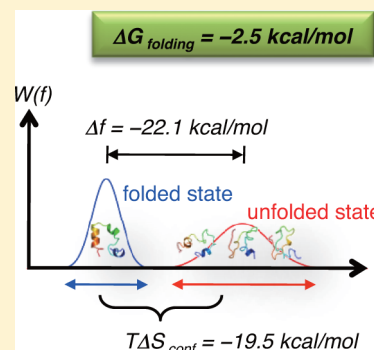


Protein Folding Thermodynamics: A New Computational Approach

Song-Ho Chong and Sihyun Ham*

Department of Chemistry, Sookmyung Women's University, Cheongpa-ro 47-gil 100, Yongsan-Ku, Seoul 140-742, Korea

ABSTRACT: Folding free energy is the fundamental thermodynamic quantity characterizing the stability of a protein. Yet, its accurate determination based on computational techniques remains a challenge in physical chemistry. A straightforward brute-force approach would be to conduct molecular dynamics simulations and to estimate the folding free energy from the equilibrium population ratio of the unfolded and folded states. However, this approach is not sensible at physiological conditions where the equilibrium population ratio is vanishingly small: it is extremely difficult to reliably obtain such a small equilibrium population ratio due to the low rate of folding/unfolding transitions. It is therefore desirable to have a computational method that solely relies on simulations independently carried out for the folded and unfolded states. Here, we present such an approach that focuses on the probability distributions of the effective energy (solvent-averaged protein potential energy) in the folded and unfolded states. We construct these probability distributions for the protein villin headpiece subdomain by performing extensive molecular dynamics simulations and carrying out solvation free energy calculations. We find that the probability distributions of the effective energy are well-described by the Gaussian distributions for both the folded and unfolded states due to the central limit theorem, which enables us to calculate the protein folding free energy in terms of the mean and the width of the distributions. The computed protein folding free energy (-2.5 kcal/mol) is in accord with the experimental result (ranging from -2.3 to -3.2 kcal/mol depending on the experimental methods).



■ INTRODUCTION

Recent advances in computing power and simulation techniques as well as refinements of force fields have developed molecular dynamics simulations into a powerful tool for computational studies of protein folding.^{1–3} Yet, an accurate estimation of folding thermodynamic quantities, in particular, the folding free energy at physiological conditions based on straightforward brute-force simulations, is still hindered by the fact that the population of the unfolded state (C_{unfolded}) at those conditions is vanishingly small relative to that (C_{folded}) of the folded state. Indeed, typical folding free energy for globular proteins at physiological conditions ranges from -5 to -15 kcal/mol,^{4,5} yielding the equilibrium population ratio of $C_{\text{unfolded}}/C_{\text{folded}} \sim 10^{-4}$ to 10^{-11} . It is practically impossible to reliably attain such a small equilibrium population ratio even with the fastest computer available today, except near the melting temperature where $C_{\text{unfolded}}/C_{\text{folded}} \sim 1$ and where the rates for the folding and unfolding transitions are enhanced.⁶

An alternative practical strategy would be to *independently* perform simulations at the folded and unfolded states and then to estimate the folding thermodynamic quantities based on protein conformations sampled in those states. This is based on the fact that the free energy is a state function, and the computation of the free energy difference between two states does not in principle require connecting pathways. In this type of approach, one has to deal with the configuration integral,^{7,8} the potential part of the partition function ($\beta^{-1} = k_B T$ with Boltzmann's constant k_B and temperature T)

$$Z = \int_X \mathrm{d}\mathbf{r} \exp[-\beta f(\mathbf{r})] \quad (X = \text{folded or unfolded}) \quad (1)$$

Here \mathbf{r} denotes the $3N$ dimensional vector representing positions of constituent N atoms, which is restricted to the configuration-space region corresponding to the state X ; $f = E_u + G_{\text{solv}}$ is the solvent-averaged protein potential energy comprising the intraprotein potential energy E_u and the solvation free energy G_{solv} , the latter arising after averaging the protein–water interaction over solvent configurations.^{7,8} The quantity f is also simply referred to as the effective energy and defines a hypersurface in the conformational space which is often called the free energy landscape.⁷

It is, in general, a formidable task to directly evaluate the configuration integral for complex macromolecules such as proteins, and introducing some simplifying approximation is inevitable. One of the most popular approximations is the quasiharmonic method in which the effective energy f is approximated by the quadratic form in terms of the atom-positional fluctuations with respect to an average protein structure.^{9–13} However, the applicability of this method to the unfolded state is questionable where, in contrast to the folded state, an average protein structure is not well-defined. Recently, an energetic approach has been developed that focuses on the fluctuations in the effective energy f instead of the atom-positional fluctuations.⁸ It has been demonstrated that the fluctuations in the effective energy f are well-described by the Gaussian distribution even for intrinsically disordered proteins which do not possess a well-defined three-dimensional

Received: January 9, 2014

Revised: April 29, 2014

Published: April 29, 2014

structure.¹⁴ In this paper, we describe in detail how the knowledge on the distribution functions of the effective energy in the folded and unfolded states can be utilized to estimate the folding thermodynamic quantities. We then apply this method to compute the folding free energy of the protein villin headpiece subdomain to demonstrate its applicability and performance.

THEORY

Our primary interest is in the standard free energy difference $\Delta F^0 = F_A^0 - F_B^0$ between two states A and B (e.g., folded and unfolded states, in which case ΔF^0 is the folding free energy) that determines the equilibrium relative population, denoted as $(C_A/C_B)_{\text{eq}}$, through the relation

$$(C_A/C_B)_{\text{eq}} = e^{-\beta\Delta F^0} \quad (2)$$

The standard free energies F_A^0 and F_B^0 refer to the free energy for the ideal solution of the same reference concentration $C^0 = M/V$, where M and V denote the number of proteins and the volume, respectively. Thus, ΔF^0 can be identified as the difference between the free energy for an ideal solution of M proteins in the state A dissolved in a volume V and the one for an ideal solution of M proteins in the state B dissolved in the same volume V .

Let us take a snapshot of the ideal solution in the state X, where here and in the following X refers to either A or B. This gives us a sample of M protein conformations in the state X characterized by the effective energies f_{i_X} ($i_X = 1, \dots, M$). The partition function for a given sample of the M effective energies is $Z(\{f_{i_X}\}) = \sum_{i_X=1}^M \exp(-\beta f_{i_X})$. Since the free energy is the average of $-k_B T \log Z(\{f_{i_X}\})$ over the samples (snapshots), the free energy difference is given by

$$-\beta\Delta F^0 = \left\langle \log \left(\sum_{i_A=1}^M e^{-\beta f_{i_A}} \right) \right\rangle_M - \left\langle \log \left(\sum_{i_B=1}^M e^{-\beta f_{i_B}} \right) \right\rangle_M \quad (3)$$

where $\langle \dots \rangle_M$ refers to the average over samples. For each sample of M effective energies $\{f_{i_X}\}$, we introduce the density of states $n_X(f)$, that is, the number of effective energies belonging to the interval $(f, f + df)$, to write

$$-\beta\Delta F^0 = \left\langle \log \int df n_A(f) e^{-\beta f} \right\rangle_M - \left\langle \log \int df n_B(f) e^{-\beta f} \right\rangle_M \quad (4)$$

At temperatures where the system is not frozen in local free energy minima (which is the case for proteins at physiological temperatures), the “quenched” average, $\langle \log \int df n_X(f) e^{-\beta f} \rangle_M$, is equal to the “annealed” average, $\log \int df \langle n_X(f) \rangle_M e^{-\beta f}$,¹⁵ yielding

$$-\beta\Delta F^0 = \log \int df \langle n_A(f) \rangle_M e^{-\beta f} - \log \int df \langle n_B(f) \rangle_M e^{-\beta f} \quad (5)$$

Since $\int df \langle n_X(f) \rangle_M = M$ by construction, one obtains after introducing the normalized probability distribution function via $W_X(f) = \langle n_X(f) \rangle_M / M$

$$-\beta\Delta F^0 = \log \int df W_A(f) e^{-\beta f} - \log \int df W_B(f) e^{-\beta f} \quad (6)$$

Thus, the explicit dependence on M drops out, and it suffices to deal with the normalized distribution function $W_X(f)$ in place of the density of states in determining the standard free energy difference.

Here, we introduce the central assumption in the present approach; that is, we assume that the distribution function $W_X(f)$ is well-approximated by the Gaussian distribution

$$W_X(f) = \frac{1}{\sqrt{2\pi\sigma_{f,X}^2}} \exp \left[-\frac{1}{2\sigma_{f,X}^2} (f - \bar{f}_X)^2 \right] \quad (7)$$

where \bar{f}_X and $\sigma_{f,X}^2$ denote the mean and the standard deviation of f in the state X, respectively. Substituting this into eq 6 yields $\Delta F^0 = (\bar{f}_A - \bar{f}_B) - \beta(\sigma_{f,A}^2 - \sigma_{f,B}^2)/2$, which we rewrite as

$$\Delta F^0 = F_A^0 - F_B^0 = \Delta \bar{f} - T\Delta S_{\text{conf}} \quad (8)$$

with

$$F_X^0 = \bar{f}_X - TS_{\text{conf},X} \quad (9)$$

by identifying $TS_{\text{conf},X} = (\beta/2)\sigma_{f,X}^2$ as the protein conformational entropy in the state X.⁸ Thus, eq 8 expresses the free energy difference ΔF^0 in terms of the difference in the well depth of the free energy landscape ($\Delta \bar{f}$) and of the difference in the extent the system explores the free energy landscape in the two states ($T\Delta S_{\text{conf}}$).

In our previous work,⁸ eqs 8 and 9 were derived by first rewriting the configuration integral (eq 1) in terms of the normalized distribution function $W_X(f)$, that is, starting from the expression $Z = \int df W_X(f) e^{-\beta f} = \overline{e^{-\beta f}}$ in which the bar denotes an average with respect to $W_X(f)$, and then by truncating the cumulant expansion of $\log \overline{e^{-\beta f}}$ at the second order, which is equivalent to assuming the Gaussian form for $W_X(f)$ (eq 7). However, the use of the normalized distribution function $W_X(f)$ in place of the (non-normalized) density of states in rewriting the configuration integral was not justified there. In the present work, the rationalization of the use of the normalized distribution function $W_X(f)$ in determining the standard free energy difference is demonstrated explicitly through the derivation of eq 6.

The approach presented here is termed the energetic approach since it is based on the statistical distribution of the effective energy $f = E_u + G_{\text{solv}}$. A major complication in this approach is that one needs to evaluate the effective energy for a sufficient number of protein conformations in order to construct the distribution function $W(f)$ whose Gaussianity can be examined. This can be done, for example, by combining molecular dynamics simulations, which generate protein conformations as well as the intraprotein energy E_u and the integral equation theory that enables the calculation of the solvation free energy G_{solv} for each of the simulated conformations (see below). This energetic approach is intimately related to the random energy model in which the probability distribution of energy is assumed to be Gaussian and the resulting entropy is expressed in terms of the energy fluctuations.^{15,16} Finally, we note that the Helmholtz free energy is essentially equal to the Gibbs free energy since the pressure–volume term is usually negligible under physiological conditions.⁷ Therefore, eqs 8 and 9 in which F (the Helmholtz free energy) is replaced by G (the Gibbs free energy) can be used in practical applications.

■ COMPUTATIONAL METHODS

We implement our energetic approach by combining the molecular dynamics simulations and the solvation free energy calculations detailed below, which will be applied to compute the folding free energy of the protein villin headpiece subdomain (HP-36).

Molecular Dynamics Simulations. Folded-State Simulations. The initial coordinates for HP-36 were taken from the NMR structure (PDB ID: 1VII).¹⁷ All-atom, explicit-water molecular dynamics simulations under neutral pH were performed with AMBER11¹⁸ using the ff99SB force field¹⁹ for protein and the TIP3P model²⁰ for water. HP-36 was placed into a cubic periodic box of the side length ~ 61 Å containing 6693 water molecules and 2 counter Cl^- ion. The system was minimized by 500 steps of steepest descent minimization and 500 steps of conjugate gradient minimization under 500 kcal/(mol Å²) harmonic restraints. This was followed by 1000 steps of steepest descent minimization and 1500 steps of conjugate gradient minimization without harmonic restraints. Canonical (NVT) ensemble simulation was then carried out for 20 ps, during which the system was gradually heated from $T = 0$ to 300 K. Subsequently, we performed constant-pressure (NPT) ensemble equilibration simulation at $T = 300$ K and $P = 1$ bar for 200 ps. We then carried out a 1 μs production run at $T = 300$ K and $P = 1$ bar. Three independent simulations were performed with different random initial velocities. The particle mesh Ewald method²¹ was used to treat long-range electrostatic interactions, whereas short-range nonbonded interactions were cut off at 10 Å. Bond lengths involving bonds to hydrogen atoms were constrained using the SHAKE algorithm.²² The time step for all the simulations was 2 fs. We used Berendsen's thermostat and barostat to control temperature and pressure with coupling constants of 1.0 and 2.0 ps, respectively.²³

Unfolded-State Simulations. Unfolded-state simulations were carried out as follows starting from heat-denatured structures.²⁴ We first heated the system after the 200 ps NPT ensemble simulation at $T = 300$ K and $P = 1$ bar mentioned above to $T = 600$ K with a constant volume and then conducted a 20 ns NVT ensemble simulation. This was followed by annealing simulations to $T = 300$ K with a 50 K interval, and we performed a 1 ns NVT ensemble simulation at each of the intervening temperatures. Subsequently, a 5 ns NPT ensemble equilibration simulation was conducted at $T = 300$ K and $P = 1$ bar. Finally, we carried out a 5 μs production run at $T = 300$ K and $P = 1$ bar. The whole procedures were repeated 10 times with different random initial velocities to generate 10 independent unfolded-state trajectories. The last 1 μs trajectory (4–5 μs) from each of the 10 independent 5 μs production runs was subjected to structural analyses and solvation free energy calculations.

Structural Analyses. From each of the three folded-state trajectories and the 10 unfolded-state trajectories of 1 μs length, we took 100 000 protein conformations with a 10 ps time interval. A hydrophobic contact is considered to be formed between a pair of residues with a sequence separation equal to or larger than four residues if the minimum distance between heavy atoms in the side-chain groups is smaller than 5.4 Å, and this was applied to the residues forming the hydrophobic core of HP-36 (hydrophobic residues Phe7, Val10, Phe11, Ala17, Phe18, Leu29, and nonpolar necks of Lys25 and Gln26).¹⁷ We calculated the secondary structure contents by using the DSSP program²⁵ in the AMBER11 distribution.¹⁸ The root mean

square deviation (rmsd)-based K -means clustering analysis was carried out to obtain representative protein conformations. We made use of the MMTSB toolset²⁶ for this analysis with the cutoff value 4.0 Å for C_α and C_β atoms in the protein.

Solvation Free Energy Calculations. To each of the simulated protein conformations, we applied the three-dimensional reference interaction site model (3D-RISM) theory^{27,28} to compute the solvation free energy. In this theory, the 3D distribution function $g_\gamma(\mathbf{r})$ of the water site γ at position \mathbf{r} around a protein is obtained by self-consistently solving the 3D-RISM equation

$$h_\gamma(\mathbf{r}) = \sum_{\gamma'} \int d\mathbf{r}' \chi_{\gamma\gamma'}(|\mathbf{r} - \mathbf{r}'|) c_{\gamma'}(\mathbf{r}') \quad (10)$$

and the closure relation

$$h_\gamma(\mathbf{r}) = \begin{cases} \exp[d_\gamma(\mathbf{r})] - 1 & \text{for } d_\gamma(\mathbf{r}) \leq 0 \\ d_\gamma(\mathbf{r}) & \text{for } d_\gamma(\mathbf{r}) > 0 \end{cases} \quad (11)$$

where $d_\gamma(\mathbf{r}) = -u_\gamma(\mathbf{r})/(k_B T) + h_\gamma(\mathbf{r}) - c_\gamma(\mathbf{r})$. Here $h_\gamma(\mathbf{r}) = g_\gamma(\mathbf{r}) - 1$ refers to the 3D total correlation function of the water site γ ; $c_\gamma(\mathbf{r})$ is the corresponding direct correlation function; $\chi_{\gamma\gamma'}(\mathbf{r})$ denotes the site–site water susceptibility function, treated as an input to the theory; and $u_\gamma(\mathbf{r})$ is the interaction potential generated by protein atoms. The same numerical procedure as described in ref 28 was employed to solve the above equations along with the susceptibility function for the TIP3P water model calculated from the dielectrically consistent RISM theory.²⁹ One can then compute the solvation free energy G_{solv} using the following analytical formula:^{27,28}

$$G_{\text{solv}} = \rho k_B T \sum_\gamma \int d\mathbf{r} \left[\frac{1}{2} h_\gamma(\mathbf{r})^2 \Theta(-h_\gamma(\mathbf{r})) - c_\gamma(\mathbf{r}) - \frac{1}{2} h_\gamma(\mathbf{r}) c_\gamma(\mathbf{r}) \right] \quad (12)$$

Here $\Theta(x)$ is the Heaviside step function, and ρ is the number density of water.

Because of the approximate nature of the closure relation, it is inevitable that the absolute value of the solvation free energy computed from the 3D-RISM theory depends on the closure relation used.^{27,28} On the other hand, relative values of the solvation free energy are reasonably accurate due to the cancellation of errors.²⁷ Since only relative values of the solvation free energy enter into the standard free energy difference (eq 8), we expect that our results to be presented below do not significantly suffer from such an inherent limitation of the 3D-RISM theory.

Statistical Analyses. Measures of Gaussianity. Our central assumption is the Gaussianity of the distribution function $W(f)$. We employed the skewness and excess kurtosis as measures of the Gaussianity, which are respectively defined by $\mu_3/\mu_2^{3/2}$ and $\mu_4/\mu_2^2 - 3$ in terms of the n th moment about the mean, $\mu_n = \int df (f - \bar{f})^n W(f)$.³⁰ The skewness is a dimensionless parameter characterizing the degree of asymmetry of a distribution around its mean, and the excess kurtosis is a dimensionless quantity measuring the peakedness or flatness of a distribution relative to a Gaussian distribution, both of which are zero if the distribution function is Gaussian.

Block-Averaging Method. We used the block-averaging approach³¹ for the error estimation of the mean effective energy \bar{f} , conformational entropy $TS_{\text{conf}} = (\beta/2)\sigma_f^2$, and $G^0 = \bar{f} - TS_{\text{conf}}$

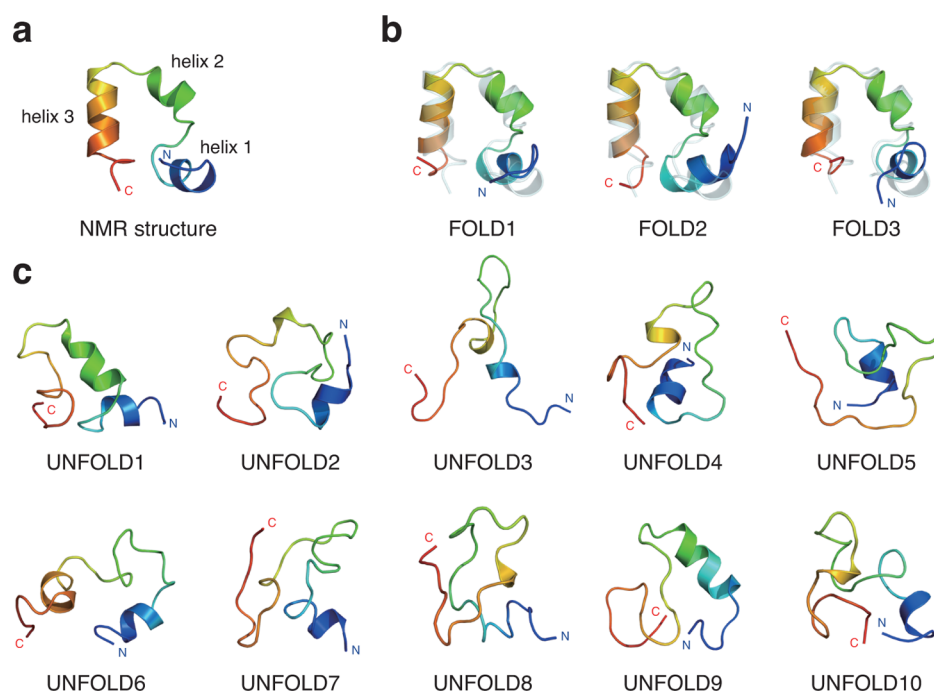


Figure 1. (a) NMR structure of HP-36 (PDB ID: 1VII).¹⁷ The protein structure is color-coded according to the sequence, ranging from blue to red at the N- and C-termini, respectively. The locations of the helix 1 (residues 4–8), helix 2 (15–18), and helix 3 (23–32) are also indicated. (b) Representative folded-state conformations from the trajectories FOLD1–FOLD3. The NMR structure is shown as a transparent gray cartoon for comparison. (c) Representative unfolded-state conformations from the trajectories UNFOLD1–UNFOLD10.

Table 1. Structural Characteristics

	C_{α} rmsd (Å) ^a	R_g (Å) ^b	native hydrophobic core contacts (%) ^c	helical contents (%) ^d		
				helix 1	helix 2	helix 3
NMR structure PDB ID: 1VII	0	9.5	100	100	100	100
folded-state trajectories						
FOLD1	1.9 ± 0.4	10.1 ± 0.3	90.9 ± 6.8	87.7	98.8	98.7
FOLD2	2.3 ± 0.6	10.3 ± 0.5	90.0 ± 8.7	80.5	99.0	86.3
FOLD3	2.0 ± 0.4	10.1 ± 0.3	93.0 ± 6.9	84.8	98.8	98.0
unfolded-state trajectories						
UNFOLD1	6.5 ± 0.9	10.5 ± 0.4	23.9 ± 7.4	39.8	28.3	12.1
UNFOLD2	6.3 ± 0.3	9.9 ± 0.3	16.1 ± 7.2	57.7	11.4	19.5
UNFOLD3	7.5 ± 0.4	10.6 ± 0.7	12.5 ± 1.5	20.3	0.0	17.1
UNFOLD4	7.2 ± 0.6	10.5 ± 1.1	21.9 ± 6.4	74.0	6.5	38.0
UNFOLD5	7.5 ± 0.9	11.7 ± 1.1	13.1 ± 8.1	60.7	0.0	0.3
UNFOLD6	7.7 ± 1.0	10.8 ± 0.8	26.8 ± 15.1	32.5	12.1	65.4
UNFOLD7	7.8 ± 0.2	9.8 ± 0.1	11.2 ± 3.9	44.3	14.8	0.0
UNFOLD8	6.9 ± 0.6	10.5 ± 0.7	13.1 ± 10.1	9.4	0.3	22.4
UNFOLD9	6.8 ± 0.6	10.2 ± 0.7	20.8 ± 12.7	34.6	10.8	10.7
UNFOLD10	6.2 ± 0.3	10.0 ± 0.4	23.3 ± 5.7	50.7	0.0	28.2

^aAverage ± standard deviation of C_{α} rmsd (excluding the terminal residues) relative to the NMR structure. ^bAverage ± standard deviation of radius of gyration. ^cAverage ± standard deviation of contents of the native hydrophobic core contacts defined by the NMR structure. ^dAverage helical contents of helix 1 (residues 4–8), helix 2 (15–18), and helix 3 (23–32).

which are denoted as θ in the following. In this approach, a simulation trajectory of length $N_{\text{sim}} = N_b \times n$ is divided into N_b blocks of length n . The average of an observable θ is calculated for each block, yielding N_b values for $\bar{\theta}_i$ with $i = 1, \dots, N_b$. For each value of n , one obtains the blocked standard error (ste) via

$$\sigma_{\text{ste}}^{\text{block}(n)}(\theta) = \left(\frac{1}{N_b(N_b - 1)} \sum_{i=1}^{N_b} [\bar{\theta}_i - \bar{\theta}]^2 \right)^{1/2} \quad (13)$$

where $\bar{\theta} = (1/N_b) \sum_{i=1}^{N_b} \bar{\theta}_i$. The standard error for θ can be estimated from the blocked standard error $\sigma_{\text{ste}}^{\text{block}(n)}(\theta)$ for large n ,³¹ where it is expected that different blocks become statistically independent and $\sigma_{\text{ste}}^{\text{block}(n)}(\theta)$ ceases to vary with n .

Bootstrap Method. For large n , the number of blocks $N_b = N_{\text{sim}}/n$ becomes small for a given length N_{sim} of a simulation trajectory. To improve the error estimate for large n , we applied the bootstrap method³² to the blocked sample $\chi = \{\bar{\theta}_1, \bar{\theta}_2, \dots, \bar{\theta}_{N_b}\}$. In this method, a new sample, called a bootstrap sample, of the same size N_b is generated from χ by sampling

folded state

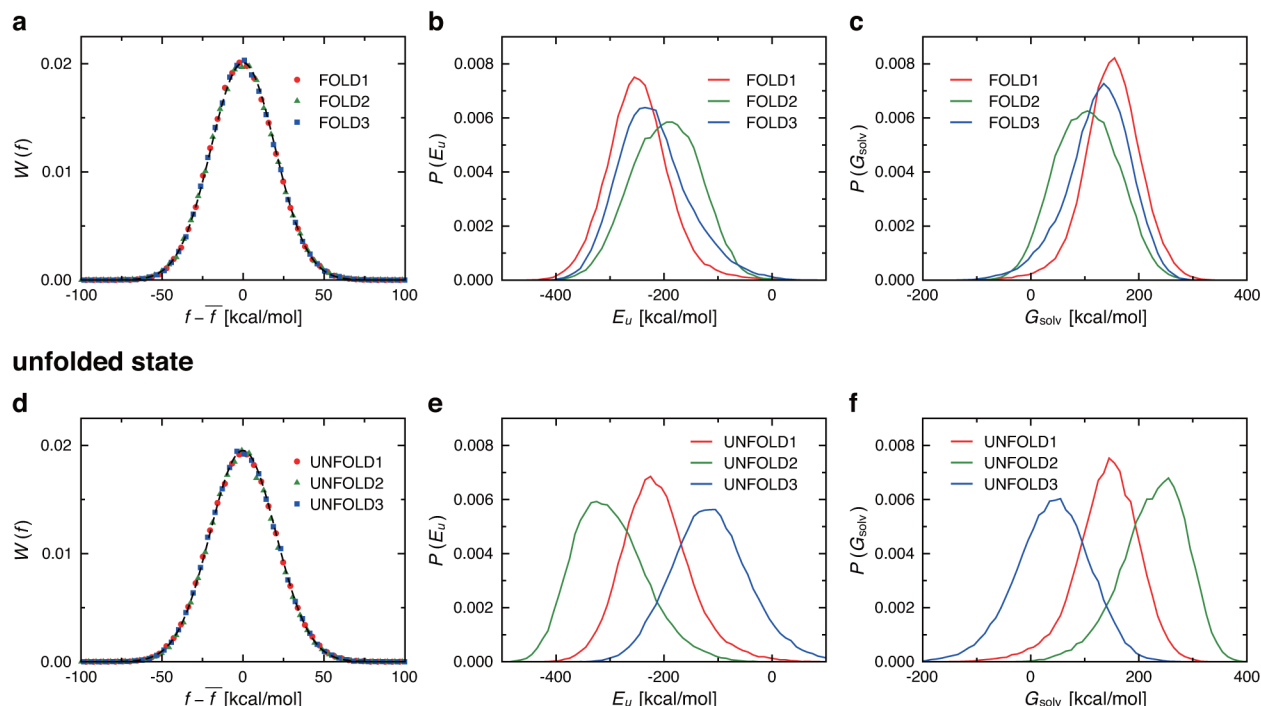


Figure 2. (a) Probability distribution function $W(f)$ of the effective energy f as a function of the deviation $f - \bar{f}$ from the mean value \bar{f} for the three independent folded-state trajectories specified by different colors (FOLD1 by red, FOLD2 by green, FOLD3 by blue). Dashed curve denotes the fit by the Gaussian distribution. (b) Probability distribution function $P(E_u)$ of the intraprotein energy E_u for the folded-state trajectories. (c) Probability distribution function $P(G_{solv})$ of the solvation free energy G_{solv} for the folded-state trajectories. (d–f) Corresponding results for representative unfolded-state trajectories (UNFOLD1–UNFOLD3).

with replacement using a random number generator, $\chi^* = \{\theta_1^*, \theta_2^*, \dots, \theta_{N_b}^*\}$, and then the average is computed, $\bar{\theta}^* = (1/N_b) \sum_{i=1}^{N_b} \theta_i^*$. This is repeated to generate a large number of bootstrap samples and their averages, $\chi_{(b)}^*$ and $\bar{\theta}_{(b)}^*$ with $b = 1, 2, \dots, B$. Finally, one computes the standard deviation of the $\bar{\theta}_{(b)}^*$ values, which is the bootstrap estimate of the standard error³²

$$\sigma_{ste}^{boot(n)}(\theta) = \left(\frac{1}{B-1} \sum_{b=1}^B [\bar{\theta}_{(b)}^* - \bar{\theta}^{boot(n)}]^2 \right)^{1/2} \quad (14)$$

where $\bar{\theta}^{boot(n)} = (1/B) \sum_{b=1}^B \bar{\theta}_{(b)}^*$. We used $B = 100$ for the bootstrapping.

Propagation of Errors. The standard error (σ_{ste}) for the difference $\Delta\theta = \theta_{folded} - \theta_{unfolded}$ of an observable θ in the folded and unfolded states can be evaluated from³³

$$\sigma_{ste}(\Delta\theta) = \sqrt{\sigma_{ste}(\theta_{folded})^2 + \sigma_{ste}(\theta_{unfolded})^2} \quad (15)$$

RESULTS

Folded-State Simulations. We illustrate the applicability and performance of the energetic approach by computing the folding free energy of the protein villin headpiece subdomain (HP-36) at the physiological condition ($T = 300$ K and $P = 1$ bar). HP-36 is the shortest naturally occurring, thermostable helical protein that autonomously folds on the microsecond time scales¹⁷ and has served as an excellent model system for studying protein folding.^{24,34–40} (Previous studies include the ones on HP-35 in which the N-terminal methionine is absent,

but the physical properties of HP-35 and HP-36 are similar.⁴¹) The NMR-derived folded structure of HP-36 (PDB ID: 1VII)¹⁷ contains three short helices surrounding a tightly packed hydrophobic core (Figure 1a).

We carried out three independent folded-state simulations of 1 μ s length each, which will be referred to as FOLD1–FOLD3 in the following. Representative protein structures in these trajectories are shown in Figure 1b. The NMR structure was stable over the time scales of the simulations: the average C_α rmsd value relative to the NMR structure (excluding the terminal residues) is ~ 2 Å; the radius of gyration (R_g) is close to that of the NMR structure, and the average contents of the native hydrophobic core contacts and of the surrounding three short helices as observed in the NMR structure are typically $\sim 90\%$ or larger (Table 1).

Unfolded-State Simulations. It is a nontrivial problem to generate unfolded protein conformations at physiological conditions since those conformations are practically unattainable just by further continuing the folded-state simulations. We initiated an unfolded-state simulation starting from a heat-denatured protein conformation,²⁴ that is, by performing a high-temperature unfolding simulation from the NMR structure. After conducting short annealing simulations, we carried out a 5 μ s production run at $T = 300$ K and $P = 1$ bar. The whole procedures were repeated to generate 10 independent 5 μ s trajectories, corresponding to an aggregated simulation time of 50 μ s. We did not observe the folding of HP-36 in either of these trajectories over the time scales of the simulations.

In the present study, the protein conformations taken from the last 1 μ s trajectory (4–5 μ s) from each 5 μ s production run

were considered to belong to the unfolded state at $T = 300$ K, regarding the earlier $4 \mu\text{s}$ trajectory as a “relaxation” process from the high- T unfolded state to the $T = 300$ K unfolded state. These 10 independent unfolded-state trajectories of $1 \mu\text{s}$ length each will be referred to as UNFOLD1–UNFOLD10 in the following. Representative protein structures in these trajectories are displayed in Figure 1c. These unfolded structures are characterized by large C_α rmsd values (~ 7 Å) to the NMR structure, barely (~ 10 to 20%) formed hydrophobic core contacts as observed in the NMR structure, and only partially (~ 0 to 70% depending on the helices and on the trajectories) formed secondary structures (Table 1).

Distribution of the Effective Energy. From each of the $1 \mu\text{s}$ length trajectories FOLD1–FOLD3 and UNFOLD1–UNFOLD10, we took 100 000 protein conformations with a 10 ps time interval. The intraprotein energies (E_u) for these conformations were directly computed during the simulations. The solvation free energies (G_{solv}) were obtained by applying the 3D-RISM theory to the simulated protein conformations. These data for E_u and G_{solv} were then combined to construct the distribution function $W(f)$ of the effective energy $f = E_u + G_{\text{solv}}$. We find that $W(f)$ curves for both the folded state (Figure 2a) and the unfolded state (Figure 2d) are very close to Gaussian. Indeed, the skewness as well as the excess kurtosis of $W(f)$ are quite small for all the folded- and unfolded-state trajectories (Table 2). This holds irrespective of the fact that

Table 2. Skewness and Excess Kurtosis of the Distribution Function of Effective Energy

	skewness	excess kurtosis
folded-state trajectories		
FOLD1	0.074	−0.007
FOLD2	0.058	0.033
FOLD3	0.063	0.034
unfolded-state trajectories		
UNFOLD1	0.049	−0.006
UNFOLD2	0.055	0.004
UNFOLD3	0.047	−0.025
UNFOLD4	0.062	0.013
UNFOLD5	0.032	0.007
UNFOLD6	0.060	0.004
UNFOLD7	0.052	−0.014
UNFOLD8	0.036	0.014
UNFOLD9	0.051	0.023
UNFOLD10	0.064	−0.011

the individual distributions for E_u and G_{solv} are rather heterogeneous and that their shapes differ in different trajectories (Figure 2b,c for the folded-state trajectories and Figure 2e,f for the unfolded-state trajectories).

The Gaussianity of $W(f)$ is essentially due to the central limit theorem because the Hamiltonian is a pairwise sum of many similar energy terms, many of them canceling each other. (The relevance of the central limit theorem for protein fluctuations has been recognized before.⁴²) In particular, the cancellation between the intraprotein interaction (E_u) and the protein–water interaction (G_{solv}) is prominent, which results in their probability distributions $P(E_u)$ and $P(G_{\text{solv}})$ being almost a mirror image of each other. (Notice from Figure 2b,c for the folded-state trajectories that the trajectory whose $P(E_u)$ is mainly populated in the lower E_u region has $P(G_{\text{solv}})$ whose peak is located in the higher G_{solv} region, and vice versa. Such

an anticorrelation between $P(E_u)$ and $P(G_{\text{solv}})$ is seen also in Figure 2e,f for the unfolded-state trajectories.) As we discussed previously,⁴³ this reflects the competition between the intraprotein interaction and the protein–water interaction. For example, when a hydrogen bond in the protein is broken and is replaced by hydrogen bond(s) with water, E_u increases whereas G_{solv} decreases, and vice versa. Thus, the large cancellation between E_u and G_{solv} leads to much milder Gaussian fluctuations in $f = E_u + G_{\text{solv}}$.

Folding Thermodynamics. The Gaussian nature of $W(f)$ curves for both the folded and unfolded states of HP-36 enables us to compute the protein folding free energy in terms of the mean \bar{f} and the width $TS_{\text{conf}} = (\beta/2)\sigma_f^2$ of the distributions (eqs 8 and 9). The \bar{f} and TS_{conf} as well as the resulting $G^0 = \bar{f} - TS_{\text{conf}}$ for all the folded- and unfolded-state trajectories are summarized in Table 3. By averaging over the independent

Table 3. Folding Thermodynamic Quantities

	\bar{f} (kcal/mol)	TS_{conf} (kcal/mol)	$G^0 = \bar{f} - TS_{\text{conf}}$ (kcal/mol)
folded-state trajectories			
FOLD1	−97.4	326.8	−424.2
FOLD2	−93.0	335.7	−428.7
FOLD3	−95.7	330.0	−425.7
average ^a	−95.4 ± 1.0	330.8 ± 2.0	−426.2 ± 1.0
unfolded-state trajectories			
UNFOLD1	−72.7	356.7	−429.4
UNFOLD2	−72.3	343.0	−415.3
UNFOLD3	−73.5	350.3	−423.8
UNFOLD4	−78.1	347.9	−426.0
UNFOLD5	−68.4	349.8	−418.2
UNFOLD6	−80.8	352.0	−432.8
UNFOLD7	−76.4	349.1	−425.5
UNFOLD8	−64.7	356.9	−421.6
UNFOLD9	−66.2	350.4	−416.6
UNFOLD10	−80.3	347.1	−427.4
average ^a	−73.3 ± 1.7	350.3 ± 1.2	−423.7 ± 1.7
	$\Delta\bar{f}$ (kcal/mol)	$T\Delta S_{\text{conf}}$ (kcal/mol)	$\Delta G^0 = \Delta\bar{f} - T\Delta S_{\text{conf}}$ (kcal/mol)
difference ^b	−22.1 ± 1.9	−19.5 ± 2.3	−2.5 ± 2.0

^aAverage ± standard error estimated with the bootstrap analysis (see text). ^bDifference between the averages for the folded and unfolded states ± the standard error estimated by eq 15.

trajectories, we obtain the folding free energy of $\Delta G^0 = -2.5$ kcal/mol, which results from a large cancellation between a favorable decrease in the effective energy, $\Delta\bar{f} = -22.1$ kcal/mol, and an unfavorable change in the protein conformational entropy, $T\Delta S_{\text{conf}} = -19.5$ kcal/mol.

To estimate standard errors in these average quantities, we performed the block averaging and bootstrap analyses (Figure 3). We observe that the error estimation for large block lengths, where the number of blocks is small, is improved with the bootstrapping. We also find that the convergence of the standard error for TS_{conf} and $G^0 = \bar{f} - TS_{\text{conf}}$ is not as good as the one for \bar{f} . Since $TS_{\text{conf}} = (\beta/2)\sigma_f^2$ is determined by the fluctuations in f , this implies that the correlation time for the fluctuations is much larger than the one for the average,³¹ which is in agreement with the previous investigation on specific heat and susceptibility.⁴⁴ The “best” estimates for the standard errors obtained with the block length of $1 \mu\text{s}$ and with the bootstrapping are presented in Table 3. The computed

folded state

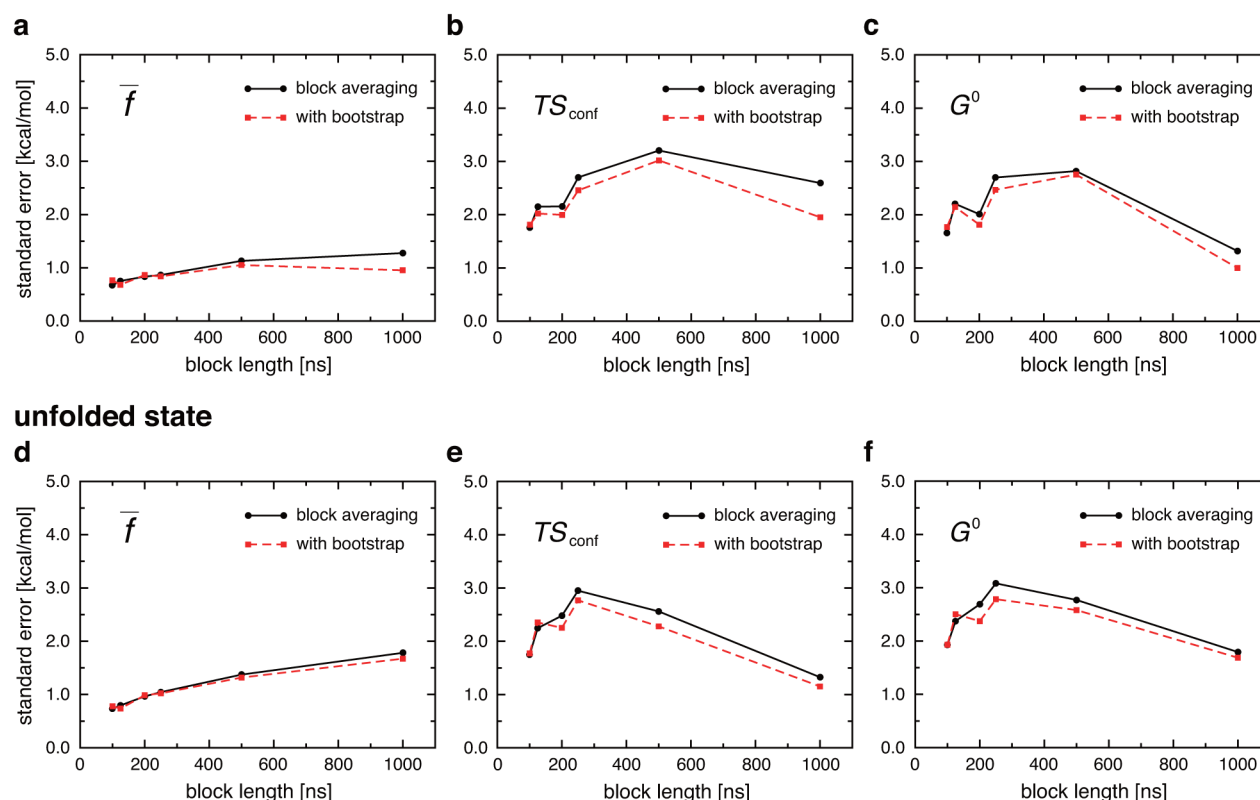


Figure 3. Standard error for (a) the mean effective energy \bar{f} , (b) protein conformational entropy TS_{conf} , and (c) $G^0 = \bar{f} - TS_{\text{conf}}$ for the folded state estimated from the block averaging (black circles) and with the bootstrap (red squares) for the block lengths of 100, 125, 200, 250, 500, and 1000 ns. (d–f) Corresponding results for the unfolded state.

folding free energy, $\Delta G^0 = -2.5$ kcal/mol, with the estimated standard error of 2.0 kcal/mol is in accord with the experimental result that ranges from -2.3 to -3.2 kcal/mol depending on the experimental methods: $\Delta G^0 = -3.1$ kcal/mol from the temperature-jump experiment,³⁴ -2.3 kcal/mol from the triplet-lifetime experiment,³⁵ -3.2 kcal/mol from the urea denaturation,³⁶ -2.4 kcal/mol from the guanidinium chloride denaturation,³⁹ and -2.9 kcal/mol from the single-molecule force spectroscopy.³⁹

DISCUSSION

Experimental folding free energy of HP-36 at physiological conditions (taken to be -3.0 kcal/mol) translates into that for 99.4% of the time the system is found in the folded state and for 0.6% of the time the system is found in the unfolded state. Such an equilibrium population ratio is not realized in our microsecond time scale simulations initiated either from the NMR structure (100% folded state) or from heat-denatured structures (100% unfolded state) due to the low rate of folding/unfolding transitions. The strong point of our energetic approach is that it nevertheless enables us to compute and analyze folding thermodynamic quantities based on such “short” simulation trajectories. Furthermore, the computed folding free energy for HP-36 is found to be in the range of experimental measurements. Of course, our numerical results might strongly depend on the choice of the force field, and this issue remains to be examined. In fact, while all the recent force fields were shown to produce a satisfactory picture of protein folding from a structural and kinetic standpoint, some

differences were observed at the level of the folding thermodynamics.⁴⁵

It is well-accepted that small values of the folding free energy result from a large cancellation between the energetic and entropic contributions.⁴⁶ Indeed, we find for HP-36 that the change in the effective energy upon folding, $\Delta \bar{f} = -22.1$ kcal/mol (Table 3), is largely offset by the decrease in the protein conformational entropy, $T\Delta S_{\text{conf}} = -19.5$ kcal/mol, yielding a small value of $\Delta G^0 = -2.5$ kcal/mol. On the other hand, due to such a large cancellation, the standard error (-2.0 kcal/mol) for the folding free energy is relatively enlarged. (Notice from Table 3 that the relative standard errors for the quantities in each of the folded and unfolded states are quite small. For example, the standard error (2.0 kcal/mol) for TS_{conf} (330.8 kcal/mol) of the folded state is only 0.6%, and the one (1.2 kcal/mol) for TS_{conf} (350.3 kcal/mol) of the unfolded state is only 0.3%.) Previous computational works on the folding thermodynamics focused on $\Delta \bar{f}$ by either ignoring $T\Delta S_{\text{conf}}$ or estimating it with empirical models.^{47–50} The significance of our approach is thus highlighted by that both $\Delta \bar{f}$ and $T\Delta S_{\text{conf}}$ can be computed on an equal footing through the distribution function $W(f)$. Our numerical result for $T\Delta S_{\text{conf}}$ (1.8 in units of cal/(mol K) per residue) is also within the range of empirical estimates (1 to 12 cal/(mol K) per residue).^{51–53}

Our computational results are in accord with the folding landscape picture^{54,55} in that the folded state is located in a much deeper minimum (characterized by the large negative $\Delta \bar{f}$) and that the extent the system explores the free energy landscape significantly decreases upon folding (reflected in the large negative $T\Delta S_{\text{conf}}$). We emphasize here a critical role of

water in determining these characteristics. For example, the protein energy (E_u) alone does not serve as a good indicator of the folded state since the average E_u for the unfolded state can be lower than that for the folded state (compare UNFOLD2 curve in Figure 2e with FOLD1–3 curves in Figure 2b). Only after considering the effective energy $f = E_u + G_{\text{solv}}$ that takes into account the solvation effect (G_{solv}) can the folded state be well-distinguished from the unfolded state and be characterized by the large negative value of Δf . (Notice from Table 3 that \bar{f} for the folded state is lower than the one for the unfolded state irrespective of the trajectories.) This confirms the well-known fact that solvent water must be handled as an integral part of biomolecules.^{56,57}

CONCLUSIONS

Developing accurate and sensible computational methods for protein folding thermodynamics is of central importance in physical chemistry. Here, we present a novel computational approach that links protein folding free energy with the probability distributions of the effective energy (solvent-averaged protein potential energy) in the folded and unfolded states. We illustrate the applicability and performance of our approach by constructing these probability distributions for the protein villin headpiece subdomain (HP-36) after conducting extensive molecular dynamics simulations and solvation free energy calculations. We find that the probability distributions of the effective energy are well-described by the Gaussian distributions for both the folded and unfolded states due to the central limit theorem, which enables us to calculate protein folding free energy in terms of the mean and the width of the distributions. The computed protein folding free energy of HP-36 (−2.5 kcal/mol) is in accord with the experimental result (ranging from −2.3 to −3.2 kcal/mol depending on the experimental methods). The novel computational approach presented here will be valuable for understanding the thermodynamics of protein folding, predicting the stability of biomolecules, and designing of biotherapeutics.

AUTHOR INFORMATION

Corresponding Author

*E-mail: sihyun@sookmyung.ac.kr. Phone: +82 2 710 9410. Fax: +82 2 2077 7321.

Notes

The authors declare no competing financial interest.

ACKNOWLEDGMENTS

This work was supported by the Basic Science Research Program through the National Research Foundation of Korea (NRF) funded by the Ministry of Education, Science and Technology (Grant Nos. 2012-0007855, 2012-0003068, and 2012R1A2A01004687). The authors would like to acknowledge the support from KISTI supercomputing center.

REFERENCES

- Freddolino, P. L.; Harrison, C. B.; Liu, Y.; Schulten, K. Challenges in Protein-Folding Simulations. *Nat. Phys.* **2010**, *6*, 751–758.
- Lindorff-Larsen, K.; Piana, S.; Dror, R. O.; Shaw, D. E. How Fast-Folding Proteins Fold. *Science* **2011**, *334*, 517–520.
- Lane, T. J.; Shukla, D.; Beauchamp, K. A.; Pande, V. S. To Milliseconds and Beyond: Challenges in the Simulation of Protein Folding. *Curr. Opin. Struct. Biol.* **2013**, *23*, 58–65.
- Pace, C. N. Conformational Stability of Globular Proteins. *Trends Biochem. Sci.* **1990**, *15*, 14–17.
- Makhatadze, G. I.; Privalov, P. L. Energetics of Protein Structure. *Adv. Protein Chem.* **1995**, *47*, 307–425.
- Piana, S.; Lindorff-Larsen, K.; Shaw, D. E. Protein Folding Kinetics and Thermodynamics from Atomistic Simulation. *Proc. Natl. Acad. Sci. U.S.A.* **2012**, *109*, 17845–17850.
- Lazaridis, T.; Karplus, M. Thermodynamics of Protein Folding: A Microscopic View. *Biophys. Chem.* **2003**, *100*, 367–395.
- Chong, S.-H.; Ham, S. Configurational Entropy of Protein: A Combined Approach Based on Molecular Simulation and Integral-Equation Theory of Liquids. *Chem. Phys. Lett.* **2011**, *504*, 225–229.
- Karplus, M.; Kushick, J. N. Method for Estimating the Configurational Entropy of Macromolecules. *Macromolecules* **1981**, *14*, 325–332.
- Levy, R. M.; Karplus, M.; Kushick, J.; Perahia, D. Evaluation of the Configurational Entropy for Proteins: Application to Molecular Dynamics Simulations of an α -Helix. *Macromolecules* **1984**, *17*, 1370–1374.
- Schlitter, J. Estimation of Absolute and Relative Entropies of Macromolecules Using the Covariance Matrix. *Chem. Phys. Lett.* **1993**, *215*, 617–621.
- Schäfer, H.; Mark, A. E.; van Gunsteren, W. F. Absolute Entropies from Molecular Dynamics Simulation Trajectories. *J. Chem. Phys.* **2000**, *113*, 7809–7813.
- Andricioaei, I.; Karplus, M. On the Calculation of Entropy from Covariance Matrices of the Atomic Fluctuations. *J. Chem. Phys.* **2001**, *115*, 6289–6292.
- Chong, S.-H.; Ham, S. Conformational Entropy of Intrinsically Disordered Protein. *J. Phys. Chem. B* **2013**, *117*, 5503–5509.
- Derrida, B. Random-Energy Model: Limit of a Family of Disordered Models. *Phys. Rev. Lett.* **1980**, *45*, 79–82.
- Elkin, M.; Andre, I.; Lukatsky, D. B. Energy Fluctuations Shape Free Energy of Nonspecific Biomolecular Interactions. *J. Stat. Phys.* **2012**, *146*, 870–877.
- McKnight, C. J.; Matsudaira, P. T.; Kim, P. S. NMR Structure of the 35-Residue Villin Headpiece Subdomain. *Nat. Struct. Biol.* **1997**, *4*, 180–184.
- Case, D. A.; Darden, T. A.; Cheatham, T. E., III; Simmerling, C. L.; Wang, J.; Duke, R. E.; Luo, R.; Merz, K. M.; Pearlman, D. A.; Crowley, M. et al. *AMBER 11*; University of California: San Francisco, CA, 2010.
- Hornak, V.; Okur, R. A. A.; Strockbine, B.; Roitberg, A.; Simmerling, C. Comparison of Multiple Amber Force Fields and Development of Improved Protein Backbone Parameters. *Proteins* **2006**, *65*, 712–725.
- Jorgensen, W. L.; Chandrasekhar, J.; Madura, J. D.; Impey, R. W.; Klein, M. L. Comparison of Simple Potential Functions for Simulating Liquid Water. *J. Chem. Phys.* **1983**, *79*, 926–935.
- Darden, T.; York, D.; Pedersen, L. Particle Mesh Ewald: An $N \log(N)$ Method for Ewald Sums in Large Systems. *J. Chem. Phys.* **1993**, *98*, 10089–10092.
- Ryckaert, J.-P.; Ciccotti, G.; Berendsen, H. J. C. Numerical Integration of the Cartesian Equations of Motion of a System with Constraints: Molecular Dynamics of n -Alkanes. *J. Comput. Phys.* **1977**, *23*, 327–341.
- Berendsen, H. J. C.; Postma, J. P. M.; van Gunsteren, W. F.; DiNola, A.; Haak, J. R. Molecular Dynamics with Coupling to an External Bath. *J. Chem. Phys.* **1984**, *81*, 3684–3690.
- Ensign, D. L.; Kasson, P. M.; Pande, V. S. Heterogeneity Even at the Speed Limit of Folding: Large-Scale Molecular Dynamics Study of a Fast-Folding Variant of the Villin Headpiece. *J. Mol. Biol.* **2007**, *374*, 806–816.
- Kabsch, W.; Sander, C. Dictionary of Protein Secondary Structure: Pattern Recognition of Hydrogen-Bonded and Geometrical Features. *Biopolymers* **1983**, *22*, 2257–2637.
- Feig, M.; Karanicolas, J.; Brooks, C. L., III MMTSB Tool Set: Enhanced Sampling and Multiscale Modeling Methods for Applications in Structural Biology. *J. Mol. Graphics Modell.* **2004**, *22*, 377–395.

- (27) Kovalenko, A. In *Molecular Theory of Solvation*; Hirata, F., Ed.; Kluwer Academic: Dordrecht, The Netherlands, 2003; p 169.
- (28) Imai, T.; Harano, Y.; Kinoshita, M.; Kovalenko, A.; Hirata, F. A Theoretical Analysis on Hydration Thermodynamics of Proteins. *J. Chem. Phys.* **2006**, *125*, 024911.
- (29) Perkyns, J.; Pettitt, B. M. A Site-Site Theory for Finite Concentration Saline Solutions. *J. Chem. Phys.* **1992**, *97*, 7656–7666.
- (30) Press, W. H.; Teukolsky, S. A.; Vetterling, W. T.; Flannery, B. P. *Numerical Recipes: The Art of Scientific Computing*, 3rd ed.; Cambridge University Press: New York, 2007.
- (31) Grossfield, A.; Zuckerman, D. M. Quantifying Uncertainty and Sampling Quality in Biomolecular Simulations. *Annu. Rep. Comput. Chem.* **2009**, *5*, 23–48.
- (32) Efron, B.; Tibshirani, R. Bootstrap Methods for Standard Errors, Confidence Intervals, and Other Measures of Statistical Accuracy. *Stat. Sci.* **1986**, *1*, 54–77.
- (33) Bevington, P. R. *Data Reduction and Error Analysis for the Physical Sciences*; McGraw-Hill: New York, 1992.
- (34) Kubelka, J.; Eaton, W. A.; Hofrichter, J. Experimental Tests of Villin Subdomain Folding Simulations. *J. Mol. Biol.* **2003**, *329*, 625–630.
- (35) Buscaglia, M.; Kubelka, J.; Eaton, W. A.; Hofrichter, J. Determination of Ultrafast Protein Folding Rates from Loop Formation Dynamics. *J. Mol. Biol.* **2005**, *347*, 657–664.
- (36) Bi, Y.; Cho, J.-H.; Kim, E.-Y.; Shan, B.; Schindelin, H.; Raleigh, D. P. Rational Design, Structural and Thermodynamic Characterization of a Hyperstable Variant of the Villin Headpiece Helical Subdomain. *Biochemistry* **2007**, *46*, 7497–7505.
- (37) Bunagan, M. R.; Gao, J.; Kelly, J. W.; Gai, F. Probing the Folding Transition State Structure of the Villin Headpiece Subdomain via Side Chain and Backbone Mutagenesis. *J. Am. Chem. Soc.* **2009**, *131*, 7470–7476.
- (38) Hu, K.-N.; Yau, W. M.; Tycko, R. Detection of a Transient Intermediate in a Rapid Protein Folding Process by Solid-State Nuclear Magnetic Resonance. *J. Am. Chem. Soc.* **2010**, *132*, 24–25.
- (39) Žoldák, G.; Stigler, J.; Pelz, B.; Li, H.; Rief, M. Ultrafast Folding Kinetics and Cooperativity of Villin Headpiece in Single-Molecule Force Spectroscopy. *Proc. Natl. Acad. Sci. U.S.A.* **2013**, *110*, 18156–18161.
- (40) Freddolino, P. L.; Schulten, K. Common Structural Transitions in Explicit-Solvent Simulations of Villin Headpiece Folding. *Biophys. J.* **2009**, *97*, 2338–2347.
- (41) McKnight, C. J.; Doering, D. S.; Matsudaira, P. T.; Kim, P. S. A Thermostable 35-Residue Subdomain within Villin Headpiece. *J. Mol. Biol.* **1996**, *260*, 126–134.
- (42) Tirion, M. M. Large Amplitude Elastic Motions in Proteins from a Single-Parameter, Atomic Analysis. *Phys. Rev. Lett.* **1996**, *77*, 1905–1908.
- (43) Chong, S.-H.; Ham, S. Impact of Chemical Heterogeneity on Protein Self-Assembly in Water. *Proc. Natl. Acad. Sci. U.S.A.* **2012**, *109*, 7376–7641.
- (44) Ferrenberg, A. M.; Landau, D. P.; Binder, K. Statistical and Systematic Errors in Monte Carlo Sampling. *J. Stat. Phys.* **1991**, *63*, 867–882.
- (45) Piana, S.; Lindorff-Larsen, K.; Shaw, D. E. How Robust Are Protein Folding Simulations with Respect to Force Field Parameterization? *Biophys. J.* **2011**, *100*, L47–L49.
- (46) Dill, K. A. Dominant Forces in Protein Folding. *Biochemistry* **1990**, *29*, 7133–7155.
- (47) Lazaridis, T.; Karplus, M. Effective Energy Function for Proteins in Solution. *Proteins* **1999**, *35*, 133–152.
- (48) Imai, T.; Harano, Y.; Kinoshita, M.; Kovalenko, A.; Hirata, F. Theoretical Analysis on Changes in Thermodynamic Quantities upon Protein Folding: Essential Role of Hydration. *J. Chem. Phys.* **2007**, *126*, 225102.
- (49) Yoshidome, T.; Kinoshita, M.; Hirota, S.; Baden, N.; Terazima, M. Thermodynamics of Apoplastocyanin Folding: Comparison between Experimental and Theoretical Results. *J. Chem. Phys.* **2008**, *128*, 225104.
- (50) Maruyama, Y.; Harano, Y. Does Water Drive Protein Folding? *Chem. Phys. Lett.* **2013**, *581*, 85–90.
- (51) Dill, K. A. Theory for the Folding and Stability of Globular Proteins. *Biochemistry* **1985**, *24*, 1501–1509.
- (52) Lee, K. H.; Xie, D.; Freire, E.; Amzel, L. M. Estimation of Changes in Side Chain Configurational Entropy in Binding and Folding: General Methods and Application to Helix Formation. *Proteins* **1994**, *20*, 68–84.
- (53) Makhataдзе, G. I.; Privalov, P. L. On the Entropy of Protein Folding. *Protein Sci.* **1996**, *5*, S07–S10.
- (54) Wolynes, P. G.; Onuchic, J. N.; Thirumalai, D. Navigating the Folding Routes. *Science* **1995**, *267*, 1619–1620.
- (55) Dill, K. A.; Chan, H. S. From Levinthal to Pathways to Funnels. *Nat. Struct. Biol.* **1997**, *4*, 10–19.
- (56) Chaplin, M. Do We Underestimate the Importance of Water in Cell Biology? *Nat. Rev. Mol. Cell Biol.* **2006**, *7*, 861–866.
- (57) Ball, P. Water as an Active Constituent in Cell Biology. *Chem. Rev.* **2008**, *108*, 74–108.