

Five hierarchical levels of sequence-structure correlation in proteins

Christopher Bystroff
Biology Department
Rensselaer Polytechnic Institute
Troy, 12180
518-276-3185., 518-276-2162.

bystrc@rpi.edu

Keywords

protein folding, hidden Markov models, motifs, pathways, contact maps, review

Abstract

A review of recent work toward modeling the protein folding pathway using a bioinformatics approach is presented. Five hierarchical statistical models have been developed that test and characterize the nucleation-condensation theory of folding for globular domains. The models are built on the principle that a non-redundant data set of structures obeys Boltzmann statistics, populating the data proportional to the free energy. A recurrent theme in the data implies an energetic selection pressure. Care must be taken to count occurrences only after accounting for redundancy at each hierarchical level. The combined models therefore represent a hierarchical pathway of energy-selected steps in folding, with the results of each model feeding into the next. The models represent five levels of structural complexity: (1) short motifs, (2) extended motifs, (3) non-local pairs of motifs, (4) three dimensional arrangements of multiple motifs, and (5) global structural homology. The models that attempt to capture the sequence-structure correlations at each level are (1) the I-sites Library of local structure motifs, (2) HMMSTR for extended motifs, (3) HMMSTR-CM for pairwise interactions between motifs, (4) SCALI for spatial arrangements of motifs, and (5) SVM-HMMSTR for global structural homology. The parallels between the statistical models and the theoretical models for folding pathways are discussed.

Access to the data used and algorithms presented in this paper are available at <http://www.bioinfo.rpi.edu/~bystrc/> or by request to bystrc@rpi.edu. HMMSTR predictions may be obtained from this web site: <http://www.bioinfo.rpi.edu/~bystrc/hmmstr/server.html>

Introduction

Proteins fold through a hierarchical accumulation of order, from short-ranged to long-ranged -- local to global. The folding pathway is somehow encoded in protein sequences. Recurrent patterns in the database of known proteins tell the story of the folding pathway and of its evolutionary history, but these are two entirely different stories. Evolution takes place on the time scale of millions of years, while protein folding happens in milliseconds. Most methods for predicting protein structure implicitly model the evolutionary process. For

example, a typical algorithm would compare two sequences and calculate a score based on the number and type of sequence differences. But we are interested in the folding process, not the evolutionary process. How can we find sequence-structure correlations in proteins so that they tell us about the folding process and not the evolutionary process?

Through a series of statistical models, we have constructed a picture of the protein folding pathway that agrees with theoretical models. Each statistical model measures the degree of a specific type of sequence-structure correlation in the database of known structures. From the perspective of statistical thermodynamics, a database correlation implies that the observed sequence pattern and structural conformation are associated energetically. A strong correlation represents a strong energetic interaction. Sequence structure correlations may therefore represent energetically stable states, such as intermediates along the protein folding pathway.

A folding pathway may be viewed as starting from folding initiation sites, which are local pieces of the chain that have a strong preference to fold into a certain structure. The chain then collapses locally, around the initiation sites – a process called “propagation.” Pairs of collapsed structures on the chain may then “condense,” or join together, subject to energetic considerations. Finally, in the last stages of the folding pathway, topological constraints predominate in dictating the packing of the preformed units along the chain (Plaxco et al., 2000; Riddle et al., 1999), since at this point in folding not all pieces of the chain can reach each other.

This is a working model for the folding pathway. Based on this model, Table 1 lists five hierarchical levels of sequence-structure correlations that should exist in all globular proteins. Each level depends on the one above it.

Table 1. Five stages of the folding pathway

Folding stage	Type of structure	Model(s)
(1) initiation	local	I-sites ¹
(2) propagation	extended local	HMMSTR ²
(3) condensation	non-local pairwise	HMMSTR-CM ³
(4) packing	non-local multibody	SCALI-HMM ⁴
(5) final	global	SVM-HMMSTR, Pfam, etc.

Global sequence-structure correlations, as detected by global sequence alignments or profile hidden Markov models, such as Pfam (Sonnhammer et al., 1998), SUPERFAMILY (Gough and Chothia, 2002), and SAM (Karplus et al., 1998), tell only the story of the last step in folding. To find the parts of the chain that define the earlier steps in folding, we have developed the models listed in the third column. A library of motifs called I-sites (Bystroff and Baker, 1997; 1998) consists of sequence-structure motifs that occur frequently in the database and which are thought to be initiation sites for protein folding. The I-sites library was extended and generalized by HMMSTR (Bystroff and Shao, 2002; Bystroff et al., 2000), a hidden Markov model of local sequence-structure correlation, to model the propagation of protein folding along the chain. HMMSTR-CM (Shao and Bystroff, 2003), a contact map prediction method using HMMSTR, goes one step further to predict non-local inter-residue contacts. Finally, SCALI (Yuan and Bystroff, 2004), a hidden Markov model for protein structure core alignment, models protein folding at the multibody-level. In this paper, we consider how these models work together to build a picture of the folding process.

Protein Folding Pathways, a Brief History

The early work of Levinthal and Anfinsen established that a protein chain folds spontaneously and reproducibly to a unique three-dimensional structure when placed in aqueous solution. Levinthal proved beyond the shadow of

¹ I-sites: initiation sites, a library of short sequence-structure motifs.

² HMMSTR: a hidden Markov model for local protein sequence-structure correlations.

³ HMMSTR-CM: a contact map prediction method using HMMSTR.

⁴ SCALI-HMM: a hidden Markov model for protein structure core alignment.

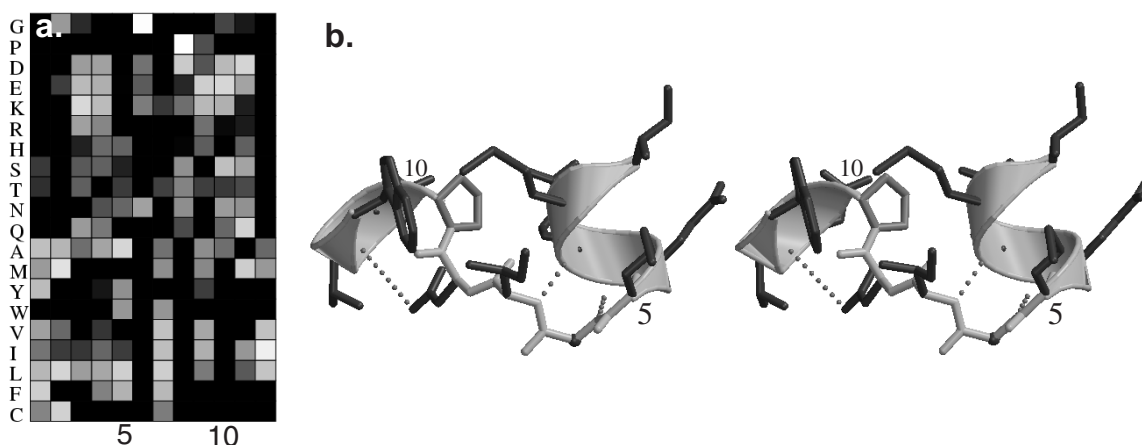


Fig. 1. a. I-sites profile for alpha-alpha corner motif. Boxes are shaded lighter in proportion to the log-likelihood ratio of each amino (Y axis) acid at each position (X axis) relative to the start of the motif. **b.** Stereo image of the alpha-alpha corner motif showing conserved H-bonds and sidechains interactions.

a doubt that the folding process cannot occur by random diffusion. Anfinsen proposed that proteins must form intermediate structures in a time-ordered sequence of events, or "pathway" (Anfinsen and Scheraga, 1975). The nature of the pathways, specifically whether they are restricted to partially native states or whether they might include non-specific interactions, such as an early collapse driven by the hydrophobic effect, was left unanswered.

Over the years, the theoretical models for folding have converged (Baldwin, 1995; Colon and Roder, 1996; Oliveberg et al., 1998; Pande et al., 1998) due to a better understanding of the structure of the "unfolded state" (Dyson and Wright, 1996; Gillespie and Shortle, 1997; Mok et al., 1999) and to a more detailed description of kinetic and equilibrium folding intermediates (Eaton et al., 1996; Gulotta et al., 2001; Houry et al., 1996). An image of the transition state of folding can now be mapped out by point mutations, or "phi-value analysis" (Fersht et al., 1992; Garbuzynskiy et al., 2004; Grantcharova et al., 2000; Gromiha and Selvaraj, 2002; Heidary and Jennings, 2002; Mateu et al., 1999; Nolting et al., 1997). The "folding funnel" model (Chan et al., 1995; Onuchic et al., 1997) has been reconciled with the "hydrophobic collapse" and "nucleation-condensation" models (Nolting and Andert, 2000) by envisioning a distorted, funicular energy landscape (Laurents and Baldwin, 1998) and a "minimally frustrated" pathway (Nymeyer et al., 2000; Shoemaker and Wolynes, 1999), where the rate limiting step is a counter-entropic search for the hole in the funnel ((Zwanzig, 1997). As such, the important role of the protein topology, especially as measured by the contact order, in determining the rate of folding is understood (Ivankov et al., 2003; Miller et al., 2002; Plaxco et al., 1998). Our set of statistical models is consistent with current thinking about folding pathways.

Systems and Methods

Knowledge-based methods for protein structure prediction assume that the frequency of an observed property in the database is a measure of its free energy, provided the database has been properly corrected for redundancy and over-counting. For example, the knowledge-based free energy of a contact between a glycine and an alanine is found by counting the frequency of finding those two amino acids in contact over a database of known protein structures. In the four models presented here, we have derived Bayesian conditional probabilities for local sequence motifs (I-sites ((Bystroff and Baker, 1998)), for sequential strings of multiple motifs (HMMSTR (Bystroff et al., 2000)), for non-local pairwise contacts between motifs (HMMSTR-CM (Shao and Bystroff, 2003)), and for non-sequential three-dimensional packing arrangements of multiple motifs (SCALI (Yuan and Bystroff, 2004)). The hierarchy of models can be roughly described as "local to global", mirroring the

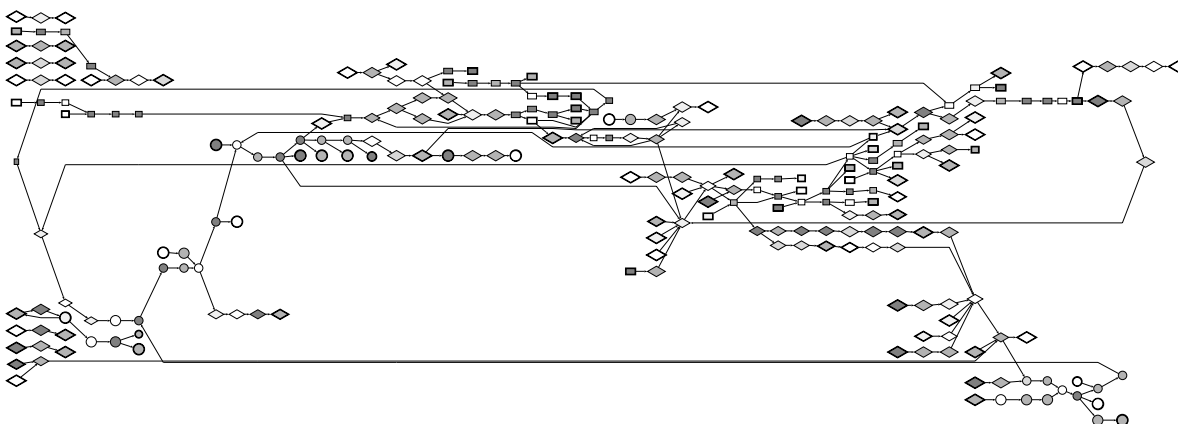


Fig. 2. HMMSTR model “R” represented as a directed graph. The symbol shape represents the secondary structure type; circles: helix; rectangles: beta sheet; diamonds: other motifs. Shading represents the amino acid preference; dark grey: non-polar; grey: polar; light-grey: proline; lightest grey: glycine; white: no preference. Only high-probability transitions are shown.

nucleation/condensation mechanism for protein folding. For each model the prediction results are an ensemble of conformational states. "Local " in this context means nearby along the chain.

Each of these models has been described elsewhere. Here we will review each model and discuss how they are related, and how the structure of the models may reflect the structure of the physical process that they are intended to represent.

The I-sites Library of Folding Initiation Site Motifs

The role of local structure motifs with regard to the initiation of folding has been discussed by Baldwin, Rooman and others (Baldwin and Rose, 1999; Efimov, 1993; Rooman et al., 1990). Recurrent local structure motifs might exist because they fold into a specific structure independent of their context, and since the structure is small and local, the folding is fast. I-sites is a library of 262 local sequence/structure motifs. A motif is expressed as a position-specific scoring matrix, a structural model, and a “confidence curve” which maps the sequence score to a probability or “confidence.” Recurrent sequence patterns of various lengths were found by first exhaustively clustering short segments of sequence families in a non-redundant database of known structures (Bystroff et al., 1996; Han and Baker, 1995; 1996; Han et al., 1997), then optimizing the sequence structure correlation using reinforcement learning (Bystroff and Baker, 1998). I-sites motifs have been used in blind prediction experiments (Bystroff and Baker, 1997; Bystroff and Shao, 2002) and have inspired several experimental studies (Jacchieri, 2000; Mendes et al., 2002; Northey et al., 2002; Skolnick and Kolinski, 2002; Steward and Thornton, 2002).

The confidence of an I-sites prediction is defined as the probability of the prediction being correct, given the sequence score. A score-to-confidence mapping was found by making predictions on a large test set of proteins that were not used to build the model. Only about one-third of all residues in all proteins are found in high-confidence (>70%) I-sites motif regions. But nearly all residues in all proteins (98%) belong to one or more of the I-sites motif structures, although many are difficult to predict.

I-sites motifs include alpha helix, helix caps, beta strands, beta hairpins and other loop structures. **Fig. 1** shows one of the I-sites motifs, the alpha-alpha corner, and the sequence pattern that predicts it. Peptide sequences that match the I-sites motifs have been shown to be structurally stable in isolation in both NMR studies (Blanco et al., 1994; Munoz et al., 1995; Viguera and Serrano, 1995; Yi et al., 1998) and in molecular dynamics simulations (Bystroff and Garde, 2003; Gnanakaran and Garcia, 2002; Krueger and Kollman, 2001). Mutations in high-

confidence I-sites motif regions are found to have dramatic effects on folding (Mok et al., 2001; Northey et al., 2002). These experiments are consistent with I-sites being early folding intermediates, or initiation sites of folding.

HMMSTR: A model for propagation

Folding initiation sites are marginally stable, but they are stable enough to provide a starting place for the propagation of structure up and down the chain. The sequence dependence of these extensions can be found by

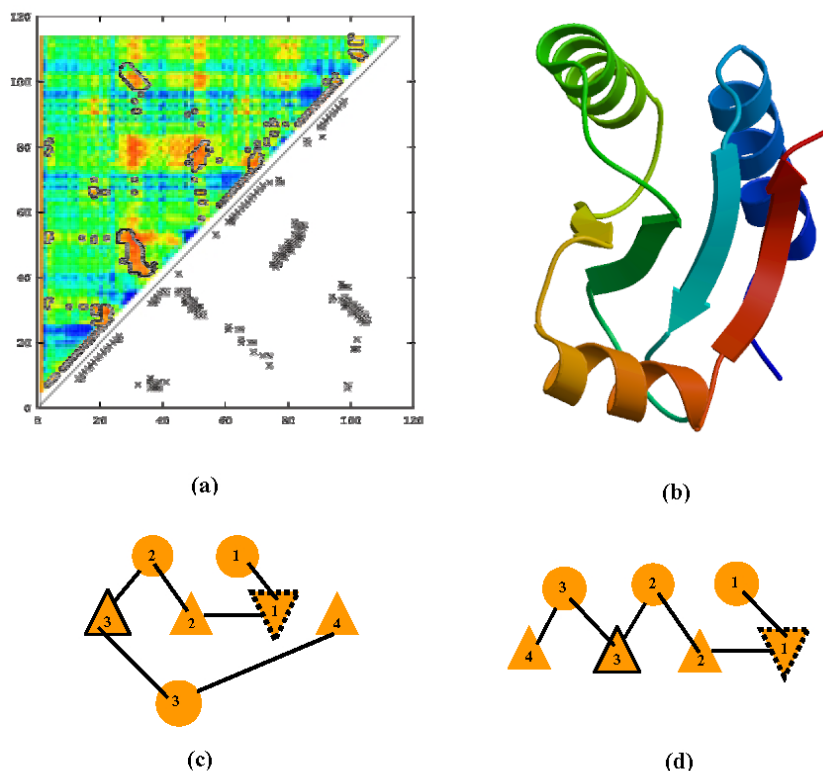


Fig. 3. HMMSTR-CM prediction of a CASP5 target.

The structure of hypothetical protein HI0073 from *H. influenzae* was successfully predicted using the HMMSTR-CM method. HI0073 (PDB code 1JOG, true structure shown in (b)) has 116 residues arranged in a three-layer all-parallel α/β sandwich. The contact potential map (a) shows that most of the true contacts were assigned favorable (red) contact potentials. However, there are also favorable regions that are non-contacts. After selecting a nucleation site, $\beta_2\alpha_2\beta_3$, contacts were assigned or erased in a 4 step pathway, as follows:

- (1) Parallel β contacts were assigned for β_2 to β_3 .
- (2) Anti-parallel β contacts were assigned for β_1 and β_2 . All other β contacts to β_2 were erased (Rule 3).
- (3) There were two choices for a right-handed crossover from β_3 to β_4 , as shown in (c) and (d). Since β_1 was more hydrophobic than β_3 , we paired β_1 and β_4 . All other β contacts to β_1 were erased, and contacts between α_2 and α_3 were erased (Rules 8, 10).
- (4) α_1 was placed on the opposite side of the sheet from α_3 , since α_3 extends across the sheet. Contacts were assigned between α_1 and α_2 and erased between α_1 and α_3 (Rule 9).

(c) The completed TOPS diagram and contact map (outlines) match the true structure. The contact map prediction has 42% contact coverage and 29% accuracy, or if we count near misses (± 1 residue), then the coverage is 75% and the accuracy is 57%. (d) The wrong choice at step (3) would give this structure.

labeling all of the I-sites motifs in the database and drawing connections between them wherever they occur adjacent to each other in the sequence. I-sites motifs that are adjacent to each other “extend” each other. For example, the amphipathic alpha helix motif extends the alpha-alpha corner motif (Fig. 1).

All adjacencies of I-sites motifs in known structures were found, counted, and the motif-motif transition probabilities were condensed into a single, non-linear hidden Markov model (HMM) called HMMSTR ("hamster") (Bystrhoff et al., 2000). The sequence preferences and the inter-motif transitions were trained on a non-redundant database of protein structures. HMMSTR models the ways that local structure can be arranged along the sequence, modeling the way an initiation site motif is likely to affect the conformation of residues that are adjacent to it in the sequence. **Fig. 2** shows the highly branched and cyclic state connectivity. Note the region containing a cycle of helix states, representing the well-known heptad repeat motif of the amphipathic alpha-helix structure.

Each state in HMMSTR represents the structure (as backbone angles) of one residue. An unbroken string of states represents a local structure motif. A branching chain of states represents two or more alternative adjacent motif structures. The result of a HMMSTR prediction is a matrix of Markov state probabilities. We may use the model to sample from this distribution, or we can choose a single structure prediction for each position by a voting procedure as described previously (Bystrhoff et al., 2000). HMMSTR predictions average 60% accuracy in predicting 8-residue fragments with RMSD < 1.4Å. HMMSTR has been used for local and secondary structure prediction (Bystrhoff et al., 2000; Rost, 2001), inter-residue contact prediction (Shao and Bystrhoff, 2003; Zaki et al., 2000), and as the source of a fragment library for Rosetta (Bystrhoff and Shao, 2002) folding simulations.

HMMSTR-CM: Pairwise condensation of motif structures

HMMSTR-CM is a model for pairwise interactions between local structure motifs. Pairwise interactions are represented as a probabilistic contact map. Contact maps are square symmetrical Boolean matrices that represent pairs of residue positions that close in space (within 8Å). A contact map may be projected into three-dimensions if it obeys certain mathematical constraints (Aszodi et al., 1997; Brunger et al., 1986; Crippen, 1988; Selvaraj and Gromiha, 2003; Vendruscolo et al., 1997). Previous contact map prediction methods have used neural nets (Fariselli and Casadio, 1999; Pollastri and Baldi, 2002), correlated mutations (Olmea and Valencia, 1997; Ortiz et al., 1998; Singer et al., 2002), and association rules (Hu et al., 2002; Zaki et al., 2000). Neural net based predictions had an average accuracy of about 21% overall (Fariselli et al., 2001), while higher accuracies were reported for local contacts (Pollastri and Baldi, 2002). Contact potential has been used for many protein structure methods (Jones et al., 1999; Pokarowski et al., 2003).

The first step in predicting a contact map is to assign a probability to each potential contact. The probability in this case is the database-derived likelihood of contact between any two local structure motifs. This implies that the local structure motif forms first, then these sub-structures condense to form larger units. The database statistics give us a free energy of interaction, similar to a binding energy. But each residue is represented not as a single motif but as a probability distribution of motifs. We may envision pairs of flickering local structures, interacting in proportion to their structural content.

The interaction potential between any two motifs is modeled as the statistical interaction potential between two corresponding Markov states. Knowledge-based Markov state “pair potentials” (actually log-likelihood ratios) were summed from the CATH database of protein domain structures (Orengo et al., 1997). Each domain was first preprocessed into HMMSTR Markov state probability distributions using the Forward/Backward algorithm (Rabiner, 1989) to get the position-dependent Markov state probability distribution of states γ (Eq. 1).

$$\gamma(i, q) = P(q | i) \quad \text{Eq. 1}$$

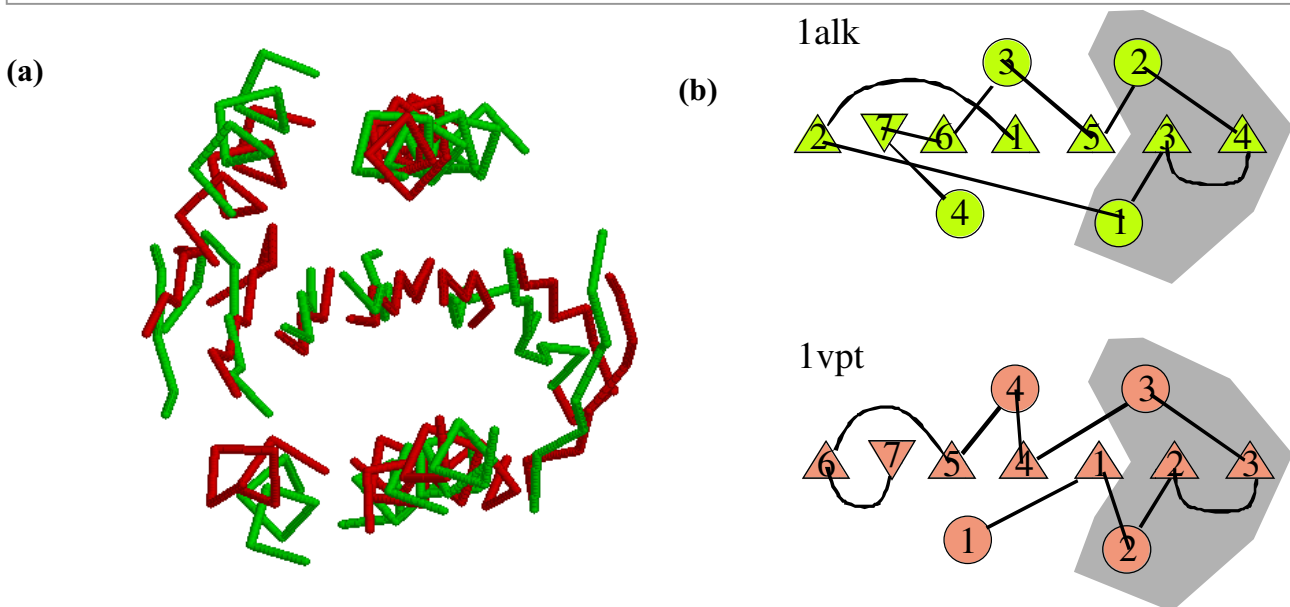
The pairwise contact potential between any two HMMSTR states p and q ($G(p, q, s)$) was calculated as the negative log of the mutual probability of these two states in contacting residues (defined as having C α -C α distance < 8Å) (Eq. 2).

$$G(p, q, s) = -\log \frac{\sum_{CATH} \sum_{i \rightarrow D_{i,i+s} < 8\text{\AA}} \gamma(i, p) \gamma(i + s, q)}{\sum_{CATH} \sum_i \gamma(i, p) \gamma(i + s, q)}, \quad \text{Eq. 2}$$

where $s=|j-i|$ is sequence separation. There is one G value for each pair of HMMSTR's 282 states and each sequence separation, from $s=4$ to 20, a total of 1037153 potential functions. (For $s > 20$, we used G for $s=20$) Using G and a target sequence, we may sum the contact potential map $E(i, j)$, which is a matrix of contact potentials between every residue pair ij in a target sequence. **Fig. 3** shows a contact potential map, E , for a protein that was one of the targets in the CASP5 prediction experiment. In this map we see patterns for super-secondary structure motifs and possible β strand pairings.

Nearly all of the true pairings are given a high score (i.e. low energy) by HMMSTR-CM, but too many high scores are given. Spatial constraints prevent many of the predicted contacts from happening. For example, at most two beta strands may pair with any one beta strand. Other rules enforce the physically possible density of contacts and mutual contacts, and the triangle inequality. Common sense rules such as these were used to extract a self-consistent set of contacts from the high scoring ones. Simple rules were sufficient to extract the correct set of contacts for some but not all of the CASP5 targets (Shao and Bystroff, 2003). However, these rules only approximate the complex topological constraints on the protein chain. A multibody model is needed to capture this higher level of organization.

Fig. 4. Non-sequential alignment of conserved core packing arrangement. Two proteins of very different overall topology, Alkaline phosphatase (1alk, green) and Vp39 from vaccinia virus (1vpt, red) share eleven superimposable secondary structure elements despite having no sequence similarity and different topologies. **(a)** Superimposed secondary structure elements. **(b)** TOPS diagrams showing topological connections. Topology is conserved only in the shaded region.



SCALI: Multibody packing arrangements of local structure motifs

HMMs have been used to predict protein structure at both the local sequence level (Bystrhoff et al., 2000) and the global level (Gough and Chothia, 2002; Karplus et al., 1998; Sonnhammer et al., 1998). But there are interesting recurrent features in proteins that are neither global nor local, specifically, three-dimensional packing arrangements in the hydrophobic core regions. Common secondary structure types have characteristic ways of arranging themselves in globular proteins (Efimov, 1994; Murzin and Finkelstein, 1988; Murzin et al., 1994; Ruczinski et al., 2002). If we take away the connections between the secondary structure units, we find that virtually all ways of packing secondary structures have already been seen in the protein database. When proteins are discovered to have a “new fold”, it is often found to be a permuted version of an “old” fold. If so, then a method for finding non-sequential (i.e. permuted) alignments would be useful for structure prediction.

Recently, we described SCALI, a new algorithm for aligning structures without sequential constraints. SCALI was compared to several structure-based alignment programs, including CE, DALI, and KENOBI, but none of these programs were able to find conserved core packing arrangements consistently when the component parts were not arranged sequentially along the chain. For example, SCALI was able to align the structures 1ALK and 1VPY and found that the positions of its secondary structure elements were superimposable with a RMSD of 4.3Å (Fig. 4) even though they are arranged very differently along the chain.

The topological constraints on folding, given the pairwise energies of interaction, may be modeled by finding all of the common ways that local structure units are arranged in space. That is, we identify the commonly recurring core packing arrangements (i.e. superimposable sets of secondary structure elements). By looking at each example protein that contains a common packing arrangement, we may trace the sequential order of the secondary structure elements by drawing connections wherever they are adjacent in the sequence. The result is a hidden Markov model. In this model a state has a specific location in space, relative to all other states. A pathway through this HMM is a structure prediction.

Recurrent core packing arrangements were found by clustering SCALI non-sequential structure-based alignments. Using a simple greedy algorithm, we found regions of proteins that occurred in multiple SCALI alignments, each protein having a different connectivity. For example, we did an all-against-all pairwise structure comparison for the 61 representative structures of the 3-layer $\alpha/\beta/\alpha$ class, a total of 1830 alignments. 56 out of 61 structures were clustered into four subclasses using a simple greedy algorithm. The structures that clustered together conserved the same core. Fig. 5 shows two of the resulting HMMs.

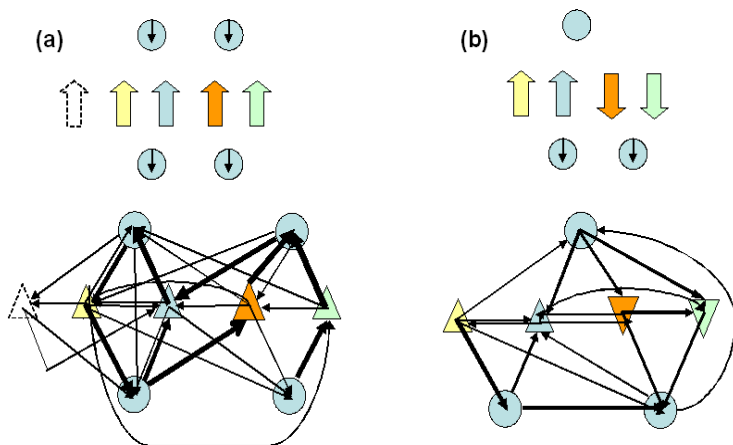


Fig. 5. Diagrammatic hidden Markov models for the two sub-classes of the 3-layer $\alpha\beta\alpha$ class of proteins based on SCALI alignments. In each subclass, the upper panel shows the topology diagram without connectivities for that core structure. (strands are arrows pointing up or down, helices are circles. Circles with arrows are helices pointing down) The lower panel is the hidden Markov model for that core, drawn as a TOPS diagram. Thicker lines indicate more frequent connections. (a) Largest sub-class 37 proteins (b) Next largest, 9 proteins.

For each SCALI-HMM, certain topological connections are observed and others are not, probably reflecting the physical constraints on secondary structure packing. Certain recurrent substructures are found in these models, for example the right-handed parallel $\beta\alpha\beta$ motif, and the shaded region in Fig. 4, a helix-hairpin-helix motif. SCALI-HMMs have been built for several recurrent classes of proteins including the “up-down bundle” alpha proteins, beta “sandwich” proteins, beta “jelly-roll” and 3-layer $\alpha\beta\alpha$ proteins. Refinement and analysis of these conserved core packing arrangements is ongoing.

SVM-HMMSTR: A Support Vector Machine for Fold Recognition Using HMMSTR

Recently a new method has been developed that uses local structure predictions from HMMSTR to characterize the proteins in the SCOP database (Hou et al., 2004). Proteins were assigned a value for each of the HMMSTR states that describes the content of that type of motif structure in the protein. Also, a new Smith-Waterman type alignment program was written that used HMMSTR states instead of amino acids, and the score for every possible alignment between SCOP (Brenner et al., 1996) proteins was calculated. These compositional and alignment-based scores were the points in a hyperdimensional space through which SVM would draw optimal

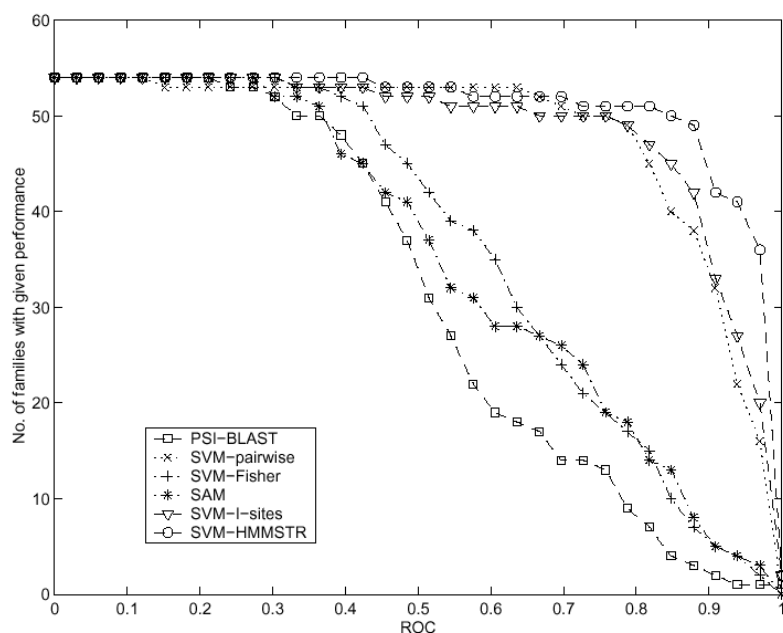


Figure 6. The Receiver-Operator Coefficient (ROC) was calculated for each of the 54 superfamilies in SCOP, for test set predictions using six different methods. Psi-BLAST (Altschul et al., 1997), SVM-pairwise, SAM (Karplus et al., 1998), SVM-Fisher (Liao and Noble, 2003), SVM I-sites (Hou et al., 2003) and SVM-HMMSTR (Hou et al., 2004). Y-axis is the number of superfamilies that had the given ROC score or better.

dividing hyperplanes that partitioned the space into 54 superfamilies.

The resulting SVM was tested on a set of proteins that were not used to train it, and the results were better than those of previous similarity measures, including other SVM-based methods, as shown in Fig. 6. The basis of SVM-HMMSTR’s success was the successful modeling of proteins at a lower level of complexity by HMMSTR. This method is still in development and better results should be possible if the input data is first filtered through one of the intermediate models, HMMSTR-CM or SCALI.

Relevance of Approaches and Results

Using database statistics in a hierarchical way, we have shown that specific recurrent themes exist at different levels of structural generality. It is reasonable to assume that an underlying physical model exists for these statistical observations. A physical model for folding has been proposed that explains the speed with which proteins fold, despite the enormous theoretical size of the conformational space to be searched. They fold fast because there is no combinatorial explosion in the number of conformational states, because at each stage of the folding pathway only a small number of options are available. Table 2 shows the added complexity at each stage of folding, in a qualitative way, as modeled by the five models described here.

Table 2. Recurrence at every level of complexity. No sparse data problem.

Step along the folding pathway:	Model	Added complexity
(1) Initiation	I-sites	~40 motifs explain 98% of proteins
(2) propagation	HMMSTR	1.1 transitions/node
(3) condensation	HMMSTR-CM	only ~5% of state contact pairs occur
(4) molten globule	SCALI	only self-avoiding paths are relevant
(5) native state	SVM-HMMSTR	Estimated total number of folds = 2000

At the stage of initiation, there are only about 40 choices. This is approximately the number of different sequence-structure patterns in I-sites, once overlap is accounted for (Bystroff and Baker, 1998; Bystroff et al., 2000).

At the second stage, propagation, there are even fewer choices. HMMSTR (Bystroff et al., 2000), the model that represents the ways that structure can propagate up and down the chain, is a sparsely branching HMM, reflecting the paucity of choices of different ways of adding structure onto structure, locally.

At the third stage of folding, condensation, the number of possibilities is not the square of the number of local motifs, because some pairs of motifs cannot physically fit together due to differences in their shape and in their surfaces. HMMSTR-CM models this by assigning a near-zero probability for about 95% of the potential pairwise interactions, since these interactions are rarely or never observed (Shao and Bystroff, 2003).

At the fourth stage of folding (we may reach a bit and call this the “molten globule” stage), the number of ways that pairwise interactions can be combined is much less than the number of secondary structure elements “choose two”, because the interaction of any two elements restricts the possibilities for the third and so on, in a process of elimination. Each secondary structure element occupies space, after all, so that its assigned location in space cannot occur twice along the chain. The number of possible tracings through the HMMs produced by SCALI (Yuan and Bystroff, 2004) is fewer than expected because only “self-avoiding” paths are physically possible.

The absence of a combinatorial explosion in the physical folding model may explain the absence of a sparse data problem in our hierarchy of statistical models. Although the database is a fixed size, we still see recurrence as we model increasing larger pieces of protein chain. Even at the global level, there are recurring themes in the database (Gough and Chothia, 2002; Orengo et al., 1994; Russell et al., 1998; Zhang and DeLisi, 1998). For example, the 8-stranded alpha-beta barrel (“TIM barrel”) seems to have independently evolved many times, as has the 7-helix topology called the “globin fold.” In many cases, there is no support for these analogous proteins having a common ancestor. Rather, they are likely to have arisen independently, by “convergent evolution.” If so, if certain topologies have been sampled many times in evolutionary history, then perhaps the total number of ways that a protein chain can fold is not so large as we have previously supposed.

Acknowledgements

This research was partially supported by NSF grant EIA-0229454.

References

- Altschul SF, Madden TL, Schaffer AA, Zhang J, Zhang Z, Miller W & Lipman DJ. (1997). Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res* 25, 3389-402.
- Anfinsen CB & Scheraga HA. (1975). Experimental and theoretical aspects of protein folding. *Adv Protein Chem* 29, 205-300.
- Aszodi A, Munro RE & Taylor WR. (1997). Distance geometry based comparative modelling. *Fold Des* 2, S3-6.
- Baldwin RL. (1995). The nature of protein folding pathways: the classical versus the view. *J Biomol NMR* 5, 103-9.
- Baldwin RL & Rose GD. (1999). Is protein folding hierarchic? I. Local structure and peptide folding. *Trends Biochem Sci* 24, 26-33.
- Blanco FJ, Rivas G & Serrano L. (1994). A short linear peptide that folds into a native stable beta-hairpin in aqueous solution. *Nat Struct Biol* 1, 584-90.
- Brenner SE, Chothia C, Hubbard TJ & Murzin AG. (1996). Understanding protein structure: using scop for fold interpretation. *Methods Enzymol* 266, 635-43.
- Brunger AT, Clore GM, Gronenborn AM & Karplus M. (1986). Three-dimensional structure of proteins determined by molecular dynamics with interproton distance restraints: application to crambin. *Proc Natl Acad Sci U S A* 83, 3801-5.
- Bystroff C & Baker D. (1997). Blind predictions of local protein structure in CASP2 targets using the I-sites library. *Proteins Suppl* 1, 167-71.
- Bystroff C & Baker D. (1998). Prediction of local structure in proteins using a library of sequence-structure motifs. *J Mol Biol* 281, 565-77.
- Bystroff C & Garde S. (2003). Helix propensities of short peptides: Molecular dynamics versus bioinformatics. *Proteins* 50, 552-62.
- Bystroff C & Shao Y. (2002). Fully automated ab initio protein structure prediction using I-SITES, HMMSTR and ROSETTA. *Bioinformatics* 18 Suppl 1, S54-61.
- Bystroff C, Simons KT, Han KF & Baker D. (1996). Local sequence-structure correlations in proteins. *Curr Opin Biotechnol* 7, 417-21.
- Bystroff C, Thorsson V & Baker D. (2000). HMMSTR: A hidden markov model for local sequence-structure correlations in proteins. *Journal of Molecular Biology* 301, 173-90.
- Chan HS, Bromberg S & Dill KA. (1995). Models of cooperativity in protein folding. *Philos Trans R Soc Lond B Biol Sci* 348, 61-70.
- Colon W & Roder H. (1996). Kinetic intermediates in the formation of the cytochrome c molten globule. *Nat Struct Biol* 3, 1019-25.
- Crippen GM, Havel, T.F. (1988). *Distance Geometry and Molecular Conformation*. Chemometrics Series, 15, John Wiley & Sons.
- Dyson HJ & Wright PE. (1996). Insights into protein folding from NMR. *Annu Rev Phys Chem* 47, 369-95.
- Eaton WA, Thompson PA, Chan CK, Hage SJ & Hofrichter J. (1996). Fast events in protein folding. *Structure* 4, 1133-9.
- Efimov AV. (1993). Standard structures in proteins. *Prog Biophys Mol Biol* 60, 201-39.
- Efimov AV. (1994). Favoured structural motifs in globular proteins. *Structure* 2, 999-1002.
- Fariselli P & Casadio R. (1999). A neural network based predictor of residue contacts in proteins. *Protein Eng* 12, 15-21.
- Fariselli P, Olmea O, Valencia A & Casadio R. (2001). Progress in predicting inter-residue contacts of proteins with neural networks and correlated mutations. *Proteins* 45, 157-62.
- Fersht AR, Matouschek A & Serrano L. (1992). The folding of an enzyme. I. Theory of protein engineering analysis of stability and pathway of protein folding. *J Mol Biol* 224, 771-82.
- Garbuzynskiy SO, Finkelstein AV & Galzitskaya OV. (2004). Outlining folding nuclei in globular proteins. *J Mol Biol* 336, 509-25.

- Gillespie JR & Shortle D. (1997). Characterization of long-range structure in the denatured state of staphylococcal nuclease. II. Distance restraints from paramagnetic relaxation and calculation of an ensemble of structures. *J Mol Biol* 268, 170-84.
- Gnanakaran S & Garcia AE. (2002). Folding of a Highly Conserved Diverging Turn Motif from the SH3 Domain. *Biophys J*.
- Gough J & Chothia C. (2002). SUPERFAMILY: HMMs representing all proteins of known structure. SCOP sequence searches, alignments and genome assignments. *Nucleic Acids Res* 30, 268-72.
- Grantcharova VP, Riddle DS & Baker D. (2000). Long-range order in the src SH3 folding transition state. *Proc Natl Acad Sci U S A* 97, 7084-9.
- Gromiha MM & Selvaraj S. (2002). Important amino acid properties for determining the transition state structures of two-state protein mutants. *FEBS Lett* 526, 129-34.
- Gulotta M, Gilmanshin R, Buscher TC, Callender RH & Dyer RB. (2001). Core formation in apomyoglobin: probing the upper reaches of the folding energy landscape. *Biochemistry* 40, 5137-43.
- Han KF & Baker D. (1995). Recurring local sequence motifs in proteins. *J Mol Biol* 251, 176-87.
- Han KF & Baker D. (1996). Global properties of the mapping between local amino acid sequence and local structure in proteins. *Proc Natl Acad Sci U S A* 93, 5814-8.
- Han KF, Bystroff C & Baker D. (1997). Three-dimensional structures and contexts associated with recurrent amino acid sequence patterns. *Protein Sci* 6, 1587-90.
- Heidary DK & Jennings PA. (2002). Three topologically equivalent core residues affect the transition state ensemble in a protein folding reaction. *J Mol Biol* 316, 789-98.
- Hou Y, Hsu W, Lee ML & Bystroff C. (2003). Efficient remote homology detection using local structure. *Bioinformatics* 19, 2294-301.
- Hou Y, Hsu W, Lee ML & Bystroff C. (2004). Remote Homology Detection Using Local Sequence-structure Correlations. *Proteins, Structure, Function and Bioinformatics* in press.
- Houry WA, Rothwarf DM & Scheraga HA. (1996). Circular dichroism evidence for the presence of burst-phase intermediates on the conformational folding pathway of ribonuclease A. *Biochemistry* 35, 10125-33.
- Hu J, Shen X, Shao Y, Bystroff C & Zaki MJ. (2002). *BIOKDD 2002, Edmonton, Canada*.
- Ivankov DN, Garbuzynskiy SO, Alm E, Plaxco KW, Baker D & Finkelstein AV. (2003). Contact order revisited: influence of protein size on the folding rate. *Protein Sci* 12, 2057-62.
- Jacchieri SG. (2000). Mining combinatorial data in protein sequences and structures. *Molecular Diversity* 5, 145-152.
- Jones DT, McGuffin LJ & Bryson K. (1999). GenTHREADER: an efficient and reliable protein fold recognition method for genomic sequences. *J Mol Biol* 287, 797-815.
- Karplus K, Barrett C & Hughey R. (1998). Hidden Markov models for detecting remote protein homologies. *Bioinformatics* 14, 846-56.
- Krueger BP & Kollman PA. (2001). Molecular dynamics simulations of a highly charged peptide from an SH3 domain: possible sequence-function relationship. *Proteins* 45, 4-15.
- Laurents DV & Baldwin RL. (1998). Protein folding: matching theory and experiment. *Biophys J* 75, 428-34.
- Liao L & Noble WS. (2003). Combining pairwise sequence similarity and support vector machines for detecting remote protein evolutionary and structural relationships. *J Comput Biol* 10, 857-68.
- Mateu MG, Sanchez Del Pino MM & Fersht AR. (1999). Mechanism of folding and assembly of a small tetrameric protein domain from tumor suppressor p53. *Nat Struct Biol* 6, 191-8.
- Mendes J, Guerois R & Serrano L. (2002). Energy estimation in protein design. *Current Opinion in Structural Biology* 12, 441-446.
- Miller EJ, Fischer KF & Marqusee S. (2002). Experimental evaluation of topological parameters determining protein-folding rates. *Proc Natl Acad Sci U S A* 99, 10359-63.
- Mok YK, Elisseeva EL, Davidson AR & Forman-Kay JD. (2001). Dramatic stabilization of an SH3 domain by a single substitution: roles of the folded and unfolded states. *J Mol Biol* 307, 913-28.

- Mok YK, Kay CM, Kay LE & Forman-Kay J. (1999). NOE data demonstrating a compact unfolded state for an SH3 domain under non-denaturing conditions. *J Mol Biol* 289, 619-38.
- Munoz V, Blanco FJ & Serrano L. (1995). The hydrophobic-staple motif and a role for loop-residues in alpha-helix stability and protein folding. *Nat Struct Biol* 2, 380-5.
- Murzin AG & Finkelstein AV. (1988). General architecture of the alpha-helical globule. *J Mol Biol* 204, 749-69.
- Murzin AG, Lesk AM & Chothia C. (1994). Principles determining the structure of beta-sheet barrels in proteins. II. The observed structures. *J Mol Biol* 236, 1382-400.
- Nolting B & Andert K. (2000). Mechanism of protein folding. *Proteins* 41, 288-98.
- Nolting B, Golbik R, Neira JL, Soler-Gonzalez AS, Schreiber G & Fersht AR. (1997). The folding pathway of a protein at high resolution from microseconds to seconds. *Proc Natl Acad Sci U S A* 94, 826-30.
- Northey JGB, Maxwell KL & Davidson AR. (2002). Protein folding kinetics beyond the Phi value: Using multiple amino acid substitutions to investigate the structure of the SH3 domain folding transition state. *Journal of Molecular Biology* 320, 389-402.
- Nymeyer H, Socci ND & Onuchic JN. (2000). Landscape approaches for determining the ensemble of folding transition states: success and failure hinge on the degree of frustration. *Proc Natl Acad Sci U S A* 97, 634-9.
- Oliveberg M, Tan YJ, Silow M & Fersht AR. (1998). The changing nature of the protein folding transition state: implications for the shape of the free-energy profile for folding. *J Mol Biol* 277, 933-43.
- Olmea O & Valencia A. (1997). Improving contact predictions by the combination of correlated mutations and other sources of sequence information. *Fold Des* 2, S25-32.
- Onuchic JN, Luthey-Schulten Z & Wolynes PG. (1997). Theory of protein folding: the energy landscape perspective. *Annu Rev Phys Chem* 48, 545-600.
- Orengo CA, Jones DT & Thornton JM. (1994). Protein superfamilies and domain superfolds. *Nature* 372, 631-4.
- Orengo CA, Michie AD, Jones S, Jones DT, Swindells MB & Thornton JM. (1997). CATH--a hierarchic classification of protein domain structures. *Structure* 5, 1093-108.
- Ortiz AR, Kolinski A & Skolnick J. (1998). Fold assembly of small proteins using monte carlo simulations driven by restraints derived from multiple sequence alignments. *J Mol Biol* 277, 419-48.
- Pande VS, Grosberg A, Tanaka T & Rokhsar DS. (1998). Pathways for protein folding: is a new view needed? *Curr Opin Struct Biol* 8, 68-79.
- Plaxco KW, Simons KT & Baker D. (1998). Contact order, transition state placement and the refolding rates of single domain proteins. *J Mol Biol* 277, 985-94.
- Plaxco KW, Simons KT, Ruczinski I & Baker D. (2000). Topology, stability, sequence, and length: defining the determinants of two-state protein folding kinetics. *Biochemistry* 39, 11177-83.
- Pokarowski P, Kolinski A & Skolnick J. (2003). A Minimal Physically Realistic Protein-Like Lattice Model: Designing an Energy Landscape that Ensures All-Or-None Folding to a Unique Native State. *Biophys J* 84, 1518-26.
- Pollastri G & Baldi P. (2002). Prediction of contact maps by GIOHMMs and recurrent neural networks using lateral propagation from all four cardinal corners. *Bioinformatics* 18 Suppl 1, S62-S70.
- Rabiner LR. (1989). A tutorial on Hidden Markov Models and selected applications in speech recognition. *Proc IEEE* 77, 257-286.
- Riddle DS, Grantcharova VP, Santiago JV, Alm E, Ruczinski I & Baker D. (1999). Experiment and theory highlight role of native state topology in SH3 folding. *Nat Struct Biol* 6, 1016-1024.
- Rooman MJ, Rodriguez J & Wodak SJ. (1990). Automatic definition of recurrent local structure motifs in proteins. *J Mol Biol* 213, 327-36.
- Rost B. (2001). Review: protein secondary structure prediction continues to rise. *J Struct Biol* 134, 204-18.
- Ruczinski I, Kooperberg C, Bonneau R & Baker D. (2002). Distributions of beta sheets in proteins with application to structure prediction. *Proteins* 48, 85-97.

- Russell RB, Sasieni PD & Sternberg MJ. (1998). Supersites within superfolds. Binding site similarity in the absence of homology. *J Mol Biol* 282, 903-18.
- Selvaraj S & Gromiha MM. (2003). Role of hydrophobic clusters and long-range contact networks in the folding of (alpha/beta)₈ barrel proteins. *Biophys J* 84, 1919-25.
- Shao Y & Bystroff C. (2003). Predicting interresidue contacts using templates and pathways. *Proteins* 53 Suppl 6, 497-502.
- Shoemaker BA & Wolynes PG. (1999). Exploring structures in protein folding funnels with free energy functionals: the denatured ensemble. *J Mol Biol* 287, 657-74.
- Singer MS, Vriend G & Bywater RP. (2002). Prediction of protein residue contacts with a PDB-derived likelihood matrix. *Protein Eng* 15, 721-5.
- Skolnick J & Kolinski A. (2002). A unified approach to the prediction of protein structure and function. In *Computational Methods for Protein Folding*, Vol. 120, pp. 131-192.
- Sonnhammer EL, Eddy SR, Birney E, Bateman A & Durbin R. (1998). Pfam: multiple sequence alignments and HMM-profiles of protein domains. *Nucleic Acids Res* 26, 320-2.
- Steward RE & Thornton JM. (2002). Prediction of strand pairing in antiparallel and parallel beta- sheets using information theory. *Proteins-Structure Function and Genetics* 48, 178-191.
- Vendruscolo M, Kussell E & Domany E. (1997). Recovery of protein structure from contact maps. *Fold Des* 2, 295-306.
- Viguera AR & Serrano L. (1995). Experimental analysis of the Schellman motif. *J Mol Biol* 251, 150-60.
- Yi Q, Bystroff C, Rajagopal P, Klevit RE & Baker D. (1998). Prediction and structural characterization of an independently folding substructure in the src SH3 domain. *J Mol Biol* 283, 293-300.
- Yuan X & Bystroff C. (2004). Non-sequential Structure-based Alignments Reveal Topology-independent Core Packing Arrangements in Proteins. *Proc Nat Acad Sci* in press.
- Zaki MJ, Shan J & Bystroff C. (2000). *Proceedings IEEE International Symposium on Bio-Informatics and Biomedical Engineering, Arlington, VA, USA*.
- Zhang C & DeLisi C. (1998). Estimating the number of protein folds. *J Mol Biol* 284, 1301-5.
- Zwanzig R. (1997). Two-state models of protein folding kinetics. *Proc Natl Acad Sci U S A* 94, 148-50.