

Folding energy landscape and network dynamics of small globular proteins

Naoto Hori^a, George Chikenji^b, R. Stephen Berry^{c,1}, and Shoji Takada^{d,e}

^aGraduate School of Science and Technology, Kobe University, Rokkodai Nada Kobe 657-8501, Japan; ^bDepartment of Computational Science and Engineering, Graduate School of Engineering, Nagoya University, Furocho, Chikusa, Nagoya 464-8603, Japan; ^cDepartment of Chemistry and The James Franck Institute, University of Chicago, Chicago, IL 60637; ^dDepartment of Biophysics, Graduate School of Science, Kyoto University, Kitashirakawa-oiwake Sakayo Kyoto 606-8502, Japan; and ^eCore Research for Evolutionary Science and Technology Japan Science and Technology Corporation, Kyoto University, Kitashirakawa-oiwake Sakayo

Contributed by R. Stephen Berry, November 14, 2008 (sent for review September 17, 2008)

The folding energy landscape of proteins has been suggested to be funnel-like with some degree of ruggedness on the slope. How complex the landscape, however, is still rather unclear. Many experiments for globular proteins suggested relative simplicity, whereas molecular simulations of shorter peptides implied more complexity. Here, by using complete conformational sampling of 2 globular proteins, protein G and src SH3 domain and 2 related random peptides, we investigated their energy landscapes, topological properties of folding networks, and folding dynamics. The projected energy surfaces of globular proteins were funneled in the vicinity of the native but also have other quite deep, accessible minima, whereas the randomized peptides have many local basins, including some leading to seriously misfolded forms. Dynamics in the denatured part of the network exhibited basin-hopping itinerancy among many conformations, whereas the protein reached relatively well-defined final stages that led to their native states. We also found that the folding network has the hierarchic nature characterized by the scale-free and the small-world properties.

contact maps | folding pathways | multiple pathways | principal coordinates

Proteins fold on large-dimensional energy landscapes through myriads of conformations. One energy-landscape theory suggests that the global shape of the landscape is primarily funnel-like with some degree of ruggedness on the slope of the funnel (1, 2). How complex/rugged the energy landscape is and how diverse the folding-pathway ensemble is are still rather controversial. Experimentally, many small fast-folding proteins exhibit single-exponential behavior, suggesting simplicity (3). For such proteins, a perfect funnel model, Go model, has been used, as an extreme of simplicity, to model folding routes, often showing modestly good agreement with experiments (4). Conversely, there exist several clear evidences of complexity in folding. Under some conditions, proteins show strange and glassy kinetics, suggesting ruggedness of the landscape (5). Some β -sheet proteins, such as β -lactoglobulin, form nonnative α -helices at early stages of folding (6, 7).

The computational approach has been the most direct to elucidate the complexity of folding energy landscapes. Methods developed in other areas, such as atomic clusters, have been applied to peptides and proteins, illustrating the multiple minima on the landscape (8–11). Recently, with background (10, 11), Krivov and Karplus developed the transition disconnectivity graph to visualize quantitatively the free-energy landscape and applied it for peptides finding a highly rugged non-funnel-like landscape with competing minima (12, 13). Caflisch and co-workers (14, 15) constructed a folding network for a designed peptide and uncovered a highly heterogeneous denatured ensemble. They both used network analyses without the data reduction to lower dimension and warned that the projection to low dimension, as is often done in conventional folding studies, can hide the complexity in the landscape. However, their analyses that required folding/unfolding trajectories were limited to

short peptides, too short to exhibit typical hydrophobic cores. These peptides were either cleaved from the wild type or human-designed. Naturally, their landscapes may be less designed and more rugged than the landscape of evolutionally designed globular proteins (see also ref. 16).

Thus, the question immediately arose of how complex and rugged the energy landscape of natural globular proteins with a typical size of hydrophobic core would be. It is, however, highly nontrivial to address this question because we need unbiased and comprehensive sampling of the energy landscape of globular proteins that should cover native basins as well as a broad range of denatured states. Here, based on the physics-based technology developed for de novo protein structure prediction, we realized it. Given only amino acid sequence information, the method we use here is able to predict native folds for small globular proteins with modest reliability and accuracy. Yet, the method is not a priori biased to the native structure, i.e., it is not a Go model, and so is able to explore nonnative parts of the landscape as well. Buttressed by these methods, we can investigate the energy landscape characteristics and network dynamics on it.

In this article, using an approximate protein model, we first obtained quite complete conformational ensembles for 2 small globular proteins, protein G and src SH3, and 2 random sequences derived from these peptides, random-G and random-S. These ensembles were then used to visualize the energy landscapes in contact map-based principal component axes. The energy landscapes of the 2 proteins were reasonably funneled in the vicinities of their native structures, whereas several misfolded minima existed. We then constructed folding networks, which were shown to have hierarchic nature characterized by the small-world and scale-free properties. Based on the master equation on the networks, we further studied folding dynamics on the unprojected space. We found basin-hopping itinerancy among many structural basins in denatured states, whereas the final stage to the native basin went through relatively specific routes.

Results

Conformational Sampling. We used a combination of the in-house-developed, coarse-grained energy function, SimFold (17–19), and a multicanonical-ensemble fragment assembly Monte Carlo method (20–22) for complete sampling of conformational spaces. This combination has been tested in 2 recent CASPs, the biennial blind tests of the protein structure prediction, showing high-level performance in de novo structure prediction of rela-

Author contributions: N.H., R.S.B. and S.T. designed research; N.H. and G.C. performed research; N.H., G.C., and S.T. analyzed data; and N.H., R.S.B., and S.T. wrote the paper.

The authors declare no conflict of interest.

¹To whom correspondence should be addressed. E-mail: berry@uchicago.edu.

This article contains supporting information online at www.pnas.org/cgi/content/full/0811560106/DCSupplemental.

© 2008 by The National Academy of Sciences of the USA

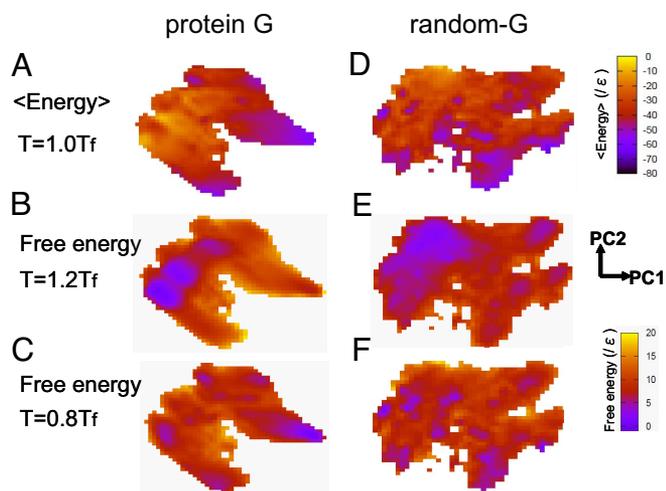


Fig. 3. Energy and free-energy landscapes drawn on the contact map-based PC1–PC2 plane. (A and D) The energy landscape for protein G (A) and for random-G (D) at $T = 1.0T_F$. (B and E) The free-energy landscape at $T = 1.2T_F$ for protein G (B) and for random-G (E). (C and F) The free-energy landscape at $T = 0.8T_F$ for protein G (C) and for random-G (F).

the N-terminal β -hairpin is located at approximately $(-3, -4)$ very far from the native structures. Even though this structure is similar to the native in the Cartesian coordinate (small RMSD), their contact maps are markedly different, which resulted in large separation of this structure from the native.

Although data classification with the contact map-based PCA was significantly better than that with the Cartesian coordinates, there were still some areas on PC1–PC2 plane (Fig. 2) where highly heterogeneous structures overlapped. Especially, structures approximately $(PC1, PC2) = (-4, -1)$ were extremely diverse, including highly extended structures and some misfolded structures as illustrated in Fig. 2. Also, near $(-1, -3)$, we found quite heterogeneous ensembles that included clustered misfolded structures as well as structures with native-like α -helix and C-terminal β -hairpin accompanied by a disordered N-terminal segment. It seems that the complete conformational space of protein G is too diverse to be well represented by a linear transformation of PCA in which only 2 PCs are used to characterize structures. Possible ways to overcome this problem may be to use (i) 3 or perhaps even 4 PCs, (ii) individual PCAs for subensembles, namely the local PCA, or (iii) nonlinear mapping as in ref. 25.

The same type of scattered plot on the PC1–PC2 plane for random-G produced a much more diverse and complex distribution of clustered structures (data not shown). The proportions, which quantify the information content on PCs (see *Materials and Methods*), of the first 2 PCs were $c1 = 30\%$ and $c2 = 5\%$ for protein G, whereas those for random-G were $c1 = 7\%$ and $c2 = 4\%$, suggesting that conformational space for random-G is more complex and diverse than that of protein G.

On the contact map-based PC axes, we calculated the energy and free-energy landscape. First, the average energy at $T = 1.0T_F$ was calculated on the PC1–PC2 plane, which is shown in Fig. 3A for protein G. Clearly, protein G had a funnel-like shape near the native with much of ruggedness in the denatured region. Random-G did not exhibit any funnel-like shape, but had multiple competing local minima (Fig. 1A and D). Free-energy surfaces, drawn in Fig. 3B and C for protein G and Fig. 3E and F for random-G, showed characteristic temperature dependences. For protein G, the free-energy surface at $T = 0.8T_F$ had its global minimum at the native basin, whereas, at $T = 1.2T_F$, extended structures, placed at the left, had lower free energies. Random-G

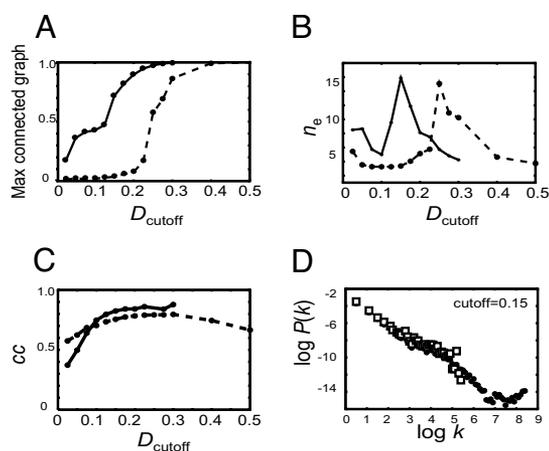


Fig. 4. Network topology analysis. (A–C) The size of the largest connected graph (A), the average number of edges that connect a pair of nodes (n_e) (B), and the cluster coefficient cc (C) are plotted as a function of the distance cutoff D_{cut} . (D) The probability of the degree k (number of structures connected from 1 structure) in log–log scale. (A–C) Solid curves indicate protein G and dashed, random-G. (D) Filled circle, protein G and open square, random-G.

also had the dominant population at extended structures at $T = 1.2T_F$. Lowering the temperature, we started to see multiple minima, which was very prominent at $T = 0.8T_F$.

Folding Network. We now proceed to a network analysis that avoids the reduction of dimension. In general, the network is defined by a combination of “node” and “edge.” As in the previous section, we used the contact map as a structure code. Nodes correspond to sampled structures or their contact maps. Introducing a distance measure D between a pair of contact maps (See *Materials and Methods* for detail), we added an edge when D for a pair is smaller than a cutoff D_{cut} . The number of edges of course increases as a function of D_{cut} (see Fig. S3). With a small D_{cut} , nodes are separated into many isolated clusters. The number of nodes (structures) involved in the largest connected graph increases with D_{cut} , and it increases very sharply (percolation transition) near a specific value of D_{cut} : For protein G, for example, near $D_{cut} = 0.15$ (Fig. 4A).

We investigated topological parameters that characterize the network. First, we estimated the average number of edges, n_e , that can connect a given pair of nodes, as a function of D_{cut} , showing a peak near 0.15 for protein G (Fig. 4B), which coincides with the onset of the percolation transition. The average number of edges n_e connecting nodes is 7 or smaller when D_{cut} is <0.2 . Second, the cluster coefficient cc was calculated. It is the probability that 2 nodes B and C that are directly connected to a node A have the edge between them. We obtained that the average cluster coefficient cc is as high as ≈ 0.7 (Fig. 4C). Networks that are characterized by relatively small n_e and large cc are called small-world networks, which were often found in social networks, web networks, and so forth (26). We also plotted the probability distribution $P(k)$ of the connectivity k , the number of edges connected to 1 node in log–log scale, finding that $P(k)$ shows a power-law dependence (Fig. 4D). This is called the scale-free property (27). The scale-free network suggests existence of small, but not negligible, number of hub-like nodes and the hierarchic nature of the network. The same analysis for random-G showed that random-G, too, had the small-world and scale-free properties.

Folding Dynamics on the Network. The folding network could express diversity of the conformation space, but it, as it was, contained neither dynamic nor energetic information. To ad-

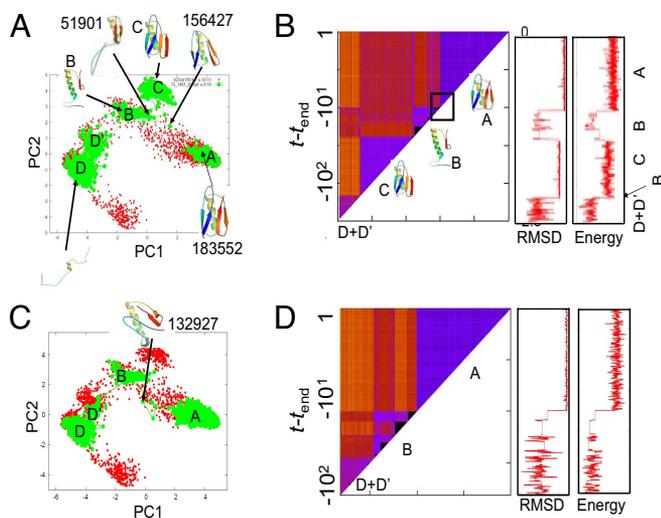


Fig. 5. Two representative folding trajectories of protein G calculated by Gillespie dynamics at $T = 0.8T_F$. (A and C) Structural propagation plotted on the PC1–PC2 plane. (B and D) (Left) Time–time distance matrix of the trajectory. The distance $D(t, t')$ between the structure at time t and that at t' was expressed at (t, t') position by color: In ascending order, from black ($D(t, t') = 0$), to blue, red, and to orange ($D(t, t') \approx 1$). Block-diagonal parts (triangles) represent clusters, which were named as D+D', C, B, and A, and representative structures there were plotted aside. The area surrounded by a square corresponds to the transition to the native. (Center) RMSD vs. time plot of the same trajectory. (Right) Energy vs. time plot of the same trajectory.

dress folding dynamics on the network, here, we assigned, for every edge of the network, transition probabilities that depend on the energy difference between the 2 nodes and that satisfy the detailed balance. The resulting master equation describes the time propagation of the probabilities $P_i(t)$ being in conformation i at time t ; the infinite time limit represents thermodynamic equilibrium. The master equation can readily be solved in 2 ways. One is to solve the probabilities $\{P_j(t)\}$ directly at an arbitrary time t by a matrix diagonalization; the other is the trajectory-based solution by Gillespie dynamics (28). The former is useful for quantitative comparison with bulk experiments, whereas the latter is powerful for elucidating the complexity of the landscape.

For protein G, we obtained 10 folding trajectories from a highly extended structure by Gillespie dynamics at the temperature $T = 0.8T_F$. The time propagation was stopped 10τ after the protein first reached the native-like structure (τ is the characteristic time of the barrierless transition defined in Eq. 2 in *Materials and Methods*). A typical folding trajectory is illustrated in Fig. 5A on the PC1–PC2 plane. At the earliest stage, the protein was highly denatured, and it transited back and forth between 2 clusters D (major) and D' (minor). After the protein departed from D+D' states, it reached the cluster B and very quickly jumped into a misfolded state C, which has quite small RMSD from the native ($\approx 3\text{--}4 \text{ \AA}$, see Fig. 5B Center) and is characterized by the flipped C-terminal β -hairpin (see cartoon in Fig. 5A). After quite a long residence in the misfolded state C, the protein went back to the state B, which is characterized by the native-like C-terminal β -hairpin packed against the central α -helix and the disordered N-terminal segment. Abruptly, then, the protein jumped into the native basin at $t = t_{\text{end}} - 10\tau$. This abrupt change occurred in only 3 steps; from the structure ID = 51901 that belongs to the cluster B (the ID number here and hereafter is the serial number for structure merely identifying each of the sampled $\approx 25,000$ structures; order does not matter) to ID = 156427, and to ID = 183552 (which is already native-like) (see Fig. 5A). Although plotting trajectories on the PC1–PC2 plane is useful to gain quick insight,

it has potential risk of hiding the highly heterogeneous nature in the denatured state because some areas on the plane contain diverse structures. To this end, we propose to use the time–time distance matrix of the trajectory, which has been used in a different field (29). Fig. 5B depicts the time–time distance matrix for the same trajectory as Fig. 5A. In Fig. 5B, the pairwise distance $D(t, t')$ between structures at time t and t' is represented by color at $(x, y) = (t, t')$. The block-diagonal region indicates that the protein resided in a cluster in the corresponding time range, and thus a move from 1 block-diagonal to another indicates basin-hopping. We note that this representation does not rely on the reduction in dimension and thus does not hide the heterogeneous nature of the trajectory. In Fig. 5B, we see 4 major block-diagonal structures (triangles). In the first triangle (at the bottom-left triangle), we recognized checkered pattern suggesting frequent transitions between (at least) 2 states, D and D'. The other 3 triangles correspond to major states, C, B, and A. Of 10 folding trajectories, 4 folded essentially through this route.

Other major routes depicted in Fig. 5C and D also folded directly from the states B to the native state A but passed through a different path. This pathway characterized by the passing of ID = 132927 was observed 4 times in 10 trajectories. Another 2 trajectories were somewhat unique. Misfolding from B to C was probabilistic and was observed in Fig. 5A but not in Fig. 5C.

We then investigated the probability-form solution for protein G and src SH3. Here, we conducted a temperature-jump refolding: Starting with thermal equilibrium at $T = 1.2T_F$, we lowered the temperature to $T = 0.8T_F$ and followed the time propagation. For protein G, between $t = 10^{-2}\tau$ and $t = 10^0\tau$, compaction started with the C-terminal β -hairpin approaching the core. At approximately $t = 10^2\tau$, the major conformation had the formed C-terminal β -hairpin packed with the central α -helix, which is the characteristic experimentally suggested for the transition state ensemble (30). By the time $t = 10^4\tau$, the population of the native basin becomes dominant with some residual population in misfolded structures. Notably, the complex itinerancy seen in Gillespie dynamics was washed out in the ensemble view of the folding.

In the case of src SH3 domain, the equilibrium ensemble $T = 1.2T_F$ contained $\approx 20\%$ of nonnative α -helix contents in the N terminus and residues between 34 and 40 as well as 20% of β -sheet contents in several segments. Between $t = \tau$ and $t = 10^2\tau$, these secondary structure contents doubled almost uniformly along the sequence (Fig. S4a and b where a representative structure is also depicted). This is perfectly consistent with experiments (31) and our earlier calculation (21). Up to $t = 10^4\tau$, the β -hairpin that contains the distal loop was formed that resulted in disappearance of the nonnative α -helix between residues 34 and 40, whereas the rest of the chains fluctuate very broadly (Fig. S4c). This structural feature is consistent with the experimental results of the ϕ -value analysis (32). The N-terminal nonnative α -helix still exists with the probability $\approx 40\%$. These all disappeared after $t = 10^4\tau$ when the final major transition to the native basin occurred (Fig. S4d).

Discussion

We analyzed folding dynamics described by the master equation on the network by 2 alternative approaches. Although the trajectory-based solution of Gillespie, once averaged over many trajectories, should converge to the probability solution, the former provided us much richer insight on the complex behavior of basin-hopping dynamics of protein G. The basin-hopping dynamics that we observed was not anticipated by the probably solution. The trajectory-based approach and the probability approach correspond to the single-molecule simulation and the ensemble computation, respectively. Complex behavior in the single-molecule observation was washed out in the ensemble view. In laboratory experiments on folding, most studies to date were ensemble-based, mostly suggesting simplicity. Recently

Instead of the solution in probability form, we can obtain the “single-molecule” version of the solution as stochastic dynamics of Gillespie (28). Ensemble sum of the Gillespie trajectories coincides with the explicit solution above.

ACKNOWLEDGMENTS. The authors thank the National Science Foundation (NSF) and Japan Society for the Promotion of Science for the U.S.–Japan

cooperative science program, which made this work possible. This work was partly supported by Grant-in-Aid for scientific research in priority areas “Chemistry of Biological Processes Created by Water and Biomolecules” from the Ministry of Education, Science, Sports, and Culture of Japan (to S.T.). In addition to the support of the NSF, R.S.B. wishes to thank the Aspen Center for Physics for its hospitality in providing the environment where the U.S. contribution to this work was completed.

1. Onuchic JN, Luthey-Schulten Z, Wolynes PG (1997) Theory of protein folding: The energy landscape perspective. *Annu Rev Phys Chem* 48:545–600.
2. Bryngelson JD, Onuchic JN, Socci ND, Wolynes PG (1995) Funnels, pathways, and the energy landscape of protein folding: A synthesis. *Proteins* 21:167–195.
3. Fersht A (1999) *Structure and Mechanism in Protein Science: A Guide to Enzyme Catalysis and Protein Folding* (Freeman, New York).
4. Takada S (1999) Go-ing for the prediction of protein folding mechanisms. *Proc Natl Acad Sci USA* 96:11698–11700.
5. Sabelko J, Ervin J, Gruebele M (1999) Observation of strange kinetics in protein folding. *Proc Natl Acad Sci USA* 96:6031–6036.
6. Hamada D, Segawa S, Goto Y (1996) Non-native alpha-helical intermediate in the refolding of beta-lactoglobulin, a predominantly beta-sheet protein. *Nat Struct Biol* 3:868–873.
7. Fernandez A, Colubri A, Berry RS (2000) Topology to geometry in protein folding: β -Lactoglobulin. *Proc Natl Acad Sci USA* 97:14062–14066.
8. Berry RS, Elmaci N, Rose JP, Vekhter B (1997) Linking topography of its potential surface with the dynamics of folding of a protein model. *Proc Natl Acad Sci USA* 94:9520–9524.
9. Despa F, Wales DJ, Berry RS (2005) Archetypal energy landscapes: Dynamical diagnosis. *J Chem Phys* 122:024103.
10. Becker OM, Karplus M (1997) The topology of multidimensional potential energy surfaces: Theory and application to peptide structure and kinetics. *J Chem Phys* 106:1495–1517.
11. Sibani P, Schon JC, Salamon P, Andersson JO (1993) Emergent hierarchical structures in complex-system dynamics. *Europhys Lett* 22:479–485.
12. Krivov SV, Karplus M (2002) Free energy disconnectivity graphs: Application to peptide models. *J Chem Phys* 117:10894–10903.
13. Krivov SV, Karplus M (2004) Hidden complexity of free energy surfaces for peptide (protein) folding. *Proc Natl Acad Sci USA* 101:14766–14770.
14. Caffisch A (2006) Network and graph analyses of folding free energy surfaces. *Curr Opin Struct Biol* 16:71–78.
15. Rao F, Caffisch A (2004) The protein folding network. *J Mol Biol* 342:299–306.
16. Ihalainen JA, et al. (2008) α -Helix folding in the presence of structural constraints. *Proc Natl Acad Sci USA* 105:9588–9593.
17. Fujitsuka Y, Chikenji G, Takada S (2006) SimFold energy function for de novo protein structure prediction: Consensus with Rosetta. *Proteins* 62:381–398.
18. Fujitsuka Y, Takada S, Luthey-Schulten ZA, Wolynes PG (2004) Optimizing physical energy functions for protein folding. *Proteins* 54:88–103.
19. Takada S, Luthey-Schulten Z, Wolynes PG (1999) Folding dynamics with nonadditive forces: A simulation study of a designed helical protein and a random heteropolymer. *J Chem Phys* 110:11616–11629.
20. Chikenji G, Fujitsuka Y, Takada S (2003) A reversible fragment assembly method for de novo protein structure prediction. *J Chem Phys* 119:6895–6903.
21. Chikenji G, Fujitsuka Y, Takada S (2004) Protein folding mechanisms and energy landscape of src SH3 domain studied by a structure prediction toolbox. *Chem Phys* 307:157–162.
22. Chikenji G, Fujitsuka Y, Takada S (2006) Shaping up the protein folding funnel by local interaction: Lesson from a structure prediction study. *Proc Natl Acad Sci USA* 103:3141–3146.
23. Vincent JJ, Tai CH, Sathyanarayana BK, Lee B (2005) Assessment of CASP6 predictions for new and nearly new fold targets. *Proteins* 61 Suppl 7:67–83.
24. Jolliffe IT (2002) *Principal Component Analysis* (Springer, New York).
25. Das P, Moll M, Stamati H, Kaviraki LE, Clementi C (2006) Low-dimensional, free-energy landscapes of protein-folding reactions by nonlinear dimensionality reduction. *Proc Natl Acad Sci USA* 103:9885–9890.
26. Watts DJ, Strogatz SH (1998) Collective dynamics of ‘small-world’ networks. *Nature* 393:440–442.
27. Barabasi AL, Albert R (1999) Emergence of scaling in random networks. *Science* 286:509–512.
28. Gillespie DT (1977) Exact stochastic simulation of coupled chemical-reactions. *J Phys Chem* 81:2340–2361.
29. Ohmine I, Saito S (1999) Water dynamics: Fluctuation, relaxation, and chemical reactions in hydrogen bond network rearrangement. *Acc Chem Res* 32:741–749.
30. McCallister EL, Alm E, Baker D (2000) Critical role of beta-hairpin formation in protein G folding. *Nat Struct Biol* 7:669–673.
31. Li J, et al. (2007) An alpha-helical burst in the src SH3 folding pathway. *Biochemistry* 46:5072–5082.
32. Riddle DS, et al. (1999) Experiment and theory highlight role of native state topology in SH3 folding. *Nat Struct Biol* 6:1016–1024.
33. Lipman EA, Schuler B, Bakajin O, Eaton WA (2003) Single-molecule measurement of protein folding kinetics. *Science* 301:1233–1235.
34. Ceconi C, Shank EA, Bustamante C, Marqusee S (2005) Direct observation of the three-state folding of a single protein molecule. *Science* 309:2057–2060.
35. Hardin C, Eastwood MP, Luthey-Schulten Z, Wolynes PG (2000) Associative memory Hamiltonians for structure prediction without homology: α -Helical proteins. *Proc Natl Acad Sci USA* 97:14235–14240.
36. Simons KT, Kooperberg C, Huang E, Baker D (1997) Assembly of protein tertiary structures from fragments with similar local sequences using simulated annealing and Bayesian scoring functions. *J Mol Biol* 268:209–225.
37. Kabsch W, Sander C (1983) Dictionary of protein secondary structure: Pattern recognition of hydrogen-bonded and geometrical features. *Biopolymers* 22:2577–2637.