

# Water in protein structure prediction

Garegin A. Papoian<sup>†‡</sup>, Johan Ulander<sup>†‡§</sup>, Michael P. Eastwood<sup>†¶</sup>, Zaida Luthey-Schulten<sup>||</sup>, and Peter G. Wolynes<sup>†,††</sup>

<sup>†</sup>Department of Chemistry and Biochemistry and Center for Theoretical Biological Physics, University of California at San Diego, 9500 Gilman Drive, La Jolla, CA 92093-0371; and <sup>||</sup>Department of Chemistry, University of Illinois at Urbana-Champaign, Urbana, IL 61801

Contributed by Peter G. Wolynes, December 4, 2003

**Proteins have evolved to use water to help guide folding. A physically motivated, nonpairwise-additive model of water-mediated interactions added to a protein structure prediction Hamiltonian yields marked improvement in the quality of structure prediction for larger proteins. Free energy profile analysis suggests that long-range water-mediated potentials guide folding and smooth the underlying folding funnel. Analyzing simulation trajectories gives direct evidence that water-mediated interactions facilitate native-like packing of supersecondary structural elements. Long-range pairing of hydrophilic groups is an integral part of protein architecture. Specific water-mediated interactions are a universal feature of biomolecular recognition landscapes in both folding and binding.**

Water is intimately involved in protein folding (1–4). That proteins denature both on heating and cooling strongly implicates the involvement of water degrees of freedom. Kauzmann (5) correctly inferred from thermodynamics the hydrophobic layering characteristic of protein structure before protein structures were determined crystallographically. The kinetics of water exclusion is often considered in discussing mechanisms of protein folding, but again it is the avoidance of water in the final folded structure that is emphasized (1). Hydrophobicity patterns have long been a dominant consideration in predicting protein structure by using sequence data (6) and are basic in synthetic protein design (7). Nevertheless, the structured character of water has not been a paramount factor in most existing algorithms for structure prediction (8). These usually rely on effective pair potentials (9) or buried surface area terms to account for the free energy of burying hydrophobic residues (10).

In this article, we hypothesize that specific water-mediated interactions help guide the folding process even before native contacts form. Using this idea we develop a bioinformatic, nonpairwise-additive interaction model accounting for water and show that it greatly improves the efficiency and accuracy of structure prediction for  $\alpha$ -helical proteins. Analysis of folding trajectories with this potential strongly implicates the guiding role of long-range water-mediated interactions. Interestingly, we find here that long-range hydrophilic interactions, as distinct from hydrophobic interactions, also take center stage.

The bioinformatic route to water-mediated potentials is difficult in several ways (for more directly physical approaches see ref. 11). Although bound water is visible in structures, localizing waters is more difficult than localizing main-chain atoms. Monomeric protein structures also have relatively few visible water-mediated interactions. Our path to a water-mediated potential started with an energy landscape analysis of protein–protein interactions and a bioinformatic survey of interfaces in dimer structures (12, 13). We found that the often-used contact potentials (9) worked well to describe hydrophobic binding interfaces; however, hydrophilic interfaces were poorly recognized (13). This finding suggests that longer-range interresidue contacts, mediated by water, play an important role in stabilizing these interfaces (13). To test this hypothesis, we derived both direct and water-mediated binding potentials (13). When these two potentials were used simultaneously (13), smooth recognition of diverse binding interfaces was achieved (in contrast to the direct contact potential). Here, we show that water-mediated

interactions play an important role not only in binding interfaces but in folding of monomeric proteins.

We use the associative memory (AM) Hamiltonian molecular dynamics model as a starting point (14–16). This Hamiltonian has two principal components: general polymer physics-based terms that are sequence independent, collectively called “backbone,” and sequence-dependent knowledge-based distance-dependent additive potentials, collectively denoted as AM/C (AM/contact). The AM part describes interactions between all pairs of residues that are separated in sequence between 3 and 12 residues. It uses a set of nonhomologous memory proteins to build a funneled energy landscape by matching fragments. The C part applies to tertiary contacts between residues separated by  $>12$  residues in sequence. All parameters in the potential have been optimized with a self-consistent procedure based on the energy landscape theory as described (15) (see *Appendix: Computational Details* and *Supporting Text*, which is published as supporting information on the PNAS web site).

The C part of the AM/C potential describes effective interactions between  $C\beta$  ( $C\alpha$  for Gly) atoms in each residue pair. It consists of three wells covering the 4.5- to 8.5-Å, 8.5- to 10.0-Å, and 10.0- to 15.0-Å distance intervals. Similarly, the potentials used in this study contain a first well for the 4.5- to 6.5-Å interval, whereas the second well is replaced by a local density-dependent potential (discussed below). They contain no third well, because it is unlikely that specific interresidue interactions are mediated by water to such a long distance (10–15 Å). There is also a residue-specific many-body burial profile potential describing coordination preferences of all 20 amino acids. The parameters for the resulting potential, which we call AM/W (W for water), were optimized by using our earlier sequence-based approach (13). We have further refined these parameters by using the self-consistent structurally based optimization scheme from energy landscape theory (15). We denote the original potential, AM/W-0, and the more refined one, AM/W-1 (see *Appendix for Computational Details* and *Supporting Text*).

For the coarse-grained models considered in our study, the definition of water-mediated contacts naturally becomes somewhat more indirect than, for example, in full-atom simulations. Because direct contacts are defined as occurring between residues that have a distance between  $C\beta$  ( $C\alpha$  for Gly) atoms of  $<6.5$  Å, a similar constraint for water-mediated contacts places them in the 6.5- to 9.5-Å distance interval. A more extensive discussion of the rationale for this choice is given in ref. 13 where the results for protein binding recognition also was found to be robust with respect to various alternative definitions of water-mediated contact range.

Abbreviations: AM, associative memory; AM/C, AM/contact; AM/W, AM/water; RMSD, rms displacement; CE, combinatorial extension; PDB, Protein Data Bank; CASP, Critical Assessment of Techniques for Protein Structure Prediction.

<sup>†</sup>G.A.P. and J.U. contributed equally to this work.

<sup>§</sup>Present address: AstraZeneca R&D Mölndal, Drug Metabolism and Pharmacokinetics and Bioanalytical Chemistry, SE-431 83 Mölndal, Sweden.

<sup>¶</sup>Present address: D. E. Shaw Research and Development, 120 West 45th Street, New York, NY 10036.

<sup>††</sup>To whom correspondence should be addressed. E-mail: pwolynes@ucsd.edu.

© 2004 by The National Academy of Sciences of the USA

To interact through water, we require that both residues are sufficiently exposed to water, or equivalently, neither residue should be buried in the protein interior (hydrophobic core). To model this we use a highly nonadditive local density-dependent potential: when either of the residues in the pair attains a local neighbor density above a critical threshold value (i.e., becomes buried), the potential switches smoothly but quickly from water mediated to protein mediated.

## Results and Discussion

**Physical Interpretation of the Interaction Potentials.** Before discussing the simulation results, we briefly analyze the main qualitative differences between the AM/C and AM/W interaction potentials. The interactions within the range of the first well and the protein-mediated interactions of the second well of the AM/W potential are qualitatively similar to their corresponding AM/C counterparts (see Fig. 1*A* and *B*). The main difference lies in the interactions between hydrophilic residues in the second-well water-mediated interactions (see Fig. 1*C*). Whereas very polar second-well interactions are destabilized on average for the AM/C potential, they are highly stabilized for the AM/W potential when two residues are in a low-density environment, i.e., when residues interact through water (Fig. 1*C*).

Although these potentials are knowledge-based in origin, examining their details gives interesting physical insight into the nature of biomolecular forces (13). Perusing of charged residue interactions in Fig. 1*B* and *C* suggests that a large desolvation penalty must be paid when a fully direct contact is formed, and, therefore, charged and highly polar residues prefer to avoid complete desolvation by interacting through one or two water layers. Even more interestingly, not only do oppositely charged residues attract each other when interacting through water, but so do residues of the same charge (Fig. 1*C*). This finding either indicates that residues of the same charge alter their mutual  $pK_a$  so only one residue is really charged (i.e., one has in fact a charged-polar interaction) or that correlated fluctuations of the counterion cloud (17) and the perturbation of the water hydrogen-bonding network bind the like-charged residues together.

**General Trends.** Given the differences among the potentials outlined above, we anticipate that the AM/W potential would significantly improve the AM/C potential results for those proteins that contain explicit water-bridged interactions in their native state. As we shall see, these water-mediated interactions also appear transiently during collapse and folding of the chain and help guide the heteropolymer into a correct topology. For each protein of 14 chosen for study (discussed below), we have carried out five distinct annealing runs ( $7.2 \times 10^5$  time steps) with each of the three potentials (AM/C, AM/W-0, and AM/W-1), starting from a randomly generated extended-coil conformation at high temperature (we have not optimized the annealing protocol for the AM/W potentials nor used other minimization techniques; ref. 18). For each run we have taken 240 snapshots at equal time intervals, monitoring the progress toward achieving a native-like conformation by using a contact overlap measure  $Q$ . Our  $Q$  measure is more stringent than the usual contact  $Q$ , because it takes into account not only the correctness of contacts that occur in the native structure but also the correctness of distances between all pairs of residues even when they are far apart in the native state. In addition to  $Q$ , when discussing various structures, we use other structural similarity measures, such as rms displacement (RMSD) and the combinatorial extension (CE) method, which makes sequence-independent alignment of two conformations (19). It is perhaps not surprising that comparing protein structures is a tricky business, involving several means of similarity as discussed (20).

The validation of any knowledge-based potential must be done on an unrelated set of test proteins because of the risk of

parameter overlearning. Nine of the 14  $\alpha$ -helical proteins used are “training” proteins for the AM/C potential, i.e., they were used to derive the parameters for the AM and C parts of the potential (15). The W part of AM/W-0 was optimized by using a sequence-based technique for an unrelated set of proteins. Thus for the AM/W-0 tertiary contact potential these nine proteins serve partially as test proteins. On the other hand, the W part of the AM/W-1 potential was refined by using the same training set of nine proteins. We emphasize these relationships to be attentive to the possibility of overlearning. It is necessary to apply the potential to an unrelated set of test proteins for confirmation.

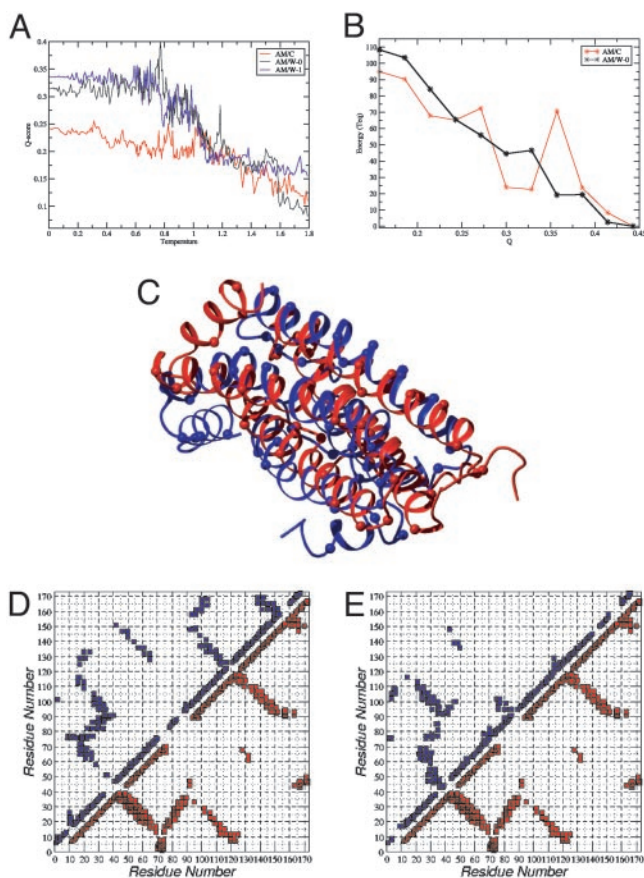
The performance of the AM/W potentials is well documented by five  $\alpha$ -helical test proteins that we discuss in detail below (see also *Supporting Text*). Two test proteins [Protein Data Bank (PDB) codes 1BG8 (21) and 1JWE (22)] were targets taken from the CASP3 (Critical Assessment of Techniques for Protein Structure Prediction) event (ref. 20; a detailed compilation of all CASP results may be found at <http://predictioncenter.llnl.gov>), and three, target T0170 [PDB code 1H40 (23)], target T172b [PDB code 1N2X (24)], and target T129a [PDB code 1IZM (A. Galkin, E. Sarikaya, C. Lehmann, A. Howard, and O. Herzberg, personal communication)], were taken from the CASP5 event (<http://predictioncenter.llnl.gov> and ref. 25). Our results compare favorably with the top CASP predictions for these proteins (<http://predictioncenter.llnl.gov>), but have nevertheless been obtained *a posteriori* (although in a fairly automatic manner) and should not be regarded as new CASP entries. Nevertheless, because numerous prediction groups participated in CASP, the CASP experiment has generated valuable statistical data that may be used to calibrate progress.

When the best  $Q$  scores obtained for each protein during all five annealing runs are compared across all 14 proteins (Fig. 2), the following trends becomes evident. First, AM/C and AM/W potentials show similar performance for small (<115 residues) training proteins. As for small test proteins, 1BG8 (21) is greatly improved by both AM/W-0 and AM/W-1 potentials, whereas T0170 is improved only by AM/W-0. The most significant trend, one that is highly desirable, comes when the largest proteins (>115 residues) are considered. A methodical improvement in the prediction of both training and test proteins is achieved by both AM/W-0 and AM/W-1 potentials, the latter showing a more uniform trend (Fig. 2). For large proteins, an improvement of 0.05–0.10 in  $Q$  is very significant, typically improving global RMSD by a few Å and significantly improving other measures of fold recognition, such as CE Z score.

**Specific Targets.** Having achieved substantial progress in protein structure prediction by using the tertiary contact potential incorporating long-range water-mediated interactions, we next investigate the cause of the improved structural recognition. We focus on three proteins: (i) PDB code 2FHA (26), a training protein, the largest one in the protein set; (ii) PDB code 1BG8 (21), a small test protein for which both AM/W-0 and AM/W-1 show very significant improvements; and (iii) CASP5 target T129a (A. Galkin, E. Sarikaya, C. Lehmann, A. Howard, and O. Herzberg, personal communication), the largest test protein in the protein set, that has two interacting domains.

Human iron storage protein, ferritin [PDB code 2FHA (26)], is the largest protein studied (172 residues). Although it was a training protein for AM/C, there was a large improvement in structure prediction when using both AM/W-0 and AM/W-1 (Fig. 3*A*). During the cooling schedule the divergence between the trajectories in nativeness occurs around  $T = 1.05$ . To further evaluate the difference between the potentials, we carried out free energy calculations as a function of  $Q$  by using the histogramming technique (16). These calculations show that the free energy minimum shifts toward more native-like structure for

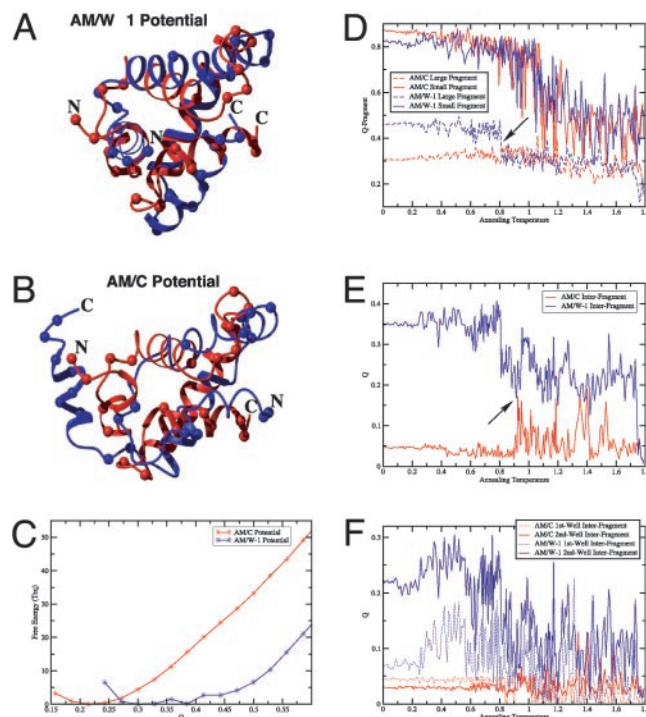




**Fig. 3.** Structure predictions for ferritin, PDB code 2FHA. (A) The best (of five for each potential)  $Q$ -score annealing trajectories are shown for three different potentials. (B) The average thermodynamic energy vs.  $Q$ . (C) Superposition of the AM/W-0 best  $Q$ -score structure (blue) and the native structure (red) is indicated. Spheres indicate charged residue  $C\alpha$  atoms. (D) The distance plot for the AM/W-0 best  $Q$ -score structure (blue, upper triangle) and the native structure (red, lower triangle). (E) Distance plot for the AM/C best  $Q$ -score structure (blue, upper triangle) and the native structure (red, lower triangle) are compared. In the AM/W-0 structure (D), only a small number of contacts are missing and a small registry shift near residue 70 occurs. In the AM/C structure (E) the C-terminal half misses on a major interhelical interface.

The small test (CASP3) protein for which we have observed large enhancement in native structure recognition is *Escherichia coli* stress-response protein HdeA [PDB code 1BG8 (21)]. The superposition of best predicted structure with AM/W-1 and the native is shown in Fig. 4A. At the overall  $Q$  score of 0.47, the CE alignment of 70 residues of total 76 residues produces an RMSD of 4.2 Å and  $Z$  score of 3.7. The global RMSD was 5.1 Å. The best AM/C prediction (Fig. 4B) again captures correctly large chunks of the structure (CE alignment of residues 7–62 produces a  $Z$  score of 3.3 and RMSD of 5.7 Å), but fails to pack them globally (overall  $Q = 0.31$ , global RMSD = 12.0 Å). Free energy calculations show that the minimum in  $F(Q)$  is shifted substantially toward the native for the AM/W potentials (Fig. 4C), rationalizing annealing results.

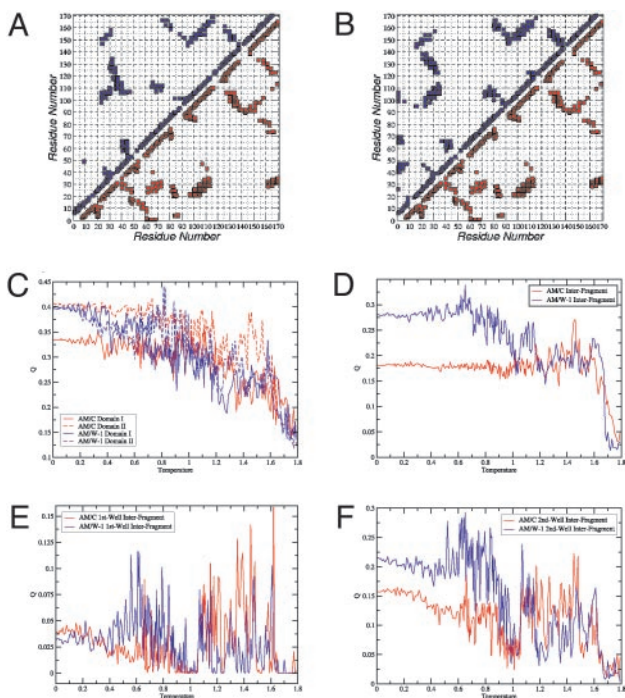
Closer examination of the HdeA sequence reveals that both N- and C-terminal 11-residue fragments are highly charged (four charged residues in the N-terminal fragment and seven charged residues in the C-terminal fragment). Fig. 4B shows that it is these terminal fragments that are packed incorrectly by the AM/C potential. For analysis, we thus partitioned the protein into two fragments: a larger N-terminal fragment consisting of residues 1–61, and a smaller C-terminal terminal fragment



**Fig. 4.** Structure predictions for 1HdeA, PDB code 1bg8 (CASP3). (A) A superposition of the best  $Q$ -score structure from the AM/W-1 potential (blue) and the native structure (red) is shown. Spheres indicate charged residue  $C\alpha$  atoms. (B) The superposition of the best  $Q$ -score structure from the AM/C potential (blue) and the native structure (red) is shown. (C) Free energy vs.  $Q$  as computed with a histogramming technique. (D) Annealing trajectories of individual fragment  $Q$  scores, large N-terminal fragment containing residues 1–61, and small C-terminal domain containing residues 62–76 are shown as a function of the instantaneous temperature through the run. (E) Annealing trajectories of interfragment  $Q$  scores are indicated. (F) Annealing trajectories of interfragment  $Q$  scores partitioned into the first-well and second-well contributions are shown.

consisting of residues 62–76. The annealing trajectories for the corresponding fragment  $Q$  scores (Fig. 4D) indicate that AM/C and AM/W-1 produce qualitatively similar fragment structures down to  $T = 0.8$ , at which temperature the larger fragment experiences a jump in the nativeness for the AM/W-1 potential. This event is immediately preceded by a jump in the interfragment  $Q$  value (Fig. 4E), suggesting that native interfragment interface formation nucleates the folding of the larger fragment. We have additionally partitioned the interfragment  $Q$  into first-well and second-well contributions (Fig. 4F). This analysis shows the major improvement in the interface recognition comes from the second-well interactions. Because protein-mediated second-well interactions are greatly diminished in AM/W-1 (see Fig. 1C) and we notice the charged nature of the C-terminal fragment (Fig. 4A), we see that it is AM/W water-mediated interactions that greatly facilitate correct packing of secondary structure elements in HdeA.

When the HdeA crystal structure was originally published, Yang *et al.* (21) could not find any sequence or structural similarity to any other known protein. Its functional role was also unknown (21). In a subsequent study (27), it was demonstrated that HdeA provides acid resistance in bacterial pathogens (HdeA is stable under extreme acidic conditions). It was suggested that in neutral pH HdeA forms a dimer (the dimer interface is formed mainly by hydrophobic residues), which dissociates to become an active monomer when pH is significantly lowered (27), the exact mechanism of dissociation being



**Fig. 5.** Structure predictions for CASP5 target protein T129a (PDB code 1IZM, structural information not yet officially released at the time of writing). (A) The distance plot for AM/W-1 best  $Q$ -score structure (blue, upper triangle) and the native structure (red, lower triangle) is shown. (B) The distance plot for AM/W-1 structure with the best sum of individual domain  $Q$  scores (blue, upper triangle) and the native structure (red, lower triangle) is shown. (C) Annealing trajectories of individual domain  $Q$  scores, N-terminal domain containing residues 1–75, and C-terminal domain containing residues 76–170 are indicated. (D) Annealing trajectories of interdomain  $Q$  scores are shown. (E) Annealing trajectories of interdomain first-well  $Q$  scores are plotted. (F) Annealing trajectories of interdomain second-well  $Q$  scores are plotted.

unclear. Gajiwala and Burley (27) hypothesized that perhaps pH-induced conformational changes of unknown nature lead to dissociation. In light of our analysis, it indeed seems plausible that a change in the protonation state of terminal fragments would lead to large structural rearrangement, perhaps causing the dimer dissociation.

The final test protein for analysis is a two-domain CASP5 target, T129a [PDB code 1IZM; the structure has not been released at the time of writing (A. Galkin, E. Sarikaya, C. Lehmann, A. Howard, and O. Herzberg, personal communication)]. The distance plot comparing the best  $Q$  (0.36) predicted structure and the crystal structure (Fig. 5A) shows that the major features of the protein fold are well captured (global RMSD was found at 8.7 Å). Interestingly enough, the same snapshot also has the best interdomain  $Q$  score for the same trajectory. However, there exist snapshots with somewhat better individual domain structures that are docked incorrectly (Fig. 5B). The individual domain II  $Q$  scores between the best AM/C and AM/W-1 trajectories are of similar quality, but the AM/W-1 potential produces more native-like structures for domain I (Fig. 5C). As in the case of HdeA, the interdomain  $Q$  (Fig. 5D) shows the most improvement for AM/W-1 compared with AM/C. Partitioning the interdomain  $Q$  into first- and second-well contributions (Fig. 5E and F), again leads to the conclusion that water-mediated interactions enhance native-like packing of supersecondary structure elements.

At a coarse-grained level, the interplay between direct contact interactions and longer-range water-mediated interactions, both guiding the folding process, suggests some new protein physics.

Direct contact potentials are crudely equivalent to surface tension between the protein and its solvent environment, whereas longer-range water-mediated interactions depend on the curvature of the protein–water interface. The complex solvation physics of polar and charged species in the presence of counterions shapes the curvature landscape. Our results imply that, at least in the cases studied, evolution has tuned both the surface tension and curvature contributions to be consistent with the principle of minimal frustration (28).

## Conclusions

In summary, specific water-mediated interactions are a universal feature of biomolecular recognition, both in folding of monomers and binding of many dimers. We have shown that the inclusion of long-range water-mediated interactions, through a nonpairwise-additive switching potential, in structure prediction Hamiltonians leads to systematically improved predictions for protein structures. Detailed analysis of annealing trajectories for the model reveals explicitly that water-mediated interactions indeed help to correctly assemble supersecondary structure elements into the global native fold. We hope that the water model presented in this article will also help advance the important ongoing efforts toward building an accurate coarse-grained representation of water for self-assembly of both biological and nonbiological systems.

## Appendix: Computational Details

**The AM/C Hamiltonian.** The AM/C Hamiltonian has been discussed at great length in the literature (14–16, 29–32). The Hamiltonian,  $H_{AM/C} = H_{bb} + H_{AM} + H_{contact}$ , consists of a general polymer physics-based backbone potential (see refs. 15, 16, and 29 for details), an AM term defining an energy funnel for residues separated by <12 residues (15, 16, 31), and a contact term that describes tertiary interactions. The contact Hamiltonian,  $H_{contact}$ , has three wells covering the 4.5- to 8.5-Å, 8.5- to 10.0-Å, and 10.0- to 15.0-Å intervals. *Supporting Text* provides additional details about the AM/C Hamiltonian.

**The AM/W Hamiltonian.** The AM/W Hamiltonian is a modification of the AM/C Hamiltonian, where the tertiary contact part of AM/C Hamiltonian is replaced by a potential based on water-mediated interactions,  $H_{AM/W} = H_{bb} + H_{AM} + H_{Rg} + H_{contact} + H_{water} + H_{burial}$ , where  $H_{bb}$  and  $H_{AM}$  are the same as in the AM/C potential,  $H_{Rg}$  is a quadratic potential that helps to collapse the chains ( $H_{Rg} = C * [R_g(\{\mathbf{r}\}) - R_g(N)]^2$ , based on work from ref. 33),  $H_{contact}$  keeps the same functional form as in AM/C, but it contains only a single, direct contact, defined between 4.5 and 6.5 Å,  $H_{water}$  is a nonpairwise additive second-well switching potential (defined below), and  $H_{burial}$  is a many-body potential indicating the burial preferences for each amino acid (defined below). The water-mediated second-well potential is,  $H_{water} = -1/2 \sum_{i,j} \Theta_{ij}^{II} (\sigma_{ij}^{wat} \gamma_{ij}^{wat} + \sigma_{ij}^{prot} \gamma_{ij}^{prot})$ , where switching functions  $\sigma_{ij}^{wat} = H(\rho_i - \rho_{trsh})H(\rho_j - \rho_{trsh})$  and  $\sigma_{ij}^{prot} = 1 - \sigma_{ij}^{wat}$  are used, that depend on local density environment of residues  $i$  and  $j$  ( $\rho_i = \sum_k \theta_{ik}^I$ ,  $\theta_{ij}^{III} = 1/4(1 + \tanh(\kappa(r_{ij} - r_{min}^{III}))) (1 + \tanh(\kappa(r_{max}^{III} - r_{ij})))$ ), and  $H(\rho_i - \rho_{trsh}) = 1/2(1 - \tanh(\kappa(\rho_i - \rho_{trsh})))$ . In these expressions  $r_{ij}$  is the distance between residues  $i$  and  $j$ ,  $r_{min}$  and  $r_{max}$  indicate the endpoints of corresponding wells (4.5–6.5 Å for the first well, 6.5–9.5 Å for the second well), and  $\kappa$  is a parameter that describes the sharpness of the switching tanh functions ( $\kappa$  was set to 5.0). The  $\sigma$  switching functions are constructed so that when the local density  $\rho$  for each residue increases beyond a threshold value of  $\rho_{trsh}$  [chosen to be 2.6 from a structural survey of the monomer database (34), see below], the  $\sigma^{wat}$  switches smoothly from 1 to 0, whereas  $\sigma^{prot}$  switches from 0 to 1.

The burial profile term,  $H_{burial}$ , is a many-body local density based on three-well potential, which indicates amino acid pref-

