

# Protein dynamics investigated by inherent structure analysis

Francesco Rao<sup>a,1</sup> and Martin Karplus<sup>a,b,1</sup>

<sup>a</sup>Laboratoire de Chimie Biophysique, Institut de Science et d'Ingénierie Supramoléculaires, Université de Strasbourg, 67000 Strasbourg, France; and  
<sup>b</sup>Department of Chemistry and Chemical Biology, Harvard University, Cambridge, MA 02138

Edited\* by Frank H. Stillinger, Princeton University, Princeton, NJ, and approved February 18, 2010 (received for review January 8, 2010)

**Molecular dynamics (MD) simulations provide essential information about the thermodynamics and dynamics of proteins. To construct the free-energy surface from equilibrium trajectories, it is necessary to group the individual snapshots in a meaningful way. The inherent structures (IS) are shown to provide an appropriate discretization of the trajectory and to avoid problems that can arise in clustering algorithms that have been employed previously. The IS-based approach is illustrated with a 30-ns room temperature "native" state MD simulation of a 10-residue peptide in a  $\beta$ -hairpin conformation. The transitions between the IS are used to construct a configuration space network from which a one-dimensional free-energy profile is obtained with the mincut method. The results demonstrate that the IS approach is useful and that even for this simple system, there exists a nontrivial organization of the native state into several valleys separated by barriers as high as 3 kcal/mol. Further, by introducing a coarse-grained network, it is demonstrated that there are multiple pathways connecting the valleys. This scenario is hidden when the snapshots of the trajectory are used directly with rmsd clustering to compute the free-energy profile. Application of the IS approach to the native state of the PDZ2 signaling domain indicates its utility for the study of biologically relevant systems.**

complex networks | conformational dynamics | energy landscapes | molecular dynamics simulations

**P**rotein function depends critically on the synergy between structure and dynamics. The dynamics often involves the interconversion of conformational states on a complex multidimensional free-energy surface. It is very difficult to study such conformational transitions experimentally at an atomic level of detail, although techniques such as single-molecule FRET (1), x-ray crystallography (2), and NMR (3) supply useful, but limited, information. Consequently, molecular dynamics simulations are playing an increasing role in determining the free-energy surface, as a supplement to the experimental studies. The energy surface is made up of many deep valleys connected by saddles (4), suggesting that protein dynamics can be divided into intravalley and intervalley motions (5). The former represent the oscillations around local minima, while the latter involve barrier crossings from one minimum to another (6). Most descriptions of the surface have been rather qualitative because computations to sample the underlying surface for such multidimensional systems have not been possible. Recently, particularly for peptides and even a few small proteins (7–9), molecular dynamics (MD) simulations that extend into the microsecond range are providing the information required for a more quantitative analysis (10). Because of the large number of degrees of freedom involved, analysis of the results based on graph theory have been found to be useful (11, 12). The essential idea is to map the calculated trajectory on a conformation space network (CSN), whose nodes represent the different conformations visited during the simulation and whose links correspond to direct transitions between the nodes (13). This approach has been successfully applied to obtain an understanding of peptide folding and biomolecular structural transitions (11, 12, 14–17), as well as to interpret electron transfer

experiments (18) and time-resolved IR measurements (19, 20). Alternative methods for determining the presence of metastable states and transition pathways combine many short trajectories to determine the kinetic connectivity (21, 22).

It has long been realized that in liquids at room temperature, thermal fluctuation can hide the underlying architecture of the energy surface. To deal with this problem, Stillinger and Weber (23) introduced the concept of "inherent structures (IS)," defined as the local minima on the potential energy surface. They are determined by calculating an MD trajectory at a given temperature and quenching the system by gentle energy minimization. All conformations that under this mapping go to the same IS define the *basin* (of attraction) of the IS. Such a partitioning has been used to study the thermodynamics of liquids, and for supercooled liquids to obtain insight into the dynamics (24–26). Early studies of proteins based on the IS concept (6, 27) demonstrated the multiminimum nature of the potential surface but were limited by the short trajectories that were accessible. Thermodynamic aspects of coarse-grained models of protein folding have been analyzed more recently (28–31) following the original prescription of Stillinger and Weber (23, 24); see also ref. 32.

In this paper, we show how the IS can be used to facilitate the analysis of MD trajectories that sample the conformation space under equilibrium conditions. The decomposition of the conformation space in terms of the IS provides a natural and simple description of a dynamical system when the timescales of the motions in a minimum and between minima are well separated. Once the IS have been determined, the IS make the mapping of the energy landscape onto a CSN essentially unique and avoid the uncertainties introduced by a purely geometrical clustering of the trajectory to define the network nodes (11, 12). Although the choice of the rmsd cutoffs for clustering, for example, is less important for the large conformational changes that occur in the transition between the unfolded and folded states of proteins (11, 33), for cases where the conformation is more restricted, as it is in the folded (native) state, the IS appear ideal for clustering the snapshots and determining the free-energy surface and associated CSN.

We study a model system, a 10-residue peptide in a  $\beta$ -hairpin conformation that is simple enough so that a full description of the dynamics with the IS approach can be achieved, while the system has a large enough number of degrees of freedom to be interesting. A 30-ns MD simulation at 300 K of the peptide in its folded ("native") state provides a full sampling of the conformation space accessible at equilibrium. Mapping the potential energy onto the free-energy surface by means of the IS-based CSN and mincut free-energy profiles (34) demonstrates the

Author contributions: F.R. and M.K. designed research; F.R. performed research; F.R. and M.K. analyzed data; and F.R. and M.K. wrote the paper.

The authors declare no conflict of interest.

\*This Direct Submission article had a prearranged editor.

<sup>1</sup>To whom correspondence may be addressed. E-mail: francesco.rao@gmail.com or marci@tammy.harvard.edu.

This article contains supporting information online at [www.pnas.org/lookup/suppl/doi:10.1073/pnas.0915087107/-DCSupplemental](http://www.pnas.org/lookup/suppl/doi:10.1073/pnas.0915087107/-DCSupplemental).

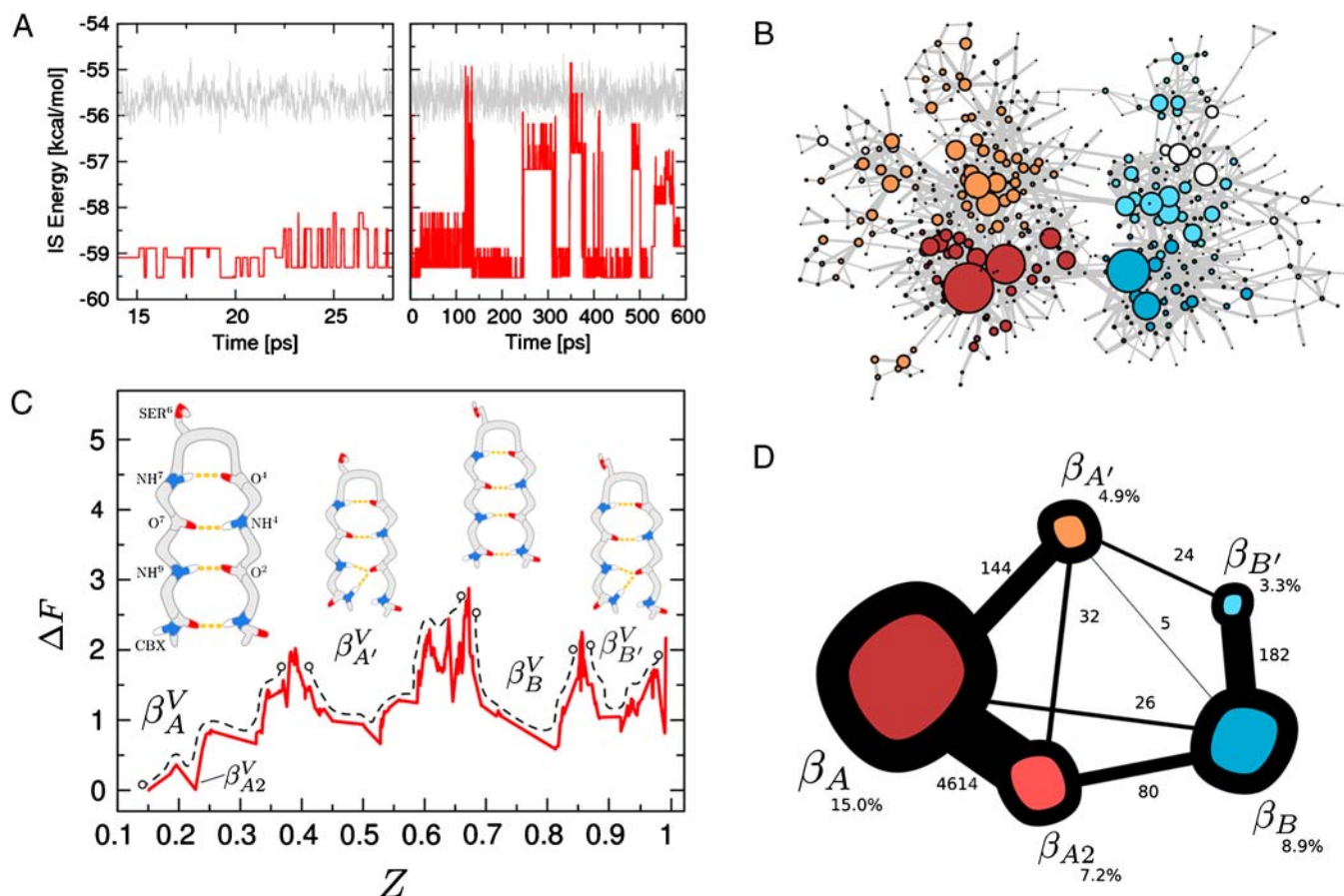
correspondence between conformational changes, energy barriers, and transition kinetics. Comparison with the standard CSN analysis based on clustering of the 300 K MD trajectory shows certain limitations of the latter. An illustrative application of the IS analysis to a PDZ2 domain demonstrates the utility of the method for studying the native state dynamics of biologically relevant systems.

## Results and Discussion

**The Conformation Space Network and Free-Energy Profile. Inherent structures and their transitions.** The application of the mapping into IS of the 30-ns long ( $1.5 \times 10^7$  snapshots) MD trajectory of the GS10 peptide results in 1,561 IS, as described in *Methods* (see Fig. S1). The set of IS defines an ensemble of *microstates* that can be used to characterize the dynamics of the peptide. Two short segments of the time series of the IS are shown in Fig. 1A. It is evident that different long-lived regions (called valleys) are sampled. The time series of the conformations at room temperature does not show this organization (gray lines in Fig. 1A). The valleys, which are described in more detail in terms of the free-energy profile (see below), include many IS basins and are characterized by different values of the potential energy. In Fig. 1A Left, a transition between two valleys ( $\beta_A^V$  and  $\beta_B^V$  in Fig. 1C) is shown at  $t \approx 22$  ps. When the system is sampling one valley, it rapidly interconverts among a small number of minima with

similar energies but slightly different conformations (e.g., all-atom rmsd between 0.3 and 0.6 Å) while, rarely, there is a transition to another valley. Fig. 1A Right shows a series of transitions between several valleys over the time scale of 600 ps. The fast transitions within a valley and the slow transitions between different ones recall the classification used in the field of super-cooled liquids between type  $\beta$ - and  $\alpha$ -transitions, respectively (see figure 3 of ref. 35 for more details).

Given the microstates defined by the IS, we construct a CSN shown in Fig. 1B. Visual inspection of the network indicates the presence of a modular organization, i.e., the presence of different groups of nodes with many links between them and fewer links to other nodes. The weights of the nodes and links have a clear physical meaning representing, respectively, the populations of the basins of attraction of the IS (i.e., proportional to their free energy) and the transition probabilities [i.e., proportional to an effective activation energy, since there could be multiple potential energy barriers between a pair of nodes (36)]. For the detailed analysis of the network, we use the procedure described in *Methods* to calculate a one-dimensional cut-based free-energy profile (CFEP) (34). In Fig. 1C the CFEP of the GS10 peptide is shown. This profile represents the free-energy surface projected on the partition function-based reaction coordinate  $Z$  (see *Methods*), relative to a given reference microstate, in this case the most populated node (arbitrarily called  $\beta_A$ ). CFEPs have proven to

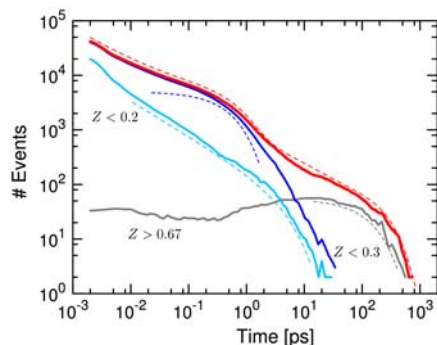


**Fig. 1.** (A) IS trajectory. Two sample windows of the IS energies of the GS10 peptide in kcal/mol and the room temperature potential energy (in arbitrary units for comparison) time series are shown in red and gray, respectively. In the left part a transition between the  $\beta_A^V$  and  $\beta_B^V$  valley is depicted, whereas in the right part a series of transitions between multiple valleys is shown. (B) CSN of the GS10 peptide. Nodes and links represent IS and MD transitions, respectively; the color code corresponds to that in D, except that white nodes are characterized by  $Z > 0.97$ . The size of the nodes and links is proportional to their populations. For clarity, only nodes that have been visited more than 200 times are shown; there are a total of 571 nodes. (C) IS CFEP of the GS10 peptide relative to the most populated microstate,  $\beta_A$ . The four most populated valleys are shown in dashed lines, and the corresponding lowest energy microstate structures are schematically sketched; only the atoms involved in the relevant hydrogen bonding and SER orientations are displayed. (D) Reaction pathways between the most important IS of the system displayed as a CCN. IS populations and average number of transitions are shown both by the numbers and by the size of the nodes and the thickness of the links, respectively (see text for comments).

be a better approach as compared with traditional free-energy profiles (where the landscape is projected onto one or more arbitrarily chosen order parameters (33, 37). The CFEPs provide a correct estimate of the height of the free-energy barriers, as well as their positions on the landscape (34, 38, 39). The cut-based free-energy profile of GS10 has a complex structure with a series of barriers and valleys. This is perhaps a somewhat surprising result for such a small peptide restricted to a well-defined  $\beta$ -hairpin structure. There are four regions that are characterized by broad valleys (dashed lines in the figure). The four valleys are labeled  $\beta_A^V$ ,  $\beta_{A'}^V$ ,  $\beta_B^V$ , and  $\beta_{B'}^V$ , and the most populated IS in each of them is  $\beta_A$  ( $E_{\text{IS}} = -59.5$  kcal/mol),  $\beta_B$  ( $E_{\text{IS}} = -59.3$  kcal/mol),  $\beta_{A'}$  ( $E_{\text{IS}} = -57.9$  kcal/mol), and  $\beta_{B'}$  ( $E_{\text{IS}} = -57.7$  kcal/mol). In addition, the first small valley of the profile located at  $Z \approx 0.22$  is labeled  $\beta_{A_2}^V$  according to the  $\beta_{A_2}$  microstate ( $E_{\text{IS}} = -59.0$  kcal/mol).

**Structural analysis.** For this system, it is possible to structurally characterize the differences between the valleys in a straightforward manner. The transition from  $\beta_A$  to  $\beta_{A'}$  is characterized by the rotation of the CBX group about the corresponding  $\psi$  dihedral angle, which disrupts the hydrogen bond between  $\text{NH}^9$  and  $\text{O}^2$  and forms a hydrogen bond between the  $\text{O}^2$  and the  $\text{NH}^{10}$  group of CBX; see the schematic peptide structures shown in Fig. 1C above the CFEP. The same atomic rearrangement characterizes the transition between  $\beta_B$  and  $\beta_{B'}$ . The transition between  $\beta_A$  and  $\beta_B$  corresponds to the rotation of the serine OH group about the dihedral angles  $\text{N-C}_\alpha^6\text{-C}_\beta\text{-O}$  ( $\chi_1$ ) and  $\text{C}_\alpha^6\text{-C}_\beta\text{-O-H}$  ( $\chi_2$ ). This transition is responsible for the highest barrier ( $\approx 3$  kcal/mol) in the CFEP at position  $Z \approx 0.67$ ; it arises primarily from the  $\chi_1$  dihedral angle barrier (see Fig. S2). The rearrangement between  $\beta_A$  and  $\beta_{A2}$  corresponds to a  $120^\circ$  rotation of the serine OH group about the  $\chi_2$  dihedral angle. Examination of the high-energy microstates that are accessed within the various valleys (see Fig. 1A) shows that they correspond to a variety of backbone distortions; see Fig. S3.

**Coarse-grained network.** A coarse CSN (CCSN) was built from the original trajectory projected on the microstates by keeping only the five IS:  $\beta_A$ ,  $\beta_{A2}$ ,  $\beta_{A'}$ ,  $\beta_B$ ,  $\beta_{B'}$ ; all others were deleted. This was done to have a simple way to capture real transitions between the valleys and eliminate the multiple passes through the transition state. The reaction pathways between the five microstates described above are shown in Fig. 1D. It is clear that there are multiple, not equally probable, pathways between the five IS, instead of a sequential pathway along the Z coordinate. This complexity is not evident from the CFEP by itself. The reaction to go



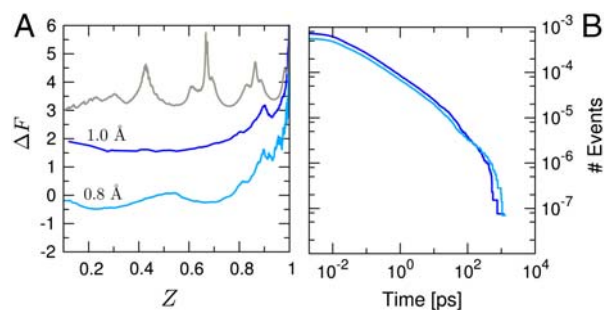
**Fig. 2.** IS FPT distributions to  $\beta_A$  focused analysis (see main text). The red, light blue, blue, and gray lines represent the full FPT, focused FPT with  $Z < 0.2$ ,  $Z < 0.3$ , and  $Z > 0.67$ , respectively (see text). Dashed lines show different functional fits. The two exponential fits have a characteristic time of 0.6 ps (dashed blue line) and 165 ps (dashed gray line). The fits have been slightly displaced for clarity. The fitting function for the full FPT distribution is  $f(t) = c_0 t^{-\alpha} e^{-t/\tau_0} + c_1 e^{-t/\tau_1} + c_2 e^{-t/\tau_2}$ , where  $c_0, c_1, c_2, \alpha, \tau_0, \tau_1$ , and  $\tau_2$  have a value of 704; 4.035; 88; 0.6; 33; 0.6; and 180, respectively.

from  $\beta_{B'}$  to  $\beta_A$  can take place by more than one possible route. The most significant pathway proceeds from  $\beta_{B'}$  to  $\beta_B$ , then to  $\beta_{A2}$ , and finally to  $\beta_A$ . Also, while the direct transition from  $\beta_B$  to  $\beta_A$  is possible, it is 3 times less probable than passing through  $\beta_{A2}$ , which acts as an intermediate that rapidly interconverts with  $\beta_A$ ; see also Fig. S2. This is not the case for the transition from  $\beta_{A'}$  to  $\beta_A$  where the pathway involving any intermediates is 4.5 times less favorable than a direct transition. Finally, we note that the transition between  $\beta_{B'}$  and  $\beta_{A'}$  is faster when the transition is direct, rather than via the intermediate  $\beta_B$ . These results demonstrate that even in the folded GS10  $\beta$ -hairpin peptide, the free-energy landscape is complex and involves multiple pathways.

**First Passage Times Analysis.** First passage time (FPT) distributions are useful for obtaining a detailed understanding of the dynamics (40). The FPT distribution obtained from the time series of the GS10 peptide is shown in Fig. 2. This plot represents the distribution of the relaxation times to the most populated IS, namely, microstate  $\beta_A$ . The distribution is broad, spanning six decades in time from femtoseconds to tens of nanoseconds. The curve shows two “bumps”: one at times of the order of 500 fs and the other one on the 100 ps time scale. From the plot the origin of this behavior is not clear. For this reason, *focused* versions of the FPT distribution have been calculated as indicated in *Methods*; i.e., they include the relaxation from only a given subset of microstates.

A FPT distribution is built considering only the IS that are in the CFEP at values of  $Z$  lower than 0.2 (see Fig. 1C), the value at which the first free-energy barrier (to  $\beta_{A2}$ ) appears. This distribution is well represented by a power law with an exponential cutoff (light blue lines in Fig. 2; the fit is shown as a dashed line). Such behavior is typical of relaxations within a single well for which the characteristic time to reach the bottom ( $\beta_A$  state) is undefined. The exponential cutoff arises from the fact that there is a typical residence time (on the order of 0.1 ps) after which the system jumps to other regions of the landscape. The CFEP, which predicts nearly barrierless transitions to  $\beta_A$ , is in agreement with this analysis.

The inclusion of the microstates up to  $Z < 0.3$  (i.e., including  $\beta_{A2}$ ) for the calculation of the FPT distribution introduces the first bump. This bump represents the fast exponential decay generated by the partial reorientation of the hydroxyl group of the SER side chain, typical of the microstates found in the  $\beta_{A2}$  region (blue curves). The time scales of these transitions overlap with the barrierless transitions generating the power-law behavior. However, the two types of transitions are microscopically different because the latter correspond to deformations of the backbone (without changing the  $\beta$ -sheet hydrogen bond pattern), while the former involves the reorientation of the hydroxyl group.



**Fig. 3.** Results of structure-based clustering of room temperature trajectory. (A) rmsd CFEF for the 1.0-Å and 0.8-Å cutoffs are shown as blue and light blue lines, respectively. The CFEF obtained by using the three most relevant degrees of freedom of the system (see text) is shown in gray. (B) FPT distributions to the most populated rmsd cluster. Blue and light blue lines represent the 1.0-Å and 0.8-Å cutoffs realizations, respectively.







16. Yang S, Roux B (2008) Src kinase conformational activation: Thermodynamics, pathways, and mechanisms. *PLOS Comput Biol* 4(3):e1000047.
17. Prada-Gracia D, Gomez-Gardenes J, Echenique P, Palo F (2009) Exploring the free energy landscape: From dynamics to networks and back. *PLOS Comput Biol* 5:e1000415.
18. Li CB, Yang H, Komatsuzaki T (2008) Multiscale complex network of protein conformational fluctuations in single-molecule time series. *Proc Natl Acad Sci USA* 105:536–541.
19. Ihalainen JA, et al. (2007) Folding and unfolding of a photoswitchable peptide from piconoseconds to microseconds. *Proc Natl Acad Sci USA* 104:5383–5388.
20. Ihalainen JA, et al. (2008) Alpha-helix folding in the presence of structural constraints. *Proc Natl Acad Sci USA* 105:9588–9593.
21. Bowman GR, Beauchamp KA, Boxer G, Pande VS (2009) Progress and challenges in the automated construction of Markov state models for full protein systems. *J Chem Phys* 131:124101.
22. Noé F, et al. (2009) Constructing the equilibrium ensemble of folding pathways from short off-equilibrium simulations. *Proc Nat Acad Sci USA* 106:19011–19016.
23. Stillinger FH, Weber TA (1982) Hidden structure in liquids. *Phys Rev A* 25:978–989.
24. Stillinger F, Weber T (1983) Dynamics of structural transitions in liquids. *Phys Rev A* 28:2408–2416.
25. Denny R, Reichman D, Bouchaud J (2003) Trap models and slow dynamics in supercooled liquids. *Phys Rev Lett* 90(2):025503.
26. Berthier L, Garrahan J (2003) Nontopographic description of inherent structure dynamics in glassformers. *J Chem Phys* 119(8):4367–4371.
27. Caves LS, Evansek JD, Karplus M (1998) Locally accessible conformations of proteins: Multiple molecular dynamics simulations of crambin. *Protein Sci* 7:649–66.
28. Baumketner A, Shea JE, Hiwatari Y (2003) Glass transition in an off-lattice protein model studied by molecular dynamics simulations. *Phys Rev E* 67:011912.
29. Nakagawa N, Peyrard M (2006) The inherent structure landscape of a protein. *Proc Nat Acad Sci USA* 103(14):5279–5284.
30. Kim J, Keyes T (2007) Inherent structure analysis of protein folding. *J Phys Chem B* 111:2647–2657.
31. Kim J, Keyes T, Straub J (2009) Relationship between protein folding thermodynamics and the energy landscape. *Phys Rev E* 79(3):030902.
32. Krivov SV, Chekmarev SF, Karplus M (2002) Potential energy surfaces and conformational transitions in biomolecules: A successive confinement approach applied to a solvated tetrapeptide. *Phys Rev Lett* 88:038101.
33. Rao F, Caflisch A (2003) Replica exchange molecular dynamics simulations of reversible folding. *J Chem Phys* 119:4035–4042.
34. Krivov SV, Karplus M (2006) One-dimensional free-energy profiles of complex systems: Progress variables that preserve the barriers. *J Phys Chem B* 110:12689–12698.
35. Debenedetti PG, Stillinger FH (2001) Supercooled liquids and the glass transition. *Nature* 410:259–267.
36. Gfeller D, de Lachapelle DM, De Los Rios P, Caldarelli G, Rao F (2007) Uncovering the topology of configuration space networks. *Phys Rev E* 76:026113.
37. Ferrara P, Caflisch A (2000) Folding simulations of a three-stranded antiparallel beta-sheet peptide. *Proc Natl Acad Sci USA* 97:10780–10785.
38. Krivov SV, Muff S, Caflisch A, Karplus M (2008) One-dimensional barrier-preserving free-energy projections of a beta-sheet miniprotein: New insights into the folding process. *J Phys Chem B* 112:8701–8714.
39. Muff S, Caflisch A (2009) Identification of the protein folding transition state from molecular dynamics trajectories. *J Chem Phys* 130:125104–125111.
40. Chekmarev SF, Krivov SV, Karplus M (2005) Folding time distributions as an approach to protein folding kinetics. *J Phys Chem B* 109:5312–5330.
41. Rao F, Settanni G, Guarnera E, Caflisch A (2005) Estimation of protein folding probability from equilibrium simulations. *J Chem Phys* 122:184901–184905.
42. Kong Y, Karplus M (2008) Signaling pathways of PDZ2 domain: A molecular dynamics interaction correlation analysis. *Proteins: Struct Funct Bioinformatics* 74:145–54.
43. Hamm P, Lim M, Hochstrasser RM (1998) Structure of the amide I band of peptides measured by femtosecond nonlinear-infrared spectroscopy. *J Phys Chem B* 102:6123–6138.
44. Kolano C, et al. (2006) Watching hydrogen-bond dynamics in a beta-turn by transient two-dimensional infrared spectroscopy. *Nature* 444:469–472.
45. Bagchi S, Falvo C, Mukamel S, Hochstrasser RM (2009) 2D-IR experiments and simulations of the coupling between amide-I and ionizable side chains in proteins: Application to the Villin headpiece. *J Phys Chem B* 113:11260–11273.
46. Brooks BR, et al. (1983) Charmm: A program for macromolecular energy, minimization, and dynamics calculations. *J Comput Chem* 4:187–217.
47. Brooks BR, et al. (2009) CHARMM: The biomolecular simulation program. *J Comput Chem* 30:1545–1614.
48. Ferrara P, Apostolakis J, Caflisch A (2002) Evaluation of a fast implicit solvent model for molecular dynamics simulations. *Proteins: Struct Funct Bioinform* 46:24–33.
49. Hartigan J (1975) *Clustering Algorithms* (Wiley, New York).
50. Seeber M, Cecchini M, Rao F, Settanni G, Caflisch A (2007) Wordom: A program for efficient analysis of molecular dynamics simulations. *Bioinformatics* 23:2625–2627.
51. Allen LR, Krivov SV, Paci E (2009) Analysis of the free-energy surface of proteins from reversible folding simulations. *PLOS Comput Biol* 5:e1000428.
52. Guarnera E, Pellarin R, Caflisch A (2009) How does a simplified-sequence protein fold? *Biophys J* 97(6):1737–1746.