

Folding Time Distributions as an Approach to Protein Folding Kinetics

Sergei F. Chekmarev,^{*,†} Sergei V. Krivov,[‡] and Martin Karplus^{*,‡,#}

Institute of Thermophysics, 630090 Novosibirsk, Russia, Laboratoire de Chimie Biophysique, ISIS, Université Louis Pasteur, 67000 Strasbourg, France, and Department of Chemistry & Chemical Biology, Harvard University, Cambridge, Massachusetts 02138

Received: July 6, 2004; In Final Form: December 16, 2004

A 27-residue lattice heteropolymer subject to Monte Carlo dynamics on a simple cubic lattice is studied over a range of temperatures. Folding time distributions are used to obtain information concerning the details of folding kinetics. The results are compared with those from methods based on mean force surfaces expressed in terms of a reduced set of variables and on a disconnectivity graph for the same system. A detailed analysis of the folding trajectories is given, and the importance of dead-end traps in determining the folding time is demonstrated. We show that the calculated folding kinetics can be modeled by a system of kinetic equations, with the essential rate constants determined from the Monte Carlo simulations and the resulting folding time distributions. The kinetic equations make possible an analysis of the variation of the importance of different channels with temperature. In particular, we show that the presence of intermediates may be masked in the folding time distributions, with the mean folding time being independent of the height of the barrier between the intermediates and collapsed globule state of the system. This and other results demonstrate that care has to be used in interpreting experimental folding data in terms of the underlying kinetics. Correspondingly, simulations are shown to have to satisfy certain requirements to obtain proper sampling of the dead-end traps.

I. Introduction

A full understanding of how proteins fold into their functional three-dimensional native structures has not been achieved, although considerable progress has been made in recent years. The problem has attracted the attention of biologists, chemists, and physicists for at least four decades.¹ The availability of an ever-increasing number of genome sequences and the realization that misfolding is an important cause of disease,² including Alzheimer's disease and spongiform encephalopathies, have greatly increased the focus on this problem. In late 1960s, Levinthal³ showed that because of the astronomical number of conformations possible for a polypeptide chain, a random search for the unique native structure would require $\sim 10^{11}$ or so years. This result suggested that the folding of a protein must involve a biased search for the native structure, rather than a random exploration of the conformation space. Early models of protein folding were concerned primarily with the structural evolution of the polypeptide chain.^{4,5} Some of these, such as the diffusion–collision⁴ and nuclear-condensation models^{5,6} have been found to be consistent with many recent experiments on protein folding; an example is provided by a set of three-helix bundle proteins.⁷ Concomitantly, what has been called a “new view” of protein folding⁸ has been developed, based on both statistical mechanical models and simulations.^{9–19} In its current state,^{20–23} the new view presumes that there exists a reduced set of variables that can describe the folding process, which involves large changes in many degrees of freedom. A variety of variables have been tried (e.g., a number of hydrogen bonds,

rmsd, radius of gyration); none appear to be satisfactory for folding, in general. For lattice folding simulations, the subject of this paper, two variables are commonly used; one determines the compactness of the conformations and serves as a “progress variable” (e.g., the total number of contacts or radius of gyration), and the other characterizes the similarity of the conformations to the native state (e.g., the number of the native contacts). Such variables have also been used for off-lattice simulations with all-atom molecular mechanics potentials.²⁴ With these variables, the effective statistical surfaces of the system, such as those concerned with the free energy, energy and entropy, can be built and used to analyze the folding trajectories. Of particular importance are the free energy surfaces, which show how folding pathways can differ and which of the pathways are favored in folding. Typically, under native conditions, these surfaces correspond to a single major basin with the energy biased toward the native state, and the less ruggedness there is on the surface, the faster is the folding. The surfaces are often described as “folding funnels” when plotted as a function of the energy change associated with the folding process.^{13,18,25} So far little is known about the detailed properties of the basins in proteins because of the lack of experimental data and the fact that simulations of protein folding with realistic models are still limited. Deviations from the single basin concept give rise to a diversity of folding scenarios. Under certain conditions (depending on the temperature or pH) the system may be trapped in off-pathway minima, which slow the folding process considerably (e.g., cytochrome *c* is a recent example of such a system at high pH²⁶), or it can exhibit “kinetic partitioning”,¹⁶ with a “ridge” on the free energy surface separating fast and slow folding trajectories (lysozyme is a protein that apparently shows such behavior²¹), or it can

* Corresponding authors, e-mail: chekmare@itp.nsc.ru (S.F.C.); marci@tammy.harvard.edu (M.K.).

[†] Institute of Thermophysics.

[‡] Université Louis Pasteur.

[#] Harvard University.

follow what appears to be a downhill folding scenario (e.g., PGK and Ub*G refolding^{27,28}).

Mean force surfaces, such as the projected free energy surface mentioned above, have been found to be useful in many studies of protein folding; reviews are given in refs 20–22. However, these reduced surfaces give an incomplete description because none of the sets of known variables are able to describe in full the folding paths of the system through the multidimensional conformation space. Consequently, alternative approaches^{29–32} that do not require projection are of interest. One of these determines all the (significant) minima and saddle points on the energy and/or free energy surface^{30,31,33} and represents them as disconnectivity graphs, in accord with the formulation of Becker and Karplus.²⁹

A source of information complementary to the various representations of the energy and free energy surfaces is provided by an analysis of time-dependent characteristics of the folding process. Experimentally, the populations of the native and/or the unfolded state are usually monitored in the course of protein folding and/or unfolding reaction. The results from such measurements are generally interpreted in terms of one (or possibly two) dimensional models. These are used to estimate the preexponential factor and the free energy barrier for the folding reaction.^{34–36} Very useful conceptualizations of this approach to understanding protein folding are given in the overviews of Socci et al.³⁷ and Gruebele.³⁸ In certain studies,^{35,39} a range for the preexponential factor was estimated from experiment, so that limits for the activation free energy could be determined. Values on the order of a few kcal/mol or less were found for fast folding proteins (e.g., the lambda-repressor^{35,36}). However, such an analysis leaves unanswered the question of whether a one-dimensional reaction coordinate, which projects many different conformations on the same point in the reduced space, results in a smoothing of the actual free energy surfaces.³³

Previous studies have shown that the time dependence of the populations may vary considerably, depending on the folding scenario. Along with the single-exponential distributions characteristic of two-state kinetics, more complex distributions (e.g., multiexponential,^{28,36,37} power law,⁴⁰ and stretched^{28,40,41} and squeezed⁴² exponential distributions) have been examined and used to approximate experimental and simulation data for some of the above-mentioned more complex folding scenarios. Certain aspects of the folding kinetics and the corresponding time-dependent distributions, in particular, can be determined from a master equation governing the state of the system on a funneled energy surface with frustration (characterized by some order parameters).^{43–45} However, a more detailed understanding can be obtained by an explicit examination of the relation between the time-dependent folding characteristics and possible folding mechanisms; it can be useful, for example, to describe the folding process in terms of the rates of transitions between the inherent coarse-grained states of the system, such as the extended state, a semicompact globule, off-pathway intermediates and the native state. The goal of the present paper is to propose such an analysis for lattice model simulations that make it possible to obtain reliable converged results. Of course, it would be more interesting to make a corresponding analysis for all-atom molecular mechanics off-lattice simulations of protein folding. However, the very large number of folding trajectories required for a definitive study of a protein have not yet been possible to obtain; some peptides with a folding transition have been studied in implicit or explicit solvent (e.g., a β -hairpin,⁴⁶ and Beta3s⁴⁷). Go models,¹⁰ which lead to nano-

second folding can also be used,^{48,49} but it has been shown that care is required in interpreting such simulations.⁵⁰

Thus, simplified models continue to be of interest due to their reasonable balance between reality and simplicity. To be useful, such models must be simple enough to permit many folding trajectories to be calculated, yet realistic enough so that they can represent a polypeptide chain. In particular, the number of possible conformations must be significantly larger than that which is sampled in a folding trajectory. Lattice models, which we deal with here, satisfy these criteria, i.e., it is possible to generate as many as 10^5 folding trajectories and the ground state for the potential is easily identified. Specifically, we used the widely studied 27-residue heteropolymer on a cubic lattice (Sali et al.¹⁴ and the references therein) with Monte Carlo dynamics employed to simulate the kinetics of folding.

Given the lattice model simulations, a simple kinetic model is developed to determine what “fingerprints” the various states of the system introduce into the time-dependent distributions. We focus on the folding time distribution, which determines the probability that the system, which starts in the unfolded state, reaches the native state in a given time. This distribution, which is proportional to the time derivative of the native state population, is more sensitive to details of the folding process than the latter. Therefore, it offers more detailed information about the folding kinetics, and it can, in principle, be measured in protein folding experiments. Related work is described in the discussion section.

The paper is organized as follows. Section II presents the model used in the simulation. We specify the energy function of the 27-bead heteropolymer, describe the Monte Carlo procedure used for the heteropolymer dynamics on the cubic lattice, and define the thermodynamic functions to be employed. Section III presents the results of the simulations and their preliminary analysis. The goal of this section is to establish a connection between the dynamics of the system and the reduced statistical surfaces, which depend on the number of total and native contacts. For this purpose, the energy, residence probability, free energy, and entropy surfaces are constructed, and the folding trajectories mapped onto these surfaces are analyzed. In Section IV, this analysis is supplemented by a free energy disconnectivity graph, which determines transitions between the inherent conformations of the heteropolymer in an explicit and accurate way. Section V describes the kinetic model used for the interpretation of the folding time distributions and discusses the general properties of the analytical solutions obtained from this model. In Section VI the rate constants for the transitions between the inherent states of the system are estimated over a wide temperature range by fitting the analytical solutions to the folding time distributions obtained in the simulations. Using the rate constants, a general discussion of the folding process is given, including the variation of the importance of different channels with the temperature. In Section VII, the statistical scatter of the calculated folding times due to the finite number of folding trajectories is discussed, and Monte Carlo simulations based on the kinetic scheme underlying the analytic model are conducted to reproduce these effects. Section VIII summarizes the results and presents some concluding remarks.

II. Specification of the System and Methodology

As the model system, we consider a 27-bead heteropolymer on a cubic lattice. Following Šali et al.,¹⁴ the energy of the system is defined as

$$E = \sum_{i < j} \Delta(r_i - r_j)(B_{ij} + B_0) \quad (1)$$

where r_i is the position of monomer i on the lattice, $\Delta(r_i - r_j)$ is equal to 1 if monomers i and j are in contact and to 0 otherwise, B_{ij} is the interaction energy between monomers i and j , and B_0 is a mean attraction energy between monomers that accounts for hydrophobic effects. For the energy matrix, we use the B_{ij} values of Dinner et al.,⁵¹ which were selected from separate Gaussian distributions corresponding to native and nonnative contacts. For native contacts, the B_{ij} have the mean $\epsilon \approx -1.67$ and standard deviation $\sigma \approx 1.0$, and for nonnative contacts, $\epsilon \approx 0.0$ and $\sigma \approx 1.0$, respectively. The value of B_0 was chosen to be -2.0 . The native (lowest energy and free energy) state of the heteropolymer is a fully compact $3 \times 3 \times 3$ -cubic structure (see Figure 1) with energy $E = -106.0917$.

To generate a folding trajectory for the heteropolymer on the lattice, the Metropolis Monte Carlo (MC) method was employed. The temperature, T , measured in the units of the energy in eq 1 (i.e., with the Boltzmann constant $k_B = 1$), ranged from 1.5 to 2.3. As in previous work,¹⁴ three types of moves were allowed; they are end flips, corner flips, and two-bead crankshaft rotations. The step (time) counter is advanced whether or not the current move is accepted; this means that the number of MC steps coincides with the number of MC trials.

The trajectories were started from an unfolded chain and terminated upon reaching the native state (Figure 1), so that the "folding time" corresponds to the first passage time. The unfolded chains were generated by adding successive beads at one end of the heteropolymer, each of the beads being randomly oriented with respect to the preceding bead in one of the three positive directions. For each temperature, 5×10^4 trajectories were run with the maximum length of the trajectories equal to 2×10^8 MC steps; this time was long enough so that all trajectories reached the native state. At fixed intervals, typically every 10^3 steps, the total number of contacts (N), the number of native contacts (N_{nat}), and the energy of the system for the given number of contacts, defined by eq 1, were calculated. Collecting these data for the ensemble of trajectories, we calculated the probability for the system to be in a state with N and N_{nat} contacts, $P(N, N_{\text{nat}})$, and the mean energy of the system at this point, $\bar{E}(N, N_{\text{nat}})$. These data were used to build the surfaces for the residence probabilities and mean energies as of function of N and N_{nat} . Since the trajectories were terminated upon reaching the native state, the functions $P(N, N_{\text{nat}})$ and $\bar{E}(N, N_{\text{nat}})$ represent nonequilibrium conditions, characterizing the kinetics of the system rather than its thermodynamics. We will show in Section IV that the function $P(N, N_{\text{nat}})$ reflects well the equilibrium residence probabilities of the globule and dead-end (DE) traps but leads to a very different value of the residence probability of the native state. However, exactly this function corresponds to the analytic model we use for the analysis of the folding time distributions (Section V and Table 1).

Having the functions $P(N, N_{\text{nat}})$ and $\bar{E}(N, N_{\text{nat}})$, we also calculated free energy of the system

$$F(N, N_{\text{nat}}) = -T \ln P(N, N_{\text{nat}}) + C \quad (2)$$

where C is the arbitrary constant, and the entropy

$$S(N, N_{\text{nat}}) = [\bar{E}(N, N_{\text{nat}}) - F(N, N_{\text{nat}})]/T \quad (3)$$

although these thermodynamic relations are of limited application in the present case, because the functions $P(N, N_{\text{nat}})$ and $\bar{E}(N, N_{\text{nat}})$ are not equilibrium through all the (N, N_{nat}) plane.

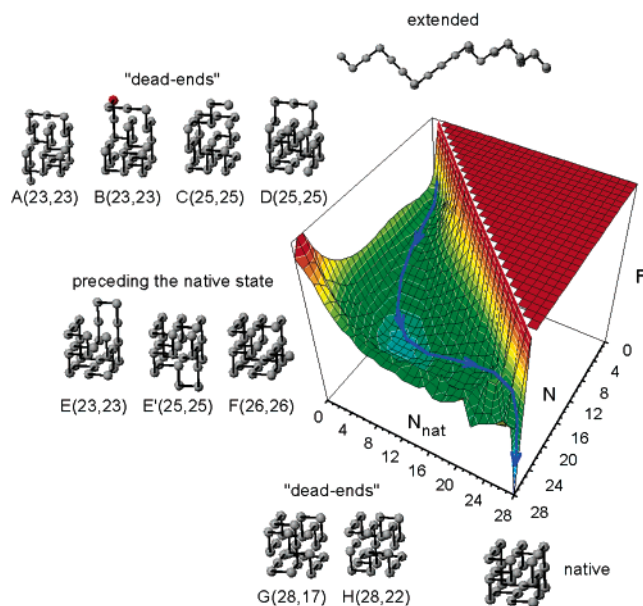


Figure 1. Schematic free energy surface for folding of the 27-bead heteropolymer. Structures A to H present characteristic configurations of the heteropolymer. In brackets following the labels of the structures there are indicated the numbers of total and native contacts, (N, N_{nat}) .

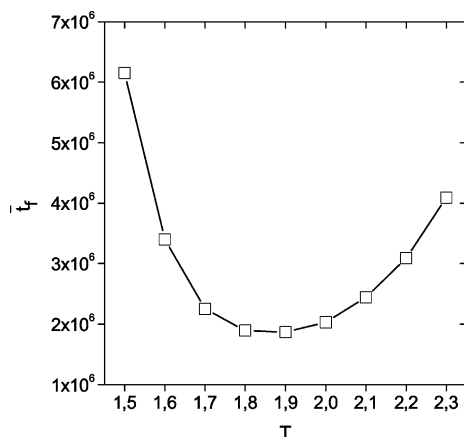
In the native state and the transition state region, connecting the native state with the globule, the first passage time sampling is essentially nonequilibrium, because the system is not allowed to either reside in the native state or explore the surroundings of the native state when it denaturates. However, the region represented by the globule and DE minima is sampled quite thoroughly, since the system spends a considerable time in these minima before it finds a way to the native state. Therefore, the functions $F(N, N_{\text{nat}})$ and $S(N, N_{\text{nat}})$, defined by eqs 2 and 3, are expected to be useful in characterization of this region of the (N, N_{nat}) plane, playing a key role in determining the first passage times, which are the goal of the present study.

As in lattice MC simulations in general, the energy $\bar{E}(N, N_{\text{nat}})$ in eq 3 is the effective potential energy, rather than the total energy; it includes, implicitly, an averaged solvation free energy. The entropy given by this equation represents the configuration entropy, which is determined by the number of configurations possible for the given values of N and N_{nat} (Appendix). According to this, the entropy of the native state should be set equal to zero, $S(28, 28) = 0$, because the native state ($N = N_{\text{nat}} = 28$) is represented by a single structure. Correspondingly, by eq 3, $F(28, 28) = \bar{E}(28, 28)$, which specifies the constant C in eq 2. We emphasize that the function $P(N, N_{\text{nat}})$ differs from the actual equilibrium function in the native state and transition state region (see Section IV) because the trajectories were terminated when the native state was reached. Thus, the above procedure serves for a specification of the arbitrary constant C in eq 2 but does not give the correct relation of the free energy and entropy of the native state with that of the globule and DE minima. In particular, in comparison with the equilibrium conditions, the free energy of the globule and DE minima is decreased approximately by $T \ln P_{\text{nat}}^{\text{eq}}$, and their entropy is increased by $\ln P_{\text{nat}}^{\text{eq}}$, where $P_{\text{nat}}^{\text{eq}}$ is the equilibrium residence probability of the system in the native state, which has not been determined. This, however, does not change the relative free energies and entropies surfaces in the part of the (N, N_{nat}) plane containing the globule and DE minima, which is of primary interest in present analysis.

TABLE 1: Fraction of the Time Spent by the System in the DEs

T	1.5	1.7	1.9	2.1	2.3
simulations ^a	0.429 (3,60,37)	0.139 (5,58,37)	0.025 (6,49,45)	0.007 (7,40,53)	0.002 (7,30,63)
kinetic model	0.486	0.164			

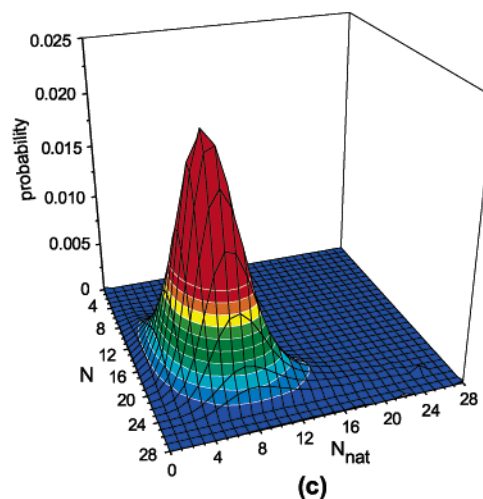
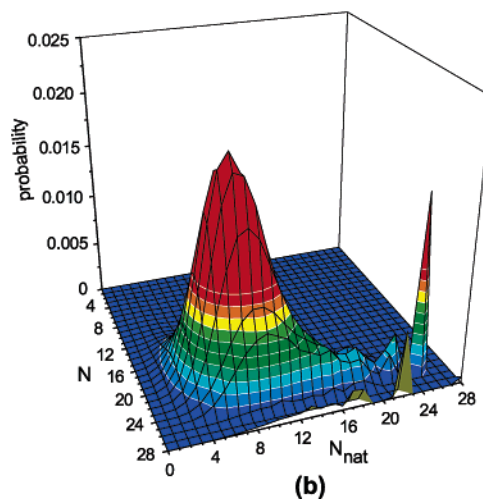
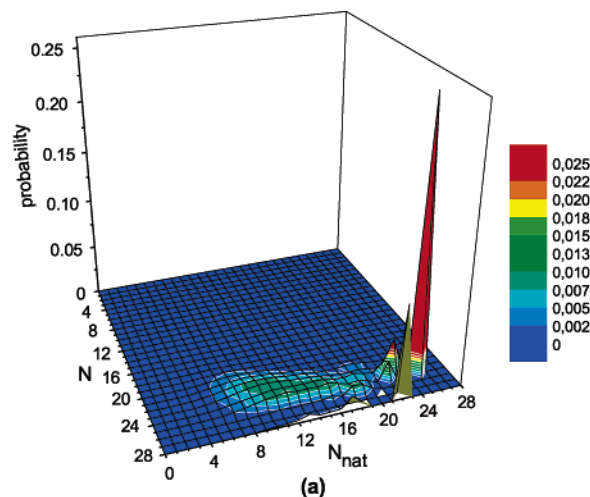
^a In the brackets, the contribution of the $N = N_{\text{nat}} = 23$, $N = N_{\text{nat}} = 25$, and $N = 28$ minima (%).

**Figure 2.** Mean folding time vs temperature. The line is to guide the eye.

We note that the surfaces calculated for the relatively small 27-mer directly from the large number of folding simulations are sufficiently accurate so that no constraints and/or special sampling methods were required, in contrast, for example, to the significantly larger (125-mer) system studied by Dinner and Karplus.⁵²

III. Simulation Results and Preliminary Analysis

Figure 2 shows how the mean folding time (MFT) \bar{t}_f , measured in terms of the number of MC steps, varies with temperature. This dependence exhibits the well-known U-shape behavior found in theoretical models and in some experiments,^{19,53} with a minimum at the optimal folding temperature T_f . For the present system, T_f is approximately 1.9, as can be seen from the figure. Figures 3 and 4 show, respectively, the residence probability surfaces (RPSs), which determine the probability of finding the system in the state with N and N_{nat} contacts, and the corresponding free energy surfaces (FESs) at three temperatures: below T_f ($T = 1.5$, panel a), at T_f ($T = 1.9$, panel b), and above T_f ($T = 2.3$, panel c). These surfaces reduce the multidimensional conformation space to the two progress variables N and N_{nat} (as shown in the figures). These two variables, although useful, do not give a full description of the folding kinetics. In particular, since each elementary move typically results in the change of several contacts, both for N and N_{nat} , the initial and final states for the move may not be, and most frequently are not, directly connected in the (N, N_{nat}) plane. Due to this, neighboring points along the folding trajectory, when it is mapped onto the FES in terms of (N, N_{nat}) , can be separated by other points in the plane. To describe this complexity, disconnectivity graphs²⁹ are needed (see Section IV). Nevertheless, the reduced surfaces give important information concerning the general character of the system and its folding behavior, as has been shown in previous studies.^{20–23} Figure 4 shows that the model protein is expected to fold roughly in two stages, corresponding to the L-shape pattern of the FESs. In the first stage, a semicompact disorganized globule is formed. It can be seen in Figure 4 that a broad basin with a large total number of contacts (approximately 15 to 25) and a moderate number of native contacts (approximately 5 to 15) is associated

**Figure 3.** Residence probability surfaces for $T = 1.5$ (a), $T = 1.9$ (b), and $T = 2.3$ (c).

with this globule. In the second stage of the folding process, this semicompact globule rearranges and makes a conforma-

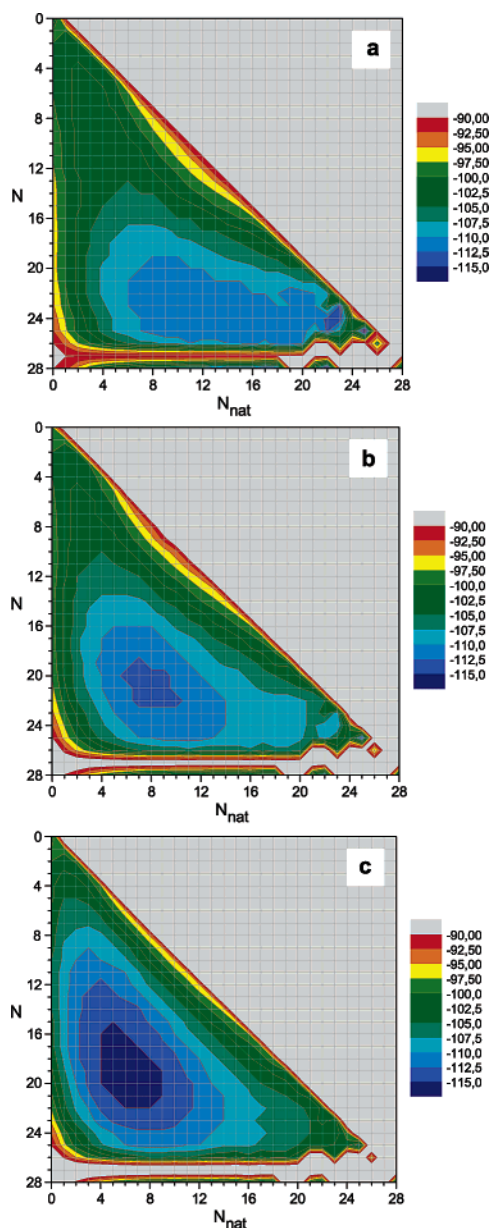


Figure 4. Free energy surfaces for $T = 1.5$ (a), $T = 1.9$ (b), and $T = 2.3$ (c).

tional search for the native structure ($N = N_{\text{nat}} = 28$). Thus, the mechanism observed for the present set of B_{ij} parameters corresponds to that found in earlier 27-mer studies (Sali et al.,¹⁴ Socci and Onuchic⁵⁴) and is similar to those that obtained with more realistic all-atom models for peptides and proteins.⁴⁶

A detailed exploration of the FESs showed that there exists a number of localized minima (basins), in addition to those corresponding to the semicompact random globule and native state. Among them, there are two pronounced minima at $N = N_{\text{nat}} = 23$ and $N = N_{\text{nat}} = 25$. Each of them contains a set of configurations corresponding to the same numbers of total and native contacts. Specifically, the minimum at $N = N_{\text{nat}} = 23$ includes nine such configurations with the energies $-89.8690 \leq E \leq -84.2973$, and that at $N = N_{\text{nat}} = 25$ three configurations, $-96.6085 \leq E \leq -95.2586$; at the temperature T_f the free energies of these minima are equal to -108.707 and -112.6342 , respectively. We note that these minima, because they consist of multiple structures, have lower free energies (but not energies) than the native state, which, according to its specification in Section II corresponds to a single structure with

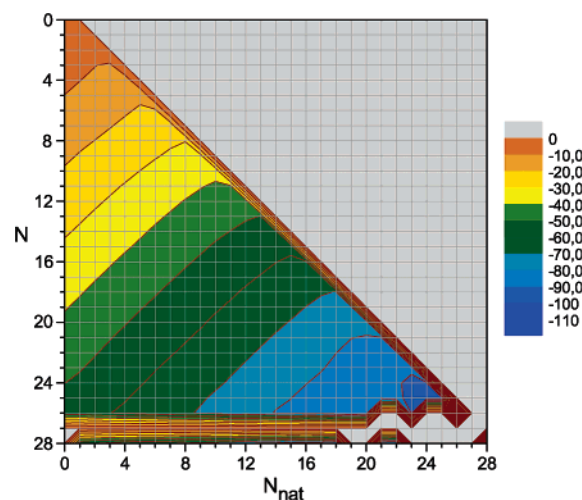


Figure 5. Mean potential energy surface, $T = 1.9$.

zero entropy. This suggests that it would be better to define the native state as a set of minima in the native basin, which have the lowest free energy (see Section IV). Some of the minima configurations, labeled A to E, are shown in Figure 1. There are also two basins associated with the fully compact states with 28 contacts, in which some contacts are nonnative. These nonnative basins, like the native basin, are separated from the rest of the surface by a “wall” along the N_{nat} coordinate at $N = 27$, which corresponds to structures that cannot be formed by a 27-bead lattice polymer because of steric restrictions. They thus correspond to one-dimensional variations in free energy as a function of N_{nat} . However, despite the wall the system is able to visit all of these $N = 28$ states because each of the elementary moves, which can transform a noncompact structure to a fully compact one, results in adding at least two contacts, i.e., the system can “jump” over the wall when passing from the globule basin to the fully compact structures. The structures E and F in Figure 1 give examples of such noncompact structures, which can transform into the fully compact ones in a single move (in the given case, it is into the native state). One of the nonnative basins contains the structures with N_{nat} increasing from 1 to 18, with the minimum of the basin at $N_{\text{nat}} = 17$; each native contact decreases the energy by approximately -1.79 , which is close to the mean value of the B_{ij} in eq 1 ($\epsilon \approx -1.67$). The other $N = 28$ basin is a narrow deep well containing configurations with N_{nat} equal to 21 and 22, with the minimum at $N_{\text{nat}} = 22$. A detailed analysis shows that each of the minima represents several structures with different energies: the minimum at $N_{\text{nat}} = 17$ contains eight structures within the energy interval $-88.3192 \leq E \leq -80.1565$, and the minimum at $N_{\text{nat}} = 22$, two structures with energies -92.6472 and -86.6908 . For each minimum, the residence probabilities of the structures determined from the MC simulation obeys the Boltzmann distribution, that is, the structures are in equilibrium, with the structures of lowest energy having the highest residence probabilities and therefore lowest free energies. Two such structures are shown at the bottom of Figure 1; structure G, with $E = -88.3192$, represents the minimum at $N_{\text{nat}} = 17$, and the structure H, with $E = -92.6472$, represents the minimum at $N_{\text{nat}} = 22$.

Additional information about the system comes from an examination of the potential energy surfaces (PESs) and entropy surfaces (ESs), Figures 5 and 7a,b,c, respectively. Figure 5 shows the PES for $T = T_f = 1.9$. This surface indicates that the mean energy of the system $\bar{E}(N, N_{\text{nat}})$ can be represented as a linear combination of the number of total and native contacts,

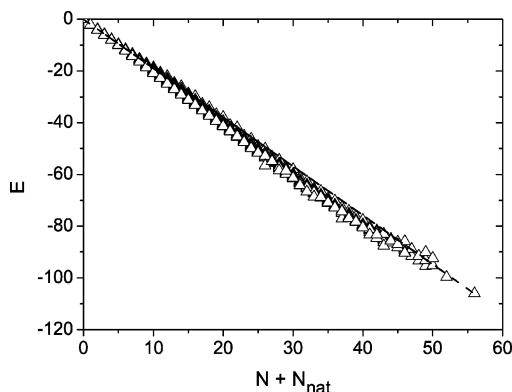


Figure 6. Approximation of the mean potential energy by eq 4.

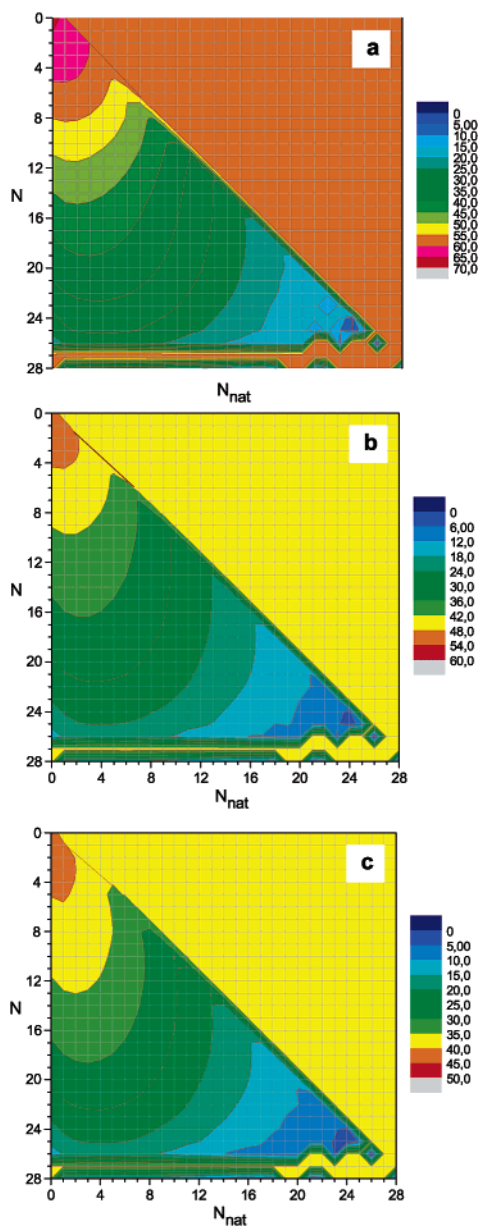


Figure 7. Entropy surfaces for $T = 1.5$ (a), $T = 1.9$ (b), and $T = 2.3$ (c).

specifically, their sum $N + N_{\text{nat}}$. Considering $\bar{E}(N, N_{\text{nat}})$ as a function of $N + N_{\text{nat}}$, we found that for entire surface, including the local minima and the native state, the results are well approximated by the equation

$$\bar{E}(N, N_{\text{nat}}) = -1.895(N + N_{\text{nat}}) = -3.79N_{\text{nat}} - 1.895N_{\text{nonnat}} \quad (4)$$

where $N_{\text{nonnat}} = N - N_{\text{nat}}$ is the number of nonnative contacts (Figure 6); the root-mean-square (standard) deviation of the simulation data from this dependence is ≈ 3.4 . From the second equality, it is seen that the coefficients of N_{nat} and N_{nonnat} are close to the values of the mean energies for the native and nonnative contacts in eq 1 (i.e., -3.67 and -2.0 , respectively); the deviation is approximately 3 and 5% for the first and second coefficient, respectively. Since a group of structures, which differ in energy, generally corresponds to a (N, N_{nat}) point (as, e.g., for the $N = N_{\text{nat}} = 23$ and $N = N_{\text{nat}} = 25$ points mentioned above), the agreement between the coefficients in eq 4 and the mean energies for the native and nonnative contacts in eq 1 suggests that structures with the same numbers of total and native contacts correspond to a random set of configurations. In this case, averaging over the structures at a given (N, N_{nat}) point should lead to eq 4 with coefficients close to the corresponding mean energies, since the energies of the native and nonnative contacts in eq 1 were selected from Gaussian distributions. This is a somewhat surprising result because one might have expected that structures with given N and N_{nat} would be correlated in some way. For example, they could be expected to be close neighbors on the MC trajectory (see also Section IV). Another consequence of the random character of the sets of configurations at the (N, N_{nat}) points is that the average PESs depend weakly on the temperature. The relative root-mean deviation of $\bar{E}(N, N_{\text{nat}})$ at $T = 1.5$ and $T = 2.3$ from that at $T = 1.9$ (Figure 5) is approximately 6 and 1%, respectively.

As has been mentioned in Section II, the entropy defined by eq 3 is the configurational entropy, which characterizes the number of configurations $Q(N, N_{\text{nat}})$ that are possible for the given values of N and N_{nat} . Specifically, according to the appendix, $S(N, N_{\text{nat}}) \sim \ln Q(N, N_{\text{nat}})$. The entropy surfaces presented in Figure 7 thus reveal how the number of configurations changes in the course of folding. In the first stage, where the semicompact globule is formed, the entropy gradually decreases as N increases, which reflects the fact that compacting the globule is accompanied by a decrease in the number of possible conformations. For each given N the entropy (and thus the number of configurations) exhibits a bell-like distribution within the allowed range of variation of N_{nat} ($N_{\text{nat}} \leq N$). A closer examination of the variation of $Q(N, N_{\text{nat}})$ with N_{nat} shows that this distribution is essentially a normal (Gaussian) distribution (not shown), which is in accord with its combinatorial origin. In the second stage, when folding corresponds to a search of the globule for the native state, the entropy continues to decrease, but here mainly as a function of increasing N_{nat} , as the system approaches the native state. We note that the entropy exhibits local minima exactly at the points where the energy surface has local minima. This indicates that the structures with minimum energy are less degenerate than the others, i.e., they correspond to particular conformations of the polymer chain.

The reduced FESs do not discriminate among the minima with fixed N and N_{nat} . However, it is of interest to determine for which minima of the system, which is moving toward the native state, has to return to the semicompact random globule basin with fewer contacts (off-pathway), and configurations representing intermediate states (on-pathway) from which the system can proceed to the native state without breaking contacts. This information can be obtained from the folding trajectories. Typical examples of the trajectories are presented in Figure 8. For each temperature, two trajectories are shown: one corresponding to a folding time t_f below the MFT \bar{t}_f for the given

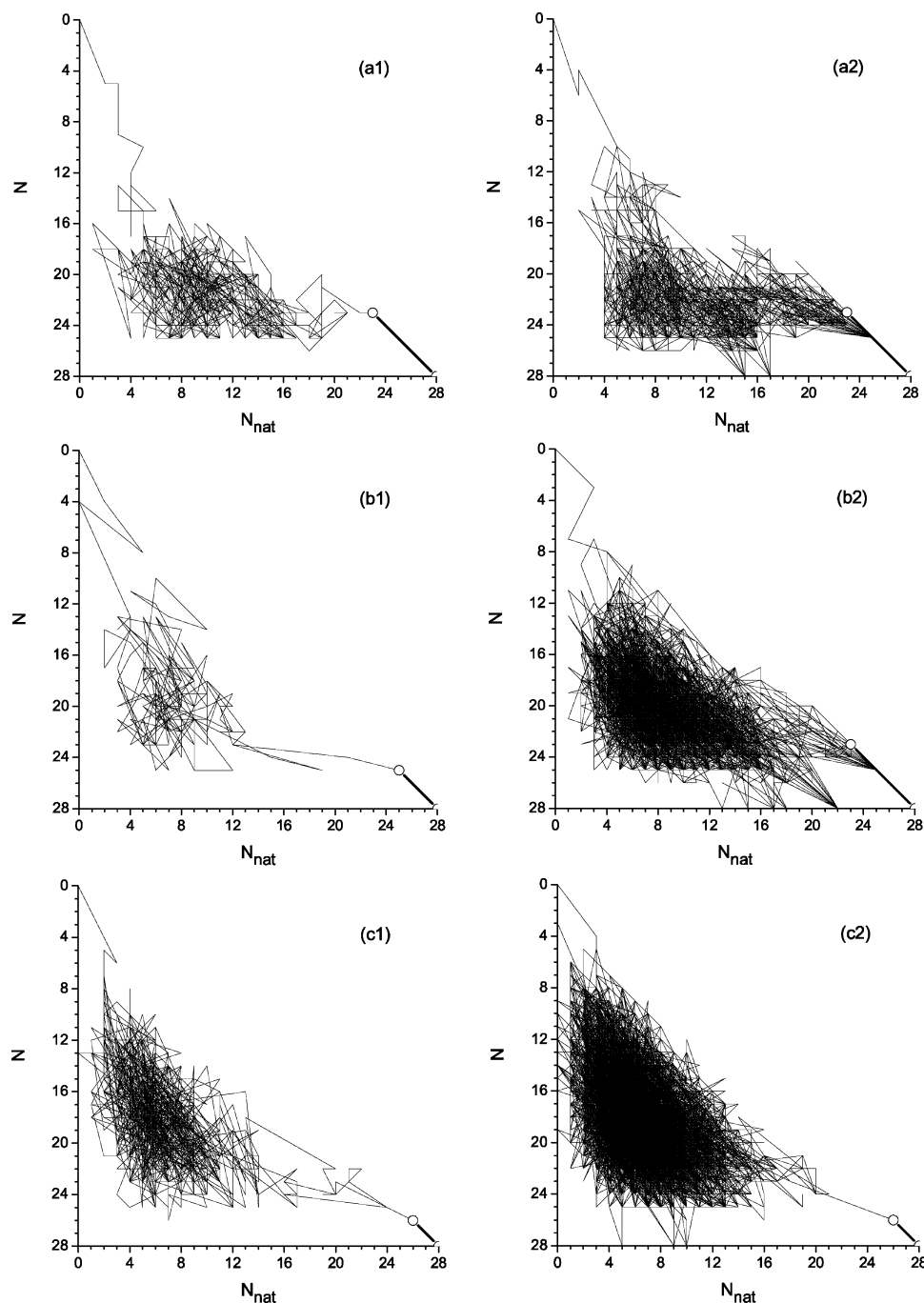


Figure 8. Typical folding trajectories for $T = 1.5$ (a), $T = 1.9$ (b), and $T = 2.3$ (c). Continuous lines connect the points separated by the interval of 10^3 steps, except for the last two points. They are indicated by circles and connected by thick lines; the latter points correspond to the native state and the configurations preceding it. For each temperature, two trajectories are shown that correspond to the folding time below and above the mean folding time $t_f(T)$. For more details see the text.

temperature (see Figure 2), and the other to one above \bar{t}_f . Specifically, panels a1 and a2 ($T = 1.5$) correspond to $t_f \approx 0.2 \bar{t}_f$ and $t_f \approx 2.5 \bar{t}_f$, panels b1 and b2 ($T = 1.9$) correspond to $t_f \approx 0.1 \bar{t}_f$ and $t_f \approx 3.6 \bar{t}_f$, and panels c1 and c2 ($T = 2.3$) correspond to $t_f \approx 0.2 \bar{t}_f$ and $t_f \approx 2.3 \bar{t}_f$. Continuous lines connect the points separated by an interval of 10^3 steps, except for the last two points. They are indicated by circles and connected by thick lines; these two points correspond to the native state and the configurations preceding it.

Figure 8 indicates that the minima corresponding to the fully compact nonnative states (G and H in Figure 1, and the other points with $N = 28$ in panels a2, b2, and c2 of Figure 8) are “dead-ends” (DEs), i.e., the system cannot proceed to the native

state without a return to the semicompact random globule basin. This is also true for some of the nonfully compact configurations with $N = N_{\text{nat}} = 23$ (such as A and B) and $N = N_{\text{nat}} = 25$ (C and D); in all cases, none of the allowed moves is able to transform these configurations into one closer (in terms of number of contacts) to the native state. Further analysis shows that some of the traps are the result of the lattice representation with a restricted move set. For configuration C, for example, a two-bead crankshaft rotation at the end of the polymer would be sufficient to pass into the native state. A similar situation occurs for configurations A and D, where a three-bead crankshaft rotation would be sufficient. Although the latter move is evidently less probable than the allowed two-bead rotation,

it would make possible a direct interconversion of configurations A and D into the native configuration and the configuration F preceding it, respectively. The configuration B is similar to configuration D, except that the monomer shown in red needs to be subjected to the corner flip, which adds two native and two total contacts. This configuration is thus a kinetic neighbor of the configuration D, although it belongs to a point of the (N, N_{nat}) plane that is not directly connected to the $N = N_{\text{nat}} = 25$ point.

Given the above, we note that Figure 8 illustrates an essential difference between the short and long trajectories. For the short trajectories (a1, b1, and c1) the system does not visit the DEs, whereas in the long trajectories (a2, b2, and c2), it spends a long time in these off-pathway minima. This is particularly evident for the DEs at $N = N_{\text{nat}} = 25$ and $N = 28$, and suggests that the longer folding times arise either from trapping the system in such minima, which is observed at low temperatures, or from repeated visiting of these minima, which is characteristic of higher temperatures (see also below).

For a more detailed insight into the folding process, a set of 10^3 trajectories, which were started at random extended configurations and continued until reaching the native state, as before, was run at the three temperatures. However, the state of the system was analyzed every 10 MC step, i.e., much more frequently than above. At the folding temperature ($T = T_f = 1.9$), approximately 90% of the trajectories were found to pass through the minimum at $N = N_{\text{nat}} = 23$, 92% through the minimum at $N = N_{\text{nat}} = 25$, and 68% through the basins for the fully compact nonnative structures, $N = 28$. None of the aforementioned minima was visited in less than 1% of the cases. Among the trajectories that visited the minimum at $N = N_{\text{nat}} = 23$, approximately 13% passed through configurations A and B, which are DEs, 63% through the configuration E, which is one of the structures preceding the native state, and 14% through the configuration, which has the defects of both configurations E and F. It requires two successive moves to reach the native state; i.e., a crankshaft rotation and end flip are required. In the case of the $N = N_{\text{nat}} = 25$ minimum, 39% of the trajectories passed through configurations C and D, which represent DEs, and 53% through a configuration E', which precedes the native state and differs from the configuration E in that the two monomers subject to the crankshaft rotation are placed, after rotation, on the distant upper edge of the cube and not on the upper face (see Figure 1).

A further analysis of the configurations preceding the native state shows that a transition to the native state occurs either through the configuration F ($N = N_{\text{nat}} = 26$, approximately 41% of events at $T = 1.9$), or through those with $N = N_{\text{nat}} = 23$ ($\approx 56\%$ of the events) or through configuration E' ($N = N_{\text{nat}} = 25$, $\approx 3\%$). Among the configurations with $N = N_{\text{nat}} = 23$, the most frequent is configuration E ($\approx 36\%$ of the events). The other configurations with $N = N_{\text{nat}} = 23$ differ from configuration E by the position of two monomers, which are subjected to the allowed crankshaft rotation. In other words, an immediate transition to the native state occurs either through an end flip (configuration F) or crankshaft rotation (configuration E and the other configurations with $N = N_{\text{nat}} = 23$ and with $N = N_{\text{nat}} = 25$ that have been just mentioned). No corner flips were observed as part of the final events in the set of 10^3 trajectories employed for the analysis.

A useful description of the folding dynamics is obtained by dividing the trajectories into three groups, depending on a set of the DEs they visited: the DEs of all types (at $N = N_{\text{nat}} = 23$, $N = N_{\text{nat}} = 25$, and $N = 28$), the DEs of two types (at $N =$

$N_{\text{nat}} = 23$ and $N = N_{\text{nat}} = 25$, etc.), or solely one of them. We found that about 10% of the trajectories visited all of the DEs. Among the trajectories that passed through two DEs, 19% correspond to the basins of the fully compact structures and the DE at $N = N_{\text{nat}} = 25$, 2% to the DEs at $N = N_{\text{nat}} = 23$ and $N = N_{\text{nat}} = 25$, and less than 1% to the fully compact nonnative structures ($N = 28$) and the DE at $N = N_{\text{nat}} = 23$. The percentage of the trajectories that pass through only one of the DEs is as follows: approximately 39% through the basins for fully compact structures, 7% through the minimum at $N = N_{\text{nat}} = 25$, and less than 1% through the DE at $N = N_{\text{nat}} = 23$. In total, among the trajectories that passed through $N = N_{\text{nat}} = 23$, $N = N_{\text{nat}} = 25$, and $N = 28$ basins, 78% visited the DEs associated with these basins, and 22% passed through the structures preceding the native state.

Variation of the temperature (from 1.5 to 2.3) does not lead to a large change in the distribution of the trajectories from that found for the folding temperature $T = T_f = 1.9$, although there is a tendency to visit a larger number of the DE minima at higher temperatures. For example, the fraction of the trajectories that passed through the fully compact nonnative structures and the DE at $N = N_{\text{nat}} = 25$ increases from 9% at $T = 1.5$ to 24% at $T = 2.3$, while the fraction of the trajectories that passed solely through the former, without visiting other DE minima, decreases from 50 to 37%, respectively. However, the fraction of the time spent by the system in the DEs decreases drastically with temperature. It can be calculated as

$$P_d = \frac{\sum_{N, N_{\text{nat}}}^{(d)} P(N, N_{\text{nat}})}{\sum_{N, N_{\text{nat}}} P(N, N_{\text{nat}})} \quad (5)$$

where $P(N, N_{\text{nat}})$ is the probability for the system to be in a state with N and N_{nat} contacts, and the superscript "d" on the summation symbol in the numerator indicates that the summation is taken over the DEs. Correspondingly, the fraction of the time spent by the system in the globule region can be estimated as

$$P_g \approx 1 - P_d \quad (6)$$

since the collapse of the unfolded state into the globule occurs very rapidly. Table 1 shows how the P_d (5) changes with temperature. The relative contribution of different DE minima to P_d also varies with temperature; it is indicated, in percent, in the brackets, with the first, second, and third figures corresponding to the minima at $N = N_{\text{nat}} = 23$, $N = N_{\text{nat}} = 25$ and $N = 28$, respectively. The increase of the MFT \bar{t}_f at low temperatures $T < T_f$ (Figure 2) results dominantly from the considerable increase of the time spent by the system in the DE minima. Thus, these minima form an essential part of the folding process, and the time spent in them before escaping contributes to the overall folding time. For more elevated temperature, $T > T_f$, \bar{t}_f increases because the system searches through a larger portion of the conformation space in the globule region (Figure 3), which is reflected by the increase of the P_g (see eq 6 and Table 1). This is in accord with the conclusion that the activation free energy is entropy dominated at higher temperatures (see Section VI).

IV. Transition Disconnectivity Graph

A more detailed view of the free energy surfaces and the resulting folding kinetics is provided by disconnectivity graphs.

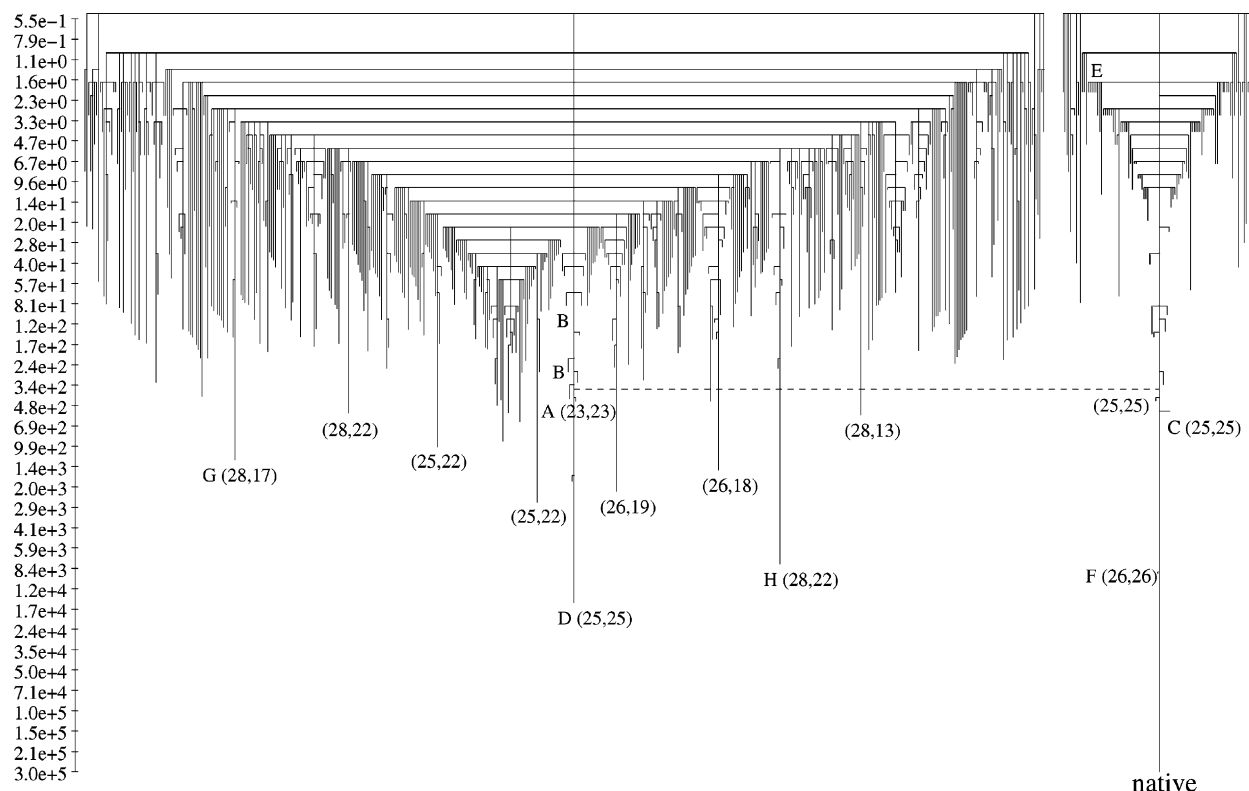


Figure 9. Semiequilibrium TRDG of the 27-bead heteropolymer at $T = 1.5$ (not all conformers with small barriers are shown to keep the size of the graph moderate). The axis shows the number of times the conformers were visited and the number of transitions the system made between different conformers. Numbers in brackets show the number of native contacts and total number of contacts for a particular conformer. Letters denote conformers according to Figure 1. Dashed line denotes a barrier between the denatured and native basins, which are explicitly separated. Barriers between individual sub-basins in the denatured basin and the native basin are considerably higher, approximately 20–50 times for the sub-basins labeled in the figure.

Such graphs were introduced by Becker and Karplus,²⁹ who used them to represent the potential energy surface landscape of complex systems with many degrees of freedom and first applied them to oligopeptides.²⁹ They have been generalized in refs 30–32 to represent free energy surfaces. In the free energy graphs, one way to calculate the free energies of the transition state ensemble between two minima is to use the transition rate between the minima. To construct such transition disconnectivity graph (TRDG), we use the algorithm proposed by Krivov and Karplus,³⁰ in which a capacitated undirected graph is introduced; in this graph the equilibrium flow reproduces the equilibrium systems dynamics. The graph consists of a set of nodes i , which corresponds to the systems states (conformations), and a set of capacitated edges, connecting the nodes with a capacity $\{c_{ij} = c_{ji}\}$ proportional to the equilibrium transition rate between corresponding states (i and j) of the systems $c_{ij} = (n_{ij} + n_{ji})/2$, where n_{ij} is the number of transition from state j to state i obtained in the MC simulation. The partition function of the free energy barrier separating any two nodes is equal to the value of minimum cut between the nodes. Correspondingly, the partition function of state i is calculated as $n_i = \sum_j n_{ij}$, i.e., it is equal to the number of times the system visited the state i . These definitions give relative values, since no overall normalization is introduced. In particular, the partition functions of the structures are proportional to the probabilities of finding the system in these structures. Knowing the free energy of the states and barriers between them allowed us to construct the disconnectivity graph following ref 29. To calculate the transition rates, semiequilibrium sampling was performed at $T = 1.5$, since at this low temperature the free energy surface is anticipated to have a larger number of the free energy minima than at a higher

temperature (Figure 4). As in Section III, each trajectory was started in a random extended conformation (Section II). However, in contrast to the simulations of Section III, the trajectory was not terminated on reaching the native state but was continued for 10^7 steps; then a new trajectory was started with an extended conformation, and so on. Without this procedure the system would spend most of its time in the native basin because the mean unfolding time $\sim 1.7 \times 10^8$ is considerably larger than the MFT $\sim 6.2 \times 10^6$. This makes investigation of the whole TRDG by pure equilibrium sampling difficult at low temperature. With the semiequilibrium procedure, the estimated relative partition function of the native state is expected to be smaller than the actual value approximately by a factor of $1.7 \times 10^8 / 6.2 \times 10^6 \approx 27$. Nevertheless, the resulting TRDG is expected to provide considerable insights concerning the free energy surface. Structures were checked every 10^4 MC steps and grouped based on their set of contacts (total and native). Thus a number of different conformers correspond to the same (N, N_{nat}) point, analogous to what was done in the previous section. The simulations were continued until 60 000 different structures (i.e., with different sets of contacts) had been visited at least once. The resulting TRDG for the whole FES is shown on the Figure 9. On the graph the denatured and native basins are explicitly separated to show the configurations that belong to the two basins. The vertical axis shows the number of times the conformers were visited (n_i), which is equivalent to their relative partition functions; it also shows the number of transition the system made between different conformers, which is equivalent to the relative partition function of the transition state ensemble (or equilibrium reaction rate) between the conformers.

The basins were separated by finding the minimum cut (mincut) between the starting configuration (all the starting configuration are the same in terms of the present analysis, since they have zero contacts) and the native configuration.³⁰ The mincut corresponds to the transition state ensemble separating the two basins, i.e., configurations on different sides of the mincut belong to the different basins.³⁰ The mincut procedure gives the following partition functions: denatured basin, $Z_d = 132426$; transition states ensemble, $Z_{d,n} = 349$; native state, $Z_n = 308992$. In contrast to the first passage time sampling of Section III, the native state partition function is a lower bound, due to the sampling procedure mentioned above. The estimated folding time is $Z_d/Z_{d,n} \times 10^4 \approx 3.8 \times 10^6$ (where 10^4 MC steps is the checking interval), which is close to the MFT obtained by direct counting ($\bar{t}_f \approx 6.2 \times 10^6$, Figure 2).

We note that the standard definition of the native state in lattice model calculations and the one used in the rest of the paper (i.e., the single conformer with $N = N_{\text{nat}} = 28$) is different from that obtained from the mincut procedure. In the latter, the native state is the whole native basin, which is separated from the denatured state by a free energy barrier. This difference, however, is not crucial for the analysis of the folding time distributions, which is the goal of the present paper, because after surmounting the free energy barrier the system rapidly attains the native state (Sali et al.¹⁴). The denatured basin consists of a number of deep sub-basins (with low enthalpy and low entropy) that have high barriers. Specifically, barriers that connect the sub-basins within the denatured basin are much higher than the barrier connecting the denatured basin with the native one. For example, the partition function of the barrier between conformer D and the denatured basin is about 20, which gives a mean transition time from the denatured basin to that conformer of about 1.1×10^8 , which is much higher than the MFT at the given temperature ($\bar{t}_f \approx 6.2 \times 10^6$).

As “dead ends” we may consider sub-basins, which are separated from the denatured basin by such a high barrier, as a result the mean transition time between the two basins is comparable to the folding time. Configurations D, G, and H, in accordance with the analysis above, can be considered as DEs, since they are in the denatured basin and separated from it by high barriers. Configuration C, as can be seen from the graph, belongs to the native basin and has a small barrier, and it thus cannot be considered as a DE trap. Configuration A is in the same sub-basin as D, as well as all configurations B. Relative residence probabilities of the DE configurations are in good correspondence with the results of the first passage time sampling. In particular, according to Figure 9, $Z_H/Z_D \approx 0.5$ and $Z_G/Z_D \approx 0.04$, and according to Figure 3, $P_H/P_D \approx 0.37$ and $P_G/P_D \approx 0.05$, where D, H, and G label the corresponding configurations. At the same time, the results for the native state differ drastically, e.g., $Z_{\text{nat}}/Z_D \approx 18.0$ (Figure 9), whereas P_{nat}/P_D is as small as 6×10^{-4} (Figure 3), since in the latter case the trajectory was terminated upon reaching the native configuration. As one can see from Figure 9, the denatured basin contains also other configurations, which can be considered to be DEs because they have relatively high partition functions and separated by high barriers. Also note that the high probability for the system to visit the F conformer before folding is explained by close position of this conformer to the native state.

V. Folding Kinetics

A. Folding Time Distribution. To study folding kinetics, we consider folding time distributions, $p_f(t)$, which present the

probability distributions of the folding times to the native state. In many studies, both of experimental and computational character (see, for example, refs 55 and 37 and 42, respectively), one deals with the time-dependent population of the native state, which is related to $p_f(t)$ by

$$n_f(t) = \int_0^t p_f(t') dt' \quad (7)$$

or with its complementary quantity, i.e., the population of the unfolded state

$$n_{\text{unfold}}(t) = 1 - n_f(t) = 1 - \int_0^t p_f(t') dt' \quad (8)$$

Due to their cumulative character, $n_f(t)$ and $n_{\text{unfold}}(t)$ are less sensitive than $p_f(t)$ to statistical noise, but at the same time they are also less sensitive to details of folding kinetics (see the discussion of eq 39 and the end of Section VII).

The folding time distributions for the various temperatures are shown in Figure 10. At $T \geq T_f$ (panels b and c), the distributions are similar and consist of a steep rise at very small times followed by an exponential decrease at long times. The first phase corresponds to the transitions from the unfolded state to the semicompact random globule, i.e., there are very few trajectories that fold directly to the native state without collapse to the semicompact globule. Note that the native state population $n_f(t)$, defined by eq 7, does not exhibit such a rise at small times but shows a relatively monotonic decrease with time. The second phase consists of folding from the globule to the native state. These results are encompassed in the analytic model given below [Section V(B)]. Since the first phase leads to a delay in the folding time (i.e., that required for the system to reach the globule state), the slope of the curve in the second phase defines the MFT $\bar{t}_f(T)$ only approximately. However, the result is very close to the true value because the delay due to the collapse is infinitesimal ($\sim 10^4$ MC steps) as compared with the transition time ($\sim 10^7$ MC steps) from the globule to the native state.

In contrast, for $T < T_f$ (panel a), the decay curve is nonexponential; at least two exponential terms are required to fit the simulated results. The first (faster) term corresponds to the short folding trajectories, which go from the globule to the native state directly (such as those on the left side of Figure 8), and the second (slower) term corresponds to the longer trajectories, which fall into the DEs when searching for a way from the globule to the native state (the right side of Figure 8). Although the trajectories of both these types exist at each temperature (both below and above T_f , Figure 8), the fraction of long trajectories is dominant for the limiting time scale at low temperatures.

Our results for the folding time distributions are consistent with the earlier study of protein-like lattice heteropolymers by Abkevich et al.,⁵⁶ who showed that the off-pathways traps are responsible for the characteristic double-exponential pattern of the distributions at low temperatures.

A characteristic feature of the folding time distributions is a large scatter of the folding times at longer times, where the statistics are essentially poorer. Although, in accord with the central limit theorem, this scatter decreases with the number of trajectories as $N_{\text{traj}}^{-1/2}$, it cannot be completely damped because it presents an intrinsic property of the exponential decay law (Section VII).

B. An Analytic Model. The results in Figure 10 for the folding kinetics obtained from the MC simulations suggest that a simple analytic model can be employed to describe the average behavior of the system. This is a useful approach, although some

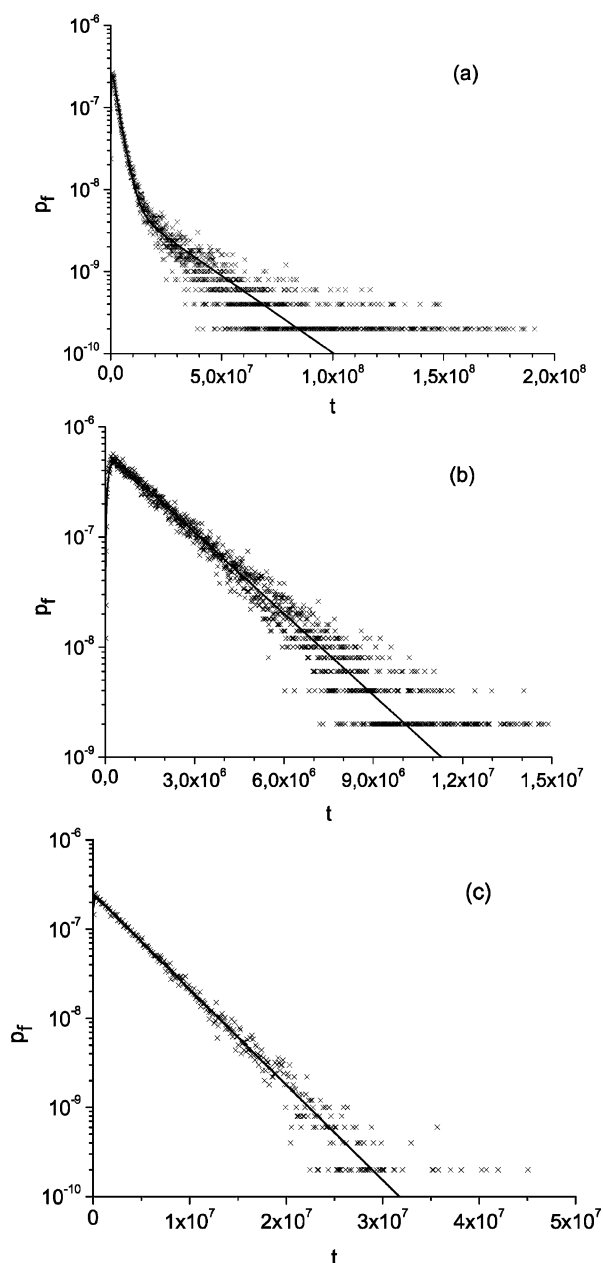


Figure 10. Folding time distributions for $T = 1.5$ (a), $T = 1.9$ (b), and $T = 2.3$ (c). Crosses correspond to the simulation data, and the curves to the eqs 14–16, and 18 (panel a) and eq 22 (panels b and c), with the values of the rate constants presented in Table 2 and Figure 12.

details of the actual folding, such as the statistical scatter of the folding times due to finite number of folding trajectories, are lost in the analytic solution. Such aspects are discussed further in Section VII based on the MC simulations.

The kinetic scheme used to describe the results is depicted in Figure 11. The index “u” corresponds to the unfolded state, “g” corresponds to the collapsed globule, “d” corresponds to a DE or off-pathway intermediate, and “f” corresponds to the native state (or the native basin; see the remark in Section IV). Arrows show the directions of the transitions. To make the model more tractable, the probability of unfolding of the compact globule was assumed to be negligible. At T equal to or below T_f , this assumption is quite accurate (see Figure 4a,b); however, at higher temperatures, such as that used in Figure 4c, unfolding of the globule can slightly alter the folding kinetics. Also, the folding time is equated to the first passage

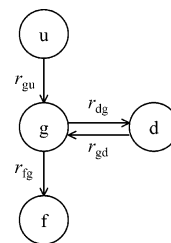


Figure 11. Schematic representation of the model system. Labels u, g, d, and f are for the unfolded chain, semicompact random globule, DE, and native state, respectively. Arrows show the directions of transitions.

time, so that native state unfolding is not included into the kinetic scheme. Under these conditions, the kinetics of the system are described by the following set of coupled linear equations

$$\frac{dn_u}{dt} = -r_{gu}n_u \quad (9)$$

$$\frac{dn_g}{dt} = r_{gu}n_u + r_{gd}n_d - (r_{fg} + r_{dg})n_g \quad (10)$$

$$\frac{dn_d}{dt} = r_{dg}n_g - r_{gd}n_d \quad (11)$$

$$\frac{dn_f}{dt} = r_{fg}n_g \quad (12)$$

where $n_\alpha = n_\alpha(t)$ is the probability of finding the system in state α at time t , and $r_{\beta\alpha}$ is the rate constant for transitions from state α to β .

We note that the applicability of this approach is not restricted by the scheme of Figure 11. It can be used for other scenarios of folding that include obligatory (on-pathway) intermediates and/or parallel pathways, which have been discussed in many papers (see, e.g., refs 20, 21, 55, and 57).

Solving the system of eqs 9–11, with the initial conditions $n_u(0) = 1$ and $n_g(0) = n_d(0) = 0$, yields for $n_g(t)$

$$n_g(t) = B_1 e^{\gamma_1 t} + B_2 e^{\gamma_2 t} + A e^{-r_{gu} t} \quad (13)$$

where

$$\gamma_{1,2} = -\frac{r_{gd} + r_{dg} + r_{fg}}{2} \pm \sqrt{\left(\frac{r_{gd} + r_{dg} + r_{fg}}{2}\right)^2 - r_{gd}r_{fg}} \quad (14)$$

$$B_{1,2} = \pm \frac{r_{gu} + A(r_{gu} + \gamma_{2,1})}{\gamma_1 - \gamma_2} \quad (15)$$

and

$$A = \frac{r_{gu}(r_{gd} - r_{gu})}{r_{gu}^2 - r_{gu}(r_{gd} + r_{dg} + r_{fg}) + r_{gd}r_{fg}} \quad (16)$$

Here subscripts 1 and 2 on γ and B correspond, respectively, to the upper and lower signs (i.e., to the plus and minus sign) in the expressions for these quantities. The folding time distribution, $p_f(t)$, is equal by eq 12 to

$$p_f(t) = \frac{dn_f}{dt} = r_{fg}n_g \quad (17)$$

which, after taking into account eq 13, becomes

$$p_f(t) = r_{fg}(B_1 e^{\gamma_1 t} + B_2 e^{\gamma_2 t} + A e^{-r_{gu} t}) \quad (18)$$

In the limiting case of the absence of the DE (off-pathway) intermediates ($r_{dg} = r_{gd} = 0$), one obtains from eqs 14–16 $\gamma_1 = B_1 = 0$ and

$$\gamma_2 = -r_{fg}, B_2 = -A = r_{gu}/(r_{gu} - r_{fg}) \quad (19)$$

and eq 18 reduces to

$$p_f(t) = \frac{r_{gu} r_{fg}}{r_{gu} - r_{fg}} [e^{-r_{fg} t} - e^{-r_{gu} t}] \quad (20)$$

If the collapse is much faster than the folding time, which is true for the present model and often the case in the folding of proteins, $r_{gu} \gg r_{fg}$. Under this condition, the second term on the right-hand side of eq 20 decays rapidly, leading to a steep rise of $p_f(t)$ at small times. The first term describes a relatively slow exponential decrease of $p_f(t)$ at longer times. This decay corresponds to the folding to the native state from the compact globule.

In the general case of the presence of the DE intermediate ($r_{dg}, r_{gd} \neq 0$), γ_1 in eq 13 becomes nonzero and negative. Also the factor B_1 , which determines the contribution of the term with γ_1 in eq 13 for $n_g(t)$, is nonzero. Due to this, an additional exponential decay term, $B_1 \exp(\gamma_1 t)$, appears in the expression for $p_f(t)$ (18). If the transition of the system from the DE to the globule is slow ($r_{gd} \ll r_{fg}$), one obtains from eqs 14–16: $\gamma_1 \approx -[r_{fg}/(r_{dg} + r_{fg})]r_{gd}$, $B_1 \approx [r_{dg}/(r_{dg} + r_{fg})^2]r_{gd}$, $\gamma_2 \approx -(r_{dg} + r_{fg})$, and $B_2 = -A \approx r_{gu}/(r_{gu} - r_{dg} - r_{fg})$, where only the linear terms in r_{gd} are kept to obtain the order of magnitude of this quantity. This shows that B_1 and γ_1 , which are proportional to r_{gd} , are small in comparison with the corresponding values of the B_2 and γ_2 . It is useful here to separate two time regimes, $t < t^*$ and $t > t^*$, where $t^* = \ln[(r_{dg} + r_{fg})^2/r_{dg}r_{gd}]/(r_{dg} + r_{gd})$, which is determined by the condition $B_1 \exp(\gamma_1 t^*) = B_2 \exp(\gamma_2 t^*)$. For short times ($t < t^*$), the term $B_2 \exp(\gamma_2 t)$ is dominant, while for longer times ($t > t^*$), the term $B_1 \exp(\gamma_1 t)$ is dominant. It thus follows that $p_f(t)$ obeys a double-exponential decay curve at a moderate and longer times. This type of behavior of $p_f(t)$ corresponds to a nonequilibrium exchange between the DE and globule, and it is characteristic of $T < T_f$ (Figure 10a). Specifically, with the rates constants of Table 2 for $T = 1.5$, $t^* \approx 1.1 \times 10^7$, as is evident in Figure 10a.

In the opposite limit, i.e., at a high rate of transition from the DE to the globule ($r_{fg} \ll r_{gd}$), eqs 14–16 give $\gamma_1 \approx -\alpha r_{fg}$, $B_1 = -A \approx \alpha r_{gu}/(r_{gu} - \alpha r_{fg})$, $\gamma_2 \approx -(r_{gd} + r_{dg})$, and $B_2 \approx r_{gu}r_{dg}/(r_{gu} - r_{gd})r_{gd}$, where

$$\alpha = 1/(1 + r_{dg}/r_{gd}) \quad (21)$$

and only terms linear in r_{fg} are kept. Assuming also that $r_{dg} \ll r_{gd}$, which is characteristic of temperatures close to and higher than T_f (see, e.g., Table 2), one finds $\gamma_1/\gamma_2 \sim r_{fg}/r_{gd} \ll 1$, and $B_1/B_2 \sim r_{gd}/r_{dg} \gg 1$. Thus, the term $B_1 \exp(\gamma_1 t)$ in eq 18 has a smaller rate of decay and a larger preexponential factor than the $B_2 \exp(\gamma_2 t)$ term. For this case, the $B_1 \exp(\gamma_1 t)$ dominates over $B_2 \exp(\gamma_2 t)$ at all times. Equation 18 then yields

$$p_f(t) \approx \frac{\alpha r_{gu} r_{fg}}{r_{gu} - \alpha r_{fg}} [e^{-\alpha r_{fg} t} - e^{-r_{gu} t}] \quad (22)$$

which corresponds to eq 20, valid in the absence of the DE, except that the rate constant r_{fg} is multiplied by the factor α , defined by eq 21. This factor is less than unity and accounts

for the slowing down of the folding process because of the presence of the DE, since an additional time is required to go to the DE and return to the globule. That the slowing down can be taken into account by a simple correction of the rate of transitions from the globule to the native state by introduction of an effective rate, $\tilde{r}_{fg} = \alpha r_{fg}$, is a consequence of the condition $r_{fg} \ll r_{gd}$, i.e., DE is in quasi-equilibrium with the globule, and so that the system “globule + DE” can be considered as an “extended” globule, characterized by the partition function $Z_g + Z_d$, where Z_g and Z_d are the partition functions of the globule and DE, respectively. Indeed, according to detailed balance, since $r_{dg}Z_g = r_{gd}Z_d$, one has

$$\alpha r_{fg} = \frac{Z_g}{Z_g + Z_d} r_{fg} = \frac{Z_{fg}}{Z_g + Z_d} = \tilde{r}_{fg}$$

where $r_{fg} = Z_{fg}/Z_g$, and Z_{fg} is the partition function of the transition state between the globule and native state.

To show more explicitly that the condition of quasi-equilibrium exchange between DE and globule plays a key role in the validity of eq 22, we assume, as the zero-order approximation, that the globule and DE are in equilibrium

$$r_{dg}n_g = r_{gd}n_d \quad (23)$$

Then $dn_d/dt = (r_{dg}/r_{gd})dn_g/dt$, with which eq 11 yields, in first order, $r_{gd}n_d - r_{dg}n_g = -(r_{dg}/r_{gd})dn_g/dt$. By substituting this expression into eq 10, the latter can be written

$$\alpha^{-1} \frac{dn_g}{dt} = r_{gu}n_u - r_{fg}n_g \quad (24)$$

where α is given by eq 21. Equations 9 and 12 do not change. Introducing new variables $\tilde{n}_g = \alpha^{-1}n_g$ and $\tilde{r}_{fg} = \alpha r_{fg}$, the system of eqs 9, 12, and (24) takes the form

$$\frac{dn_u}{dt} = -r_{gu}n_u$$

$$\frac{d\tilde{n}_g}{dt} = r_{gu}n_u - \tilde{r}_{fg}\tilde{n}_g$$

$$\frac{dn_f}{dt} = \tilde{r}_{fg}\tilde{n}_g$$

which corresponds to the system of eqs 9–11 in the absence of DE ($r_{dg} = r_{gd} = 0$). The desired solution is then obtained from eq 20 as

$$p_f(t) = \frac{r_{gu}\tilde{r}_{fg}}{r_{gu} - \tilde{r}_{fg}} [e^{-\tilde{r}_{fg} t} - e^{-r_{gu} t}] = \frac{\alpha r_{gu} r_{fg}}{r_{gu} - \alpha r_{fg}} [e^{-\alpha r_{fg} t} - e^{-r_{gu} t}]$$

which coincides with eq 22.

The behavior of $p_f(t)$, given by eq 22, corresponds to that observed for $T \geq T_f$ (Figure 10b,c). It thus follows that for $T \geq T_f$ the DE intermediate, if it exists, does not leave a distinct “fingerprint” on the folding time distribution, in contrast to the result at $T < T_f$ (Figure 10a), that is, the form of the function $p_f(t)$ is exactly the same as in the absence of the DE. Importantly, this means that the single-exponential behavior of a folding time distribution at longer time scales only indicates that there are no off-pathway intermediates that are not in equilibrium with the globule, rather than that there are no intermediates. Clearly, this has consequences for the experimental analysis of folding time distributions.

TABLE 2: Fitted Values of the Rate Constants^a

<i>T</i>	1.5	1.7	1.9	2.1	2.3
r_{gu}	3.3×10^{-6}	5.4×10^{-6}	$1.1 \times 10^{-5\dagger} (0.9 \times 10^{-5})$	$1.5 \times 10^{-5\dagger} (1.5 \times 10^{-5})$	$2.6 \times 10^{-5\dagger} (2.6 \times 10^{-5})$
r_{dg}	5.1×10^{-8}	1.1×10^{-7}			
r_{gd}	5.1×10^{-8}	5.2×10^{-7}			
r_{dg}/r_{gd}	1.0	2.1×10^{-1}	$2.5 \times 10^{-2\dagger}$	$7.3 \times 10^{-3\dagger}$	$1.7 \times 10^{-3\dagger}$
r_{fg}	3.2×10^{-7}	6.1×10^{-7}	$5.9 \times 10^{-7\dagger} (5.9 \times 10^{-7})$	$4.3 \times 10^{-7\dagger} (4.2 \times 10^{-7})$	$2.5 \times 10^{-7\dagger} (2.5 \times 10^{-7})$

^a For values with a dagger (†) or double dagger (‡), see Section VI for an explanation.

Along with the folding time distribution, $p_f(t)$, the model describes the time-dependent populations of the globule, unfolded state, DE and the native state, which are given by the functions $n_g(t)$, $n_u(t)$, $n_d(t)$, and $n_f(t)$, respectively. The function $n_g(t)$ is determined by eq 13 [or eq 22, where $n_g(t) = p_f(t)/r_{fg}$], and $n_u(t)$, $n_d(t)$, and $n_f(t)$ are obtained by solving eqs 9, 11, and 17 with the initial conditions $n_u(0) = 1$, $n_d(0) = 0$, and $n_f(0) = 0$, respectively. The resulting expressions are

$$n_u(t) = e^{-r_{gu}t} \quad (25)$$

$$n_d(t) = r_{dg} \int_0^t n_g(t') e^{r_{gd}(t'-t)} dt' \quad (26)$$

and

$$n_f(t) = r_{fg}[(B_1/\gamma_1)(e^{\gamma_1 t} - 1) + (B_2/\gamma_2)(e^{\gamma_2 t} - 1) - (A/r_{gu})(e^{-r_{gu}t} - 1)]$$

where $\gamma_{1,2}$, $B_{1,2}$, and A are given by eqs 14–16.

Integrating eqs 13, 25, and 26 over the time, we can calculate the fractions of the time spent by the system in the unfolded, globule and DE states

$$P_\alpha = \frac{\int_0^\infty n_\alpha(t) dt}{\sum_\alpha \int_0^\infty n_\alpha(t) dt} \quad (27)$$

where α stands for “u”, “g”, and “d”; the contribution of the native state to the denominator of eq 27, which is $\int_0^\infty p_f(t) dt = 1$, is neglected because it is small in comparison with the other terms. The fractions P_α are essentially the same as those defined by eqs 5 and 6, which were introduced for the analysis of the simulation results, and present a coarse grained characterization of the nonequilibrium residence probability surfaces (Figure 3). The only difference between the residence probabilities of Figure 3 and these fractions is that the former are associated with individual points on the (N, N_{nat}) plane and the latter with the collections of these points, which represent characteristic regions of the surface, i.e., the unfolded, globule and DE states.

C. Mean Folding Time (MFT). The MFT is defined as $\bar{t}_f = \int_0^\infty t p_f(t) dt$. Using eq 18 for $p_f(t)$, \bar{t}_f can be written as

$$\bar{t}_f = \frac{1}{r_{fg}} \left(1 + \frac{r_{dg}}{r_{gd}} \right) + \frac{1}{r_{gu}} \quad (28)$$

This equation shows explicitly that the collapse of the unfolded state into the globule, characterized by the mean time

$$\bar{t}_u = 1/r_{gu} \quad (29)$$

leads to a delay in the folding time.

Interestingly, the MFT, given by eq 28, does not depend on the degree of equilibration between the DE and globule. The contribution of the DE to the MFT is determined by the ratio

of the rate constants r_{dg}/r_{gd} . By detailed balance, as already mentioned, r_{dg}/r_{gd} is equal to the ratio of the partition functions of the globule and DE regions, Z_d/Z_g , and does not depend on the height of the barrier between these regions. It follows that in the case of a very high barrier, when the probability to fall into the DE is negligible, the MFT is the same as in the case of a low barrier, when the globule and DE regions are in equilibrium due to a fast exchange between them. This is an important result, which makes clear the limitation of experimental results for obtaining insight into the folding process.

The folding trajectories can be divided into two classes: trajectories that go from the globule to the native state directly and those that pass through the DE. The trajectories that go directly are short, while those that go through the DE are long (see Figure 8 and the comments concerning it). According to the kinetic model, Figure 11, the mean lifetimes of the system in the globule and DE regions are equal, respectively, to

$$\bar{t}_g = 1/(r_{fg} + r_{dg}) \quad (30)$$

and

$$\bar{t}_d = 1/r_{gd} \quad (31)$$

and the fractions of the trajectories passing from the globule to the native state directly and via the DE, respectively, are

$$f_g = r_{fg}/(r_{fg} + r_{dg}) \quad (32)$$

and

$$f_d = r_{dg}/(r_{fg} + r_{dg}) = 1 - f_g \quad (33)$$

The average length and fraction of the direct trajectories are $\bar{t}_1 = \bar{t}_g + \bar{t}_u$ and $f_1 = f_g$, where \bar{t}_u , \bar{t}_g , and f_g are given by eqs 29, 30, and 32, respectively, and those of the indirect trajectories, are \bar{t}_2 and $f_2 = f_d$, where f_d is given by eq 33, and \bar{t}_2 includes the collapse time \bar{t}_u (eq 29) and also the time spent by the system in both the globule and DE regions, possibly with repeated visits to the DE. The MFT can be written as

$$\bar{t}_f = \bar{t}_1 f_1 + \bar{t}_2 f_2 \quad (34)$$

Comparing this equation with eq 28, one finds for \bar{t}_2

$$\bar{t}_2 = \frac{r_{gu}(r_{fg} + r_{dg})^2(r_{gd} + r_{dg}) - r_{fg}^2 r_{gd}(r_{gu} + r_{fg} + r_{dg})}{r_{gu} r_{fg} r_{dg} r_{gd}(r_{fg} + r_{dg})} \quad (35)$$

We consider now the case of rare visits of the system to the DE, $r_{dg}/r_{fg} \ll 1$, and, for simplicity neglect the collapse time ($r_{gu}/r_{fg} \gg 1$). Here $f_1 \approx 1$, $\bar{t}_1 \approx 1/r_{fg}$, $f_2 \approx r_{dg}/r_{fg}$, and $\bar{t}_2 \approx 1/r_{gd} = \bar{t}_d$. The contributions of the direct and indirect trajectories to the MFT are then $\bar{t}_1 f_1 \approx 1/r_{fg}$ and $\bar{t}_2 f_2 \approx (r_{dg}/r_{gd})/r_{fg} = (r_{dg}/r_{fg}) \bar{t}_d$, respectively, and the MFT, \bar{t}_f , given by eq 34, becomes

$$\bar{t}_f = \frac{1}{r_{fg}} + \frac{r_{dg}}{r_{fg}} \bar{t}_d = \frac{1}{r_{fg}} \left(1 + \frac{r_{dg}}{r_{gd}} \right) \quad (36)$$

in accord with eq 28 (after neglecting the collapse time \bar{t}_0). It is seen that, although the fraction of the indirect trajectories is small ($f_2 \approx r_{dg}/r_{fg} \ll 1$), their contribution to the MFT can be dominant, if the statistical weight of the DE region is large in comparison with that of the globule region, $r_{dg}/r_{gd} = Z_d/Z_g \gg 1$. It follows that a correct estimation of the MFT, or, more generally, a complete analysis of the folding process, requires that the simulated trajectories be sufficiently long, so that $t_{traj} \gg 1/r_{dg}$ to allow the system to visit the DE. In addition, the total number of trajectories must be sufficiently large that $N_{traj} \gg 1/f_2 = r_{fg}/r_{dg}$, to ensure a proper sampling of the DE region.

In the opposite limit of very frequent visits to the DE, $r_{dg}/r_{fg} \gg 1$, one obtains $f_1 \approx r_{fg}/r_{dg}$, $t_1 \approx 1/r_{dg}$, $f_2 \approx 1$, and $t_2 \approx (r_{gd} + r_{dg})/r_{fg}r_{gd} = (r_{gd} + r_{dg})t_d/r_{fg}$; the collapse time is neglected as previously, here assuming that $r_{gu}/r_{dg} \gg 1$. In this case, the corresponding contribution of the direct and indirect trajectories to the MFT is $t_1f_1 \approx r_{fg}/r_{dg}^2$ and $t_2f_2 \approx (r_{gd} + r_{dg})t_d/r_{fg} = (1 + r_{dg}/r_{gd})/r_{fg}$, respectively, with the former being negligible in comparison with the latter. Then eq 34 for the MFT leads to eq 36, as demonstrated above. The difference between these two cases (rare versus frequent visits to the DE) is that in the former case the factor r_{dg}/r_{fg} in the right-hand side of eq 36 is less than unity and thus not every folding trajectory visits the DE; in the latter case, this factor is larger than unity and thus essentially every trajectory repeatedly visits the DE (on the average, r_{dg}/r_{fg} times).

VI. Estimation of the Rate Constants from the Simulation Data

To estimate the rate constants for the folding model, we compared the analytical solutions of eqs 18 and 22 with the folding time distributions found in the computer simulations (such as in Figure 10). For this purpose we used the maximum likelihood approach. Specifically, the logarithm of the likelihood function $\Phi = \sum_i p_f^{\text{sim}}(t_i) \ln[p_f^{\text{sim}}(t_i)/p^{\text{th}}_f(t_i)] \geq 0$ was minimized with respect to the values of the rate constants; here the labels “sim” and “th” correspond to the computer simulation data and analytical solutions, respectively, and the t_i ($i = 1, \dots, 2000$) are the centers of the time intervals in which $p_f^{\text{sim}}(t_i)$ were stored. In the case of the general solution, given by eq 18, the minimization was conducted with respect to the rate constants r_{gu} , r_{dg} , r_{gd} , and r_{fg} , and in the case of the quasi-equilibrium exchange between the globule and DE intermediate, given by eq 22, with respect to r_{gu} and αr_{fg} , where the factor α is defined by eq 21.

It should be emphasized that the model does not distinguish among the off-pathway minima contributing to the simulations but consider them as a single effective DE; correspondingly, the transition rates r_{dg} and r_{gd} are effective rates that determine transitions to/from all these off-pathway minima. Moreover, according to the consideration in Section V(B), the slower decay term in eq 18 is associated with the off-pathway minima that are not in equilibrium with the globule. Correspondingly, the rate constants r_{dg} and r_{gd} represent the transitions between an extended globule, which incorporates the off-pathway minima in equilibrium with the globule, and an effective DE, which includes all off-pathway minima that are not in equilibrium with the globule. The rate constants r_{gu} and r_{fg} thus determine the transitions from the unfolded state to the extended globule and from the extended globule to the native state, respectively. If the individual DEs were treated separately, a master equation approach would be required, but it would add little to the simpler picture used here.

The results of the fitting procedure are shown in Table 2 and Figures 10 and 12: Table 2 and Figure 12 present the values

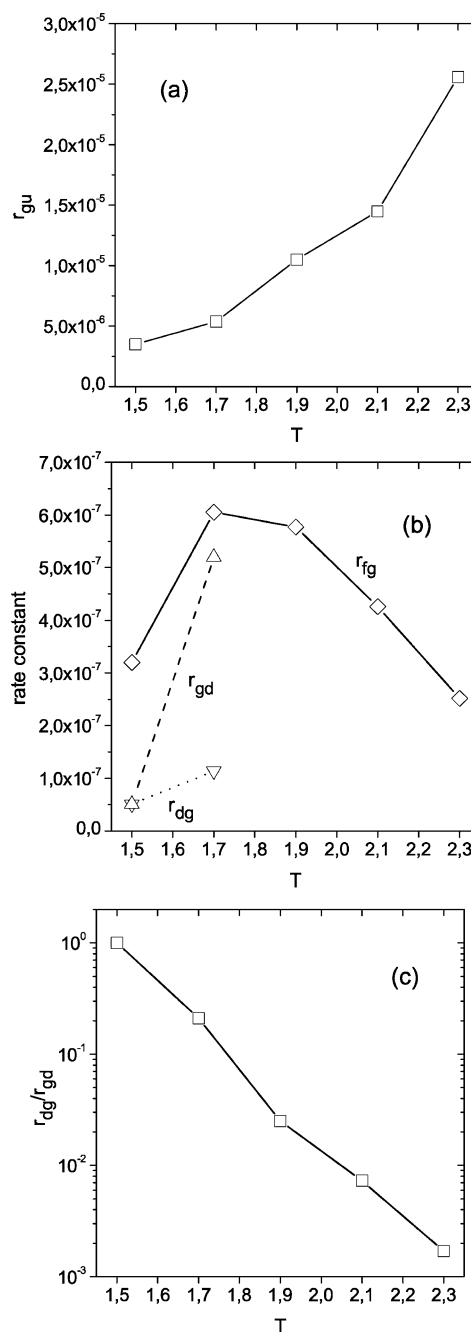


Figure 12. Fitted values of the rate constants: (a) r_{gu} , (b) r_{dg} , r_{gd} and r_{fg} , and (c) r_{dg}/r_{gd} . The lines are to guide the eye.

of the rate constants, and Figure 10 shows how the analytical solutions with these values of the rate constants approximate the folding time distributions obtained in the simulations. To find the rate constants for $T < T_f = 1.9$ (i.e., at $T = 1.5$ and $T = 1.7$), where the decay curves are not single-exponential (Figure 10a), the general solution, given by eqs 14–16 and 18, was employed. In this case, as can be seen from Table 2 and Figure 12, r_{fg} is faster than r_{gd} , so that quasi-equilibrium is not established, and the individual rate constants involving the DE (r_{dg} and r_{gd}) play a role. For $T \geq T_f$ we tried both the general solution, eq 18, and the solution for the quasi-equilibrium exchange between the DE intermediate and globule, eq 22; in Table 2 the results obtained with the latter equation are labeled with the superscript \dagger . Comparing the two sets of results, we see that the values of the rate constants for the collapse, r_{gu} , and for the transition from the globule to the native state, r_{fg} ,

are robust and agree well for eq 18 and eq 22. When eq 22 is used, the product αr_{fg} is determined rather than rate constant r_{fg} ; consequently, this product was divided by the factor $\alpha = (1 + r_{dg}/r_{gd})^{-1}$ (21), which was obtained from the ratio r_{dg}/r_{gd} , estimated from the residence probabilities (see below). In fact, this correction is negligible, because α is very close to unity, as evident from the values of the r_{dg}/r_{gd} given in Table 2. The values of r_{gu} and r_{fg} obtained in the present study are in good correspondence with those of Sali et al.;¹⁴ the unfolded chain collapsed to a semicompact globule in approximately $1/r_{gu} \sim 10^4$ MC steps, and the subsequent search for the native state required approximately $1/r_{fg} \sim 10^7$ steps.

As was mentioned in Section V(B), the process of globule unfolding can affect the kinetics for $T > T_f$, although, to simplify the analytical solution, we left it out of consideration. The resulting rate constant r_{gu} may be not quite accurate for $T > T_f$. However, this constant can be considered as an effective constant, which takes into account the globule unfolding rate, characterized by a rate constant r_{ug} . In the presence of globule unfolding, eqs 9 and 10 would be written as

$$\frac{dn_u}{dt} = r_{ug}n_g - r_{gu}n_u \quad (37)$$

and

$$\frac{dn_g}{dt} = r_{gu}n_u - (r_{ug} + r_{fg})n_g \quad (38)$$

The terms corresponding to the exchange between the DE and globule in the latter equation are omitted, because at $T > T_f$ the DE and globule are in quasi-equilibrium. In the zero-order, the term $r_{ug}n_g$ can be neglected in eq 38, but it should be kept in eq 37. Further, since the maximum unfolding rate is attained in the region where $n_g(t)$ is maximal, dn_g/dt in eq 38 is close to zero; this is essentially the steady-state approximation. Equation 38 then gives $r_{gu}n_u \approx r_{fg}n_g$, or $n_g \approx (r_{gu}/r_{fg})n_u$. With the latter expression for n_g , eq 37 becomes

$$\frac{dn_u}{dt} = -r_{gu} \left(1 - \frac{r_{ug}}{r_{fg}} \right) n_u$$

which shows, as expected, that the globule unfolding slows down the collapse of the unfolded chain, and that this effect becomes significant if the unfolding rate r_{ug} is comparable with the rate of the globule folding to the native state, r_{fg} .

The rate constants for the transitions between the globule and DE, r_{dg} and r_{gd} , are successfully determined with eq 18 only for $T < T_f$, in accord with the discussion in Section V. Also in accord with that discussion, under quasi-equilibrium conditions, when eq 22 is valid, only the complex αr_{fg} , with α given by eq 21, can be determined from the data for $T \geq T_f$. Estimates for r_{dg}/r_{gd} can be obtained from the simulation results of Table 1. We note that for $T > T_f$ the DE and globule are in quasi-equilibrium, and thus, according to detailed balance, $r_{gd}P_d \approx r_{dg}P_g$. Since the time spent by the system in the unfolded state and DE regions is negligible in comparison with the time spent in the globule region (see Figure 3c), $P_g \approx 1$. Then, $r_{dg}/r_{gd} \approx P_d$, which gives from Table 1 the following values of r_{dg}/r_{gd} : 2.5×10^{-2} , 7.3×10^{-3} and 1.7×10^{-3} at $T = 1.9$, $T = 2.1$, and $T = 2.3$, respectively. These values of r_{dg}/r_{gd} are also included in Table 2, where they are labeled by the superscript ‡, and they are shown in Figure 12c.

The same procedure of fitting of the rate constants is applicable to the time-dependent populations of the unfolded

state, $n_{\text{unfold}}(t)$ (8), which should be compared with the corresponding function obtained in the simulations. At $T = 1.5$ the best fit is obtained with

$$r_{gu} = 5.7 \times 10^{-6}, r_{dg} = 4.2 \times 10^{-8}, r_{gd} = 4.0 \times 10^{-8}, \text{ and } r_{fg} = 2.9 \times 10^{-7} \quad (39)$$

These values are slightly different from those in Table 2, with the larger difference being obtained for r_{gu} ; this is responsible for the steep rise in $p_f(t)$ at small times that is not visible in the $n_f(t)$ (7) and $n_{\text{unfold}}(t)$ (8) functions.

According to eq 33 and the data of Table 2, the fraction of the indirect trajectories (i.e., the trajectories that visit the DE before reaching the native state), $f_d = r_{dg}/(r_{fg} + r_{dg})$ is equal to ≈ 0.14 at $T = 1.5$. At the same time, direct counting of the trajectories that visited the off-pathway minima (Section III) show that the fraction of such trajectories at this temperature is as large as ≈ 0.76 . This suggests, in accord with the above analysis, that most of the DEs listed in Figure 1 are actually in equilibrium with the globule, and thus constitute part of the extended globule. Configuration D (Figure 1) with $N = N_{\text{nat}} = 25$ is the best candidate for a DE that is not in equilibrium with the globule because the lifetime of the system in this configuration ($\approx 1.3 \times 10^6$) is at least 10 times longer than those of the other DEs configurations, including the fully compact nonnative configurations. At $T = 1.5$, the fraction of the trajectories that visited the minimum at $N = N_{\text{nat}} = 25$ is approximately equal to 0.26; among these trajectories $\approx 40\%$ visited configuration D, and $\approx 60\%$ configuration C, whose lifetime is about 4 order of magnitude less than that for configuration D. According to the transition disconnectivity graph (Figure 9), configuration C contributes to the extended native state basin rather than corresponds to a DE trap. Thus, only the trajectories that visited configuration D are not in equilibrium with the globule, and the above-mentioned value of the fraction of the trajectories is reduced to $0.26 \times 0.4 \approx 0.11$, which is in reasonable agreement with the previously given estimate for f_d on the basis of the data of Table 2 (≈ 0.14). As mentioned in Section III, the fractions of the trajectories passing through different off-pathway minima do not vary much with temperature. The rate constants of Table 2 support this conclusion. For example, at $T = 1.7$ eq 33 gives $f_d \approx 0.15$, which is also in good correspondence with results of the direct counting of the fraction of the trajectories that visited configuration D, which is approximately equal to 0.12.

Having the rate constants, we can also calculate the fraction of the time spent by the system in the DEs, predicted by the kinetic model, eq 27, and compare it with the corresponding values of P_d found in the simulations (Table 1). This is straightforward for $T < T_f$, where all rate constants are known. Correspondingly, eqs 13, 25, and 26 are employed to define $n_g(t)$, $n_u(t)$, and $n_d(t)$ in eq 27. Results of the calculations are presented in Table 1 and show good agreement with the simulation results.

As is seen from Table 2 and Figure 12, the transition rate from the globule to the native state, r_{fg} , exhibits a maximum near $T = T_f$ (see below). This explains the origin of the minimum in the MFT (Figure 2). Indeed, since at $T \approx T_f$ both r_{fg}/r_{gu} and r_{dg}/r_{gd} are small, the general expression for the MFT, given by eq 28, simplifies approximately to

$$\bar{t}_f \approx 1/r_{fg} \quad (40)$$

Actually, the minimum in t_f is slightly shifted to higher

temperature, as compared to the maximum in r_{fg} (the former is located at $T = T_f \approx 1.9$, and the latter at $T = T_{fg, \min} \approx 1.7$). Differentiating eq 28 with respect to T and equating the right-hand side of the resulting equation to zero, one finds

$$\frac{dr_{fg}}{dT} = -\frac{1}{\gamma} \left(\frac{r_{fg}}{r_{gu}} \right)^2 \frac{dr_{gu}}{dT} + \frac{r_{fg}}{\gamma} \frac{d\gamma}{dT}$$

where $\gamma = (1 + r_{dg}/r_{gd})$. Since $dr_{gu}/dT > 0$ and $d(r_{dg}/r_{gd})/dT < 0$ at $T = T_f$ (Figure 12), it follows that the minimum in t_f is attained in the temperature range when the dr_{fg}/dT is negative, that is, at $T > T_{fg, \min}$.

VII. Simulation of the Statistical Effects in the Folding Time Distributions: Monte Carlo Simulations

As has been mentioned in Section V, the system of kinetic equations 9–11 describes the average behavior of the system and does not account for the statistical scatter of the folding times at long time scales, as evident in Figure 10. However, this shortcoming is characteristic of eqs 9–11 and not of the kinetic scheme (Figure 11) that underlies these equations. The scatter arises from the finite number of folding trajectories and presents an intrinsic property of the exponential decay law.

To see that, consider a system randomly escaping from a certain state, with the mean lifetime in this state being equal to τ . Then the lifetime distribution (i.e., the probability density for the system to live for a given time t) obeys the Poisson distribution

$$p = \tau^{-1} \exp(-t/\tau)$$

The variation of this equation gives

$$\delta t = -\tau^2 \exp(t/\tau) \delta p$$

which shows that the probability intervals of a constant length, δp , are mapped onto the exponentially increasing time intervals δt , and thus the number of the events corresponding to the current time t decreases with t as $\sim Np \sim N \exp(-t/\tau)$, where N is the total number of events. According to the central limit theorem, the current time variance σ_t^2 is then $\sim 1/Np$, so that the root-mean-square deviation in the lifetimes is

$$\sigma_t \sim N^{-1/2} p^{-1/2} \sim N^{-1/2} \exp[t/(2\tau)] \quad (41)$$

This equation shows that for every given t , the scatter decreases as the number of the events increases. However, at larger N less probable events come into existence, which do not occur at smaller N because of low probabilities of these events. As a result, the scatter shifts to longer times rather than is completely damped.

To mimic the scatter of the folding times within the kinetic scheme (Figure 11), we use MC simulations; we note that they have no relation to the MC simulations used for the lattice model. For this purpose, we assume that the lifetime of the system in each of the characteristic states (i.e., the unfolded chain, globule, and DE regions) obeys the Poisson distribution, so that

$$t_{\beta\alpha} = -1/r_{\beta\alpha} \ln P \quad (42)$$

where $r_{\beta\alpha}$ is the rate constant for the transition from state α to β , $t_{\beta\alpha}$ is the corresponding waiting time, and P is a random

variable uniformly distributed between zero and one. Given the values of the rate constants, the kinetic scheme of Figure 11 can be implemented: Starting in the unfolded state, the transition time to the globule, t_{gu} , is calculated according to eq 42. In the globule state, two quantities, t_{fg} and t_{dg} , are selected randomly based on eq 42. For $t_{fg} < t_{dg}$, the system passes directly into the native state, and the trajectory is terminated with the folding time

$$t_f = t_{gu} + t_{fg}$$

Otherwise, that is for $t_{fg} \geq t_{dg}$, the system visits the DE. In this case, the waiting time t_{gd} is selected, as given by eq 42, and the system is returned to the globule region. Then, the times t_{fg} and t_{dg} are selected again, and a decision is taken where the system should go, either to the native state or to the DE again. This process is continued until the system reaches the native state, and the folding time is calculated as

$$t_f = t_{gu} + n(t_{dg} + t_{gd}) + t_{fg}$$

where n is the number of times the system visited the DE.

Figure 13 shows the results of the MC simulations of the model equations for different numbers of folding trajectories. The temperature is equal to 1.5, and for the rate constants of the transitions between the states of the system (r_{gu} , r_{dg} , r_{gd} , and r_{fg}), the values from Table 2 are used for this temperature. Panels a, b, and c correspond to a successive increase of the number of folding trajectories, 5×10^4 to 5×10^5 and to 5×10^6 . The results are in accord with the previous analysis of eq 41, that is, at a given t the scatter of the data decreases with the number of trajectories (approximately as $N^{-1/2}$), but it is not completely damped, shifting to longer times. Panel a also presents the comparison with the lattice simulation results for the same number of trajectories (Figure 10a, 5×10^4 trajectories), which shows that the MC simulation of the model equations reproduces the scatter of folding times, observed in the lattice simulations, quite well.

We also performed MC simulations for the model equations with the rate constants obtained by fitting the time-dependent populations of the unfolded states (see eq 39). To compare the results with those with the rate constants in Table 2, the root-mean-square (rms) deviation from the corresponding folding time distribution obtained in the lattice simulation at $T = 1.5$ was calculated. It was found that for the rate constants in Table 2 the rms is $\approx 5 \times 10^{-10}$, whereas in the case of eq 39 it is about five times larger, $\approx 2.5 \times 10^{-9}$.

VIII. Concluding Discussion

Folding time distributions provide information about the folding reaction that is difficult to obtain from mean force surfaces. Being directly related to the folding kinetics, they offer an opportunity of characterization of the folding process in the terms of the rate constants of transitions between the inherent states of the system. In this paper, we describe the results of a comprehensive study of the folding reaction for a 27-residue heteropolymer on a cubic lattice with the Monte Carlo method used to simulate the folding dynamics and suggest a simple kinetic model to describe the folding kinetics. This model system was chosen to make it possible to analyze the interrelation between the system properties and the kinetics, rather than to mimic a specific protein or family of proteins. Because the system reproduces many characteristic features of the folding

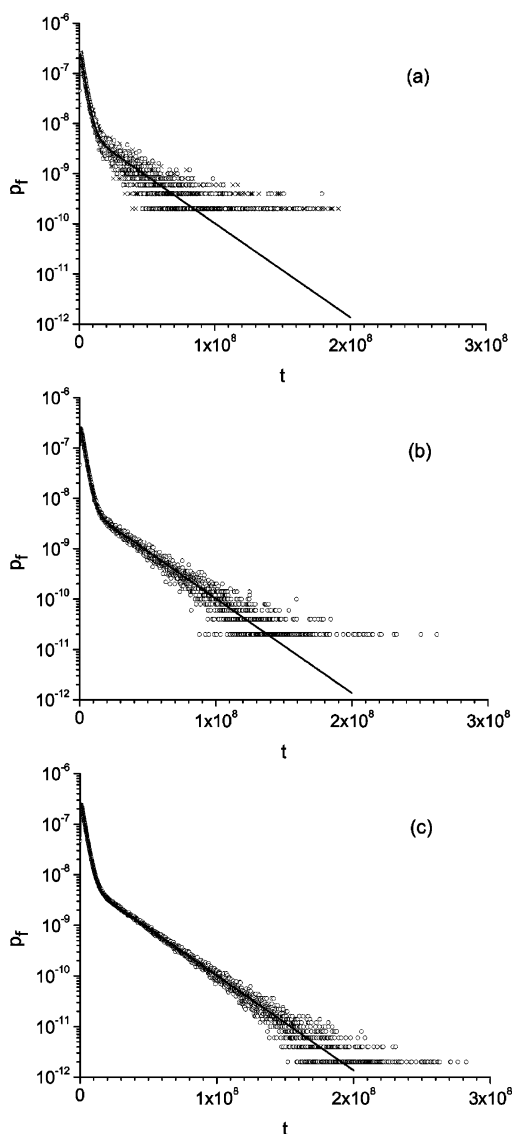


Figure 13. Folding time distributions, $T = 1.5$. Circles correspond to the MC simulations of the model equations for different numbers of the trajectories: 5×10^4 (a), 5×10^5 (b), and 5×10^6 (c). Crosses in panel a show the results of the lattice MC simulations from Figure 10a for the same number of trajectories (5×10^4). The solid line presents the analytical solution, eqs 14–16 and 18, with the rate constants from Table 2.

process and at the same time is a simple enough to allow a comprehensive study, it has been and continues to be a widely used model for protein folding. We undertook a detailed examination of the mean force surfaces (free energy, potential energy, and entropy surfaces) at various temperatures. As in previous papers, we found that, starting with the unfolded state, the system collapsed rapidly to a random globule, from which it folds to the native state. Mapping folding trajectories onto these surfaces, we determined the essential off-pathway minima, or DEs, which play an important role in the folding process, particularly at low temperatures. The minima correspond either to fully compact nonnative structures or nearly compact structures that are unable to convert into the native structure without breaking some contacts (i.e., without returning to the globule state). The role of the DEs was confirmed by constructing the transition disconnectivity graph, in which the free energy barriers between the minima are calculated from the transition rates between the minima. This graph shows the complexity of

the surface, which is often missed when a projection on a small number (usually two) progress variables are used. For example, when the two variables (i.e., N and N_{nat}) are used for the potential energy and entropy surfaces, the results represent smoothed surfaces which miss the full complexity of the multidimensional surfaces obtained with the transition disconnectivity graph. By this we mean that the denatured state has a number of free energy basins, separated by high barriers. We have found that the folding trajectories can be divided into two classes: trajectories that pass directly from the globule region to the native state (i.e., without visiting the off-pathway minima), and trajectories that visit the off-pathway minima, possibly repeatedly. In general, the former trajectories correspond to short folding times, while the latter have longer folding times. At each temperature studied (below, at, and above the folding temperature), both types of the trajectories are observed, but the fraction of “direct” trajectories decreases as the temperature is elevated, i.e., more trajectories visit the DE traps.

The folding time distributions have a steep rise at short times and decay exponentially at long times for temperatures equal to or above the folding temperature. At lower temperatures, by contrast, the folding time distributions exhibit a double-exponential decay at longer times. To interpret the folding time distributions from the simulations, we introduced a simple kinetic model, which described the transitions between the characteristic states of the system, i.e., the unfolded chain, the semicompact random globule, DE, and the native state. The analytical solution obtained for this model showed that the steep rise at short times corresponds to the transitions from an unfolded chain to a semicompact random globule, and the decay at longer times to those from the globule to the native state. The double-exponential decay curve has significant contributions from the two types of trajectories mentioned above; the rapidly decreasing exponential term corresponds to the folding trajectories that do not visit the DEs, and the slowly decreasing term to the trajectories that do visit the DEs, possibly returning repeatedly to the globule state before reaching the native state.

It should be noted that although some of the DE minima are the result of the lattice representation with a restricted move set (Section III), the origin of these minima does not invalidate the analysis of the interrelation between the free energy surface landscape and kinetics of the system, which is the goal of the present paper.

The analytical solution was used to obtain results of interest that are confirmed by the simulations. If the DEs exist but are in equilibrium with the globule (due to a fast rate of return from the DE to globule), the folding time distribution is single-exponential at longer times, which is characteristic of folding in the absence of the DE. In this case, the globule and DE can be considered to constitute an extended globule, and the transition from the globule to the native state is characterized by a single, effective rate constant, which is renormalized to take into account the time delay due to the interconversion between the DE and globule. This shows that a single-exponential decay of the folding time distribution at intermediate and longer time scales does not prove that off-pathway minima do not contribute. Another striking result is that the mean folding time is independent of the height of the barrier between the globule and DE regions, with the contribution of these regions completely determined by the ratio of their partition functions.

Kinetics that are well approximated by double-exponential decay curves were observed in recent folding experiments with PGK and Ub*G^{27,28} and with the λ -repressor.^{35,36} In the first case the fast phase was associated with on-pathway intermedi-

ates, and in the second with motions on the “molecular” time scale. Our results suggest an alternative explanation; i.e., that in both cases there are off-pathway intermediates not in equilibrium with the globular state, and that the fast phase corresponds to the folding trajectories which do not visit these intermediates.

Given the kinetic model and the folding time distributions, the rate constants and principal channels of transitions were obtained from the kinetic model. The rate constants for a wide temperature range showed how the significance of different channels varied with the temperature, and which of the off-pathway minima are in equilibrium with the globule state and which are not.

A random (frustrated) lattice 27-mer with hydrophobicity, was recently studied by Socci et al.,³⁷ who used the fraction of unfolded structures $n_{\text{unfold}}(t)$, defined in eq 8, to analyze the folding kinetics. At temperatures higher than the folding temperature, this system exhibits two-state behavior with single-exponential kinetics. The results contrast with those for the minimally frustrated sequence, which shows downhill folding with multiexponential kinetics. As is seen from Figure 17 of the paper by Socci et al., where the folding trajectory is mapped onto the free energy surface, the system frequently visits fully compact nonnative structures, which play the role of off-pathway minima (DEs). The single-exponential kinetics found in the calculations, indicates that, at the given temperatures, the DEs are in equilibrium with the globule, and a more complex, e.g., double-exponential kinetics, can be expected for this system at lower temperatures.

We also showed that a large scatter of the folding times at long times, where the statistics are poor, represents an intrinsic property of the exponential decay law. Monte Carlo simulations based on the kinetic scheme underlying the analytic model successfully reproduce this statistical effect.

The analysis of folding time distributions presented here has demonstrated their importance as an approach to understanding the mechanism of folding. It suggests that measurements of these distributions both in ensemble and single molecule experiments would be of interest. Single molecule experiments⁵⁸ could be particularly useful because they provide a way, in principle, of determining whether nonexponential folding time distributions arise from homogeneous or inhomogeneous processes. This has its analogue in early work on the rebinding of CO to myoglobin after photolysis.^{59,60}

Acknowledgment. This work was supported by a grant from the INTAS (#2001-2126). S.F.C. also acknowledges support from the RFBR (#02-03-32048) and SB RAS (#119). The work done at Harvard University was supported in part by the National Institutes of Health.

Appendix: Configuration Entropy

Equation 3 can be written as

$$TS(N, N_{\text{nat}}) = \frac{\sum E(N, N_{\text{nat}}) \exp[-E(N, N_{\text{nat}})/T]}{\sum \exp[-E(N, N_{\text{nat}})/T]} + T \ln \sum \exp[-E(N, N_{\text{nat}})/T] \quad (43)$$

where $E(N, N_{\text{nat}})$ is the energy of a structure with N total and N_{nat} native contacts, and the sums are over the structures which are characterized by the given values of N and N_{nat} . We group the structures with close energies and write $\sum \exp[-E(N, N_{\text{nat}})/T] = \sum_k Q_k(N, N_{\text{nat}}) \exp[-E_k(N, N_{\text{nat}})/T]$, where $Q_k(N, N_{\text{nat}})$ is

the number of the structures in group k , which is characterized by energy $E_k(N, N_{\text{nat}})$. Since the energy of the system, given by eq 1, is mainly determined by the numbers of total and native contacts (which, in particular, resulted in good correlation of the mean energy \bar{E} with the values of N and N_{nat} ; see the discussion of eq 4), the distribution of $Q_k(N, N_{\text{nat}})$ over k typically has a sharp maximum at certain $k = m$. Therefore, we have

$$\sum \exp[-E(N, N_{\text{nat}})/T] = \sum_k Q_k(N, N_{\text{nat}}) \exp[-E_k(N, N_{\text{nat}})/T] \approx Q_m(N, N_{\text{nat}}) \exp[-E_m(N, N_{\text{nat}})/T]$$

and

$$\sum E(N, N_{\text{nat}}) \exp[-E(N, N_{\text{nat}})/T] = \sum_k E_k(N, N_{\text{nat}}) Q_k(N, N_{\text{nat}}) \exp[-E_k(N, N_{\text{nat}})/T] \approx E_m(N, N_{\text{nat}}) Q_m(N, N_{\text{nat}}) \exp[-E_m(N, N_{\text{nat}})/T]$$

Consequently, eq 43 reduces to

$$S(N, N_{\text{nat}}) \approx \ln Q_m(N, N_{\text{nat}}) \sim \ln Q(N, N_{\text{nat}})$$

where $Q(N, N_{\text{nat}})$ is the number of the structures with the given N and N_{nat} .

References and Notes

- (1) Dinner, A. A.; Sali, A.; Smith, L.; Dobson, C. M.; Karplus, M. *Trends Biol. Sci.* **2000**, 25, 331–339.
- (2) Dobson, C. M. *Nat. Rev.* **2003**, 2, 154–160.
- (3) Levinthal, C. *J. Chim. Phys.* **1968**, 65, 44–45.
- (4) Karplus, M.; Weaver, D. L. *Nature* **1976**, 260, 404–406.
- (5) Dill, K. A. *Biochemistry* **1985**, 24, 1501–1509.
- (6) Abkevich, V. I.; Gutin, A. M.; Shakhnovich, E. I. *Biochemistry* **1994**, 33, 10026–10036.
- (7) Islam, S. A.; Karplus, M.; Weaver, D. L. *J. Mol. Biol.* **2002**, 318, 199–215.
- (8) Baldwin, R. L. *Nature* **1994**, 369, 183–184.
- (9) Pitsyn, O. B.; Rushin, A. A. *Biophys. Chem.* **1975**, 3, 1–20.
- (10) Gō, N. *Annu. Rev. Biophys. Bioeng.* **1983**, 12, 183–210.
- (11) Harrison, S. C.; Durbin, R. *Proc. Natl. Acad. Sci. U.S.A.* **1985**, 82, 4028–4030.
- (12) Bryngelson, J. D.; Wolynes, P. G. *J. Phys. Chem.* **1989**, 93, 6902–6915.
- (13) Leopold, P. E.; Montal, M.; Onuchic, J. N. *Proc. Natl. Acad. Sci. U.S.A.* **1992**, 89, 8721–8725.
- (14) Sali, A.; Shakhnovich, E.; Karplus, M. *J. Mol. Biol.* **1994**, 235, 1614–1636; Sali, A.; Shakhnovich, E.; Karplus, M. *Nature* **1994**, 369, 248–251.
- (15) Pande, V. S.; Grosberg, A. Y.; Tanaka, T. *Proc. Natl. Acad. Sci. U.S.A.* **1994**, 91, 12972–12975.
- (16) Guo, Z. Y.; Thirumalai, D. *Biopolymers* **1995**, 36, 83–102.
- (17) Zwanzig, R. *Proc. Natl. Acad. Sci. U.S.A.* **1995**, 92, 9801–9804.
- (18) Dill, K. A.; Chan, H. S. *Nat. Struct. Biol.* **1997**, 4, 10–19.
- (19) Karplus, M. *Fold. Des.* **1997**, 2, 569–576.
- (20) Onuchic, J. N.; Luthey-Schulten, Z.; Wolynes, P. G. *Annu. Rev. Phys. Chem.* **1997**, 48, 545–600.
- (21) Dobson, C. M.; Sali, A.; Karplus, M. *Angew. Chem., Int. Ed.* **1998**, 37, 869–893.
- (22) Shea, J.-E.; Brooks, C. L., III. *Annu. Rev. Phys. Chem.* **2001**, 52, 499–535.
- (23) Plotkin, S. S.; Onuchic, J. N. *Q. Rev. Biophys.* **2002**, 35, 111–167; **2002**, 35, 205–286.
- (24) Lazaridis, T.; Karplus, M. *Science* **1997**, 278, 1928–1931.
- (25) Wolynes, P. G.; Onuchic, J. N.; Thirumalai, D. *Science* **1995**, 267, 1619–1620.
- (26) Colon, W.; Wakem, L. P.; Sherman, F.; Roder, H. *Biochemistry* **1997**, 36, 12535–12541.
- (27) Sabelko, J.; Ervin, J.; Gruebele, M. *Proc. Nat. Acad. Sci. U.S.A.* **1999**, 96, 6031–6036.

- (28) Osváth, S.; Sabelko, J.; Gruebele, M. *J. Mol. Biol.* **2003**, *333*, 187–199.
- (29) Becker, O. M.; Karplus, M. *J. Chem. Phys.* **1997**, *106*, 1495–1517.
- (30) Krivov, S. V.; Karplus, M. *J. Chem. Phys.* **2002**, *117*, 10894–10903.
- (31) Evans, D. A.; Wales, D. J. *J. Chem. Phys.* **2003**, *118*, 3891–3897.
- (32) Gavrilov, A. V.; Chekmarev, S. F. In *Bioinformatics of Genome Regulation and Structure*; Kolchanov, N., Hofstaedt, R., Eds.; Kluwer Academic Publishers: Boston, 2004; pp 171–178.
- (33) Krivov, S. V.; Karplus, M. *Proc. Nat. Acad. Sci. U.S.A.* **2004**, *101*, 14766–14770.
- (34) Munoz, V.; Thompson, P.; Hofrichter, J. A.; Eaton, W. A. *Nature* **1997**, *390*, 196–199.
- (35) Yang, W. Y.; Gruebele, M. *Nature* **2003**, *423*, 193–197.
- (36) Yang, W. Y.; Gruebele, M. *Biophys. J.* **2004**, *87*, 596–608.
- (37) Socci, N. D.; Onuchic, J.; Wolynes, P. G. *Proteins* **1998**, *32*, 136–158.
- (38) Gruebele, M. *Curr. Opin. Struct. Biol.* **2002**, *12*, 161–168.
- (39) Schuler, B.; Lipman, E. A.; Eaton, W. A. *Nature* **2002**, *419*, 743–747.
- (40) Metzler, R.; Klafter, J.; Jortner, J.; Volk, M. *Chem. Phys. Lett.* **1998**, *293*, 477–484.
- (41) Skorobogatiy, M.; Guo, H.; Zuckermann, M. *J. Chem. Phys.* **1998**, *109*, 2528–2535.
- (42) Nakamura, H. K.; Sasai, M.; Takano, M. *Proteins* **2004**, *55*, 99–106.
- (43) Saven, J. G.; Wang, G.; Wolynes, P. G. *J. Chem. Phys.* **1994**, *101*, 11037–11043.
- (44) Wang, G.; Saven, J. G.; Wolynes, P. G. *J. Chem. Phys.* **1996**, *105*, 11276–11284.
- (45) Wang, G.; Onuchic, J.; Wolynes, P. G. *Phys. Rev. Lett.* **1996**, *25*, 4861–4864.
- (46) Zhou, R. H. *Proteins* **2003**, *53*, 148–161.
- (47) Cavalli, A.; Haberthür, U.; Paci, E.; Caflisch, A. *Protein Sci.* **2003**, *12*, 1801–1803.
- (48) Clementi, C.; Nymeyer, H.; Onuchic, J. N. *J. Mol. Biol.* **2000**, *298*, 937–953.
- (49) Li, L.; Shakhnovich, E. I. *Proc. Natl. Acad. Sci. U.S.A.* **2001**, *98*, 13014–13018.
- (50) Paci, E.; Caflisch, A.; Vendruscolo, M. (to be published).
- (51) Dinner, A. A.; Šali, A.; Karplus, M. *Proc. Natl. Acad. Sci. U.S.A.* **1996**, *93*, 8356–8361.
- (52) Dinner, A. A.; Karplus, M. *Nat. Struct. Biol.* **1998**, *5*, 236–241.
- (53) Oliveberg, M.; Yan, Y.-J.; Fersht, A. R. *Proc. Nat. Acad. Sci. U.S.A.* **1995**, *92*, 8926–8929.
- (54) Socci, N. D.; Onuchic, J. N. *J. Chem. Phys.* **1994**, *101*, 1519–1528.
- (55) Gruebele, M. *Annu. Rev. Phys. Chem.* **1999**, *50*, 485–516.
- (56) Abkevich, V. I.; Gutin, A. M.; Shakhnovich, E. I. *J. Chem. Phys.* **1994**, *101*, 6052–6062.
- (57) Zhou, Y.; Karplus, M. *Nature* **1999**, *401*, 400–403.
- (58) Fernandez, J. M.; Hongbin, L. *Science* **2004**, *303*, 1674–1678.
- (59) Austin, R. H.; Beeson, K. W.; Eisenstein, L.; Frauenfelder, H.; Gunsalus, I. C. *Biochemistry* **1975**, *14*, 5355–5372.
- (60) Petrich, J. W.; Lambry, J.-C.; Kuczera, K.; Karplus, M.; Poyart, C.; Martin, J.-L. *Biochemistry* **1991**, *30*, 3975–3987.