

where we have estimated the first term as $1/\beta(\sigma_w^{(b)})^2 \approx \beta k_w^\beta$. That is, assumed the width of each biased window's distribution is dominated by its force constant in the canonical ensemble, a strong force condition. We estimated the second term by setting $R - \mu_w = (\Delta R_w^e/2)$, its (maximal) boundary value. The two terms in eq 36, the first of which is $\text{var}[\beta f_w]$ and the second $\text{var}[\ln(\tilde{\rho}_w^{(b)} \mu_w (\sigma_w^{(b)})^2 | R)]$, are of the same order. This estimate shows that if neighbor window separations are large compared to the corresponding window widths, then the $\text{var}[\beta f_w]$ contributions will be significant compared with the variance arising from the parameters of the windows.

2.4. Variance of $p_w(R)$. To estimate $\text{var}[p_w(R)]$, with the $p_w(R)$ defined in eq 5, using propagation of error, view the parameters in the propagation of error as the average number of counts $\overline{\rho_w^{(b)}(R)}$ in each ΔR interval of the histograms of the window densities. Error propagation on $\text{var}[p_w(R)]$ then yields

$$\begin{aligned} \frac{\text{var}[p_w(R)]}{p_w^2(R)} &= \left(\frac{\text{var}[\rho_w^{(b)}(R)]}{(\overline{\rho_w^{(b)}(R)})^2} + \frac{\sum_{w'} \text{var}[\rho_{w'}^{(b)}(R)]}{(\sum_{w'} \overline{\rho_{w'}^{(b)}(R)})^2} \right. \\ &\quad \left. - \frac{2\text{var}[\rho_w^{(b)}(R)]}{(\overline{\rho_w^{(b)}(R)})(\sum_{w'} \overline{\rho_{w'}^{(b)}(R)})} \right) \end{aligned} \quad (37)$$

assuming that correlations between the $\rho_w^{(b)}(R)$ densities in the different windows are negligible.

The i.i.d. standard density variance estimate with histogram width ΔR that satisfies $\overline{\rho_w^{(b)}(R)} \Delta R \ll 1$ ^{7,29} is

$$\text{var}[\rho_w^{(b)}(R)] = \frac{\overline{\rho_w^{(b)}(R)}}{N_w \Delta R} \quad (38)$$

and yields

$$\begin{aligned} \text{var}[p_w(R)] &= p_w^2(R) \left(\frac{1}{N_w \Delta R} \right) \left(\frac{1}{\overline{\rho_w^{(b)}(R)}} - \frac{1}{\sum_{w'} \overline{\rho_{w'}^{(b)}(R)}} \right) \end{aligned} \quad (39)$$

An instructive alternative form is

$$\begin{aligned} \text{var}[p_w(R)] &= \left(\frac{1}{N_w} \right) \left(1/\Delta R \sum_{w'} \overline{\rho_{w'}^{(b)}(R)} \right) \{ p_w(R) (1 - p_w(R)) \} \end{aligned} \quad (40)$$

Around $R = \mu_w$, the variance is small because there is a lot of data. Away from there, the variance increases in accord with less data, and still further away the $\text{var}[P_w(R)] \rightarrow 0$ because there is no more data. Note that $p_w(R)(1 - p_w(R))$ tends to cancel because $p_w(R)$ and $p_w^2(R)$ are quite similar, especially around $R = \mu_w$.

The $\text{var}[\beta\Delta\text{PMF}_w(R)]$ in eq 34 and $\text{var}[P_w(R)]$ in eq 40 expressions, combined in eq 20 as local contributions weighted by either $H_w(R)$ or $p_w^2(R)$ to provide $\text{var}[\beta\Delta\text{PMF}(R)]$, are the main results of this PMF(R) variance analysis. Both contributions to the variance are proportional to $1/N_w$, showing that the PMF standard deviation has the characteristic statistical error scaling.

3. MODEL AND SOLUTION METHOD

3.1. Model. A model^{25,26} that mimics a protein with a slow mode in some complex coordinate arising from a large number of coupled two-state degrees of freedom (e.g., from dihedral states) consists of N “particles” with coordinates $\mathbf{y} = (y_1, y_2, \dots, y_N)$ interacting with potential

$$U(\mathbf{y}) = \sum_{i=1}^N \left[-(1-\theta)\frac{1}{2}y_i^2 + \frac{1}{4}y_i^4 \right] + \frac{\theta}{2} \sum_{i=1}^N (y_{i+1} - y_i)^2; \quad y_{N+1} = y_1 \quad (41)$$

The one particle potentials present double wells with minima at $y_{\pm\theta} = \pm(1-\theta)^{1/2}$ and a barrier at $y_0 = 0$, with barrier height $U_B = 1/4(1-\theta)^2$, for $0 < \theta < 1$. The total potential for window simulations

$$V_w(\mathbf{y}) = U(\mathbf{y}) + u_w(\mathbf{y}) \quad (42)$$

is obtained by adding the restraint potential

$$u_w(\mathbf{y}) = \frac{k_w}{2}(R(\mathbf{y}) - R_w^e)^2 \quad (43)$$

with $R(\mathbf{y}) = \sum_{i=1}^N (y_i/N)$ a collective reaction coordinate, and R_w^e and k_w ($w = 1, 2, \dots, W$) the window restraint positions and force constants, respectively.

3.2. Langevin Equation. The window trajectories are generated with use of a (dimensionless) Langevin equation (LE)^{30,31}

$$dy_i/dt = -\partial V_w(\mathbf{y})/\partial y_i + \sqrt{2D}\eta_i(t) \quad (44)$$

where the Gaussian, delta correlated, white noise $\eta_i(t)$ satisfies

$$\langle \eta_i(t) \eta_j(t') \rangle = 2D\delta_{ij}\delta(t - t') \quad (45)$$

The LE is guaranteed to produce the correct equilibrium distribution function, as readily shown by conversion to the equivalent Fokker–Planck equation.³⁰ Dynamical information is not available from this extreme overdamped LE; the “time” is just a way to parametrize the evolution to the equilibrium distribution, $P_{\text{eq}}(\mathbf{y}) \sim \exp(-V(\mathbf{y})/D)$. Solution of the Langevin equation is carried out with a stochastic version of the second-order Runge–Kutta algorithm.³² A step-size $h = 0.01$ was used for all trajectory generation, and data was saved every 100 steps to ensure that succeeding trajectory points are independent, or every 10 steps for the investigation of statistical inefficiency (see Appendix A). The trajectories are aged sufficiently to generate equilibrium sampling. The pseudorandom numbers for generating the noise are obtained with a “Mersenne Twister”³³ that has a very long period. The Gaussian sampling is carried

out by a simple, fast method that provides the correct mean and variance, as set by the value of D .^{32,34}

3.3. Expected Potential of Mean Force. In the LE, the temperature is defined by the value of D . For N particles, and a θ value, in the absence of the window potentials, $u_w(\mathbf{y})$, initiating a trajectory from one set of minima of the double well potentials will, for large D , readily sample both sides of the double well potentials and results in a potential of mean force for the reaction coordinate R that consists of a single well. As the temperature is reduced, going over the barriers becomes increasingly difficult, requiring exponentially longer trajectories to provide equilibrium sampling. The PMF will display a barrier around $R = 0$ and two minima to either side of $R = 0$; a double well PMF. The addition of suitably chosen window potentials will provide reasonable sampling in all relevant regions of the reaction coordinate and, with the use of WHAM to combine the data and unbiased the window trajectories, permit construction of the PMF. For $N = 10$ and $\theta = 0.5$, a temperature of $D = 0.15$ results in a double well PMF. For these parameters, relatively long trajectories are required to adequately sample repeated transitions between the two sides, while the use of restraint windows and the WHAM procedure leads to accurate PMFs with 50 000 steps/window.

3.4. The Reaction Coordinate. That the collective coordinate R is an appropriate reaction coordinate was verified by carrying out a Principal Component Analysis (PCA)^{35–37} on a nonwindow trajectory (potential $U(\mathbf{y})$) that is sufficiently long to adequately sample configuration space, as verified by obtaining essentially the same PMF as found by the window method. A PCA constructs the covariance matrix of the particle fluctuations from their respective average positions and diagonalizes this matrix. It finds linear combinations of the particle coordinates that successively account for decreasing amounts of the total mean square fluctuation over the trajectory (the sum of the eigenvalues of the covariance matrix). If the eigenvalue of the first mode is much larger than those of the remaining modes, then there is a slow mode separated from the other modes. The first PCA mode eigenvalue evaluated for the conditions used here is well separated from the other modes, occupying 77% of the total mean square fluctuation. The first mode’s trajectory has a Pearson correlation coefficient of 0.97 with the collective coordinate, R , trajectory. Thus, as verified by the PCA, R reports on a slow motion of the system, and is the analog of some slow conformational coordinate in a protein trajectory.

3.5. The Potential of Mean Force. The model potential equations of motion were integrated for $N = 10$, $\theta = 0.5$, and $D = 0.15$ using 11 and 7 windows with R_w^e covering the range -1.5 to 1.5 in 0.3 unit increments using 11 windows and the same range with 0.45 increments using 7 windows, with $k_w = 25$ for all windows. Each window’s trajectory consisted of 50 000 steps, after the equilibration period. Histograms of the reaction coordinate for each window are displayed in Figure 1. They are quite Gaussian, as is verified by evaluating their third and fourth cumulants that should be strictly zero. Successive windows do overlap as required for WHAM derived PMFs. The data for 7 windows is similar with of course smaller overlaps of the histograms.

The 11 window WHAM based PMF is displayed in Figure 2. (The PMF from 7 windows is essentially the same as that shown in Figure 2, with the important limitation that the range over which it can be determined is limited to about $(-1.2 < R < 1.2)$.) The line labeled βPMF is generated with WHAM

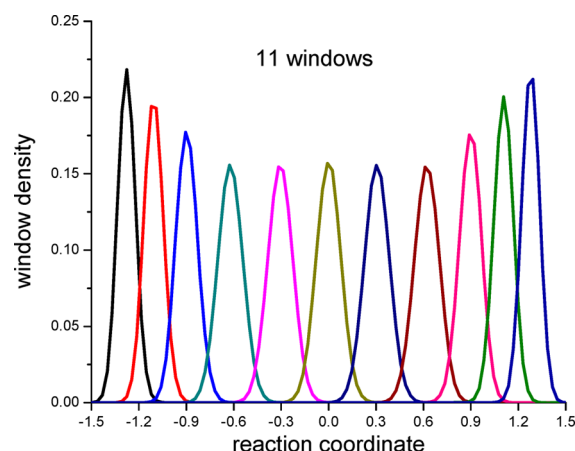


Figure 1. Histograms of the reaction coordinate for the 11 windows. All are quite Gaussian as monitored by their moments. The end ones have peak positions moved inward from their respective R_w equilibrium positions.

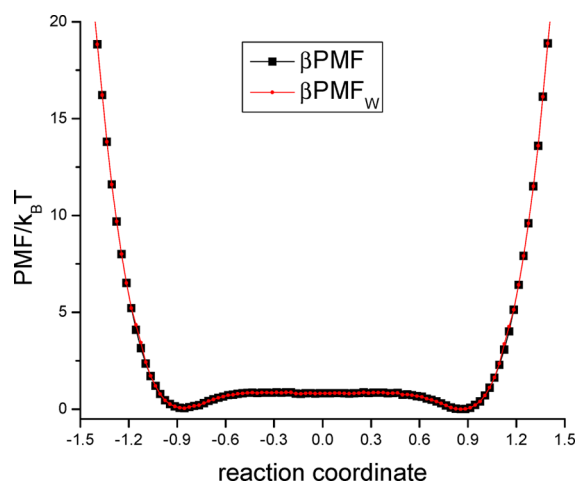


Figure 2. WHAM-obtained PMF in temperature units for 11 windows. The $\beta\text{PMF}(R)$ uses all the windows in WHAM, as conventionally done. The $\beta\text{PMF}_w(R)$ is constructed with WHAM using each pair of neighboring windows (see text). To graph resolution, they provide the same PMF.

pooling the data from all 11 windows. The line labeled βPMF_w uses only each succeeding neighbor pair to construct the PMF. To do so, the two part strategy is used. The free energies f_w are obtained from eq 6 using only neighbor pairs data. Then, the PMF is generated using eqs 7–9, for each local range defined in eq 19, on the desired scale of resolution. To the resolution of the graph, they provide identical PMFs. As long as window overlaps between next-nearest and further neighbor windows are negligible, the two methods should provide essentially the same PMF.

4. MODEL VARIANCE

The expression for the PMF variance in eqs 20 and 21 has two origins: from the variance of the $p_w(R)$ defined in eq 40 and from the variance of the local PMFs, $\beta\Delta\text{PMF}_w(R)$, defined in eq 35. Because the variances have been defined “point wise”, that is, over slices of size ΔR , there will be a dependence on this resolution of the PMF. The difference in variance definitions in eqs 20 and 21 lead in to somewhat different results as we now show.

4.1. Point Wise Variance. For the 11 window data, we display the two contributions to the variances in Figures 3 and

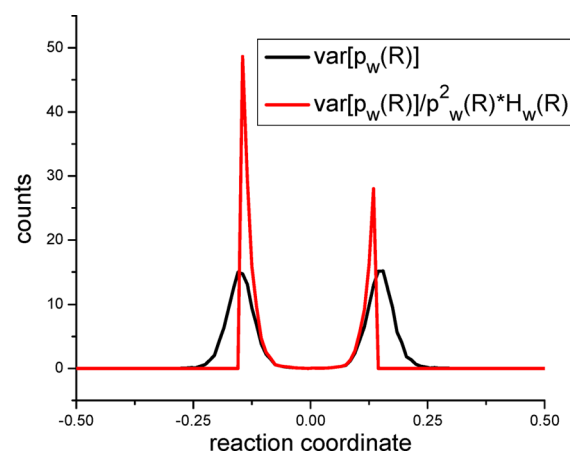


Figure 3. Variances of the window probability contribution for window 6 as a function of the reaction coordinate evaluated using a cutoff $H_w(R)$ or using the $p_w^2(R)$ (see eq 20) as ways to localize the variance for each window.

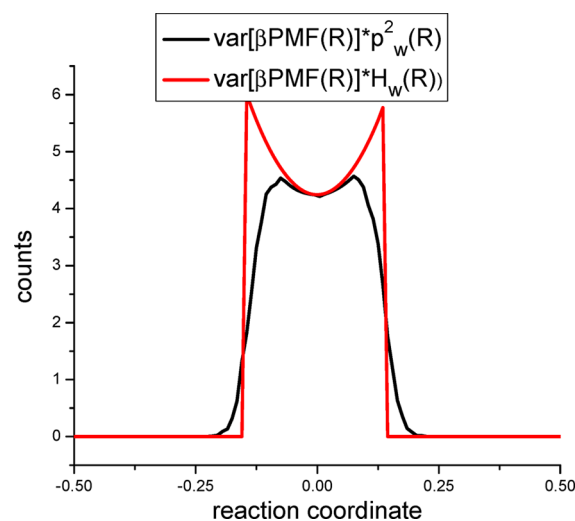


Figure 4. Variances of the PMF contribution for window 6 as a function of the reaction coordinate evaluated using a cutoff $H_w(R)$ or using $p_w^2(R)$ (see eq 20) as ways to localize the variance for each window.

4. Only the middle window data is displayed; all others are similar in shape for both variance contributions, with amplitudes that fall off toward the beginning and end windows. These and all the following plots do not include the number of points in the windows, the N_w . The plots are for the two ways of localizing the data as given in eqs 17 and 19–21. Note that the hard cutoff $H_w(R)$ version requires sharp boundary values for R , and the boundary resolution is limited by the size of the grid spacing ΔR , here $\Delta R = 0.01$. The structure of $\text{var}[P_w(R)]$ given in eq 40 reflects the behavior of the window histograms $\rho_w^{(b)}(R)$ with their essentially Gaussian behavior. As noted after eq 40, the variance is small around the peak of $\rho_w^{(b)}(R)$, increases in the wings, and decreases when there is no more data. As expected, the $\text{var}[P_w(R)]$ depends on the choice of ΔR , increasing in magnitude when the grid spacing is finer. For the

$\text{var}[\beta\Delta\text{PMF}_w(R)]$ contribution, the results are essentially independent of ΔR , as is clear from eq 34. Its variance contribution maximizes around the restraint position and has contributions from the $\text{var}[\beta f_w]$ and the variances of the Gaussian parameters.

4.2. Averaged Variance. The point wise variances vary strongly as functions of the $\rho_w^{(b)}(R)$ biased densities. To obtain results on the scale of variation of the PMF(R) it is sensible to integrate over the scale of the window separations. Figure 5

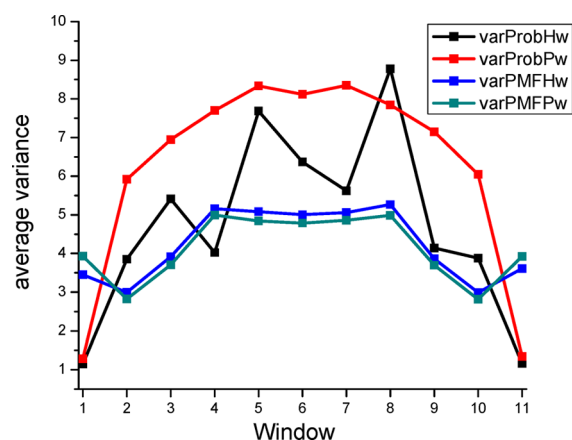


Figure 5. Variance contributions from $p_w(R)$ (labeled as Prob) and $\beta\text{PMF}_w(R)$ (labeled as PMF) averaged over the local ranges selected using $H_w(R)$ (labeled with Hw) and $p_w^2(R)$ (labeled with Pw) for the 11 windows. The lines are just a guide for the eye. The grid spacing for the point wise variances is $\Delta R = 0.01$. The two contributions to the overall variance are of comparable magnitude.

shows these averaged variances for the 11 windows. The variance for $p_w(R)$ using $H_w(R)$ as a local cutoff is noisier, reflecting that found for the point wise data. The scale of the two contributions to the variance is similar for this choice of grid spacing ΔR . Of course, if a coarser grid spacing is used, the contribution from $\text{var}[p_w(R)]$ will be reduced. Thus, this contribution to the variance can be minimized by using a coarser grid with the caution that the PMFs' variation with R may be lost. Figure 6 shows the dependence of $\text{var}[p_w(R)]$ on

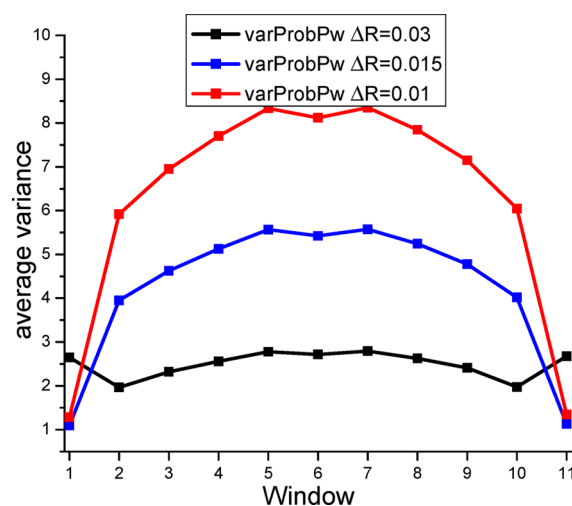


Figure 6. Variance contributions from $p_w(R)$ averaged over the local ranges selected using $p_w^2(R)$ for the 11 windows for several grid spacings, ΔR , for the point wise variance evaluation.

the grid spacing. As the number of bins is decreased, the averaged variances decrease in proportion. The end windows are an exception: the cutoff factor for their ranges is not equivalent to those of the interior windows. The variances are roughly constant across the interior windows, more so as the grid spacing is increased. The corresponding plots for the $\beta\Delta\text{PMF}_w(R)$ variance contribution shows hardly any effect of grid spacing (data not shown).

The point wise variance data for 7 windows mirrors that for the 11 windows. The variance contributions for $p_w(R)$ are scaled upward, as results from the reduced density where neighbor windows overlap. The $\beta\Delta\text{PMF}_w(R)$ contribution variance is not as greatly affected for the interior windows. These results are reflected in the averaged variances. Figure 7

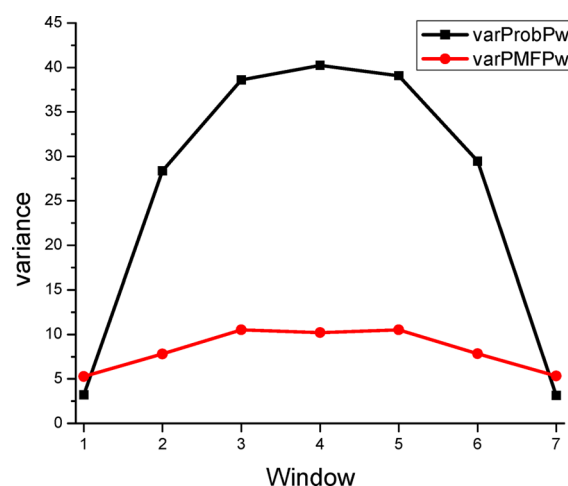


Figure 7. Variance contributions from $p_w(R)$ and $\beta\Delta\text{PMF}_w(R)$ averaged over the local ranges selected using $p_w^2(R)$ (labeled with Pw) for the 7 windows. The grid spacing for the point wise variances is $\Delta R = 0.01$.

plots these averaged variances. The $p_w(R)$ averaged variance contributions are considerably larger than those with 11 windows, and dominates those from $\beta\Delta\text{PMF}_w(R)$. The dependence on grid spacing mirrors that for 11 windows.

4.3. Dependence of Averaged Variance on the Number of Samples. The data for the above plots were generated with 50 000 samples per window. In principle, since the sampling points are i.i.d., similar results should be obtained with fewer samples per window. For the $\text{var}[\beta\Delta\text{PMF}_w(R)]$ contribution, as long as parameters μ_w and $\sigma_w^{(b)}$ are adequately evaluated, the corresponding variance should not change (excluding the N_w factor). For the $\text{var}[p_w(R)]$ contribution, there will be some effect that arises mainly from the point wise variance in regions where the windows overlap and the lack of data becomes serious (see Figure 3). In Figures 8 and 9, we plot these variances for 50 000, 12 500, and 3125 samples per window. There is sufficient data that even the least amount used provides essentially identical results. The exception is for the extreme windows of the $\text{var}[\beta\Delta\text{PMF}_w(R)]$ that as noted above is not on an equivalent footing with the interior windows.

The corresponding plots for the 7 window data are shown in Figures 10 and 11. The $\text{var}[\beta\Delta\text{PMF}_w(R)]$ is hardly affected by the reduction in point density but the variance contributions from $p_w(R)$ show an effect. Redoing the same procedure for a completely independent seven window run shows that the

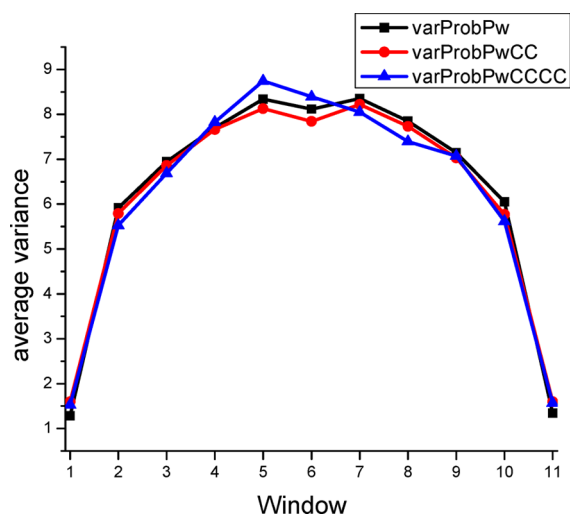


Figure 8. Variance contributions from $p_w(R)$ averaged over the local ranges selected using $p_w^2(R)$ for the 11 windows as a function of the number of samples per window. The 50 000 samples/window is compared with the CC (CCCC) line for 12 500 (3125) samples per window.

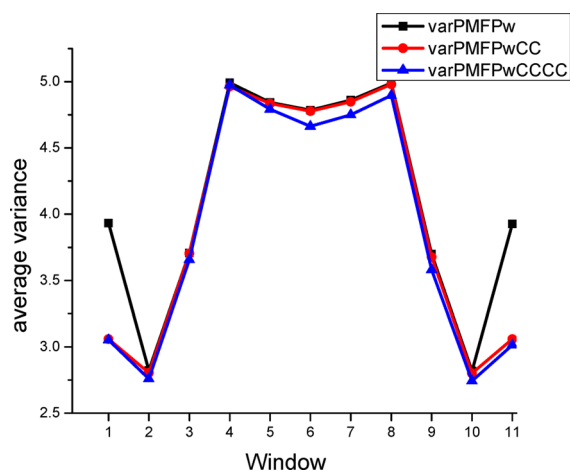


Figure 9. The $\text{var}[\beta\Delta\text{PMF}_w(R)]$ averaged over the local ranges selected using $p_w^2(R)$ for the 11 windows as a function of the number of samples per window. The 50 000 samples/window is compared with the CC (CCCC) line for 12 500 (3125) samples per window.

result for $\text{var}[\beta\Delta\text{PMF}_w(R)]$ is essentially the same as shown in Figure 11 but the variance contributions from $p_w(R)$ for 3125 samples per window are not available.

5. DISCUSSION

In this work, an analytic expression for the variance of a potential of mean force generated with the use of WHAM has been generated based on a number of assumptions. A local approximation eq 13 that partitions the variance into contributions from each window's data range is central to the analysis. Because the neighbor window densities, $\rho_w^{(b)}(R)$, must have some overlap but more remote window densities essentially do not overlap, this local approach is appropriate, as illustrated in Figure 2. In order to make umbrella sampling efficient, one wants to minimize the number of windows to simulate; there would be little point in having further than neighbor windows overlap. Thus, this condition is consistent

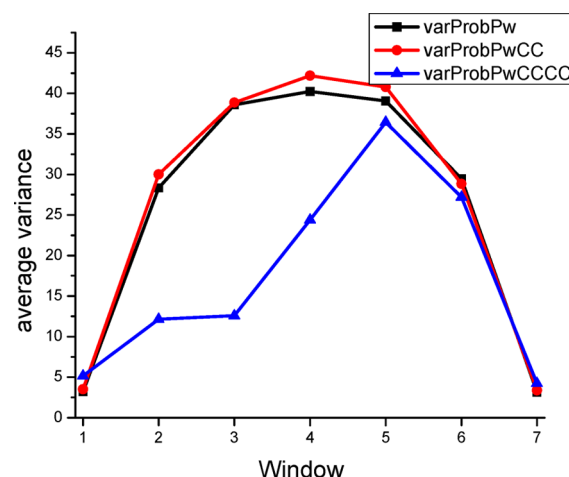


Figure 10. Variance contributions from $p_w(R)$ averaged over the local ranges selected using $p_w^2(R)$ for the 7 windows as a function of the number of samples per window. The 50 000 samples/window is compared with the CC (CCCC) line for 12 500 (3125) samples per window.

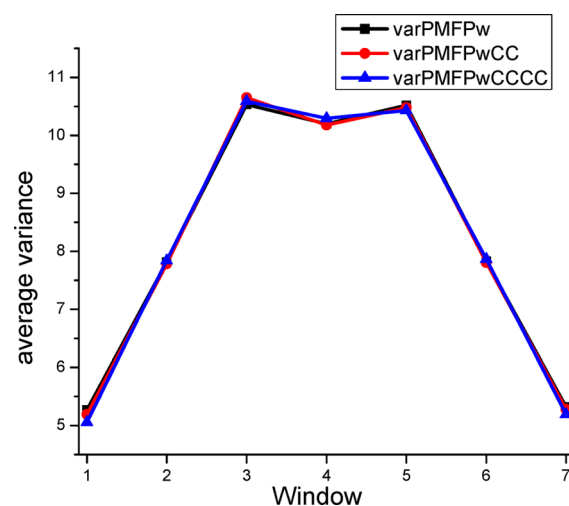


Figure 11. The $\text{var}[\beta\Delta\text{PMF}_w(R)]$ averaged over the local ranges selected using $p_w^2(R)$ for the 7 windows as a function of the number of samples per window. The 50 000 samples/window is compared with the CC (CCCC) line for 12 500 (3125) samples per window.

with umbrella sampling or related multiple restrained window methods.

The variance expression is given by either eq 20 or eq 21, with the latter more numerically stable. It has two sources: from the $p_w(R)$ and from the $\text{PMF}_w(R)$, the local PMF variance defined in eq 34. The contribution from the local PMF is determined by the means and widths of the assumed Gaussian distributions of the sampled biased window densities, and by the window free energies. Obtaining Gaussians for the window densities is largely governed by the sizes of the restraint force constants, k_w . If the restraints are not sufficiently strong, they will not confine the system to the desired values of the reaction coordinate. To the extent that an average window position $\mu_w = \bar{R}_w$ will be held reasonably close to the imposed restraint position, R_w^c , the corresponding k_w will be large enough that the Gaussian assumption tends to be robust. The $\text{var}[f_w]$ can be obtained by relating them to the $\text{var}[\mu_w]$, providing a completely analytic expression for the $\text{var}[\text{PMF}_w(R)]$ con-

tributions, as given in eq 35. The size of $\text{var}[f_w]$ is controlled by the ratio of the window separations ΔR_w^e and the window widths $2\sigma_w^{(b)}$ and the ratio of mechanical k_w and thermal k_w^b force constants. For a large enough k_w , the thermal width of the window is determined by this force constant; thus, $k_w^b \rightarrow k_w$. As successive windows separate more, relative to the window widths, the $\text{var}[f_w]$ increases, reflecting that the f_w are coupled together, in principle, through all the windows though, again, in practice, mainly through neighbor window pairs. The Gaussian parameters' contribution is minimal at $R = \mu_w$ where the data is maximal, and increases from there. Note that this local expression is defined over its range, so the variance in eq 35 does not increase without bound. These $\text{var}[\text{PMF}_w(R)]$ contributions to the total variance are essentially independent of the grid spacing, as the analytic expression suggests, and is verified by the numerical results.

The analysis of the contribution from $\text{var}[p_w(R)]$ relies on asserting that the correlations between the $\rho_w^b(R)$ densities in the different windows can be neglected. This is a reasonable assumption because each window's trajectory is thermally generated. It should be an even better assumption for distance REM based simulations.¹⁴ The point wise $\text{var}[p_w(R)]$ has a characteristic dependence on its R range, with a near-zero value around the peak of $\rho_w^b(R)$ and a rapid increase to peak where the neighbor window overlaps occur. There is an explicit dependence on the grid spacing ΔR and this variance can be reduced by increasing ΔR . Whether this contribution is needed relative to that from the local PMF variance then becomes a matter of what resolution of the PMF is desired. As a practical matter, even for a rapidly changing region of the PMF, most likely ΔR can be set sufficiently big to suppress this contribution. If N_w is decreased, the contribution from the $\text{var}[p_w(R)]$ will become less well predicted, as noted in section 4.3. Thus, other things being equal, a larger ΔR will be required to decrease its contribution to the overall variance. If the $\text{var}[p_w(R)]$ are considered as negligible, then $\text{var}[\beta\Delta\text{PMF}_w(R)] \approx \sum_w \text{var}[\beta\Delta\text{PMF}_w(R)]p_w^2(R)$, an expression analogous to what Kastner and Thiel²¹ present for the mean force corresponding to the $\text{PMF}_w(R)$. Working with the mean force does not incorporate the $\text{var}[f_w]$ contributions.

From the analysis of the $\text{var}[\beta\Delta\text{PMF}_w(R)]$ and $\text{var}[\beta p_w(R)]$, recommendations for setting the width and separation of windows can be made. The thermal width of a biased window, for a sufficiently large force constant where the restraint potential dominates, scales as $\langle R^2 \rangle = (1/\beta k)$. For typical situations at ambient temperature where $1/\beta = 0.6$ kcal/mol, a force constant of 6 kcal/mol/Å² gives $(\sigma_w^{(b)}) \sim 0.33$ Å. To keep the $\text{var}[\beta\Delta\text{PMF}_w(R)]$ contribution that scales as $(\Delta R_w^e/2)^2/(\sigma_w^{(b)})^2$ below a factor of about 10 permits window separations of ~ 2 Å, according to the estimates in eq 36. With 1000 snapshots per window of simulation data, that would produce a standard deviation of $(10/1000)^{1/2} = 0.1$ in temperature scaled units, which is ~ 0.25 kcal/mol for the standard deviation of $\Delta\text{PMF}_w(R)$. Doubling the window separation increases the standard deviation to ~ 0.5 kcal/mol that now is essentially the thermal energy. Of course, if more data is available, then, with the overall $(1/N_w)^{1/2}$ scaling, stronger force constants and/or larger window separations can be used. The other, $\text{var}[p_w(R)]$, contribution is controlled by the number of data points in the ΔR slices. If the window half-width is $(\sigma_w^{(b)}) \sim 0.33$ Å, then $\Delta R \sim 0.03$ Å would produce about the same standard deviation contribution as that from the $\Delta\text{PMF}_w(R)$ contribution. If window separations of ~ 2 Å are being used, and presumably

the scale of variation of the potential of mean force will be similar to the window separation, it would be sensible to pick a $\Delta R > 0.3$ Å that is sufficiently larger to make this contribution to the overall variance negligible.

The other approximation made to generate an analytic result is that the sampling is from an independent identically distributed process entailing two requirements. First, the trajectory has to be sufficiently aged that equilibrium sampling is achieved. Statistical tests do exist for checking this aspect and have been used to analyze MD trajectory data.^{20,38} For the model potential used here, after a transient period, the sampling can be considered as equilibrium sampling because the barriers in the potential are uniform, versus complex systems where there is no a priori estimate of the distribution of barrier heights. Second, the use of extreme overdamped Langevin equation equations of motion is an effective strategy to generate less correlated trajectories, compared with MD algorithms where the mechanical forces are more dominant. Though, it should be noted that for more complex, solvated protein systems the exploration of configuration space is slow in the extreme overdamped LE limit. The block averaging method used to obtain the statistical inefficiency (see Appendix A) showed that by outputting every 100th integration step the sampling in each window is essentially i.i.d., so that the data is appropriate for the variance analysis.

In summary, an expression for the variance of a PMF along its reaction coordinate has been generated that is inversely proportional to the number of samples in the windows. There are two contributions with different dependences on the reaction coordinate grid spacing, ΔR . The variance from the window probabilities $p_w(R)$ that depends on ΔR can be minimized by choosing a sufficiently large ΔR . The variance from the local $\text{PMF}_w(R)$ then will most likely dominate, as it is essentially independent of ΔR . The variance of the $\text{PMF}(R)$ can readily be evaluated from the biased window densities and their corresponding means and standard deviations.

■ APPENDIX A

Because the analytic variance expressions that are developed here rely on sampling from an independent, identically distributed (i.i.d.) random variable, the generated data is tested to make sure sample points are separated sufficiently to be i.i.d.

Block averaging³⁹ is the basis of an efficient way to ascertain the degree of correlation in serial data, here the reaction coordinate trajectory data for each window. It avoids the estimation of a correlation time that may not be numerically reliable for the given data, and it uses all the data in a computationally efficient manner. It is also useful for the generation of a lower bound on how many blocks are needed before accurate estimates of the degree of correlation can be obtained.

Consider a random variable y with probability distribution $p(y)$, average $\langle y \rangle = \int y p(y) dy$, and standard deviation $\sigma = (\int (y - \langle y \rangle)^2 p(y) dy)^{1/2}$. Let \bar{y}_n denote the sample average

$$\bar{y}_n = \frac{1}{n} \sum_{i=1}^n y_i \quad (\text{A1})$$

If the sampling is from an i.i.d. random variable, then $\sigma(n)$, the sample variance of the \bar{y}_n , obeys²⁷

$$\sigma(n) = \sigma/\sqrt{n} \quad (\text{A2})$$

To test if the data is i.i.d., serially divide the n samples into n_b blocks of length n/n_b with the number of blocks satisfying $n_b = 2^b$ ($b = 0, 1, 2, \dots, B-1$) with $n_{B-1} = 2^{B-1} = n$. Let

$$\bar{y}_b(n_b) = \left(\frac{1}{n/n_b} \right) \sum_{i=b(n/n_b)+1}^{(b+1)(n/n_b)} y_i \quad (\text{A3})$$

be the sample average of the data in block b of length n/n_b . Block averaging proceeds by computing the variance of the blocked sample averages.³⁹

$$\text{var}(n_b) = \left(\frac{1}{n_b} \right) \sum_{b=0}^{n_b-1} \bar{y}_b^2(n_b) - \left(\left(\frac{1}{n_b} \right) \sum_{b=0}^{n_b-1} \bar{y}_b(n_b) \right)^2 \quad (\text{A4})$$

With this definition of block averages, $\bar{y}_b(n_b)$, note that

$$\text{var}(n_b = 1) = 0$$

$$\begin{aligned} \text{var}(n_b = n) &= \frac{1}{n} \sum_{b=0}^{n_b-1} y_{b+1}^2 - \left(\frac{1}{n} \sum_{b=0}^{n_b-1} y_{b+1} \right)^2 \\ &\equiv \text{var} \end{aligned} \quad (\text{A5})$$

with var denoting the variance of the length n trajectory data. (The block average order here is opposite to that of Flyvbjerg and Petersen.³⁹) For this division into n_b blocks, the sample standard deviation of the block averages is

$$\sigma(n_b) = \sqrt{\text{var}(n_b)/(n_b - 1)} \quad (\text{A6})$$

Note that in the following $\sigma(n_b) = (\text{var}(n_b)/n_b)^{1/2}$ is actually evaluated, for display convenience; the distinction is of no import for practical values of n_b . To study the dependence of the standard deviation on n_b , it is convenient to define

$$\sigma_a \equiv \sigma(n_b) \sqrt{n_b/n} \quad (\text{A7})$$

If σ_a were independent of n_b ($n_b = n, n/2, \dots, 1$), each σ_a would be constant with value

$$\sigma_a = \sigma(n) = \sigma/\sqrt{n} \quad (\text{A8})$$

If the sample points are correlated, then one does need to do block averaging to construct i.i.d. variables. A convenient measure of the degree of correlation is given by the statistical inefficiency (SI)^{40,41}

$$\text{SI} = \lim_{n/n_b \rightarrow \infty} \frac{(n/n_b) \text{var}(n_b)}{\text{var}} \quad (\text{A9})$$

where $\text{var}(n_b)$ is given in eq A4 and var is the variance for the entire trajectory. For a trajectory of a particular length, the SI provides the effective length of the trajectory. Or, said otherwise, it indicates how much longer a trajectory would have to be to produce i.i.d. sampling. Thus, the SI provides the factor to reduce n in the $\sigma(n)$ formula to obtain an accurate estimate of the standard deviation. Note that the SI increases with increasing block size, n/n_b , until it is large enough to provide i.i.d. points. The SI becomes inaccurate for a small number of blocks because the error in $\text{var}(n_b)$ itself, which, for block sizes sufficiently large that the central limit theorem applies, is proportional to $(\text{var}(n_b)/n_b)^{1/2}$, increases.^{39,40}

To assess how accurate $\sigma(n_b)$ is as a function of block size, we first generated a set of variances from the normal distribution $N(0,1)$, using the Box Mueller method⁴¹ and the c++ `rand()` random number generator (using the clock to define each new seed). Figure 12 displays $\sigma_a \equiv \sigma(n_b)(n_b/n)^{1/2}$

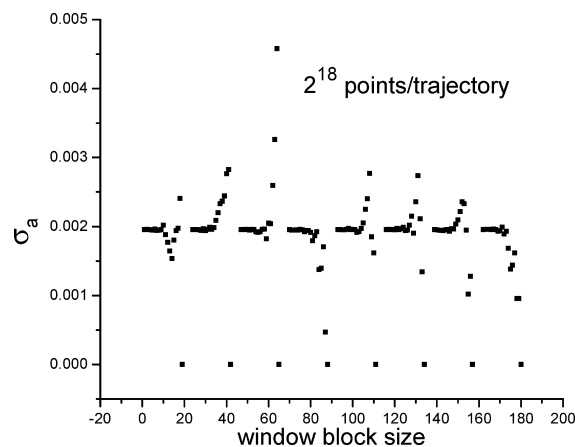


Figure 12. Approximated standard deviations $\sigma_a \equiv \sigma(n_b)(n_b/n)^{1/2}$ for sampling from $N(0,1)$ for eight independent “trajectories”. The number of points in each block increases from 1 to n from left to right for each of the trajectories. (The block size indicator in the figure is just a running number.) If the number of blocks is too small, the corresponding variance cannot be accurately obtained. In the above data, $n_b \geq \sim 256$ provides an accurate estimation the variance.

versus increasing n/n_b , the number of points in each block, for $2^{18} = 262\,144$ points, for eight independent runs. If σ_a were independent of n_b ($n_b = n, n/2, \dots, 1$), each σ_a would be $(1/2^{18})^{1/2} = 0.001953$. The stable value of $\sigma_a \approx 0.00196$ is obtained when $n_b \geq \sim 256$. As long as there are order 1000 data points that can be considered as independent, the sample standard deviation $\sigma(n_b)$ should be accurately predicted from the data.

For the data generated from the model potential, we first block average the window reaction coordinates when they are sampled every 100th integration step. For this sampling rate, Figure 13, a plot of σ_a , defined in eq A7, for the 11 windows shows that the variances reach plateau values about when the window block size also is ~ 256 . The inner windows show some slight non-i.i.d. behavior. The corresponding statistical inefficiency displayed in Figure 14 shows that $\text{SI} \sim 1$ with slightly larger values for some of the “interior” windows. As noted above, the SI will not be accurate for large block sizes.

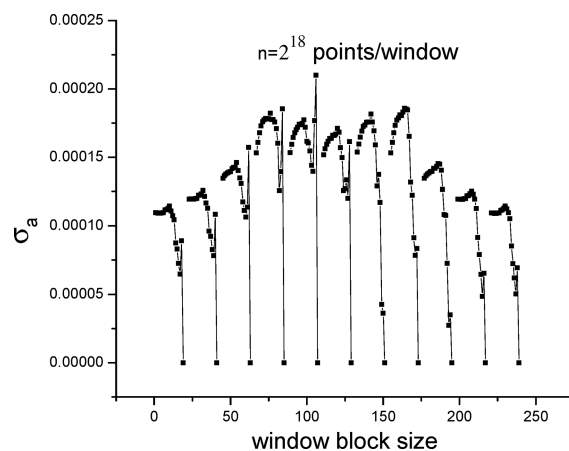


Figure 13. Approximated standard deviations $\sigma_a \equiv \sigma(n_b)(n_b/n)^{1/2}$ of each window's reaction coordinate trajectory for the model potential data. The number of points in each block increases from 1 to n from left to right for each of the trajectories. The data sampling is every 100th integration step. The σ_a approach their plateau values for block sizes $n_b \geq \sim 256$.

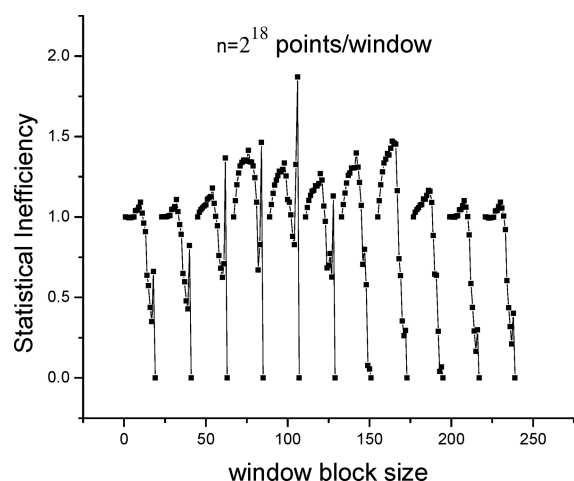


Figure 14. Statistical inefficiency of each window's reaction coordinate trajectory. The number of points in each block increases from 1 to n from left to right for each of the trajectories. The data sampling is every 100th integration step. The statistical inefficiencies are close to unity.

If the window data is block averaged when it is sampled every 10th integration step (for the same size integration step), a plot of the statistical inefficiency in Figure 15 shows that the data is

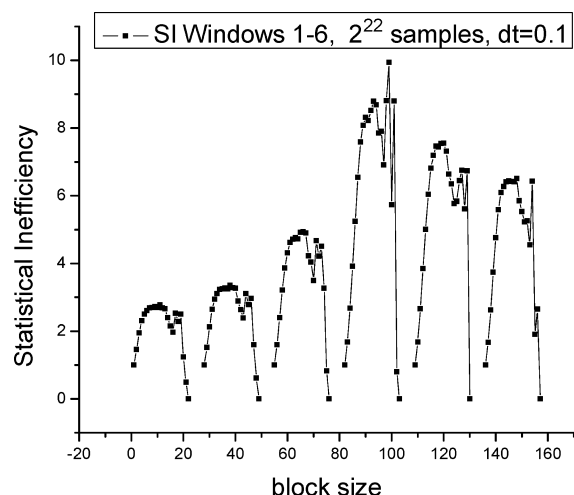


Figure 15. Statistical inefficiency of reaction coordinate trajectories for windows 1–6 as a function of decreasing block size, n_b . The data sampling is every 10th integration step.

correlated. For clarity, the results in Figure 15 are shown only for windows 1–6; by symmetry, the window 7–11 results are similar. Thus, sampling the reaction coordinate at this rate will violate the i.i.d. condition. The analysis in section 4 is carried out using data sampled every 100th integration step.

AUTHOR INFORMATION

Corresponding Author

*E-mail: cukier@chemistry.msu.edu. Phone: 517-355-9715 ext 263.

Notes

The authors declare no competing financial interest.

ACKNOWLEDGMENTS

The author thanks Professor M. Morillo of the Universidad de Sevilla for introducing him to the model and its method of solution used to generate the window trajectories.

REFERENCES

- (1) Kirkwood, J. G. Statistical Mechanics of Fluid Mixtures. *J. Chem. Phys.* **1935**, *3*, 300–313.
- (2) Torrie, G. M.; Valleau, J. P. Monte Carlo Free Energy Estimates Using Non-Boltzmann Sampling: Application to the Sub-Critical Lennard-Jones Fluid. *Chem. Phys. Lett.* **1974**, *28*, 578–581.
- (3) Frenkel, D.; Smit, B. *Understanding Molecular Simulation: From Algorithms to Applications*; Academic: San Diego, 1996.
- (4) Ferrenberg, A. M.; Swendsen, R. H. Optimized Monte-Carlo Data-Analysis. *Phys. Rev. Lett.* **1989**, *63*, 1195–1198.
- (5) Kumar, S.; Bouzida, D.; Swendsen, R. H.; Kollman, P. A.; Rosenberg, J. M. The Weighted Histogram Analysis Method for Free-Energy Calculations on Biomolecules 0.1. The Method. *J. Comput. Chem.* **1992**, *13*, 1011–1021.
- (6) Swendsen, R. H.; Wang, J. S. Replica Monte-Carlo Simulation of Spin-Glasses. *Phys. Rev. Lett.* **1986**, *57*, 2607–2609.
- (7) Souaille, M.; Roux, B. Extension to the Weighted Histogram Analysis Method: Combining Umbrella Sampling with Free Energy Calculations. *Comput. Phys. Commun.* **2001**, *135*, 40–57.
- (8) Sugita, Y.; Okamoto, Y. Replica-Exchange Molecular Dynamics Method for Protein Folding. *Chem. Phys. Lett.* **1999**, *314*, 141–151.
- (9) Woods, C. J.; Essex, J. W.; King, M. A. The Development of Replica-Exchange-Based Free-Energy Methods. *J. Phys. Chem. B* **2003**, *107*, 13703–13710.
- (10) Hansmann, U. H. E.; Okamoto, Y. Numerical Comparisons of Three Recently Proposed Algorithms In the Protein Folding Problem. *J. Comput. Chem.* **1997**, *18*, 920–933.
- (11) Lyubartsev, A.; Laaksonen, A. Parallel Molecular Dynamics Simulations of Biomolecular Systems. *Appl. Parallel Comput.* **1998**, *1541*, 296–303.
- (12) Lyubartsev, A. P.; Martsinovski, A. A.; Shevkunov, S. V.; Vorontsov-Velyaminov, P. N. New Approach to Monte-Carlo Calculation of the Free-Energy - Method of Expanded Ensembles. *J. Chem. Phys.* **1992**, *96*, 1776–1783.
- (13) Chodera, J. D.; Swope, W. C.; Pitner, J. W.; Seok, C.; Dill, K. A. Use of the Weighted Histogram Analysis Method for the Analysis of Simulated and Parallel Tempering Simulations. *J. Chem. Theory Comput.* **2007**, *3*, 26–41.
- (14) Su, L.; Cukier, R. I. Hamiltonian and Distance Replica Exchange Method Studies of Met-Enkephalin. *J. Phys. Chem. B* **2007**, *111*, 12310–12321.
- (15) Lou, H. F.; Cukier, R. I. Molecular Dynamics of Apo-Adenylate Kinase: A Distance Replica Exchange Method for the Free Energy of Conformational Fluctuations. *J. Phys. Chem. B* **2006**, *110*, 24121–24137.
- (16) Efron, B.; Tibshirani, R. *An Introduction to the Bootstrap*; Chapman & Hall: New York, 1994.
- (17) Fasnacht, M.; Swendsen, R. H.; Rosenberg, J. M. Adaptive Integration Method for Monte Carlo Simulations. *Phys. Rev. E* **2004**, *69*, 056704.
- (18) Bennett, C. H. Efficient Estimation of Free-Energy Differences from Monte-Carlo Data. *J. Comput. Phys.* **1976**, *22*, 245–268.
- (19) Shirts, M. R.; Pande, V. S. Comparison of Efficiency and Bias of Free Energies Computed by Exponential Averaging, The Bennett Acceptance Ratio, and Thermodynamic Integration. *J. Chem. Phys.* **2005**, *122*, 144107.
- (20) Kastner, J.; Thiel, W. Bridging the Gap between Thermodynamic Integration and Umbrella Sampling Provides a Novel Analysis Method: “Umbrella Integration”. *J. Chem. Phys.* **2005**, *123*, 144104.
- (21) Kastner, J.; Thiel, W. Analysis of the Statistical Error in Umbrella Sampling Simulations by Umbrella Integration. *J. Chem. Phys.* **2006**, *124*, 234106.

- (22) Bereau, T.; Swendsen, R. H. Optimized Convergence for Multiple Histogram Analysis. *J. Comput. Phys.* **2009**, *228*, 6119–6129.
- (23) Zhu, F. Q.; Hummer, G. Convergence and Error Estimation in Free Energy Calculations Using the Weighted Histogram Analysis Method. *J. Comput. Chem.* **2012**, *33*, 453–465.
- (24) Gallicchio, E.; Andrec, M.; Felts, A. K.; Levy, R. M. Temperature Weighted Histogram Analysis Method, Replica Exchange, and Transition Paths. *J. Phys. Chem. B* **2005**, *109*, 6722–6731.
- (25) Desai, R. C.; Zwanzig, R. Statistical-Mechanics of a Non-Linear Stochastic-Model. *J. Stat. Phys.* **1978**, *19*, 1–24.
- (26) Kometani, K.; Shimizu, H. Study of Self-Organizing Processes of Nonlinear Stochastic Variables. *J. Stat. Phys.* **1975**, *13*, 473–490.
- (27) Cowan, G. *Statistical Data Analysis*; Oxford University Press: Oxford, 1998.
- (28) Cho, E.; Cho, M. Variance of Sample Variance. *Proceedings of the 2008 Joint Statistical Meetings, Section On Survey Research Methods*; American Statistical Association, Washington, DC; pp 1291–1293.
- (29) Muller-Krumbhaar; Binder, K. Dynamic Properties of The Monte Carlo Method In Statistical Mechanics. *J. Stat. Phys.* **1973**, *8*, 1–24.
- (30) Risken, H. *The Fokker-Planck Equation: Methods of Solution and Applications*; Springer-Verlag: Berlin, 1984.
- (31) Zwanzig, R. *Nonequilibrium Statistical Mechanics*; Oxford University: Oxford, 2001.
- (32) Greiner, A.; Strittmatter, W.; Honerkamp, J. Numerical-Integration of Stochastic Differential-Equations. *J. Stat. Phys.* **1988**, *51*, 95–108.
- (33) Matsumoto, M.; Nishimura, T. Mersenne Twister MT19937, <http://www.Math.Sci.Hiroshima-U.Ac.Jp/~M-Mat/MT/Emt.html>.
- (34) Paul, W.; Yoon, D. Y. Stochastic Phase-Space Dynamics with Constraints for Molecular-Systems. *Phys. Rev. E* **1995**, *52*, 2076–2083.
- (35) Amadei, A.; Linssen, A. B. M.; Berendsen, H. J. C. Essential Dynamics of Proteins. *Proteins: Struct., Funct., Genet.* **1993**, *17*, 412–425.
- (36) García, A. E. Large-Amplitude Nonlinear Motions in Proteins. *Phys. Rev. Lett.* **1992**, *68*, 2696–2699.
- (37) Jolliffe, I. T. *Principal Component Analysis*, second ed.; Springer Science: New York, 2004.
- (38) Yang, W.; Bitetti-Putzer, R.; Karplus, M. Free Energy Simulations: Use of Reverse Cumulative Averaging to Determine the Equilibrated Region and the Time Required for Convergence. *J. Chem. Phys.* **2004**, *120*, 2618–2628.
- (39) Flyvbjerg, H.; Petersen, H. G. Error-Estimates on Averages of Correlated Data. *J. Chem. Phys.* **1989**, *91*, 461–466.
- (40) Friedberg, R.; Cameron, J. E. Test Of Monte-Carlo Method - Fast Simulation of a Small Ising Lattice. *J. Chem. Phys.* **1970**, *52*, 6049–6058.
- (41) Allen, M. P.; Tildesley, D. J. *Computer Simulation of Liquids*; Clarendon Press: Oxford, 1987.