# Statistical mechanics of a correlated energy landscape model for protein folding funnels

Steven S. Plotkin, Jin Wang, and Peter G. Wolynes
*Department of Physics and School of Chemical Sciences, University of Illinois, Urbana, Illinois 61801*

In heteropolymers, energetic correlations exist due to polymeric constraints and the locality of interactions. Pair correlations in conjunction with the *a priori* specification of the existence of a particularly low energy state provide a method of introducing the aspect of minimal frustration to the energy landscapes of random heteropolymers. The resulting funneled landscape exhibits both a phase transition from a molten globule to a folded state, and the heteropolymeric glass transition in the globular state. We model the folding transition in the self-averaging regime, which together with a simple theory of collapse allows us to depict folding as a double-well free energy surface in terms of suitable reaction coordinates. Observed trends in barrier positions and heights with protein sequence length and thermodynamic conditions are discussed within the context of the model. We also discuss the new physics which arises from the introduction of explicitly cooperative many-body interactions, as might arise from sidechain packing and nonadditive hydrophobic forces. © *1997 American Institute of Physics.* [S0021-9606(97)52406-8]

## I. INTRODUCTION

Molecular scientists view protein folding as a complex chemical reaction. Another fruitful analogy from statistical physics is that folding resembles a phase transition in a finite system. A new view of the folding process combines these two ideas along with the notion that a statistical characterization of the numerous possible protein configurations is sufficient for understanding folding kinetics in many regimes.

The resulting energy landscape theory of folding acknowledges that the energy surface of a protein is rough, containing many local minima like the landscape of a spin glass. On the other hand, in order to fold rapidly to a stable structure there must also be guiding forces that stabilize the native structure substantially more than other local minima on the landscape. This is the principle of minimum frustration.[1] The energy landscape can be said then to resemble a ''funnel.''[2] Folding rates then depend on the statistics of the energy states as they become more similar to the native state at the bottom of the funnel.

One powerful way of investigating protein energy landscapes has been the simulation of ''minimalist'' models. These models are not fully atomistic, but caricature the protein as a series of beads on a chain either embedded in a continuum[3] or on a lattice.[4] A correspondence, in the sense of phase transition theory, between these models and real proteins has been set up using energy landscape ideas.[5] Many issues remain to be settled however in understanding how these model landscapes and folding mechanisms change as the system under study becomes larger and as one introduces greater complexity into the modeling of this correspondence, as for example, by explicitly incorporating many-body forces and extra degrees of freedom. Simulations become cumbersome for such surveys, and an analytical understanding is desirable.

Analytical approaches to the energy landscape of proteins have used much of the mathematical techniques used to treat spin glasses[6] and regular magnetic systems.[7] The polymeric nature of the problem must also be taken into account. Mean field theories based on replica techniques[8] and variational methods[9] have been very useful, but are more difficult to make physically intuitive than the straightforward approach of the random energy model,[10] which flexibly takes into account many of the types of partial order expected in biopolymers.[11] Recently we have generalized the latter approach to take into account correlations in the landscape of finite-sized random heteropolymers.[12] This treatment used the formalism of the generalized random energy model (GREM) analyzed by Derrida and Gardner.[13] In this paper, we extend that analysis to take into account the minimum frustration principle and thereby treat proteinlike, partially nonrandom heteropolymers.

There are various ways of introducing the aspect of minimum frustration to analytical models with rugged landscapes. One way recognizes that many empirical potentials actually are obtained by a statistical analysis of a database, and when the database is finite, there is automatically an aspect of minimal frustration for any member of that database. Thus the so-called ''associative memory'' Hamiltonian models[14] have coexisting funnel-like and rugged features in their landscape. Other methods of introducing minimal frustration model the process of evolution as giving a Boltzmann distribution over sequences for an energy gap between a fixed target structure and unrelated ones.[15] All of the above approaches can be straightforwardly handled with replica-based analyses. Here we show that the GREM analyses can be applied to minimally frustrated systems merely by requiring the energy of a given state to be specified as having a particularly low value (i.e., less than the putative ground state value). Minimally frustrated, funneled landscapes are just a special case of the general correlated landscape studied earlier.

A convenient aspect of the correlated landscape model is that it allows the treatment of the polymer physics in a very direct way, using simple statistical thermodynamics in the tradition of Flory.[16] Here we will show how the interplay of collapse and topological ordering can be studied. In order to do this we introduce a simple ''core-halo'' model to take into account the spatially inhomogeneous density. We will also discuss the role of many-body forces in folding. Explicitly cooperative many-body forces have often been involved in the thinking about protein structure formation. Hydrophobic forces are often modeled as involving buried surface area. Such an energy term is not pairwise additive but involves three or more interacting bodies. Sidechain packing involves objects fitting into holes created by more than one other part of the chain, thus the elimination of sidechains from the model can yield an energy function for backbone units with explicit nonadditivity. These many-body forces can be treated quite easily by the GREM, and we will see that they can make qualitative changes in the funnel topography.

To illustrate the methods here, we construct two-dimensional free energy surfaces for the folding funnel of minimally frustrated polymers. These explicitly show the coupling between density and topological similarity in folding. We pay special attention to the location of the transition state ensemble and discuss how this varies with system size, cooperatively of interactions, and thermodynamic conditions. In the case of the 27-mer on a lattice, a detailed fit to the lattice simulation data[4] is possible. Although delicate cancellations of energetic and entropic terms are involved in the overall free energy, plausible parameters fit the data.

The trends we see in the present calculations are in rough agreement with experimental information on the nature and location of the transition state ensemble,[17,18] although the theory suggests that fluctuation mechanisms, in the form of independently folding units (foldons) [A. R. Panchenko et al., Proc. Natl. Acad. Sci. USA **93**, 2008 (1996)] become more important at larger $N$. We intend later to return to the experimental comparison, especially taking into account more structural details within the protein.

The organization of this paper is as follows: In Sec. II we introduce a theory of the free energy at constant density, and in this context investigate the effects of cooperative interactions on the transition state ensemble and corresponding free energy barrier. In Sec. III we detail a simple theory coupling collapse with topological similarity, and resulting in the core-halo model described there. In Sec. IV we apply this collapse theory to obtain the free energy in terms of density and topological order, now coupled via the core-halo model. In the same section we compare our model of the minimally frustrated heteropolymer with lattice simulations of the 27-mer. In terms of the categorization of Bryngelson et al.[2] these free energy surfaces depict scenarios described as type I or type IIa folding. We then study the quantitative aspects of the barrier as a function of the magnitude of three-body effects. The dependence of position and height of the barrier as a function of sequence length is studied, as well as the effects of increasing the stability gap. Finally, we study

the denaturation curve as determined by the constant and variable density models. In Sec. VI we discuss the results and conclude with some remarks.

## II. A THEORY OF THE FREE ENERGY

In this section, we show how the existence of a particularly low energy configuration, together with energetic correlations for similar states, leads to a model for the folding transition and corresponding free energy surface in protein-like heteropolymers. This ansatz for the correlated energy landscape corresponds to the introduction of minimal frustration in a random energy landscape, where the order parameter here (which will function as a reaction coordinate for the folding transition) counts the number of native contacts or hydrogen bonds.

We start by assuming a simple ''ball and chain'' model for a protein which is readily comparable with simulations, e.g., of the 27-mer, which is widely believed to capture many of the quantitative aspects of folding (Sec. IV). Proteins with significant secondary structure have an effectively reduced number of interacting units as may be described by a ball and chain model. Properties of both, when appropriately scaled by critical state variables such as the folding temperature $T_F$, glass temperature $T_G$, and collapse temperature $T_C$, will obey a law of corresponding states.[5] Thus the behavior describing a complicated real protein can be validly described by an order parameter applied to a minimal ball and chain model in the same universality class.

For a 27-mer on a three-dimensional cubic lattice, there are 28 contacts in the most collapsed (cubic) structure. For concreteness we take such a maximally compact structure to be the configuration of our ground state, the generalization to a less compact ground state being straightforward in the context of the model to be described. For a collapsed polymer of sequence length $N$, the number of pair contacts per monomer, $z_N$, is a combination of a bulk term, a surface term, and a lattice correction[19] [see Eq. (A4)]. The effect of the surface on the number of contacts is quite important even for large macromolecules, as $z_N$ approaches its bulk value of 2 contacts per monomer rather slowly, as $\sim 2-3N^{-1/3}$.

To describe states that are not completely collapsed, we introduce the packing fraction $\eta \cong N\sigma/R_g^3$ as a measure of the density of the polymer, where $\sigma$ is the volume per monomer and $R_g$ is the radius of gyration of the whole protein. So for less dense states the total number of contacts is reduced from its collapsed value $Nz_N$, to $Nz_N\eta$.

In the spirit of the lattice model we have in mind for concreteness, we introduce a simple contact Hamiltonian to determine the energy of the system,

$$\mathcal{H} = \sum_{i<j} \epsilon_{ij}\sigma_{ij}, \tag{2.1}$$

where $\sigma_{ij}=1$ when there is a contact made between monomers $\{ij\}$ in the chain, and $\sigma_{ij}=0$ otherwise. Here contact means that two monomers $\{ij\}$, nonconsecutive in sequence along the backbone chain, are adjacent in space at neighboring lattice sites. $\epsilon_{ij}$ is a random variable so that, at constant

density, the total energies of the various configurations, each the sum of many $\epsilon_{ij}$, are approximately Gaussianly distributed by the central limit theorem, with mean energy at a given density $\eta$ given by $\bar{E}_\eta = N z_N \eta \bar{\epsilon}$, where $\bar{\epsilon}$ is simply defined as the mean energy per contact and $N z_N \eta$ is again the total number of contacts, and variance $\Delta E_\eta^2 = N z_N \eta \epsilon^2$, where $\epsilon^2$ is the effective width of the energy distribution per contact.

Suppose there exists a configurational state $n$ of energy $E_n$ (which will later become the "native" state). Then if the Hamiltonian for our system is defined as in Eq. (2.1), we can find the probability that configuration $a$ has energy $E_a$, given that $a$ has an overlap $Q_{an}$ with $n$,[12] where $Q_{an} \equiv Q$ is the number of contacts that state $a$ has in common with $n$, divided by the total number of contacts $N z_N \eta$,

$$Q = \frac{1}{N z_N \eta} \sum \sigma_{ij}^a \sigma_{ij}^n. \qquad (2.2)$$

Since this analysis is at constant density both $a$ and $n$ have $N z_N \eta$ contacts. This probability is obtained directly from the Hamiltonian (2.1) by averaging over Gaussian distributions of contact energies $\epsilon_{ij}$,

$$\frac{P_{an}(E_a, Q, E_n)}{P_n(E_n)} = \frac{\langle \delta[E_a - \mathcal{H}(\{\sigma_{ij}^a\})] \delta[E_n - \mathcal{H}(\{\sigma_{ij}^n\})] \rangle}{\langle \delta[E_n - \mathcal{H}(\{\sigma_{ij}^n\})] \rangle}, \qquad (2.3)$$

where $\{\sigma_{ij}^a\}$ is the set of contacts in configuration $a$. The conditional probability distribution is simply a Gaussian with a $Q$ dependent mean and variance,

$$\frac{P_{an}(E_a, Q, E_n)}{P_n(E_n)} \sim \exp\left( -\frac{[(E_a - \bar{E}) - Q(E_n - \bar{E})]^2}{2 N z_N \eta \epsilon^2 (1 - Q^2)} \right). \qquad (2.4)$$

When $Q = 1$, states $a$ and $n$ are identical and must then have the same energy, which Eq. (2.4) imposes by becoming delta function, and when $Q = 0$ states $a$ and $n$ are uncorrelated and then Eq. (2.4) becomes the Gaussian distribution of the random energy model for the energy of state $a$. Expression (2.4) holds for all states of the same density as $n$, e.g., all collapsed states if $n$ is the native state (the degree of collapse must be a somewhat coarse-grained description to avoid fluctuations due to lattice effects coupled with finite size).

Previously a theory was developed of the configurational entropy[20] $S_\eta(Q)$ as a function of similarity $Q$ with a given state, at constant density $\eta$.[12] The results of this theory are summarized in Appendix A. Given $S_\eta(Q)$ and the conditional probability distribution (2.4), the average number of states of energy $E$ and overlap $Q$ with state $n$, all at density $\eta$, is

$$\langle n_\eta(E, Q, E_n) \rangle = e^{S_\eta(Q)} \frac{P(E, Q, E_n)}{P(E_n)}$$

$$\sim \exp N \left\{ s_\eta(Q) - \frac{1}{2(1 - Q^2)} \times \left( \frac{(E - \bar{E}) - Q(E_n - \bar{E})}{N J_\eta} \right)^2 \right\}, \qquad (2.5)$$

where $J_\eta^2 \equiv z_N \eta \epsilon^2$ and $s_\eta(Q) \equiv S_\eta(Q)/N$. Equation (2.5) is still Gaussian with a large number of states provided $E - \bar{E}$ is within a band of energies having $E_c - \bar{E} = Q(E_n - \bar{E}) \pm N J_\eta \sqrt{2(1 - Q^2) s_\eta(Q)}$ as upper and lower bounds. There is a negligibly small number of states with energies above or below this range [where the exponent changes sign in Eq. (2.5)].

At temperature[21] $T$, the Boltzmann factor $e^{-E/T}/z$ weighting each state shifts the number distribution of energies so that the maximum of the thermally weighted distribution can be interpreted as the most probable (thermodynamic) energy at that temperature

$$E_\eta(T, Q, E_n) = \bar{E} + Q(E_n - \bar{E}) - \frac{N z_N \eta \epsilon^2 (1 - Q^2)}{T}. \qquad (2.6)$$

The above expression for the most probable energy is useful provided the distribution (2.5) is a good measure of the actual number of states at $E$ and $Q$, the condition for which is that the fluctuations in the number of states be much smaller than the number of states itself. To this end, we make here the simplifying assumption that in each "stratum" defined by the set of states which have an overlap $Q$ with the native state, the states themselves are not further correlated with each other, i.e., $P(E_a, Q, E_b, Q | E_n) = P(E_a, Q, E_n) \times P(E_b, Q, E_n)$, so that in each stratum of the reaction coordinate $Q$, the set of states is modeled by a random energy model. Then since the number of states $n_\eta(E, Q, E_n)$ counts a collection of random uncorrelated variables—large when $E > E_c$—the relative fluctuations $\sqrt{\langle (n - \langle n \rangle)^2 \rangle}/\langle n \rangle$ are $\sim \langle n \rangle^{-1/2}$ and are thus negligible. So $n(E, Q, E_n) \approx \langle n(E, Q, E_n) \rangle$, and we can evaluate the exponent in the number of states (2.5) at the the most probable energy (2.6) as an accurate measure of the ($Q$-dependent) thermodynamic entropy at temperature $T$,

$$S_\eta(T, Q, E_n) = S_\eta(Q) - \frac{N z_N \eta \epsilon^2 (1 - Q^2)}{2 T^2}. \qquad (2.7)$$

The assumption of a REM at each stratum of $Q$ is clearly a first approximation to a more accurate correlational scheme. The generalization to treat each stratum itself as a GREM as in our earlier work is nevertheless straightforward, since our earlier work suggested only quantitative changes, which we will not pursue here. If two configurations $a$ and $b$ have an overlap $Q$ with state $n$ and thus are correlated to $n$ energetically, they are certainly correlated to each other, particularly for large overlaps where the number of shared contacts is large. Using the REM scheme at each stratum is more accurate for small $Q$ and breaks down to some extent for large $Q$. In the ultrametric scheme of the GREM, states $a$

and $b$ have an overlap $q_{ab} \geqslant Q$, which is more accurate for large overlap than at small $Q$, since at small $Q$ states $a$ and $b$ need not share *any* bonds and still can both have overlap $Q$ with $n$. One can also further correlate the energy landscape of states by stratifying with respect to $q_{ab} = q$ and so on, resulting in a hierarchy of overlaps and correlations best treated using renormalization group ideas.

All of the states in each stratum defined by $Q$ are still correlated to state $n$, and their statistics are correspondingly modified. For smaller values of $Q$, most of the states have zero overlap with each other essentially because their configurational entropy is largest when their is no topological constraint between the states. Microscopic ultrametricity is broken in that $q_{ab}$, the overlap between any two states $(a,b)$ in the stratum, is less than $Q$. As $Q$ increases, there is a crossover to a regime where microscopic ultrametricity becomes a more accurate description. We assume here that this happens typically after $Q \approx 0.5$, where there must be some overlap between states $a$ and $b$. In this regime the form of the thermodynamic functions is modified by the replacement of $(1-Q^2)$ in expressions (2.6), (2.7), (2.9), etc., by $(1-Q)$. For a derivation of these formulas in the ultrametric regime, see Appendix B.

Just as the number of states (2.5) has a characteristic energy for which it is exponentially small, the REM entropy for a stratum at $Q$ (2.7) vanishes at a characteristic temperature

$$\frac{T_g(Q)}{\epsilon} = \sqrt{\frac{z_N \eta (1-Q^2)}{2 s_\eta(Q)}}, \tag{2.8}$$

which signals the trapping of the polymer into a low energy conformational state within the stratum characterized by $Q$.

If $T_g(Q)$ is a monotonically decreasing function of $Q$, as the temperature is lowered the polymer will gradually be thermodynamically confined in its conformational search to smaller and smaller basins of states. The basin around the native state is the largest basin with the lowest ground state, and hence is the first basin within which to be confined. Its characteristic size at temperature $T$ is just the number of states within overlap $Q_0(T)$, where $Q_0(T)$ is the value of overlap $Q$ that gives $T_g(Q_0) = T$ in Eq. (2.8). Thus there is now no longer a single glass temperature at which ergodic confinement suddenly occurs, as in the REM, but there is a continuum of basin sizes to be localized within at corresponding glass temperatures for those basins.

If $T_g(Q)$ has a single maximum at say $Q^*$, the glass transition is characterized by a sudden REM-like freezing to a basin of configurations whose size is determined by $Q^*$. The range of glass temperatures will turn out to be lower than the temperature at which a folding transition occurs (see Fig. 3), so that this model predicts a proteinlike heteropolymer whose folded state is stable by several $k_B T$ at temperatures where freezing becomes important. A replica-symmetric analysis of the free energy is therefore sufficient to describe the folding transition to such deep native states that are minimally frustrated.

From the thermodynamic expressions for the energy (2.6) and entropy (2.7) with the mean energy at density $\eta$, we can write down the free energy per monomer above the glass temperature as the sum of four terms,

$$\frac{F_\eta}{N}(T, Q, E_n) = z_N \eta \bar{\epsilon} + Q z_N \delta \epsilon_n - T s_\eta(Q) - \frac{z_N \eta \epsilon^2}{2T}(1-Q^2), \tag{2.9}$$

where $z_N \delta \epsilon_n = z_N(\epsilon_n - \bar{\epsilon}) = (E_n - \bar{E})/N$ is the extra energy for each bond beyond the mean homopolymeric attraction energy (the energy ''gap'' between an average molten globule structure and the minimally frustrated one), times the number of bonds per monomer, and $s_\eta(Q) = S_\eta(Q)/N$ is the entropy per monomer described in Appendix A.

The first term in Eq. (2.9) multiplied by $N$ is just the homopolymeric attraction energy between all the monomers, for a polymer of density $\eta$. It depends only on the degree of collapse, and not on how many contacts are native contacts. The second term is the average extra bias energy if a contact is native, times the average number of native contacts per monomer. The third term measures the equilibrium bias toward larger configurational entropy at smaller values of the reaction coordinate $Q$. The last term accounts for the diversity of energy states that exist on a rough energy landscape, the variance of which lowers thermodynamically the energy more than the entropy, and so lowers the equilibrium free energy.

For a special surface in $(\delta \epsilon_n, \epsilon, T)$ space, expression (2.9) has a double minimum structure in the reaction coordinate $Q$, with one entropic minimum at low $Q$ corresponding to the ''molten globule'' state, separated by a barrier from an energetic minimum at high $Q$ corresponding to a ''folded'' state. For a given temperature, values of $\delta \epsilon_n$ and $\epsilon^2$ can be obtained which are reasonably close to the values obtained by a more accurate analysis which includes the coupling of density with topology, but we will not examine the constant density case in much detail for reasons discussed below, except to make the following remarks: (1) The true coupling between density and $Q$ constraints need not be strong to obtain a double-well free energy structure. (2) For monomeric units with pair interactions, at constant density, the molten globule and folded minima are not at $Q=0$ and 1, respectively. The position of the molten globule state is near the maximum of the entropy of the system, which is at $Q \cong 0.1$ for the 27-mer due to the interplay of confinement effects and the combinatorial mixing entropy inherent in the ''coarse-grained'' description $Q$.[12] The native minimum shifts to $Q=1$ when many-body interactions are introduced (see the next section). (3) The barrier height, at position $Q^0 \cong 0.25$ for the 27-mer with proteinlike parameters $(T_F/T_G \cong 2)$, is small $(\Delta F^0 \approx k_B T_F)$, due to the effective cancellation of entropy loss by negative energy gain, as the system moves toward the native state (This cancellation is reduced when many-body forces are taken into account). (4) When a linear form for the entropy is used in Eq. (2.9), e.g., $s(Q) = s_0(1-Q)$ instead of the more accurate $s(Q)$ obtained in Ref. 12 the double minimum structure disappears and is replaced by a single minimum near $T_F$ at $Q \approx 1/2$, with the

$Q=0$ and $Q=1$ states becoming free energy maxima. So folding is downhill or spinodal-like in this approximation.

## A. Effects of cooperative interactions

As the interactions between monomeric segments become more explicitly cooperative, the energetic correlations between states become significant only at greater similarity, with the system approaching the REM limit for $\infty$-body interactions, where the statistical energy landscape assumes a rough ''golf-course'' topography with a steep funnel close to the native state.

In the presence of $m$-body interactions, the homopolymer collapse energy scales as a higher power of density ($\sim \overline{\epsilon z}^{m-1}$). For even moderate $m \sim \mathcal{O}(1)$ a first-order phase transition to collapsed states results, which effectively confines all reaction paths in the coordinate $Q$ between molten globule and folded states to those where the density is constant and $\approx 1$. So within this constant density approximation we can investigate the nature of the folding transition as a function of the cooperativity of the interactions, and see how the correlated landscape simplifies to the REM in the limit of $m$-body interactions with large $m$.

In the presence of $m$-body interactions, the $Q$ dependence in the pair energy distribution (2.4) scales with $Q$ as $Q^{m-1}$, when $Q$ is defined as in Eq. (2.2), and the terms $\sigma_{ijk\cdots m}$ in the modified Hamiltonian factorize into pair interaction terms $\sigma_{ij}\sigma_{jk}\cdots$ through a suitable decomposition law such as in the superposition approximation in the theory of fluids.[22] Using this modified pair distribution along with the collapsed homopolymeric state as our zero point energy, the free energy (2.9) becomes

$$\frac{F}{N}(T,Q,E_n) = -Ts_1(Q) - Q^{m-1} z_N |\delta\epsilon_n|$$
$$- \frac{z_N \epsilon^2}{2T}(1 - Q^{2(m-1)}), \qquad (2.10)$$

where $s_1(Q)$ is the entropy as a function of constraint $Q$ for a fully collapsed polymer. For pure three-body interactions and higher, the globule and folded states are very nearly at $Q \cong 0$ and $Q \cong 1$, respectively [see Fig. 1(a)]. To the extent that this approximation is good, we can equate the free energies of the molten globule (at $Q = Q_{MG} \cong 0$) and folded ($Q = 1$) structures and obtain an $m$-independent folding temperature (note again that this is not a good approximation for pair interactions),

$$T_F = \frac{z_N |\delta\epsilon_n|}{2s_0}\left(1 + \sqrt{1 - \frac{2s_0\epsilon^2}{z_N \delta\epsilon_n^2}}\right), \qquad (2.11)$$

where $s_0$ is the maximum of the entropy as a function of constraint $Q$ (essentially the log of the total number of configurations).

From expression (2.11) we can obtain a first approximation to the constraint on the magnitude of the gap energy $\delta\epsilon_n$ in order to have a global folding transition (rather than merely a local glass transition) to the low energy state in question. The condition that the square root term in Eq. (2.11) be real gives the minimum gap for global foldability in terms of the roughness $\epsilon$,

$$\frac{\delta\epsilon_n^{(c)}}{\epsilon} = \sqrt{\frac{2s_0}{z_N}} \approx \sqrt{2}, \qquad (2.12)$$

where the minimum folding temperature is then $k_B T_F^{(c)} \approx \delta\epsilon_n^{(c)}/2$ (or equivalently, one can obtain the maximum roughness for foldability as $\approx 1/\sqrt{2}$ of a given gap energy). For typical proteins (with folding temperatures at $\approx 330$ K) gap energies are (at least) $\approx 1$ kcal/mol (lattice unit). Note that Eq. (2.12) is precisely the same result, as it should be, to that obtained previously[23] in the context of finding optimal folding energy functions, by requiring the quantity $T_F/T_G > 1$, where the glass temperature $T_G = \sqrt{z_N \epsilon^2/(2s_0)}$ is evaluated at the molten globule overlap $Q = Q_g$.

Evaluating $F(T,Q,E_n)/N$ [Eq. (2.10)] with proteinlike energetic parameters at the folding temperature $T_F$, we obtain free energy curves as in Fig. 1(a), plotted for illustrative examples with $m=3$ and $m=12$, for a 27-mer lattice protein.

Note that the transition state ensemble (the collection of states at $Q = Q^*$ where the free energy is a maximum) becomes more and more native like (and thus the ensemble becomes smaller and smaller, eventually going to 1 state in the REM) as the energy correlations become more short-ranged in $Q$ (i.e., as $m$ increases)—see Fig. 1(b). The corresponding free energy barrier then grows with $m$ as the energetic bias ($\sim Q^{m-1}$) overcomes the entropic barrier only much closer to the native state, and the barrier becomes more and more entropic and less energetic [see Fig. 1(c)]. There are less kinetic paths to the native state through the transition state ensemble.

As was already mentioned, the above analysis was for a polymer of constant collapse density. However, experimental evidence of folding, as well as numerical evidence for lattice models, suggest a coupling of density with nativeness, with energetically favorable nativelike states typically being denser. So to this end we now investigate in detail a simple theory coupling collapse density $\eta$ with nativeness $Q$, assuming a native ($Q=1$) state which is completely collapsed ($\eta=1$). Including this effect in Eq. (2.9) will complete our simple model of the folding funnel topography in two reaction-coordinate dimensions.

## III. A SIMPLE THEORY OF COLLAPSE

The GREM theory for random heteropolymers developed by us earlier investigates the interplay between entropy loss and energetic roughness as a function of similarity to any given reference state, all at fixed density. However for exceptional reference states such as the ground state of a well-packed protein, the density is not independent of configurational similarity, so a theory of the coupling of density $\eta$ with topological similarity $Q$ must also be developed.

We wish to obtain the polymer density $\eta$ as a function of both the fraction of total contacts $\overline{z}/z_N$ and fraction of native contacts $Q$. At low degrees of nativeness, a good approxima-

(a)



(b)



(a)



(b)

FIG. 2. (a) A model of the partially native protein can be pictured as a frozen native core surrounded by a halo of non-native polymer of variable density. (b) The halo density $\eta_H$ as a function of the fraction of the total contacts $\bar{z}/z_N$.
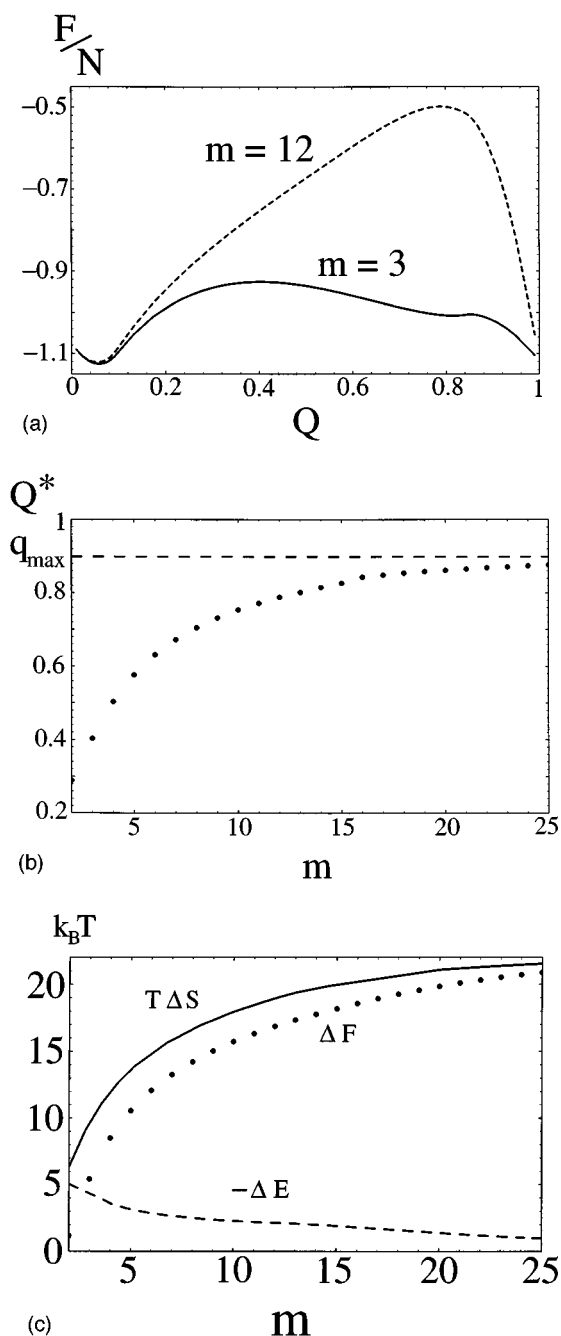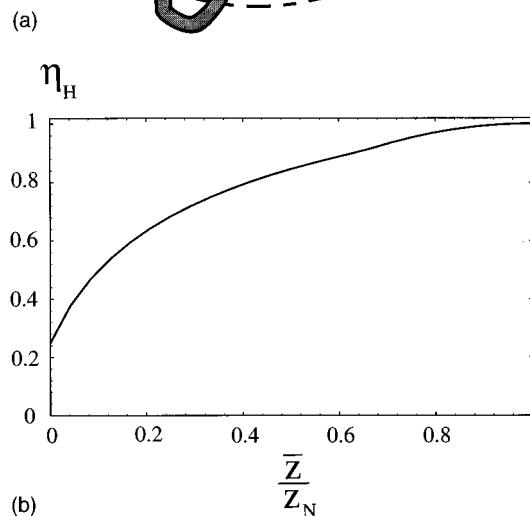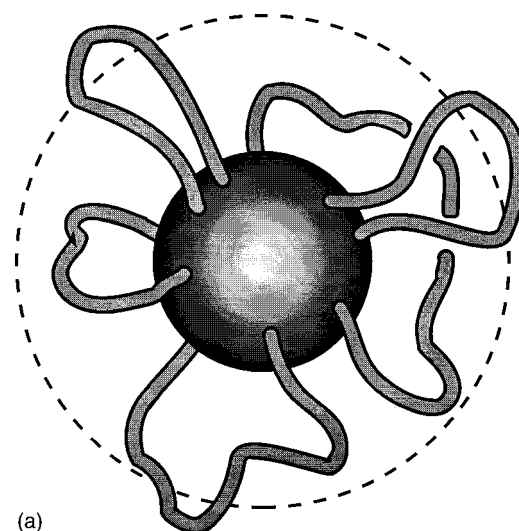
FIG. 1. (a) Free energy per monomer $F/N$ for a 27-mer, in units of $k_B T_F$ as a function of $Q$, at constant density $\eta=1$, with $s_0=0.8$, for proteinlike energetic parameters $(\delta\epsilon_n,\epsilon)=(-2.28, 1.55)$. For these parameters $T_F \approx |\delta\epsilon_n|$. For illustrative purposes, two values of $m$-body interactions are chosen: (*solid line*) pure three-body interactions; (*dashed line*) pure twelve-body interactions. Note the trends in height and position of the barrier, and note how in the $m=12$ case the free energy curve is essentially $-T$ times the entropy curve $s(Q)$ of Fig. 10 with $\eta=1$, until $Q$ is very large. (b) Position of the transition state ensemble $Q^*$ along the reaction coordinate $Q$ as a function of the explicit cooperativity in pure $m$-body forces, $m$. The fact that the asymptotic limit $Q_{max}$ is less than one is due to the finite size of the system, so that $Q$, the fraction of native contacts, is not a continuous parameter. (c) The free energy barrier height $\Delta F$ in units of $k_B T_F$ as a function of the explicit cooperativity of the $m$-body force, $m$. The barrier height rises to the limit of $S(Q=0)$ as $m \to \infty$, when it becomes completely entropic. Also shown are the energetic (dashed) and entropic (solid) contributions to the barrier.

tion to the density can be obtained by assuming homogeneous collapse, or $\eta=\bar{z}/z_N$. In this approximation the density is a function of total contacts only, irrespective of nativeness.

At higher degrees of nativeness, we adopt a simple model consisting of a native ''core'' region of density $\eta\cong 1$ surrounded by a typically less dense ''halo'' region ($\eta_H \lesssim 1$) consisting of dangling loops and ends [see Fig. 2(a)]. We then seek the functional form of the halo density $\eta_H(Q,\bar{z})$.

We make the approximation that at constant $\bar{z}$, the halo density is approximately constant. i.e., contacts will increase the halo density by reducing the effective loop size in the halo, which is determined by total contacts only, irrespective of nativeness. Note that the *total* density includes both core and halo density, and is of course $Q$ dependent.

We determine the function $\eta_H(\bar{z}/z_N)$ by first constructing a theory of loop density for a given length. Then we obtain the loop length as a function of total contacts by calculating it along the line $\bar{z}=z_N Q$, and using the $Q$ dependence of free loop length from the high $Q$ entropy theory of

Appendix A. The high $Q$ theory gives $l_E(Q)$ and $\langle l_{\text{melted}}(Q) \rangle$ along the line $Q = \bar{z}/z_N$. Using the fact that $\eta_H$ should be a function of total contacts only gives the density for all $Q$ and $\bar{z}$. The high $Q$ expressions for loop sizes will be a good approximation since the molten globule and folded states are largely collapsed, so that $Q = \bar{z}/z_N$ is in the strongly constrained regime.

To estimate the packing fraction of a polymer end $\eta_H^E(l_E)$, consider a chain of sequence length $l_E$ with intrinsic volume $l_E b^3$, confined to a half-plane making a self-avoiding walk. The characteristic volume of space it occupies is given by $\frac{1}{2} R_{\text{rms}}^3 = \frac{1}{2} l_E^{9/5} b^3$, so its packing fraction is $\eta_H^E \cong 2 l_E^{-4/5}$. For a loop of length $l$ its characteristic volume is approximately $\frac{1}{2}(l/2)^{9/5} b^3$, giving a denser packing fraction $\eta_H^L(l) \cong 2 \times 2^{9/5} l^{-4/5}$. Next, we average the density over the total number of loops and ends,

$$\langle \eta_H \rangle = \frac{\Sigma n_i \eta_i}{\Sigma n_i} = 2 \frac{Nf 2^{9/5} l^{-4/5} + 2 l_E^{-4/5}}{Nf + 2}. \tag{3.1}$$

Equation (3.1) is the halo density when the loops and ends have no cross-links or bonds, i.e., along the line $\bar{z} = z_N Q$. The quantities $l(Q) = \langle l_{\text{melted}}(Q) \rangle$, $l_E(Q)$, and $Nf(Q)$ are taken from the high $Q$ entropy analysis (see Ref. 12 and Appendix A). Putting these values into Eq. (3.1) gives the halo density $\eta_H(Q)$ along the line $Q = \bar{z}/z_N$. Then we use the independence of loop density on the nativeness of contacts made so that $\eta_H(Q, \bar{z}/z_N) = \eta_H(\bar{z}/z_N)$.

Figure 2(b) gives a plot of $\eta_H(\bar{z}/z_N)$. The value at $\bar{z} = 0$ is the density of an end of length $N/2$ ($\approx 0.24$ for a 27-mer). The true packing fraction should be roughly $1/2$ of this, however this artifact of the theory has little effect on the folding transition, which involves states with $\bar{z}/z_N$ typically larger than $\cong 0.6$.

We can now reinvestigate the glass transition temperature as a function of both $Q$ and $\bar{z}$ through the insertion of the halo density $\eta_H(\bar{z}/z_N)$ into Eq. (2.8). This gives the regions in the space of these reaction coordinates where the dynamics tends to become glassy if $T_g(Q, \bar{z})$ is comparable to $T_F$ (see Fig. 3). We can see from Fig. 3 that $T_F/T_G$ grows during the folding process, with a corresponding slowing down of the dynamics. $T_F \cong T_G$ when $Q \cong 0.85$. At these high values of $Q$ the dynamics is glassy. At the transition state, $(Q^*, \bar{z}^*) \approx (0.50, 0.92)$, $T_F/T_G \cong 1.5$ [these coordinates are determined in Sec. IV, see Fig. 4(b)]. $T_F/T_G \approx 2.3$ in the molten-globule phase at $(Q_{\text{MG}}, \bar{z}_{\text{MG}}) = (0.14, 0.67)$. This value is larger than that obtained from simulations: $T_F/T_G(Q_{\text{MG}}, \bar{z}_{\text{MG}}) \cong 1.6$. Values closer to 1.6 are easy to obtain by adjusting the energetic parameters, however these new parameters move the transition state to lower $Q$ values. In any event, both theory and simulations justify a replica-symmetric treatment of the folding transition for minimally frustrated polymers, particularly regarding the characterization of the barrier. Thermodynamic quantities are self-averaging at $T_F$, with the exception of very nativelike states where the free energy becomes strongly sequence dependent for a finite size polymer.
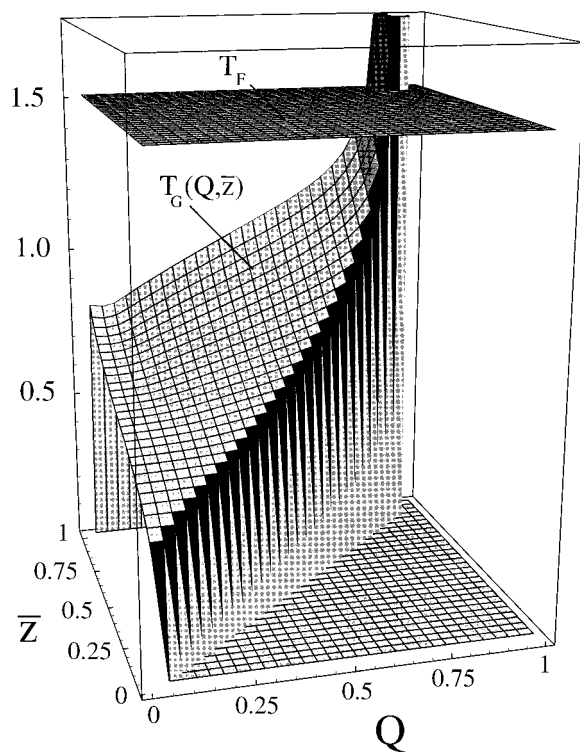


FIG. 3. The folding temperature $T_F$ and glass transition temperature $T_G$ as a function of the fraction of native contacts $Q$ and the total contacts per monomer $\bar{z}$. The folding temperature is above the glass temperature (2.8) for most values of $Q$ and $\bar{z}$, for proteinlike energetic parameters used in fitting the theory to simulations ($\epsilon \cong 1.1$ and $\delta \epsilon_n \cong -2.1$).

The scalar $T_F/T_G(Q, \bar{z})$ is a rather simple indication of self-averaging, and a more rigorous method to determine the degree of self-averaging would be to follow the calculations by Derrida and Toulouse[24] of the moments of the probability distribution of $Y = \Sigma_j W_j^2$, measuring the sample to sample fluctuations of the sum of weights of the free energy valleys, and generalize them to finding the probability distribution of $Y(Q, \bar{z})$.

## IV. THE DENSITY-COUPLED FREE ENERGY

In this section we obtain the free energy in terms of the reaction coordinates $Q$ and $\bar{z}$ through the introduction of the halo density (3.1).

The halo density $\eta_H(\bar{z}/z_N)$ from Eq. (3.1) will appear in the roughness term of Eq. (2.9) since this term arises as a result of non-native interactions which contribute to the total variance of state energies. The entropic term in Eq. (2.9) contains the configurational entropy $s_\eta(Q)$ at $Q$ and density $\eta$. The most accurate values of barrier position and height are obtained by inserting in $s_\eta(Q)$ an interpolated form of density, between homogeneous collapse valid at low $Q$, and the high $Q$ core-halo formula (3.1),

$$\eta_H^{\text{full}}(Q, \bar{z}) = (1-Q) \frac{\bar{z}}{z_N} + Q \eta_H\left(\frac{\bar{z}}{z_N}\right).$$
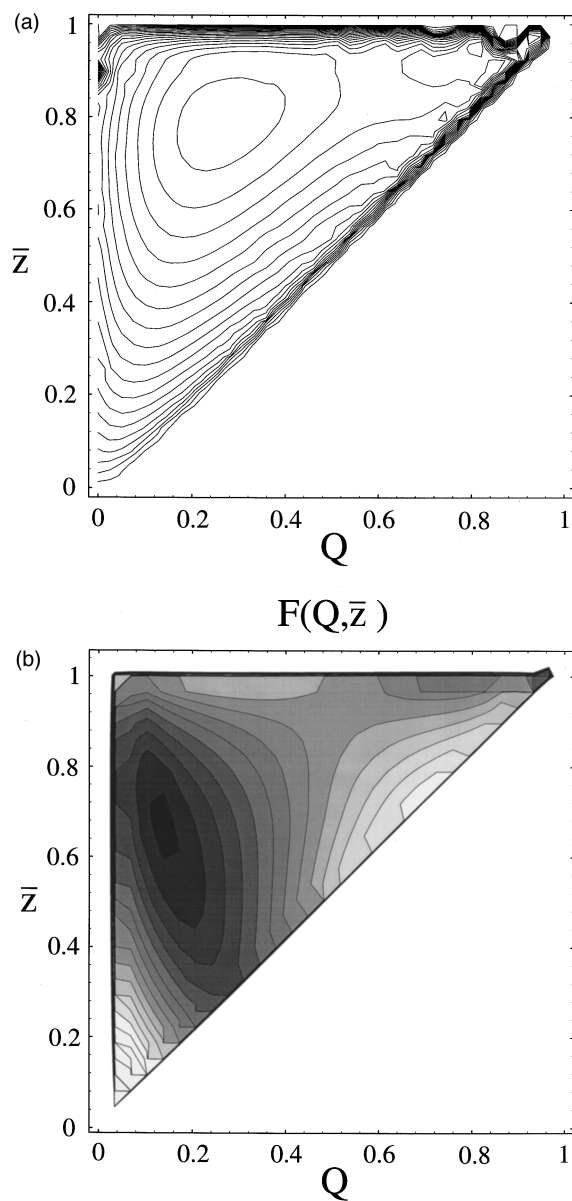
$$F(Q,\bar{z})$$



FIG. 4. (a) The free energy vs $Q$ and $\bar{z}$, at the folding temperature $T_F$, from simulations (Ref. 5). The native minimum is in the upmost right-hand corner. (b) Free energy surface at $T_F$ for the 27-mer, obtained from Eq. (4.1) with the parameters $(\epsilon, \bar{\epsilon}, \delta\epsilon_n, T_F) = (1.1, -1.27, -2.1, 1.51)$. The surface has a double well structure (darker is deeper) with a transition state ensemble at $Q^* \cong 0.50$, and barrier height $\cong 3.0 k_B T$.

This gives more weight to the two behaviors in their respective regimes: mean-field uniform density at weak constraint or low $Q$, and core-halo behavior at strong constraint.

The total density $\eta_{\text{tot}} = \bar{z}/z_N$ appears in the homopolymeric energy since this term is a function only of the number of contacts, irrespective of whether they were native or not. The extra gap energy defined with respect to fully collapsed states in Eq. (2.9) is an energetic contribution added to each native bond formed, independent of $\bar{z}$, up to the limit $Q z_N = \bar{z}$, where the gap term in Eq. (2.9) becomes simply $\bar{z}\delta\epsilon_n$.

These substitutions in Eq. (2.9) describe a free energy

surface as a function of the reaction coordinates $Q$, the fraction of native contacts, and $\bar{z}$, the total contacts per monomer,

$$\frac{F}{N}(T,Q,\bar{z}|E_n) = -\bar{z}|\bar{\epsilon}| - Q z_N|\delta\epsilon_n| - Ts(Q,\bar{z})$$

$$- \frac{z_N \eta_H \epsilon^2}{2T}(1 - Q^2), \qquad (4.1)$$

where $s(Q,\bar{z}) = s(Q, \eta_H^{\text{full}}(Q,\bar{z}))$. The first term is an equilibrium bias toward states that simply have more contacts and depends only on $\bar{z}$, whereas the second term is a bias toward states with greater nativeness and depends only on $Q$, although the maximum value of this bias $\bar{z}\delta\epsilon_n$ does depend on $\bar{z}$. The entropic term biases the free energy minimum toward both small values of $Q$ and $\bar{z}$ where the entropy is largest. The free energy bias due to landscape roughness is largest when there are many non-native contacts ($\bar{z}$ is large and $Q$ is small), which means that the protein can find itself in non-native low energy states due to the randomness of those non-native interactions.

To model the protein behavior at the folding temperature, the temperature $T$ is held fixed at a value $T_F$ described below, and the other energetic parameters ($\bar{\epsilon}$, $\delta\epsilon_n$, and $\epsilon$) are adjusted so as to give the free energy a double well structure with folded and unfolded minima of equal depth.

### A. Comparison with a simulation

The 27-mer lattice model protein has been simulated for polymer sequences designed to show minimal frustration.[4,25,26] The system we are interested in is modeled by a contact Hamiltonian as in Eq. (2.1), but now the beads representing the monomers are of three different kinds with respect to their energies of interaction. If like monomers are in contact, they have an energy $\epsilon_{ij} = -3$, otherwise $\epsilon_{ij} = -1$, where the interaction energy is in arbitrary units of order $k_B T$. Specific sequences are modeled to have a fully collapsed native state with a specific set of 28 contacts and a ground state energy of $-3 \times 28$.

In the thermodynamic limit, the discrete interaction energies used in the simulation give a Gaussian distribution for the total energy of the system by the central limit theorem, whose mean and width naturally depends on the fraction of native contacts.

If we call $\bar{Z}$ the total number of contacts of any kind, the energy at $Q$ and $\bar{Z}$ is determined simply by the energies of these native and non-native contacts above, while the entropy at high temperatures is the log of the number of states satisfying the constraints of $\bar{Z}$ total contacts and $\mu$ native contacts. However, the temperature range where folding occurs is well below the temperature of homopolymeric collapse, and so the polymer can be considered to be largely collapsed. This can be seen either by direct computation or by computing the entropy, defined through

$$S(Q,\overline{z},T) = -\sum_i p_i \log p_i$$

$$= -\sum_i \left(\frac{e^{-E_i/T}}{Z_p}\right) \log\left(\frac{e^{-E_i/T}}{Z_p}\right), \qquad (4.2)$$

where $Z_p$ is the (partial) partition function, the sum being over all of the states consistent with the constraints characterized by $\mu$ and $\overline{Z}$ above.

Onuchic et al.[5] have obtained the free energy as $F = E - TS$ for the 27-mer, which mimics the landscape of a small helical protein, as a surface plot versus the total number of contacts per monomer $\overline{z} = \overline{Z}/N$, and $Q = $ (total number of native bonds)/28, see Fig. 4(a). The largest value of $Q$ for a given $\overline{Z}$ is $\overline{Z}/28$, because there cannot be more native contacts than there are total contacts, hence the allowable region is the upper left-hand side of the surface plot. The surface plot in Fig. 4(a) has a double minimum structure at a specific (folding) temperature $T_F = 1.51$ on the energy scale where $\epsilon_{ij} = \{-3, -1\}$ described above. The free energy barrier of $\approx 2k_B T_F$ is small compared with the entropic barrier of the system ($\sim 14 k_B T_F$). The transition ensemble at reaction coordinates $(Q^*, \overline{z}^*) \cong (0.54, 0.88)$, consists of about $\exp Ns(Q^*, \overline{z}^*) \cong 2000$ thermally occupied states and $\sim 10^5$ configurational states.

There are four energetic parameters in the free energy theory ($\epsilon$, $\overline{\epsilon}$, $\delta\epsilon_n$, and $k_B T_F$), and three parameters in the simulation [$\epsilon$(like units), $\epsilon$(unlike units), and $k_B T_F$], plus the roughness parameter, which is implicitly evaluated through the diversity of energies consistent with overlap $Q$. Minimal frustration in the lattice simulation is implicit in the sequence design, in that the ground state is topologically consistent with all the pair interactions between like monomers. Energetic correlations are implicit from minimal frustration, since states at similarity $Q$ to the native state are also low in energy.

We should note however that the gap energy in these simulations is functionally somewhat different than our theoretical model in that contacts between like monomers are always favored whether native or not, and in the theory only true native contacts have explicit contributions to the energy gap.

We do not undertake here a comparison of simulations at all parameter values with theory. Rather, we compare simulations and theory only for the 27-mer, with parameters chosen to be proteinlike according to the corresponding states principle analysis of Onuchic et al.[5] The scheme for comparison between the simulations and theory for the 27-mer is to hold $T_F$ fixed at the simulational value of 1.5, and then determine the remaining three energetic parameters ($\epsilon, \delta\epsilon_n, \overline{\epsilon}$) by the conditions of folding equilibrium [double-well structure of $F(Q,\overline{z})$], and a barrier position and height consistent with simulations and experiments.

The result of this is shown in Fig. 4(b), which shows the free energy surface at $T_F$ obtained from the parameters $(\epsilon, \overline{\epsilon}, \delta\epsilon_n, T_F) = (1.1, -1.27, -2.1, 1.51)$. These values are very compatible with those of the simulations. The gap to roughness ratio for this minimal model is $|\delta\epsilon_n|/\epsilon \cong 1.9$, satis-

fying the conditions for global foldability. The system has a double-well structure with a weakly first-order transition between a semicollapsed globule, at $(Q_{MG}, \overline{z}_{MG}) \cong (0.14, 0.67)$, and a fully collapsed, near-native folded state at $(Q_F, \overline{z}_F) \cong (0.98, z_N \cong 1.03)$.

In what follows let $N_C = $ the number of native residues in the core, and let $N_H = $ the number of non-native residues in the halo. The molten globule states have a halo density of $\eta_H \cong 0.91$, and there are $N_C \cong N z_N Q^*/z_{NQ} \sim 8$ native residues in the various $\sim N!/N_C! N_H! \approx 2 \times 10^6$ cores consisting of $Q_{MG} N z_N \cong 4$ native contacts. The total density is given by $\eta = N \eta_H / (N_C \eta_H + N_H)$ which is $\cong 0.93$ in the molten globule state. The folded state has a core with $\cong 27$ contacts, containing $\cong 26$ monomers (almost all) and at density $\eta_C = 1$ and a collapsed halo of about 1 monomer at density $\eta_H = 1$. The folded state is fully collapsed. It is energetically favored over the molten globule by about a $k_B T$ per monomer $(E_f - E_{mg} \cong -28.3 k_B T)$ and thus less entropic $[T_F(S_f - S_{mg}) \cong -28.3 k_B T]$.

The core residues in the transition state ensemble at $(Q^*, \overline{z}^*) \cong (0.50, 0.92)$ contain approximately $N z_N Q^*/z_{NQ} \cong 16$ monomers. The 11 remaining monomers in the dangling loops and ends are nearly collapsed, with $\eta_H \cong 0.99$, so the total density is very nearly nearly one. The transition state ensemble in the theory consists of $\approx 1600$ thermally occupied states and $\approx 4 \times 10^5$ configurational states. Its thermal entropy is $\cong 7.4 k_B$. The folding free energy barrier $\Delta F \equiv F(Q^*, \overline{z}^*) - F(Q_{MG}, \overline{z}_{MG})$ is $\cong 3.0 k_B T$. The energetic gain from Eq. (2.6) is $\cong -17.7 k_B T$, and the entropic barrier from Eq. (2.7) is $\cong 20.7 k_B T$. The full free energy barrier ($\cong 3.0 k_B T$) arises from the delicate incomplete cancellation of entropic losses with energetic gains.

One aspect of lattice simulations also present in our theoretical model is an essentially native folded minimum which persists up to high temperatures.[26,27] In simulations this is a lattice effect. The very few near native states available on a lattice create an entropic barrier to escape from the native state. This barrier has led to some confusion in identifying the position of the folding transition state.[27] The native state has a glass temperature $T_G(Q=1)$, which in the GREM formalism is determined by the ratio of energetic gains to entropic losses as the native state is approached. The discrete nature of lattice models gives the native state an effective high glass temperature above any simulation temperature. In our analytical model, the collective nature of melting (i.e., $l_C$ and $l_{EC}$ are $>0$, see Appendix A) leads to a similar gap in the density of states along the $Q$ coordinate. This also causes a weak barrier ($\cong 0.4 k_B T$) between a near-folded local minimum and the native folded minimum [see Fig. (4b)]. Again this weak barrier is a result of discrete quantities in the simulations and theoretical model.

More elaborate theories should incorporate a local nativeness parameter $Q(\mathbf{x})$ which varies in space, allowing for rigid as well as fluctuating regions of the protein.

## B. Explicit three-body effects

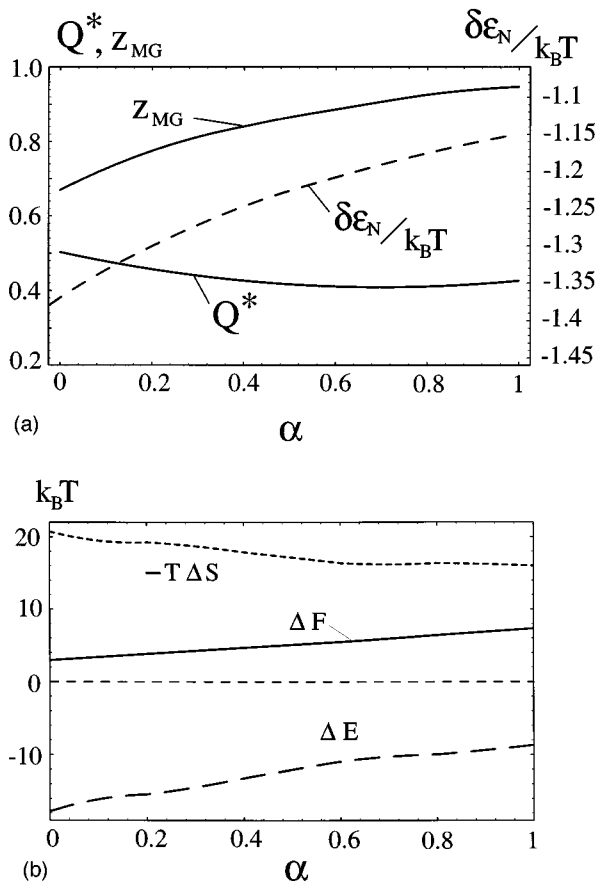It is interesting to investigate the effects of explicit many-body cooperativity on the folding funnel by introduc-

FIG. 5. (a) (*Left axis*) Positions of the barrier $Q^*$ and total bonds per monomer in the molten globule ($\bar{z}_{MG}$), as a function of the three-body coefficient $\alpha$. (*Right axis*) Decrease with $\alpha$ in the necessary energy gap to maintain equilibrium at $T_F$. Plots are for the 27-mer with fitted energy parameters (Sec. IV A). (b) Free energy barrier $\Delta F = F(Q^*,\bar{z}^*) - F(Q_{MG},\bar{z}_{MG})$ in units of $k_B T_F$, and its energetic and entropic contributions, for the 27-mer, as a function of $\alpha$. The barrier grows moderately with $\alpha$, i.e., the energetic drop decreases faster than entropic losses, due to the collective nature of the energetic interactions.

ing a three-body interaction in addition to the pair interactions already present. Models with such partially explicit cooperativity mimic the idea that only formed secondary structure units can couple, and have been introduced in lattice models by Kolinski *et al.*[4] Three-body interactions enter into the energetic contributions of Eq. (4.1) as an additional $Q^2$ term in the bias and roughness, and $\bar{z}^2$ term in the homopolymer attraction, so that those terms in the free energy become

$$-[(1-\alpha)\bar{z}+\alpha\bar{z}^2]|\bar{\epsilon}| - [(1-\alpha)Q+\alpha Q^2]z_N|\delta\epsilon_n|$$

$$-\frac{z_N\eta_H(\bar{z}/z_N)\epsilon^2}{2T}\{1-[(1-\alpha)Q+\alpha Q^2]^2\}, \qquad (4.3)$$

where $\alpha$ is a measure of the amount of three-body force present. In the model defined by Eqs. (4.1) and (4.3), proteins with more three-body forces need not be as strongly optimized, and so the magnitude of the gap is a decreasing function of $\alpha$ at fixed $T_F$, $\bar{\epsilon}$, and $\epsilon$ [see Fig. 5(a)]. With this correction included, we find the barrier position $Q^*(\alpha)$ to be

a weakly decreasing function of $\alpha$ [Fig. 5(a)] (although as described above, $Q^*$ is not independent of $m$, the order of the $m$-body interactions). We eventually expect this trend to reverse for larger $m$ as in Fig. 1(b). The position of the folded state $Q_F$ remains near native and $Q_{MG}$ is also roughly constant. However as homopolymer attraction becomes more collective (increasing $\alpha$), $\bar{z}_{MG}$ increases and the molten globule state becomes denser [Fig. 5(a)]. The transition becomes more first-order-like with increasing $\alpha$, as the trend in energetic loss is not as great as entropic loss, so that the free energy barrier increases with $\alpha$ [Fig. 5(b)].

## C. Dependence of the barrier on sequence length

It is simple in our theory to vary the polymer sequence length. One recalculates $s(Q)$ at constant density for a larger chain[12] and inserts this, along with the density $\eta_H$ [Eq. (3.1)] at the larger value of $N$, into the free energy (4.1). To model the barrier at $T_F$, one must scale the temperature with $N$ since in our model larger proteins fold at higher temperatures [see Fig. 6(a)], e.g., in Eq. (2.11) $T_F \propto z_N$. The barrier position $Q^*$ mildly decreases with $N$ [Fig. 6(b)].[28] Plotted along with the theoretical curve are two experimental measurements of the barrier position. The square represents the measurement for truncated $\lambda$ repressor,[17] a ~80 residue protein fragment with largely helical structure. The corresponding states analysis[5] shows that the formation of helical secondary structure within the $\lambda$ repressor makes it entropically similar to the lattice 27-mer. Also plotted in Fig. 6 (circle) is the experimental barrier measurement for Cytochrome C,[29] a 104 residue helical protein which is entropically similar to the 64-mer lattice model. The simple proposed model has the same decreasing trend in the position of the transition state ensemble that is observed experimentally, but decreases much slower. This suggests that a local nucleation description may become more appropriate as $N$ increases, rather than the homogeneous mean field theory proposed here. However the question as to whether energetic heterogeneity induces a specific nucleus[30] rather than an ensemble of nuclei with correspondingly many kinetic paths, is an open issue. In Appendix C, we show for thoroughness that experimental plots of folding rate versus equilibrium constant used in the experiments above are indeed a measure of the position of the transition state ensemble.

Figure 6(c) shows the roughly linear trend of the barrier height with $N$. This mean-field result applies for small $N$; as $N$ grows, fluctuation mechanisms begin to dominate the scaling behavior, and may reproduce the sublinear scaling with $N$ seen in lattice studies[31] and by scaling arguments.[32]

The overall folding time results from a combination of thermodynamic barrier crossing dealt with here, and kinetic diffusion between locally stable basins.[26,33] Experimental measurements from the folding rate at $T_F$ can thus lead to an estimate for the reconfiguration time at $Q^*$. For example, rearranging Eq. (C2), the folding rate over the reconfiguration rate is given by

$$\ln(k_F \overline{t(Q^*)}) = -\frac{F(Q^*)-F_u}{k_B T}. \qquad (4.4)$$
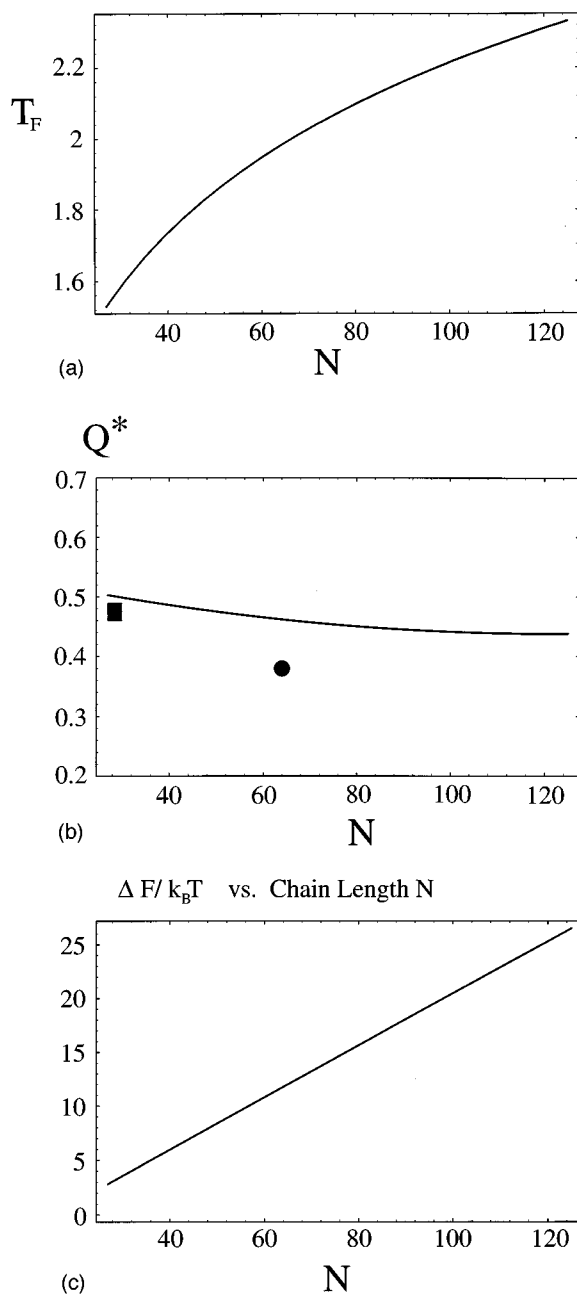
FIG. 6. (a) The folding temperature $T_F$ is an increasing function of polymer sequence length $N$. (b) Position of the barrier $Q^*$ as a function of sequence length $N$. The solid line is the theory as determined by Eq. (4.1), and the points marked are experimental results (see the text). (c) Free energy barrier height $\Delta F$ in units of $k_B T_F$, as a function of sequence length $N$.



FIG. 7. (a) Two plots of the free energy vs $Q$ for the 27-mer with the fitted parameters [Fig. 4(b)]. The upper curve is the free energy with the fitted stability gap $\delta\epsilon_n = -2.1$, and $\delta\epsilon_n = -2.5$ in the lower curve. These one-dimensional plots are the most folding probable paths (minimum free energy along coordinate $Q$ determined by $dF/d\bar{z}=0$) on the 2D surface plot of Fig. 4. From the figure we can see that as $|\delta\epsilon_n|$ increases the folding becomes downhill. Folding becomes purely downhill with no barrier at $|\delta\epsilon_n| \cong 2.75$. Because of the stability of the molten globule position, the barrier shifts slightly to lower $Q^*$ (from 0.51 to 0.47), and decreases in height (from about $3k_B T$ to $0.6k_B T$). (b) Position of the barrier $Q^*$ as a function of magnitude of the energy gap $|\delta\epsilon_n|$ (in units of $k_B T$), for the fitted 27-mer described in (a). $Q^*$ weakly decreases until $|\delta\epsilon_n|/k_B T \cong 1.75$, and then rapidly merges with the position $Q_{MG}$ of the molten globule at $|\delta\epsilon_n|/k_B T \cong 1.82$ as the barrier vanishes. (c) Free energy barrier in units of $k_B T$ vs magnitude of stability gap $|\delta\epsilon_n|$. The short dashed line is the entropic contribution to the barrier, and the long dashed line is minus the energetic contribution. These two terms merge and become zero at $|\delta\epsilon_n|/k_B T \cong 1.82$ where the barrier vanishes.

For the $\lambda$ repressor[17] at the folding midpoint, the folding rate $k_F$ is about 400 s$^{-1}$. Using the barrier height $F(Q^*) - F_u \cong 3k_B T$ for the corresponding 27-mer gives a re-configuration time $t(Q^*) \approx 10^{-4}$ s. Since the Rouse–Zimm time is typically in the $\mu$s range, this suggests that configurational diffusion in the transition region is typically activated. Of course there are many issues involved in the local dynamics and structure of folding proteins, which make precise comparisons with specific examples difficult. The num-
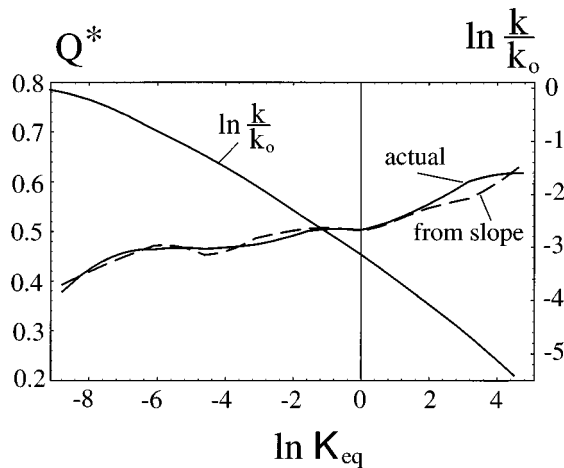
FIG. 8. (*Right axis*) Plot of the logarithm of the folding rate vs the logarithm of the unfolding equilibrium constant. (*Left axis*) Reading the slope of the folding rate gives a measure of the position of the barrier $Q^*$. Also plotted is the actual value of $Q^*$ directly calculated from the free energy curves. The values compare well for most values of the gap where the free energy has a double-well structure.



FIG. 9. (*Solid line*) Probability to be in the unfolded state vs temperature (in simulation units) for the two-order parameter model used in fitting the 27-mer simulations. With this model $T_F = 1.51$ and $\Delta T / T_F \cong 0.14$. (*Dashed line*) Same probability for the one-order parameter model using Eq. (4.5). Here $T_F = 2.14$ and $\Delta T / T_F \cong 0.22$.

bers we quote should be interpreted as estimates showing the reasonableness of the current parametrizations.

## D. Dependence of the barrier on the stability gap, at fixed temperature and roughness

As the stability gap is increased at fixed temperature, folding approaches a downhill process, with a folded global equilibrium state [see Fig. 7(a)]. We can see from Fig. 7(a) that the barrier position and height are decreasing functions of stability gap, with true downhill folding (zero barrier) occurring when $|\delta\epsilon_n|/\epsilon \cong 2.5$ or $|\delta\epsilon_n|/T_F s_0 \cong 1.4$ for the 27-mer [see Figs. 7(b) and 7(c)]. At folding equilibrium, $|\delta\epsilon_n|/T_F s_0 \cong 1.1$. Thus, achieving downhill folding requires a considerable change of stability—an estimate for a 60-aa protein (27-mer lattice model) would be an excess stability of $\approx 12 k_B T_F$.

We can apply the equations of Appendix C to changes of the transition state free energy by modifying stability. Figure 8 shows a plot of the log of a normalized folding rate $\ln(k_F/k_0)$ vs the log of the unfolding equilibrium constant $\ln K_{eq}$, whose slope is a measure of the barrier position $Q^*$. The increasing magnitude of slope with increasing $\ln K_{eq}$ means that the barrier position is shifting toward the native state as the gap decreases. Also shown is a comparison between the position of the barrier $Q^*$ calculated directly from the theory, and $Q^*$ as derived from the slope of $\ln(k_F/k_0)$ using Eq. (C5). The linear free energy relation works well for the range of parameters having a double-well free energy surface.

## E. Denaturation with increasing temperature

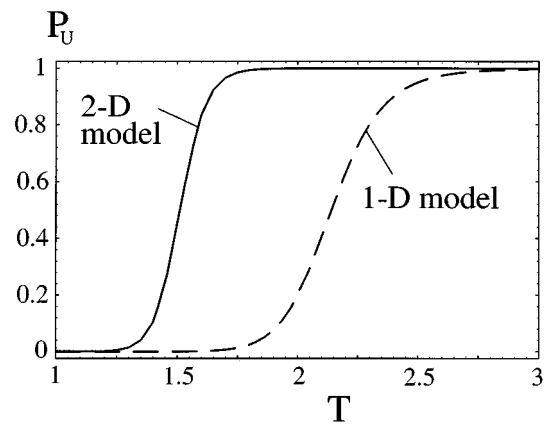The probability $P_u$ for the protein to be in the unfolded globule state at temperature $T$ is

$$P_u = \left[ 1 + \exp\left[ -\frac{1}{T}(F_f - F_u) \right] \right]^{-1},$$

where $F_u$ and $F_f$ are the free energies at temperature $T$ of the unfolded and folded minima (at $T_F$, $F_f = F_u$ and $P_u = 1/2$). This can be used to obtain denaturation curves as a function of temperature. For illustration, we make the simplifying assumptions that both the folded and globule states are collapsed, making $P_u$ independent of $\bar{\epsilon}$, and that the folded and globule states occur approximately at $Q_F = 1$ and $Q_{MG} = 0$. As the temperature is lowered, the molten globule freezes into a low energy configuration at $T_g = \epsilon\sqrt{z_N(1 - Q_{mg}^2)/(2s(Q_{mg}))} \cong \epsilon\sqrt{z_N/(2s_0)}$ (see Fig. 3), and the expression for $P_u$ becomes one of equilibrium between two temperature independent states with the corresponding ''Shottky'' form of the energy and specific heat:

$$P_u = \left[ 1 + \exp(-Ns_0)\exp Nz_N\left( \frac{|\delta\epsilon_n|}{T} - \frac{\epsilon^2}{2T^2} \right) \right]^{-1}, \quad T_g < T$$

$$= \left[ 1 + \exp\frac{Nz_N}{T}\left( |\delta\epsilon_n| - \epsilon\sqrt{\frac{2s_0}{z_N}} \right) \right]^{-1}, \quad T < T_g. \quad (4.5)$$

The condition that $T_F/T_G \geq 1$ gives Eq. (2.12). Using the glass temperature of the globule state, this is equivalent to

$$T_g \geq \frac{\epsilon^2}{\delta\epsilon_n},$$

which is the temperature where the high $T$ expression for $P_u$ in 4.5 has a minimum. Hence cold denaturation will not be seen in the constant density model (as it would if there were no glass transition), and $P_u$ will always decrease to zero at low temperatures.

In the limit of large $T$, Eq. (4.5) becomes $\approx 1/(1 + \exp -Ns_0) \approx 1$, indicating denaturation. At small $T$ Eq. (4.5) tends to zero as $\exp -(\text{const.} \times N/T)$.

Allowing density to vary modifies the denaturation behavior. In Fig. 9, denaturation curves are plotted for the vari-

able density model with fitted parameters $[(\epsilon,\bar{\epsilon},\delta\epsilon_n,T_F) =(1.1, -1.27, -2.1, 1.51)]$, and for the one-dimensional (1D) case for a completely collapsed ($\eta=1$) protein. In the 1D case, the same energetic parameters $[(\delta\epsilon_n,\epsilon)=(-2.1, 1.1)]$ are used in Eq. (4.5) but the molten globule entropy is the fully collapsed value ($S_{mg} \cong 27 \times 0.88$). $T_F$ in the 1D case increases to $\cong 2.14$, in units where $T_F=1.51$ as in the simulation.

If we define the width $\Delta T$ of the transition between 10% and 90% denatured, the ratios of widths to folding temperatures $\Delta T/T_F$ are about 0.22 and 0.14 for the 1D and 2D cases, respectively. Allowing density to vary sharpens the transition. The values lie between those obtained in lattice studies of the 27-mer,[25,34] where $\Delta T/T_F \approx 0.3$ for foldable sequences, and from measurements of the thermal denaturation of $\lambda$ repressor,[35] where $\Delta T/T_F \approx 0.05$. This suggests that many-body forces may play a role in the stabilization of the native state for laboratory proteins.

## V. CONCLUSION

In this paper we have shown that if the energy of a given configuration of a random heteropolymer is known to be lower than expected for the ground state of a completely random sequence (i.e., the protein is minimally frustrated), then correlations in the energies of similar configurations lead to a funneled landscape topography. The interplay of entropic loss and energetic loss as the system approaches the native state results in a free energy surface with weakly two-state behavior between a dense globule of large entropy, and a rigid folded state with nearly all native contacts. The weak first-order transition is characterized by a free energy barrier which functions as a "bottleneck" in the folding process.

The barrier is small compared with the total thermal energy of the system—on the order of a few $k_B T$ for smaller proteins of sequence length about 60 amino acids. For these small proteins the model predicts a position of the barrier $Q^*$ about halfway between the unfolded and native states ($Q^* \cong 1/2$). For larger proteins, the barrier height rises linearly, and its position $Q^*$ moves away from the native state toward the molten globule ensemble roughly as $1/z_N$, due essentially to the fact that the entropy decrease per contact is independent of $N$ initially. Experimental measurements of the barrier for fast folding proteins are consistent with this predicted shift in position with increasing sequence length, but with a shift somewhat greater in magnitude.

The unfolded and transition states are not single configurations but ensembles of many configurations. The transition state ensemble according to the theory consists of about 1600 thermally accessible states for a small protein such as the $\lambda$ repressor. There are about $\approx 4 \times 10^5$ configurational states in the transition state ensemble, less than the $\sim 10^7$ ways of choosing 16 core residues in the 27-mer minimal model, so that many but not all nuclei are sampled. The multitude of states at $Q^*$ seen in the theory is in harmony with a picture of a transition state ensemble of generally delocalized nuclei, a subject investigated recently by various authors.[36-38]

A simple theory of collapse was introduced to couple protein density with nativeness. This resulted in a density contraction in the process of folding. During folding, a dense inner native core forms, which grows while possibly interchanging some native contacts with others upon completion of folding. This core is surrounded by a halo of non-native polymer which shrinks and condenses in the folding process, as topological constraints upon folding make dangling loops shorter and denser. The folded free energy minimum is essentially native in the model when parameters are chosen to fit simulated free energy curves for the 27-mer.

Explicitly cooperative interactions were shown to enhance the first-order nature of the transition through an increase in the size of the barrier, and a shift toward more nativelike transition state ensembles (i.e., at higher $Q^*$). For the constant density scenario the barrier becomes almost entirely entropic when the order $m$ of the $m$-body interactions becomes large, and the transition state ensemble becomes correspondingly more nativelike. In the energy landscape picture, as explicit cooperativity increases, the protein folding funnel disappears, and the landscape tends toward a golf-course topography with energetic correlations less effective and more short range in $Q$ space. The correlation of stability gaps and $T_F/T_G$ ratios with kinetic foldability is true only for fixed $m$ much less than $N$.

A full treatment of the barrier as a function of the three energetic parameters ($\epsilon,\bar{\epsilon},\delta\epsilon_n$) plus temperature $T$ would involve the analysis of a multidimensional surface defining folding equilibrium in the space of these parameters. We shall return to this issue in the future, but we have deferred it for now in favor of the simpler analysis of seeking trends in the position and height of the barrier as a function of individual parameters such as $\delta\epsilon_n$ and $T$.

## APPENDIX A

Here we summarize the derivation of the configurational entropy $S(Q,\eta)$, as a function of the topological constraints and density $\eta$. For a complete derivation see Ref. 12.

The entropy of an unconstrained polymer is given by

$$S_0(\eta)=N\left[\ln\frac{\nu}{e}-\left(\frac{1-\eta}{\eta}\right)\ln(1-\eta)\right], \tag{A1}$$

where $\nu=$ the number of configurations per monomer (six for three-dimensional cubic lattice models).

For a weakly constrained polymer (low values of $Q$), the entropy loss from the unconstrained state can be decomposed into three terms;

$$S_{low}(Q,\eta)=S_0(\eta)+\Delta S_B(Q,\eta)+\Delta S_{mix}(Q,\eta)$$
$$+\Delta S_{AB}(Q,\eta). \tag{A2}$$

The first term $\Delta S_B(Q,\eta)$ is the reduction in searchable phase space due to $\mu = QN z_N \eta$ cross-links in the polymer chain, first derived by Flory.[16] For polymers in three dimensions this is

$$\Delta S_B(Q,\eta) = \tfrac{3}{2} N Q z_N \eta \, \ln(C Q z_N \eta), \tag{A3}$$

where $z_N$, the coordination number for a chain of length $N$, is given by

$$z_N \approx \frac{1}{N} \, \mathrm{Int}[2N - 3(N+1)^{2/3} + 3] \tag{A4}$$

($\mathrm{Int}[...]$ means take the integer part), and

$$C = \frac{3}{4\pi e} \left( \frac{\Delta\tau}{b^3} \right)^{2/3}, \tag{A5}$$

where $(\Delta\tau/b^3)^{1/3}$ is the ratio of the bond radius to the persistence length.

There are many ways $QN z_N \eta$ cross-links can be formed from the $N z_N \eta$ total contacts in each molten globule state, at least for weakly constrained polymers. This results in a reduction of the entropy lost due to a combinatoric or mixing entropy given by

$$S_{\mathrm{mix}}(Q,\eta) = -N z_N \eta [Q \ln Q + (1-Q)\ln(1-Q)]. \tag{A6}$$

There is an additional reduction in the searchable phase space because of the $(N z_N \eta - QN z_N \eta)$ native contacts that *cannot* be formed because the overlap can be no larger than $Q$:

$$\Delta S_{AB}(Q,\eta) = \frac{N}{C} \int_{CQ_{z_N\eta}}^{Cz_N\eta} dx \, \ln(1 - x^{3/2}) = -\frac{N}{2C} \left[ 3 C z_N \eta (1-Q) - 2 C z_N \eta (\ln[1 - (C z_N \eta)^{3/2}] - Q \, \ln[1 - (Q C z_N \eta)^{3/2}]) \right.$$

$$+ \ln\left[ \left( \frac{1 - \sqrt{C z_N \eta}}{1 - \sqrt{Q C z_N \eta}} \right)^2 \left( \frac{1 + Q C z_N \eta + \sqrt{Q C z_N \eta}}{1 + C z_N \eta + \sqrt{C z_N \eta}} \right) \right] + 2\sqrt{3} \left( \arctan \frac{1}{\sqrt{3}} [1 + 2\sqrt{C z_N \eta}] \right.$$

$$\left. \left. - \arctan \frac{1}{\sqrt{3}} [1 + 2\sqrt{Q C z_N \eta}] \right) \right], \tag{A7}$$

where $C$ is given in expression (A5).

In Ref. 12, surface effects were also considered as reducing the conformational search when the rms loop size was comparable to the size of the globule. This somewhat modifies the previous formulas for small values of $Q$.

When $QN z_N \eta \cong N$, there is about 1 cross-link per monomer, $S_{\mathrm{low}}(Q,\eta) \cong 0$, and the low $Q$ formula is no longer valid. At some point before this, configurations having fluctuations in the mean field contact pattern begin to dominate the free energy. It then becomes more accurate to switch the description of entropy loss from that due to a dilute ''gas'' of contacts, to an atomistic description ascribing entropy to lengths of chain melted out from the frozen $(Q=1)$ three-dimensional native structure, and the combinatorics of these pieces of melted chain.

In what follows, let $Nf$ be the total number of melted or unconstrained pieces in the chain of any length, let $l_E$ be the sequence length of the melted ends of the chain, and $l_c$ and $l_{EC}$ be the minimum or critical lengths of the melted pieces or ends, respectively. We can characterize a state by the number distribution of melted and frozen pieces of length $l$, $\{n_l\}$, and $\{m_l\}$, respectively, and the probability distribution for an end of the chain to have length $l$, $\{p_l\}$. It follows that

$$\sum_{l_c}^{N} l n_l = N_M, \tag{A8}$$

$$\sum_{1}^{N} l m_l = N_F, \tag{A9}$$

$$\sum_{l_{EC}}^{N} l p_l = l_E, \tag{A10}$$

where $N_M$ and $N_F$ are the numbers of melted (free) and frozen (constrained) monomers, respectively. It follows that $\Sigma l n_l + \Sigma l m_l + 2\Sigma l p_l = N$. Furthermore, if there are $N_F$ frozen monomers, there are $z_N \eta N_F = QN z_N \eta$ frozen bonds, so, e.g., $\Sigma l n_l = N(1-Q) - 2l_E$. Since $\Sigma n_l = \Sigma m_l = Nf$, the average melted loop length at $Q$ is given by

$$\langle l_{\mathrm{melted}}(Q) \rangle = \frac{\Sigma_{l_c}^{N} l n_l}{\Sigma_{l_c}^{N} n_l} = \frac{1-Q}{f(Q)} - \frac{2 l_E(Q)}{N f(Q)}. \tag{A11}$$

This expression will be useful in modeling the polymer density as a function of collapse.

In terms of the macroscopic parameters $Q$, $f$, and $l_E$, the entropy in the strongly constrained regime is given by[12]

$$\frac{S(Q,f,l_E)}{N} = \left[\ln\frac{\nu}{e} - \left(\frac{1-\eta}{\eta}\right)\ln(1-\eta)\right]\left[1 - Q - f(l_c-1) - \frac{2(l_{EC}-1)}{N}\right] + Q\ln Q - (Q-f)\ln(Q-f) - 2f\ln f$$

$$+ \left(1 - Q - f(l_c-1) - 2\frac{l_E}{N}\right)\ln\left(1 - Q - f(l_c-1) - 2\frac{l_E}{N}\right) - \left(1 - Q - fl_c - 2\frac{l_E}{N}\right)\ln\left(1 - Q - fl_c - 2\frac{l_E}{N}\right)$$

$$+ 2\left(\frac{l_E - (l_{EC}-1)}{N}\right)\ln\left(\frac{l_E - (l_{EC}-1)}{N}\right) - 2\left(\frac{l_E - l_{EC}}{N}\right)\ln\left(\frac{l_E - l_{EC}}{N}\right) - 2\frac{\ln N}{N}. \tag{A12}$$

Maximizing $S(Q,f,l_E)$ with respect to $f$ and $l_E$ gives

$$f(Q,l_E) = \frac{1 - Q - 2l_E/N}{l_c + l_E - l_{EC}} \tag{A13}$$

and

$$(l_E - l_{EC})^{l_c}\left[Q(l_c + l_E - l_{EC}) + Q - 1 + 2\frac{l_E}{N}\right]$$

$$- \left(1 - Q - 2\frac{l_E}{N}\right)(l_E - l_{EC} + 1)^{l_c-1}\mu^{l_c-1} = 0. \tag{A14}$$

The solution to Eq. (A14) is numeric for $l_c \geqslant 3$. These equations put into expression (A12) gives the entropy $S_{high}(Q,\eta)$ in the strongly constrained (high $Q$) regime.

The total entropy may be numerically approximated by an interpolation between the low $Q$ and high $Q$ expressions,

$$S_{tot}(Q,\eta) = (1-Q)S_{low}(Q,\eta) + QS_{high}(Q,\eta), \tag{A15}$$

which is plotted in Fig. 10 for a 27-mer with $l_c=3$, and $l_{EC}=1.5$. Equation (A15) together with Eq. (3.1) gives $S(Q,\bar{z})$ used in Eq. (4.1).
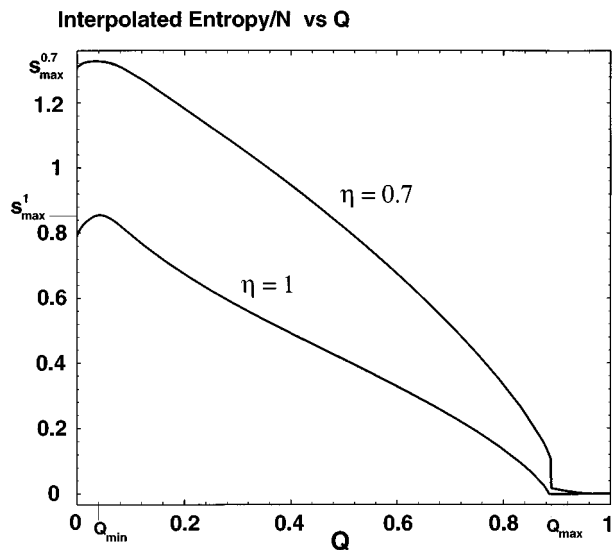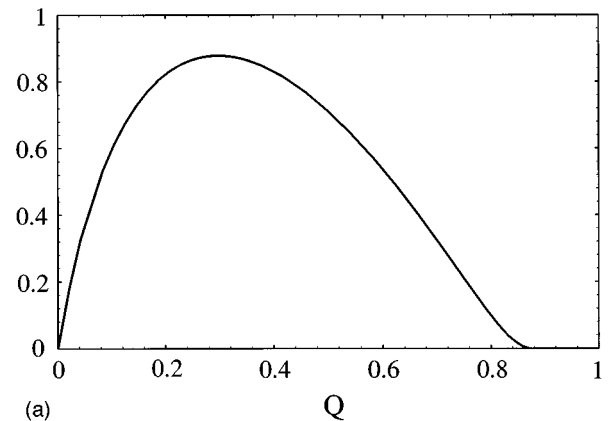
Equations (A13) and (A14) determine $f(Q)$ and $l_E(Q)$, from which we can obtain $\langle l_{melted}(Q)\rangle$ from Eq. (A11). These quantities, plotted in Fig. 11 for parameters appropriate for a 27-mer, will be used in calculating the density as a function of total contacts. In the high $Q$ analysis the total polymer length in the halo is a linearly decreasing function of $Q$.
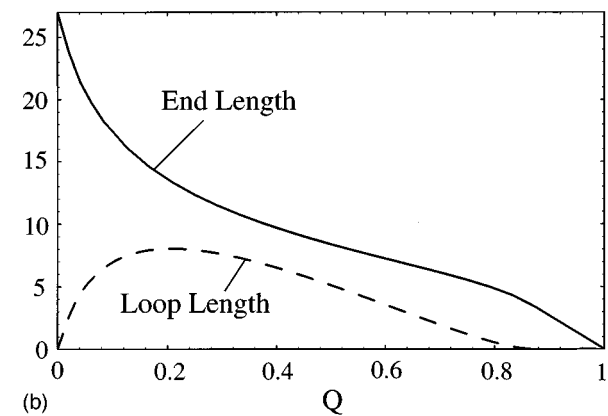
## APPENDIX B

Here we derive the form of the thermodynamic functions in the ultrametric approximation, valid for high $Q$. We wish to find the free energy relative to the state $n$ with energy $E_n$,



FIG. 10. Interpolated entropy from Eq. (A15) for a 27-mer with $l_C=3$, and $l_{EC}=1.5$, for two densities: $\eta=1$ and $\eta=0.7$. The maximum value of $S(Q)$ at $Q_{min}\cong 0.04$ is the entropy at the statistically most likely overlap between any two states for the 27-mer. The value of $Q=Q_{max}<1$ where the entropy vanishes is due to the finite chain length that must be collectively melted out from the frozen structure.

FIG. 11. (a) Number of melted pieces (free loops) $Nf(Q)$ vs $Q$ for a 27-mer with $l_c=4.5$, $l_{EC}=1.5$. (b) Average free end sequence length in the polymer $2l_E(Q)$ (solid line), and mean melted loop length in the polymer at $Q$, $Nf(Q)\langle l_{melted}(Q)\rangle$ (dashed line) for the same 27-mer. These parameters are used in deriving the density function $\eta(Q=\bar{z}/z_N)$ [Fig. 2(b)].
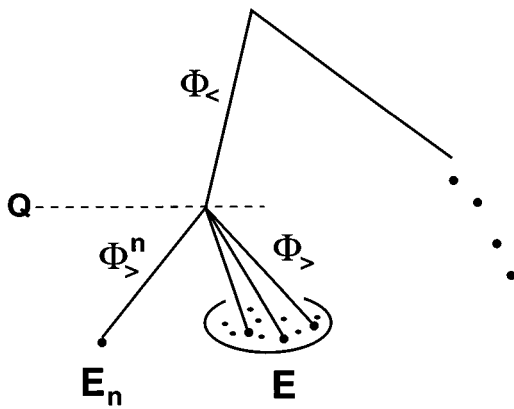
FIG. 12. Diagram of the hierarchy used in a microscopically ultrametric model.
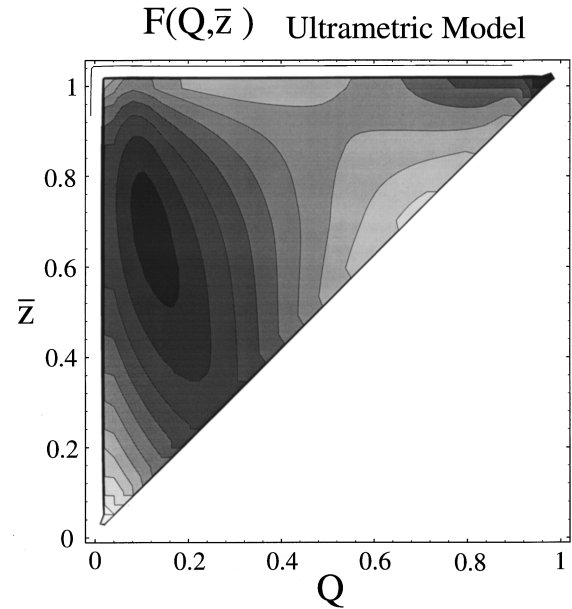


$F(Q,\bar{z})$  Ultrametric Model

FIG. 13. Free energy surface at $T_F$ for an ultrametric 27-mer, obtained from Eq. (B5) with the parameters $(\epsilon,\bar{\epsilon},\delta\epsilon_n,T_F)=(1.1,-1.17,-2.1,1.54)$. The surface has a double-well structure with a transition state ensemble at $Q^*\cong0.52$.

as a function of $Q$, if all the states are organized ultrametrically as in Fig. 12. Call the energetic contribution common to states with overlap $Q$ $\phi_<$, and call the different contributions $\phi_>^n$ for state $n$, and $\phi_>$ for the group (stratum) of states at $Q$ ($\phi_>$ varies from state to state). As in Fig. 12,

$$E_n=\phi_<+\phi_>^n, \quad E=\phi_<+\phi_>.$$

Without the constraint $\phi_<+\phi_>^n=E_n$, the distributions of $\phi_<$ and $\phi_>$ are simple Gaussians of widths $Q\Delta E^2$ and $(1-Q)\Delta E^2$, where $\Delta E^2=Nz_N\eta\epsilon^2$. These are just the number of bonds contributing to $\phi_<$ and $\phi_>$ times the individual variance of a bond $\epsilon^2$. Native contacts are energetically heterogeneous and also contribute to the total width. We will let $\bar{\epsilon}=0$ for simplicity here since it has no bearing on the calculation.

For the state $n$ of energy $E_n$ however, $\phi_<$ is chosen from the conditional probability distribution $P_Q(\phi_<|E_n)$, where

$$P_Q(\phi_<|E_n)=\frac{P(\phi_<,E_n)}{P(E_n)}$$

$$\approx\frac{\exp\left(-\dfrac{\phi_<^2}{2\Delta E^2Q}\right)\exp\left(-\dfrac{(E_n-\phi_<)^2}{2\Delta E^2(1-Q)}\right)}{\exp\left(-\dfrac{E_n^2}{2\Delta E^2}\right)}$$

$$\approx\exp\left(-\frac{(\phi_<-QE_n)^2}{2\Delta E^2Q(1-Q)}\right). \quad (B1)$$

[which approaches $\delta(\phi_<)$ and $\delta(\phi_<-E_n)$ when $Q\to0$ and $Q\to1$ respectively]. Supposing we have found a state $n$ with $E_n$ and a given contribution $\phi_<$, the number of states $N(E,q,E_n)$ having overlap $Q$ with it and energy $E$ is given by $\exp[S_\eta(Q)]P(\phi_>)$ or

$$\ln N(E,Q,E_n)=S_\eta(Q)-\frac{(E-\phi_<)^2}{2\Delta E^2(1-Q)} \quad (B2)$$

with $\phi_<$ chosen from the distribution (B.1). Using $1/T=\partial\ln N/\partial E$, we obtain the free energy relative to the state $n$ as a function of $Q$,

$$F(Q,\phi_<)=\phi_<-TS_\eta(Q)-\frac{\Delta E^2}{T}(1-Q). \quad (B3)$$

Since $\phi_<$ is chosen from the distribution (B1), the free energy function $F$ in Eq. (B3) is chosen from the distribution

$$P_Q(F)\approx\exp-\frac{1}{2\Delta E^2Q(1-Q)}\left(F-QE_n+TS_\eta(Q)\right.$$

$$\left.+\frac{\Delta E^2}{T}(1-Q)\right)^2, \quad (B4)$$

with the mean and most probable free energy function $F^*(Q,T,E_n)$ given by

$$F^*(Q,T,E_n)=QE_n-TS_\eta(Q)-\frac{\Delta E^2}{2T}(1-Q) \quad (B5)$$

for temperatures above the glass temperature $T_G$ (the glass temperatures in this model have been investigated elsewhere[13,33]).

So the modification of the free energy from Eq. (2.9) is simply to replace $(1-Q^2)$ by $(1-Q)$ in the roughness term. The entropy and energy follow straightforwardly from the derivatives of $F^*(Q,T,E_n)$.

Including homopolymer attraction in Eq. (B5) yields a free energy surface in coordinates $(Q,\bar{z})$ as before (see Fig. 13). For the same energy parameters as in Sec. IV A, the surface is very similar to the nonultrametric case. The folding temperature is slightly higher $[(\epsilon,\bar{\epsilon},\delta\epsilon_n,T_F)=(1.1,-1.17,-2.0,1.54)]$ presumably because the ultrametric landscape is smoother. Equivalently, ultrametric polymers need not be as strongly optimized. The barrier is slightly

larger and more native [$\Delta F \cong 4.4 k_B T$ and $(Q^*, \bar{z}^*) = (0.52, 0.93)$], and its height scales roughly linearly with $N$ as before, up to $\cong 32.0 k_B T$ for a 125-mer.

## APPENDIX C

For small, $Q$ dependent changes in the free energy (4.1), e.g., changes in temperature, we can approximate the change in the free energy at the position of the barrier $\delta F(Q^*)$ as a linear interpolation between the free energy changes of the unfolded and folded minima $\delta F_u$ and $\delta F_f$, here estimated to be at $Q_u \cong 0$ and $Q_f \cong 1$, respectively,

$$\delta F(Q^*) \cong Q^* \delta F_f + (1 - Q^*) \delta F_u . \qquad (C1)$$

Furthermore let us approximate the kinetic folding time by the thermodynamic folding time[38]

$$\tau = \overline{t(Q^*)} e^{(F(Q^*) - F_u)/k_B T} \qquad (C2)$$

where $Q^*$ is the position of the barrier and $\overline{t(Q^*)}$ is the lifetime of the microstates of the transition ensemble. Then the log of the folding rate $k_F$ is $\propto F_u - F(Q^*)$. The equilibrium constant for the unfolding transition $K_{eq}$ is the probability to be in the unfolded minimum over the probability to be in the folded minimum, and so $\ln K_{eq} \propto F_f - F_u$. If we plot $\ln k_F$ vs $\ln K_{eq}$, the assumption of a linear free energy relation (C1) and a stable barrier position results in a linear dependence of rate upon equilibrium constant, with slope

$$\frac{\delta \ln(k_F/k_0)}{\delta \ln K_{eq}} \cong \frac{\delta[F_u - F(Q^*)]}{\delta[F_f - F_u]} = \frac{\delta F_u - \delta F(Q^*)}{\delta F_f - \delta F_u} = -Q^* \qquad (C3)$$

so that experimental slopes of folding rates versus unfolding equilibrium constants are indeed a measure of the position of the barrier in our theory.

If the unfolded and folded states are not assumed to be at $Q = 0$ and $Q = 1$, respectively, Eqs. (C1) and (C3) are modified by

$$\delta F(Q^*) \cong \left( \frac{Q^* - Q_U}{Q_F - Q_U} \right) \delta F_f + \left( \frac{Q_F - Q^*}{Q_F - Q_U} \right) \delta F_u \qquad (C4)$$

and

$$\frac{\delta[F_u - F(Q^*)]}{\delta[F_f - F_u]} = -\left( \frac{Q^* - Q_U}{Q_F - Q_U} \right), \qquad (C5)$$

where $Q_U$ and $Q_F$ are the respective positions of the unfolded and folded states. So one can obtain the barrier position from the slope of a plot of $\ln k_F$ vs $\ln K_{eq}$, given the positions of the unfolded and folded states (see Fig. 8).

[1] J. Bryngelson and P. G. Wolynes, Proc. Natl. Acad. Sci. USA **84**, 7524 (1987).

[2] J. Bryngelson, J. O. Onuchic, N. D. Socci, and P. G. Wolynes, Proteins: Struct., Funct., Genet. **21**, 167–195 (1995).

[3] M. S. Friedrichs and P. G. Wolynes, Tet. Comp. Meth. **3**, 175 (1990). D. Thirumalai and Z. Guo, Biopolymers **35**, 137 (1995).

[4] A. Kolinski, A. Godzik, and J. Skolnick, J. Chem. Phys. **98**, 7420 (1993);

N. D. Socci and J. N. Onuchic, *ibid*. **101**, 1519 (1994); K. A. Dill, S. Bromberg, K. Yue, K. M. Fiebig, D. P. Yee, P. D. Thomas, and H. S. Chan, Protein Sci. **4**, 561 (1995); A. Sali, E. Shakhnovich, and M. Karplus, J. Mol. Biol. **235**, 1614 (1994).

[5] J. N. Onuchic, P. G. Wolynes, Z. Luthey-Schulten, and N. D. Socci, Proc. Natl. Acad. Sci. USA **92**, 3626 (1995).

[6] M. Mezard, G. Parisi, and M. A. Virasoro, *Spin Glass Theory and Beyond* (World Scientific, Singapore, 1986). K. H. Fischer and J. A. Hertz, *Spin Glasses* (Cambridge University Press, Cambridge, 1991).

[7] L. D. Landau and E. M. Lifshitz, *Statistical Physics, Part 1*, 3rd ed., edited by J. B. Sykes, M. J. Kearsley trans., Course of Theoretical Physics Vol. 5 (Pergamon, Oxford, 1980); N. Goldenfeld, *Lectures on Phase Transitions and the Renormalization Group* (Addison-Wesley, New York, 1992).

[8] T. Garel and H. Orland, Europhys. Lett. **6**, 307 (1988); E. I. Shakhnovich and A. M. Gutin, *ibid*. **8**, 327 (1989); Biophys. Chem. **34**, 187 (1989).

[9] M. Sasai and P. G. Wolynes, Phys. Rev. Lett. **65**, 2740 (1990).

[10] B. Derrida, Phys. Rev. B **24**, 2613 (1981).

[11] Z. Luthey-Schulten, B. E. Ramirez, and P. G. Wolynes, J. Phys. Chem. **99**, 2177 (1995); J. G. Saven and P. G. Wolynes, J. Mol. Biol. **257**, 199 (1996).

[12] S. S. Plotkin, J. Wang, and P. G. Wolynes, Phys. Rev. E **53**, 6271 (1996).

[13] B. Derrida and E. Gardner, J. Phys C **19**, 2253 (1986).

[14] M. S. Friedrichs and P. G. Wolynes, Science **246**, 371 (1989); M. S. Friedrichs, R. A. Goldstein, and P. G. Wolynes, J. Mol. Biol. **222**, 1013 (1991).

[15] S. Ramanathan and E. Shakhnovich, Phys. Rev. E **50**, 1303 (1994); V. S. Pande, A. Y. Grosberg, and T. Tanaka, J. Phys. II France **4**, 1711 (1994).

[16] P. J. Flory, J. Am. Chem. Soc., **78**, 5222 (1956).

[17] G. S. Huang and T. G. Oas, Proc. Natl. Acad. Sci. USA **92**, 6878 (1995).

[18] S. E. Jackson and A. R. Fersht, Biochemistry **30**, 10428 (1991).

[19] J. F. Douglas and T. Ishinabe, Phys. Rev. E **51**, 1791 (1995).

[20] Unless otherwise indicated, all entropies have dimensions of Boltzmann's constant $k_B$.

[21] Unless otherwise indicated, all temperatures have units of energy, with Boltzmann's constant $k_B$ included in the definition of $T$.

[22] J. G. Kirkwood, J. Chem. Phys. **3**, 300 (1935).

[23] R. A. Goldstein, Z. A. Luthey-Schulten, and Peter G. Wolynes, Proc. Natl. Acad. Sci. USA **89**, 4918 (1992).

[24] B. Derrida and G. Toulouse, J. Phys. Lett. **46**, 223 (1985).

[25] E. I. Shakhnovich and A. M. Gutin, Proc. Natl. Acad. Sci. USA **90**, 7195 (1993).

[26] N. D. Socci, J. N. Onuchic, and P. G. Wolynes J. Chem. Phys. **104**, 5860 (1996).

[27] A. Sali, E. Shakhnovich, and M. Karplus, Nature 248 (1994).

[28] An explanation for this is that in larger polymers, entropy loss due to topological constraints is more dramatic in $Q$ because a smaller fraction of total native contacts is necessary to constrain the polymer. That is, as $N$ increases, the number of bonds per monomer in the fully constrained state ($Q = 1$) approaches the bulk limit of 2 [see Eq. (A4)], while only one bond is needed to constrain a monomer. So this pushes the position $Q^*$ of the barrier in, roughly as $1/Z_N$. Density coupling quantitatively modifies the above picture in that as $N$ grows, there is a slower decrease in barrier position $Q^*$.

[29] T. Pascher, J. P. Chesick, J. R. Winkler, and H. B. Gray, Science **271**, 1558 (1996).

[30] V. I. Abkevich, A. M. Gutin, and E. I. Shakhnovich, Biochemistry **33**, 10026 (1994).

[31] A. M. Gutin, V. I. Abkevich, and E. I. Shakhnovich, Phys. Rev. Lett. **77**, 5433 (1996).

[32] D. Thirumalai, J. Phys. 1 France **5**, 1457 (1995).

[33] J. Wang, S. S. Plotkin, and P. G. Wolynes, J. Phys 1 France (in press).

[34] N. D. Socci and J. N. Onuchic, J. Chem. Phys. **101**, 1519 (1994).

[35] G. S. Huang and T. G. Oas, Biochemistry **34**, 3884 (1995).

[36] L. S. Itzhaki, D. E. Otzen, and A. R. Fersht J. Mol. Biol. **254**, 260 (1995).

[37] E. M. Boczko and C. L. Brooks, Science **269**, 393 (1995).

[38] J. N. Onuchic, N. D. Socci, Z. Luthey-Schulten, and P. G. Wolynes, Folding & Design **1**, 441 (1996).