

Geometry and symmetry prescript the free-energy landscape of proteins

Trinh Xuan Hoang*, Antonio Trovato†, Flavio Seno†, Jayanth R. Banavar‡§, and Amos Maritan†§¶

*Institute of Physics, National Centre for Natural Science and Technology, 46 Nguyen Van Ngoc, Hanoi, Vietnam; †Italian National Institute for Materials Physics and Dipartimento di Fisica "G. Galilei," Università di Padova, Via Marzolo 8, 35131 Padua, Italy; ‡Department of Physics, 104 Davey Laboratory, Pennsylvania State University, University Park, PA 16802; and §Abdus Salam International Center for Theoretical Physics, Strada Costiera 11, 34014 Trieste, Italy

Communicated by Morrel H. Cohen, Rutgers, The State University of New Jersey, Piscataway, NJ, April 8, 2004 (received for review October 9, 2003)

We present a simple physical model that demonstrates that the native-state folds of proteins can emerge on the basis of considerations of geometry and symmetry. We show that the inherent anisotropy of a chain molecule, the geometrical and energetic constraints placed by the hydrogen bonds and sterics, and hydrophobicity are sufficient to yield a free-energy landscape with broad minima even for a homopolymer. These minima correspond to marginally compact structures comprising the menu of folds that proteins choose from to house their native states in. Our results provide a general framework for understanding the common characteristics of globular proteins.

Protein folding (1–5) is complex because of the sheer size of protein molecules, the twenty types of constituent amino acids with distinct side chains, and the essential role played by the environment. Nevertheless, proteins fold into a limited number (6, 7) of evolutionarily conserved structures (8, 9). It is a familiar, yet remarkable, consequence of symmetry and geometry that ordinary matter crystallizes in a limited number of distinct forms. Indeed, crystalline structures transcend the specifics of the various entities housed in them. Here, we ask the analogous question (10): is the menu of protein folds also determined by geometry and symmetry?

We show that a simple model that encapsulates a few general attributes common to all polypeptide chains, such as steric constraints (11–13), hydrogen bonding (14–16), and hydrophobicity (17), gives rise to the emergent free-energy landscape of globular proteins. The relatively few minima in the resulting landscape correspond to putative marginally compact native-state structures of proteins, which are assemblies of helices, hairpins, and planar sheets. A superior fit (18, 19) of a given protein or sequence of amino acids to one of these predetermined folds dictates the choice of the topology of its native-state structure. Instead of each sequence shaping its own free energy landscape, we find that the overarching principles of geometry and symmetry determine the menu of possible folds that the sequence can choose from.

Following Bernal (20), the protein problem can be divided into two distinct steps: first, analogous to the elucidation of crystal structures, one must identify the essential features that account for the common characteristics of all proteins; second, one must understand what makes one protein different from another. Guided by recent work (21, 22) that has shown that a faithful description of a chain molecule is a tube and using information from known protein native-state structures, our focus, in this paper, is on the first step: we demonstrate that the native-state folds of proteins emerge from considerations of symmetry and geometry within the context of a simple model.

We model a protein as a chain of identical amino acids, represented by their C_α atoms, lying along the axis of a self-avoiding flexible tube. The preferential parallel placement of nearby tube segments approximately mimics the effects of the anisotropic interaction of hydrogen bonds whereas the space needed for the clash-free packing of side chains is approximately captured by the non-zero tube thickness (21, 22). Here, we

carefully incorporate these key geometrical features by means of an extensive statistical analysis of experimentally determined native-state structures in the Protein Data Bank (PDB).

A tube description places constraints on the radii of circles drawn through both local and nonlocal triplets of C_α positions of a protein native structure (22, 23). Furthermore, when one deals with a chain molecule, the tube picture underscores the crucial importance of knowing the context that an amino acid is in within the chain. The standard coarse-grained approach considers the locations of interacting amino acid pairs. Here, instead, we incorporate the strongly directional hydrogen bonding between a pair of amino acids, through an analysis of the PDB to determine the constraints on the mutual orientation of the local coordinate systems defined from a knowledge of the locations of the C_α atoms (see *Methods* and Fig. 1). The geometrical constraints associated with the tube and hydrogen bonds that we consider here are representative of the typical aspecific behavior of the interacting amino acids.

There are two other ingredients in the model: a local bending penalty, which is related to the steric hindrance of the amino acid side chains, and a pair-wise interaction of the standard type mediated by the water (17). Even though these two properties clearly depend on the specific amino acids involved in the interaction, here, we choose to study the phase diagram of a homo-peptide chain by varying its overall hydrophobicity and local bending penalty, while keeping them constant along the chain. This is the simplest and most general way to assess their relevance in shaping the free-energy landscape.

Methods

Tube Geometry. The protein backbone is modeled as a chain of C_α atoms (Fig. 2*a*) with a fixed distance of 3.8 Å between successive atoms along the chain, an excellent assumption for all but non-cis proline amino acids (24). The geometry imposed by chemistry dictates that the bond angle associated with three consecutive C_α atoms is between 82° and 148°.

Self-avoiding conformations of the tube whose axis is the protein backbone are identified by considering all triplets of C_α atoms and drawing circles through them and ensuring that none of their radii is smaller than the tube radius (25) (Fig. 2*a*). At the local level, the three-body constraint ensures that a flexible tube cannot have a radius of curvature any smaller than the tube thickness, to prevent sharp corners, whereas, at the nonlocal level, it does not permit any self-intersections. There is an inherent local anisotropy due to the special direction singled out by consecutive atoms along the chain, which enforces a preference for parallel alignment of neighboring tube segments in a compact conformation.

The backbone of C_α atoms is treated as a flexible tube of radius 2.5 Å, a constraint imposed on all (local and nonlocal)

§To whom correspondence should be addressed. E-mail: banavar@psu.edu or maritan@pd.infn.it.

© 2004 by The National Academy of Sciences of the USA

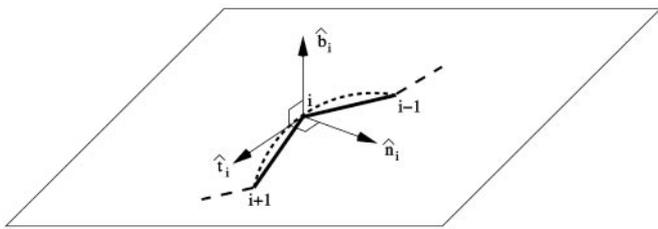


Fig. 1. Sketch of the local coordinate system. For each C_α atom i (except the first and the last one), the axes of a right-handed local coordinate system are defined as follows. The tangent vector \hat{t}_i is parallel to the segment joining $i-1$ with $i+i$. The normal vector \hat{n}_i joins i to the center of the circle passing through $i-1$, i , and $i+1$, and it is perpendicular to \hat{t}_i , \hat{t}_{i+1} and \hat{n}_i along with the three contiguous C_α atoms lie in a plane shown in the figure. The binormal vector \hat{b}_i is perpendicular to this plane. The vectors \hat{t}_i , \hat{n}_i , and \hat{b}_i are normalized to unit length.

three-body radii, an assumption validated for protein native structures (23). It is interesting to note that recent observations of residual dipolar couplings in short peptides (26) in the denatured state have demonstrated their stiffness and their anisotropic deformability; the building blocks of proteins are relatively stiff segments with strong directional preferences.

Sterics. Steric constraints require that no two nonadjacent C_α atoms are allowed to be at a distance closer than 4 Å. Ramachandran and Sasisekharan (11) showed that steric considerations based on a hard sphere model lead to clustering of the backbone dihedral angles in two distinct α and β regions for non-glycyl and non-prolyl residues. The two backbone geometries that allow for systematic and extensive hydrogen bonding

(14–16) are the α -helix and the β -sheet, obtained by a repetition of the backbone dihedral angles from the two regions, respectively (13). Short chains rich in alanine residues, which are a good approximation to a stretch of the backbone, can adopt a helical conformation in water (see refs. 27–32 for a detailed discussion of experimental conditions that would lead to a helical conformation). However, when one has more heterogeneous side chains, the helix backbone could sterically clash with some side chain conformers, resulting in a loss of conformational entropy (33). When the price in side chain entropy is too large, an extended backbone conformation results, pushing the segment toward a β -strand structure (13). These steric constraints are approximately imposed through an energy penalty (denoted by e_R) when the local radius of curvature is between 2.5 Å and 3.2 Å. (The magnitude of the penalty does not depend on the specific value of the radius of curvature, provided it is between these values.) There is no cost when the local radius exceeds 3.2 Å. Note that the tube constraint does not permit any local radius of curvature to take on a value less than the tube radius, 2.5 Å.

Hydrogen Bonds. We do not allow more than two hydrogen bonds to form at a given C_α location. In our representation of the protein backbone, local hydrogen bonds form between C_α atoms separated by two residues along the sequence with an energy defined to be -1 unit whereas nonlocal hydrogen bonds are those that form between C_α atoms separated by more than three residues along the sequence with an energy of -0.7 . This energy difference is based on experimental findings that the local bonds provide more stability to a protein than do the nonlocal hydrogen bonds (34). Cooperativity effects (35, 36) are taken into account by adding an energy of -0.3 units when consecutive hydrogen bonds along the sequence are formed. There is some latitude in the choice of the values of these

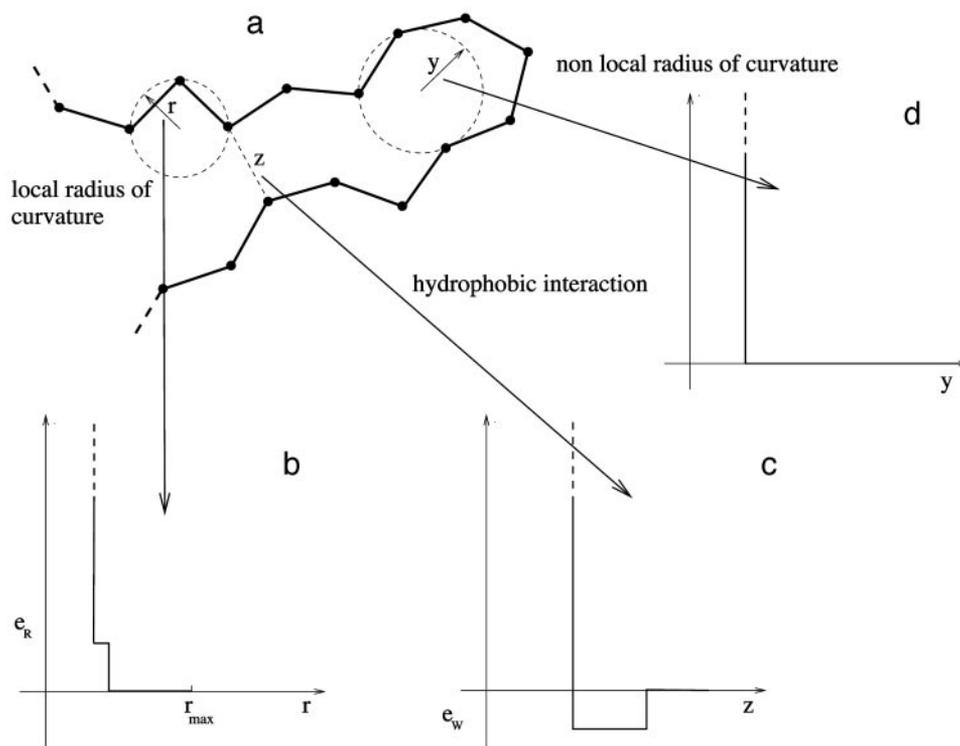


Fig. 2. Sketch of a portion of a protein chain. (a) The black spheres represent the C_α atoms of the amino acids. The local radius of curvature r is defined as the radius of the circle passing through three consecutive atoms and is constrained to lie between 2.5 Å and 7.9 Å (r_{\max}). A penalty e_R is imposed when $2.5 \leq r \leq 3.2$ (see b). The hydrophobic interaction, e_W , is operative when two atoms separated by more than two along the sequence are within 7.5 Å of each other (see c). Note that two nonadjacent atoms cannot be closer than 4 Å. A flexible tube is characterized by the constraint that none of the three-body radii is less than the tube thickness, chosen here to be 2.5 Å (see b and d).

energy parameters. The results that we present are robust to changes (at least of the order of 20%) in these parameters.

Geometrical Constraints Due to Hydrogen Bonding. Three noncolinear consecutive atoms ($i - 1$, i , and $i + 1$) of the chain define a plane. At atom i (special care is needed to adapt these rules to atoms at the C and N termini), one may define a tangent vector (along the direction joining the $i - 1$ and $i + 1$ atoms) and a normal vector (along the direction joining the i th atom and the center of the circle passing through the three atoms), which together define a plane. One then defines a binormal vector \vec{b}_i perpendicular to the plane with the tangent, normal, and binormal forming a right-handed local coordinate system (Fig. 1). This coordinate system defines the context of an amino acid within a chain, a feature that plays a crucial role in the tube picture. For hydrogen bond formation between atom i and j , the distance between these atoms ought to be between 4.7 Å and 5.6 Å (4.1 Å and 5.3 Å) for the local (nonlocal) case. A study of protein native state structures reveals an overall nearly parallel alignment of the axes defined by three vectors: the binormal vectors at i and j and the vector \vec{r}_{ij} joining the i and j atoms. A hydrogen bond is allowed to form only when the binormal axes are constrained to be within 37° of each other whereas the angle between the binormal axes and that defined by \vec{r}_{ij} ought to be <20°. Additionally, for the cooperative formation of nonlocal hydrogen bonds, one requires that the corresponding binormal vectors of successive C_α atoms make an angle >90°. The first and the last residues of the chain are special cases because their binormal vectors are not defined. In order for such residues to form a hydrogen bond (with each other or with other internal residues in the chain), it is required that the angle between the associated ending peptide link and the connecting vector to the other residue participating in the hydrogen bond is between 70° and 110°. As in real protein structures, when helices are formed, they are constrained to be right-handed. This constraint is enforced by requiring that the backbone chirality associated with each local hydrogen bond is positive. The chirality is defined as the sign of $(\vec{r}_{i,i+1} \times \vec{r}_{i+1,i+2}) \cdot \vec{r}_{i+2,i+3}$.

Our approach for the derivation of the geometrical constraints imposed by hydrogen bonds is similar to that carried out at the level of an all-atom description of the protein chain (37). For the simpler C_α atom-based description, hydrogen bond energy functions have been introduced previously (38, 39) but without any input from a statistical analysis of protein structures.

Hydrophobic Interactions. The hydrophobic (hydrophilic) effects mediated by the water are captured through a relatively weak interaction e_W (either attractive or repulsive) between C_α atoms that are within 7.5 Å of each other (Fig. 2c). Note that hydrogen bonds can easily be formed between the amino acid residues in an extended conformation and the water molecules. Within our model, the intra-chain hydrogen bond interaction introduces an effective attraction, because water molecules are not explicitly present. The hydrophobicity scale is thus renormalized (e.g., even when e_W is weakly positive, there could be an effective attraction resulting in structured conformations such as a single helix or a planar sheet). A negative e_W is, in any case, crucial for promoting the assembly of secondary motifs in native tertiary arrangements. The properties of the model are summarized in Table 1.

Results and Discussion

Fig. 3 shows the ground state phase diagram obtained from Monte-Carlo computer simulations using the simulated annealing technique (40). [The solvent-mediated energy, e_W , and the local radius of curvature energy penalty, e_R , (see *Methods* for a description of the energy parameters) are measured in units of the local hydrogen bond energy.] When e_W is sufficiently repul-

Table 1. Summary of all geometrical and energetical parameters involved in the model definition

Parameter	Constraint
Tube approximation*	$R_{ijk} \geq 2.5\text{Å}, \forall i < j < k$
Local radius of curvature	$2.5\text{Å} \leq R_{i-1,i,i+1} \leq 7.9\text{Å}, \forall 1 < i < N^{\dagger}$
Self-avoidance	$r_{ij} \geq 4\text{Å}, \forall i < j - 1$
Amino acid specific?	No
Local hydrogen bond [‡]	$j = i + 3$
C_α - C_α distance	$4.7\text{Å} \leq r_{ij} \leq 5.6\text{Å}$
Binormal-binormal correlation [§]	$ \vec{b}_i \cdot \vec{b}_j > 0.8$
Binormal-connecting vector [§]	$ \vec{b}_i \cdot \vec{c}_{ij} > 0.94, \vec{b}_j \cdot \vec{c}_{ij} > 0.94$
Chirality	$(\vec{r}_{i,i+1} \times \vec{r}_{i+1,i+2}) \cdot \vec{r}_{i+2,i+3} > 0$
Energy	-1
Amino acid specific?	No
Nonlocal hydrogen bond [‡]	$j > i + 4$
C_α - C_α distance	$4.1\text{Å} \leq r_{ij} \leq 5.3\text{Å}$
Binormal-binormal correlation [§]	$ \vec{b}_i \cdot \vec{b}_j > 0.8$
Binormal-connecting vector [§]	$ \vec{b}_i \cdot \vec{c}_{ij} > 0.94, \vec{b}_j \cdot \vec{c}_{ij} > 0.94$
Energy	-0.7
Amino acid specific?	No
Cooperative hydrogen bonds	between (i, j) and $(i \pm 1, j \pm 1)$
β -sheet zig-zag pattern ^{§**}	$\vec{b}_i \cdot \vec{b}_{i\pm 1} < 0, \vec{b}_j \cdot \vec{b}_{j\pm 1} < 0$
Energy per pair	-0.3
Amino acid specific?	No
Bending rigidity	$R_{i-1,i,i+1} \leq 3.2\text{Å}$
Energy	e_R
Amino acid specific?	Yes (for a heteropolymer)
Hydrophobic contact	$j > i + 2$
C_α - C_α distance	$r_{ij} \leq 7.5\text{Å}$
Energy	e_W
Amino acid specific?	Yes (for a heteropolymer)

All geometrical properties have been derived by means of a thorough analysis of PDB native structures.

* R_{ijk} is the radius of a circle drawn through the C_α positions of i, j , and k .

[†] N is the number of residues.

[‡]Each residue is constrained to form no more than two hydrogen bonds (except the residues located at the chain termini, which form at most one hydrogen bond).

[§]Applied only when the corresponding binormal vectors exist.

^{||}For $i = 1$ and/or $j = N$, this is replaced by the constraint that the connecting vector is making an angle between 70° and 110° with the extremal peptide links.

^{||}The connecting vector, $\vec{c}_{ij} = \vec{r}_{ij}/r_{ij}$ is a unit vector joining i and j .

^{**}Applied when at least one of the two cooperative hydrogen bonds is nonlocal.

sive (hydrophilic) (and $e_R > 0.3$ in the phase diagram), one obtains a swollen phase with very few contacts between the C_α atoms. When e_W is sufficiently attractive, one finds a very compact, globular phase with featureless ground states with a high number of contacts.

Between these two phases (and in the vicinity of the swollen phase), a marginally compact phase emerges (the interactions barely stabilize the ordered phase) with distinct structures including a single helix, a bundle of two helices, a helix formed by β -strands, a β -hairpin, three-stranded β -sheets with two distinct topologies, and a β -barrel-like conformation. Strikingly, these structures are the stable ground states in different parts of the phase diagram. Furthermore, other conformations, closely resembling distinct supersecondary arrangements observed in proteins (6), such as the β - α - β motif, are found to be competitive local minima whose stability can be enhanced by sequence design (for example, nonuniform values of curvature energy penalties for single amino acids and hydrophobic interactions for amino acid pairs). Fig. 4 shows a compendium of various structures

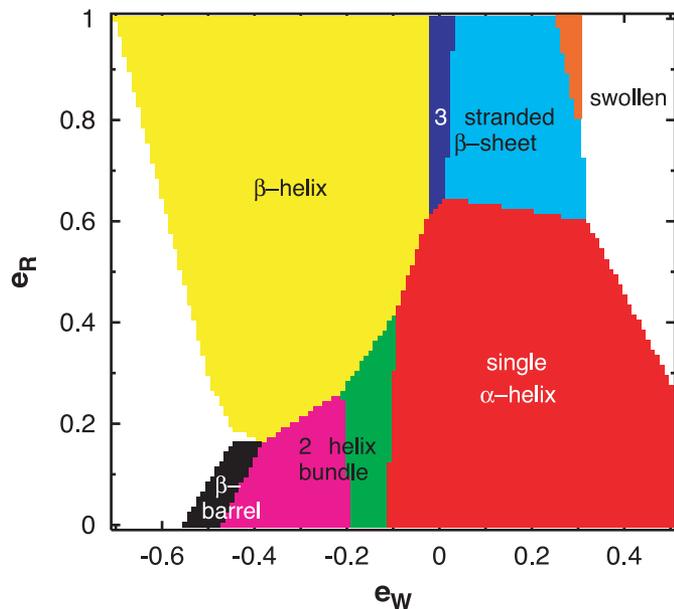


Fig. 3. Phase diagram of ground state conformations. The ground state conformations were obtained by using Monte-Carlo simulations of chains of 24 C_α atoms. e_R and e_W denote the local radius of curvature energy penalty and the solvent-mediated interaction energy, respectively. Over 600 distinct local minima were obtained in different parts of parameter space, starting from a random conformation and successively distorting the chain with pivot and crankshaft moves commonly used in stochastic chain dynamics (43). A Metropolis Monte-Carlo procedure is used with a thermal weight $\exp(-E/T)$, where E is the energy of the conformation and the temperature T is set initially at a high value and then decreased gradually to zero. In the orange phase, the ground state is a two-stranded β -hairpin. Two distinct topologies of a three-stranded β -sheet (dark and light blue phases) are found corresponding to conformations shown in conformations *i* and *j* in Fig. 4, respectively. The helix bundle shown in conformation *b* in Fig. 4 is the ground state in the green phase whereas the ground state conformation in the magenta phase has a slightly different arrangement of helices. The white region in the left of the phase diagram has large attractive values of e_W , and the ground state conformations are compact globular structures with a crystalline order induced by hard sphere packing considerations (44) and not by hydrogen bonding (conformation *l* in Fig. 4).

obtained in our studies, including for comparison a generic compact conformation of a conventional polymer chain (with no tube geometry or hydrogen bonds), which neither is made up of helices or sheets nor possesses the significant advantages of protein structures. Although there is a remarkable similarity between the structures that we obtain and protein folds, our simplified coarse-grained model is not as accurate as an all-atom representation of the polypeptide chain in capturing features such as the packing of amino acid side chains.

The fact that different putative native structures are found to be competing minima for the same homopolymeric chain clearly establishes that the free-energy landscape of proteins is presculpted by means of the few ingredients used in our model. At the same time, relatively small changes in the parameters e_W and e_R lead to significant differences in the emergent ground state structure, underscoring the sensitive role played by chemical heterogeneity in selecting from the menu of native state folds.

Fig. 5*a* is a contour plot of the free energy at a temperature higher than the folding transition temperature (identified by the specific heat peak) for the parameter values $e_W = -0.08$ and $e_R = 0.3$ for which the ground state is an α -helix (Fig. 3). The free energy landscape has just one minimum corresponding to the denatured phase whose typical conformations are somewhat compact but featureless. The contour plot at the folding tran-

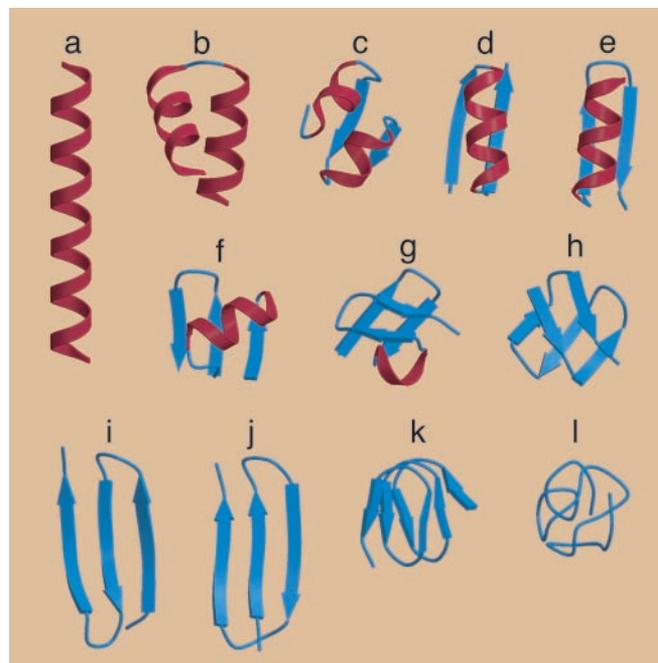


Fig. 4. MOLSCRIPT representation of the most common structures obtained in our simulations. Helices and strands are assigned when local or nonlocal hydrogen bonds are formed according to the described rules. Conformations *a*, *b*, *h*, *i*, *j*, and *k* are the stable ground states in different parts of the parameter space shown in Fig. 3. Conformations *c*, *d*, *e*, *f*, and *g* are competitive local minima. Conformation *l* is that of a generic compact polymer chain, obtained by switching off hydrogen bonds, the tube constraint, and curvature energy penalty, and is obtained on maximizing the total number of hydrophobic contacts.

sition temperature (Fig. 5*b*) has three local minima corresponding to an α -helix, a three-stranded β -sheet, and the denatured state. At lower temperatures, the α -helix is increasingly favored and the β -sheet is never the global free-energy minimum. Many protein-folding experiments show that, for small globular proteins at the transition temperature, only two states (folded and unfolded) are populated. The appearance in the present framework of multiple states for a homopolymer chain suggests that two-state folders might have been evolutionarily selected by sequence design favoring the native-state conformation over competing folds in the presculpted landscape.

Such a design is indeed straightforward within our model. For example, the α - β - α motif shown in Fig. 4*d* (which is a local energy minimum for a homopolymer) can be stabilized into a global energy minimum for the sequence HPHHHPPPHHP-PHHPPPHHHPP, with $e_W = -0.4$ for HH contacts and $e_W = 0$ for other contacts, and $e_R = 0.3$ for all residues.

It is interesting to note that lattice models of compact homopolymers yield large amounts of secondary structure (41); local radius of curvature constraints are built into lattice models. However, an all-atom study of polyalanine has shown that compactness alone is insufficient to obtain secondary structures (42). Even a simple tube subject to an attractive self-interaction favoring compaction has a tendency to form helices, hairpins and sheets when the ratio of the tube thickness to the range of attractive interaction is tuned properly (22). Our work here underscores the importance of hydrogen bonds in stabilizing both helices and sheets simultaneously (without any need for adjustment of the tube thickness), allowing the formation of tertiary arrangements of secondary motifs. Indeed, the fine tuning of the hydrogen bond and the hydrophobic interaction is of paramount importance in the selection of the marginally

