

Bayesian Segmental Models with Multiple Sequence Alignment Profiles for Protein Secondary Structure and Contact Map Prediction

Wei Chu, Zoubin Ghahramani, Alexei Podtelezhnikov and David L. Wild,

Abstract

In this paper, we develop a segmental semi-Markov model (SSMM) for protein secondary structure prediction which incorporates multiple sequence alignment profiles with the purpose of improving the predictive performance. The segmental model is a generalization of the hidden Markov model where a hidden state generates segments of various length and secondary structure type. A novel parameterized model is proposed in the likelihood function that explicitly represents multiple sequence alignment profiles to capture the segmental conformation. Numerical results on benchmark data sets show that incorporating the profiles results in substantial improvements and the generalization performance is promising. By incorporating the information from long range interactions in β -sheets, this model is also capable of carrying out inference on contact maps. This is an important advantage of probabilistic generative models over the traditional discriminative approach to protein secondary structure prediction.

Index Terms

Bayesian segmental semi-Markov models, generative models, protein secondary structure, contact maps, multiple sequence alignment profiles, parametric models.

I. INTRODUCTION

PROTEIN secondary structure prediction remains an important step on the way to full tertiary structure prediction in both fold recognition (threading) and ab-initio methods, as well as providing useful information for the design of site directed mutagenesis experiments to elucidate protein function. A variety of approaches have been proposed to derive the secondary structure of a protein from its amino acid sequence as a classification problem. Beginning with the seminal work of Qian and Sejnowski [1], many

Wei Chu and Zoubin Ghahramani are with Gatsby Computational Neuroscience Unit, University College London. Alexei Podtelezhnikov and David L. Wild are with the Keck Graduate Institute of Applied Life Sciences, Claremont, CA.

of these methods have utilized neural networks. A major improvement in the prediction accuracy of these methods was made by Rost and Sander [2], who proposed a prediction scheme using multi-layered neural networks, known as PHD. The key novel aspect of this work was the use of evolutionary information in the form of profiles derived from multiple sequence alignments instead of training the networks on single sequences. Another kind of alignment profile, position-specific scoring matrices (PSSM) derived by the iterative search procedure PSI-BLAST [3], has been used in neural network prediction methods to achieve further improvements in accuracy [4] [5].

All the above approaches treat the secondary structure prediction problem as a supervised discriminative classification problem. An alternative approach is to treat the problem from the perspective of generative models. One of the first applications of hidden Markov models (HMMs) to the secondary structure prediction problem was described by Delcher et al. [6]. Generalized HMMs with explicit state duration, also known as segmental semi-Markov models, have been widely applied in the field of gene identification [7] [8] [9] [10]. Recently, Schmidler [11] presented an interesting statistical generative model for protein structure prediction, based on a segmental semi-Markov model (SSMM) [12] for sequence-structure relationships. The SSMM is a generalization of hidden Markov models that allows each hidden state to generate a variable length sequence of observations. One advantage of such a probabilistic framework is that it is possible to incorporate varied sources of sequence information using a joint sequence-structure probability distribution based on structural segments. Secondary structure prediction can then be formulated as a general Bayesian inference problem. However, the secondary structure prediction accuracy of the SSMM as described by Schmidler [11] falls short of the best contemporary discriminative methods. Incorporation of multiple alignment profiles into the model might be a plausible way to improve the performance. In this paper, we propose a novel parameterized model as the likelihood function for the SSMM to exploit the information provided by the profiles. Moreover, we incorporate the long range interaction information in β -sheets into the modelling. We describe a Markov Chain Monte Carlo sampling scheme to perform inference in this model, and then demonstrate the capability of the parametric SSMM to carry out inference on β -sheet contact maps in the Bayesian segmental framework. This ability to infer contact maps is one of the advantages of a probabilistic modelling approach over the traditional discriminative approach to protein secondary structure prediction.

The paper is organized as follows. We describe the Bayesian framework of the SSMM in section II. In section III we extend the model to incorporate long range interactions, and point out the capability to infer contact maps. In section IV we discuss the issue of parameter estimation in detail. In section V

we describe a general sampling scheme for prediction. In section VI we present the results of numerical experiments, and conclude in section VII.

II. BAYESIAN MODELLING FRAMEWORK

The key concept underlying our modelling approach is the notion of proteins as collections of local structural fragments, or segments, which may be shared by unrelated proteins - an approach which has gained increasing currency in the protein modelling community in recent years [13], [14]. The modelling framework we adopt is that of the segmental semi-Markov model (SSMM) [12] - a generalization of hidden Markov models that allows each hidden state to generate a variable length sequence of the observations.

The observation sequence includes both a residue sequence and a multiple alignment profile for each protein chain, and is denoted as $O = [O_1, O_2, \dots, O_i, \dots, O_n]$. The associated secondary structure can be fully specified in terms of segment locations and segment types. The segment locations can be identified by the positions of the last residue of these segments, denoted as $e = [e_1, e_2, \dots, e_m]$, where m is the number of segments. We use three secondary structure types. The set of secondary structure types is denoted as $\mathcal{T} = \{H, E, C\}$ where H is used for α -helix, E for β -strand and C for Coil. The sequence of segment types can be denoted as $T = [T_1, T_2, \dots, T_i, \dots, T_m]$ with $T_i \in \mathcal{T} \forall i$. In Figure 1, we present an illustration for the specification of the secondary structure of an observed sequence. Based on a set of protein chains with known secondary structure, we learn an explicit probabilistic model for sequence-structure relationships in the form of a segmental semi-Markov model.

In this approach, the segment types are regarded as a set of discrete variables, known as states. Each of the segment types possesses an underlying generator, which generates a variable-length sequence of observations, i.e. a segment. A schematic depiction of the SSMM is presented in Figure 2 from the perspective of generative models. The variables (m, e, T) describe the secondary structure segmentation of the sequence. In a Bayesian framework the secondary structure prediction problem then consists of computing the posterior probability, $\mathcal{P}(m, e, T|O)$ for an observed sequence O . For this purpose we need to define the prior probability $\mathcal{P}(m, e, T)$ and the likelihood $\mathcal{P}(O|m, e, T)$. This Bayesian framework is described in more detail in the following sections.

A. Multiple alignment profiles

In our model, each of the primary sequences of amino acid residues we are given, denoted as $R = [R_1, R_2, \dots, R_i, \dots, R_n]$ with $R_i \in \mathcal{A}$ where $1 \leq i \leq n$ and \mathcal{A} is the set of 20 amino acids, is associated with a profile derived by multiple sequence alignment [15] or PSI-BLAST [3].

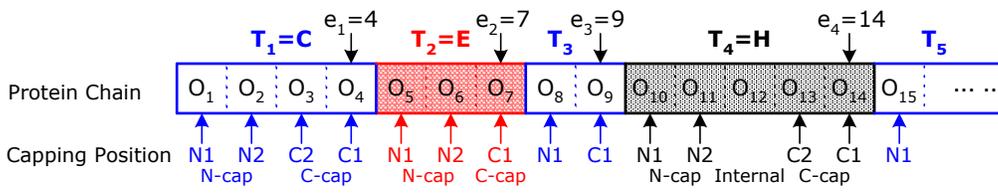


Fig. 1. Presentation of the secondary structure of a protein chain in terms of segments. The square blocks denote our observations of these amino acid residues. The rectangular blocks with solid borders denote the segments. The model represents the segment type $T = [C, E, C, H, \dots]$ and the segmental endpoints $e = [4, 7, 9, 14, \dots]$. Capping positions specify the N- and C-terminal positions within a segment. Here, both the N-capping and C-capping length are fixed at 2, and then $\{N1, N2, \text{Internal}, C2, C1\}$ are used to indicate the capping positions within a segment.

- **Multiple Sequence Alignment Profiles:** for a sequence of amino acid residues, we employ the techniques of pairwise sequence comparison to search a non-redundant protein sequence database for several other sequences which are similar enough at the sequence level to be evolutionarily related. These homologs are then aligned using standard multiple sequence alignment techniques [15]. Ideally, a row of aligned residues occupy similar structural positions and all diverge from a common ancestral residue. By counting the number of occurrences of each amino acid at each location, we obtain an alignment profile. Formally, the alignment profile $M = [M_1, M_2, \dots, M_i, \dots, M_n]$ is a sequence of 20×1 vectors, where M_i contains the occurrence counts for the 20 amino acids at location i .
- **Profiles from PSI-BLAST:** PSI-BLAST [3] is a gapped-version of BLAST that uses an effective scheme for weighting the contribution of different numbers of specific residues at each position in the sequence via a intermediate sequence profile, known as position-specific score matrix (PSSM). Jones [4] explored the idea of using this PSSM as a direct input to a secondary structure prediction method rather than extracting the homologous sequences and then producing an multiple sequence alignment. The PSSM from PSI-BLAST is a matrix with $20 \times n$ elements, where n is the length of the sequence, and each element represents the log-likelihood of the particular amino acid substitution at that position. The profile matrix elements can be mapped to relative occurrence counting by using the standard logistic function: $\frac{1}{1+\exp(-x)}$.

B. Prior Distribution

The prior distribution for the variables describing secondary structure $\mathcal{P}(m, e, T)$ is factored as

$$\mathcal{P}(m, e, T) = \mathcal{P}(m)\mathcal{P}(e, T|m) = \mathcal{P}(m) \prod_{i=1}^m \mathcal{P}(e_i|e_{i-1}, T_i)\mathcal{P}(T_i|T_{i-1}). \quad (1)$$

The segment type depends on the nearest previous neighbour in the sequence through the state transition probabilities $\mathcal{P}(T_i|T_{i-1})$, which are specified by a 3×3 transition matrix. $\mathcal{P}(e_i|e_{i-1}, T_i)$, more exactly

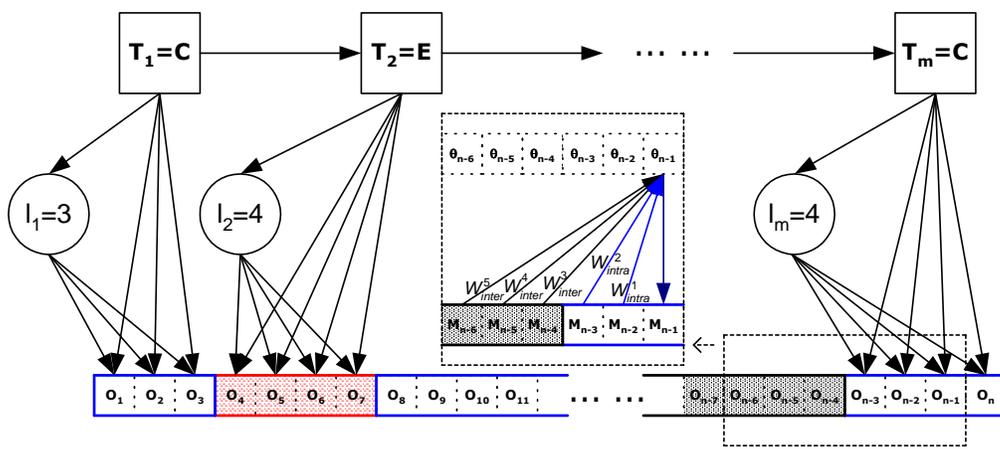


Fig. 2. The segmental semi-Markov model illustrated as generative processes. A variable-length segment of observations associated with random length l_i is generated by the state T_i . The observations within a segment need not be fully correlated, while there might be dependencies between the residues in adjacent segments. The dashed rectangle denotes the dependency window with length 5 for the observation O_{n-1} . In the enlarged dependency window, θ_{n-1} is a vector of latent variables that defines the multinomial distribution in which we observe M_{n-1} , while θ_{n-1} is assumed to be dependent on M_{n-6}, \dots, M_{n-2} .

$\mathcal{P}(l_i|T_i)$ where $l_i = e_i - e_{i-1}$, is the segmental length distribution of the type T_i .¹ Note that the prior on length implicitly defines a prior on the number of segments m for a sequence of a given length. A uniform prior can be assigned for m , i.e. $\mathcal{P}(m) \propto 1$, as this does not have much effect on inference.

C. Likelihood Function

The likelihood is the probability of observing the sequence of alignment profiles given the set of random variables $\{m, e, T\}$. Generally, the probability of the observations can be evaluated as a product of the segments specified by $\{m, e, T\}$:

$$\mathcal{P}(O|m, e, T) = \prod_{i=1}^m \mathcal{P}(S_i|S_{-i}, T_i) \quad (2)$$

where $S_i = O_{[e_{i-1}+1:e_i]} = [O_{e_{i-1}+1}, O_{e_{i-1}+2}, \dots, O_{e_i}]$ is the i -th segment, and $S_{-i} = [S_1, S_2, \dots, S_{i-1}]$. The likelihood function $\mathcal{P}(S_i|S_{-i}, T_i)$ for each segment can be further factorized as a product of the conditional probabilities of individual observations,

$$\mathcal{P}(S_i|S_{-i}, T_i) = \prod_{k=e_{i-1}+1}^{e_i} \mathcal{P}(O_k|O_{[1:k-1]}, T_i) \quad (3)$$

where O_k is the pair of $\{R_k, M_k\}$. R_k is a column vector with 20 elements in which only one element is 1, indicating the amino acid type of the k -th residue, while others are 0, and M_k is the count vector obtained

¹ $e_0 = 0$ is introduced as an auxiliary variable.

from the alignment profile. The likelihood function for each residue should be capable of capturing the core features of the segmental composition in the protein structure.

Schmidler [16] proposed a helical segment model with lookup tables to capture helical capping signals [17] and the hydrophobicity dependency [18] in segmental residues, where the number of free parameters is exponential with the length of dependency window. To overcome this drawback, Chu et al. [19] proposed an extended sigmoid belief network with parameterization for likelihood evaluation. However, these methods were designed to use the primary sequence only, and their secondary structure prediction accuracy still falls short of the best contemporary methods. Incorporation of multiple alignment profiles into the model might be a plausible way to improve the performance.

The plausibility of a protein structure should be evaluated from various perspectives, such as segmental dependency [18], helical capping signals [17] and steric restrictions [20] etc. It is hard to incorporate all the relevant perspectives by using a single model as the likelihood function. Therefore we adopt the concept of a “product of experts” [21] for likelihood evaluation. In the present work, we introduce two experts for the segmental dependency and the helical capping signals respectively. One is a novel parameterized model for segmental dependency that explicitly represents the multiple sequence alignment profile; another is a set of discrete distributions that captures helical capping signals. The conditional probabilities of individual observations can be evaluated in the form of a product of experts, i.e.

$$\mathcal{P}(O_k|O_{[1:k-1]}, T_i) = \mathcal{P}(M_k|M_{[1:k-1]}, T_i)\mathcal{P}(R_k|R_{[1:k-1]}, T_i) \quad (4)$$

More details are given in the following.

1) *An Expert for Segmental Dependency:* The existence of correlated side chain mutations in α -helices has been well studied [18] [22]. These correlations in nonadjacent sequence positions are induced by their spatial proximity in the folded protein molecule and provide an important source of information about the underlying structure. We propose a novel parametric model to capture the segmental dependency by exploiting the information in the multiple sequence alignment profile.

a) *Multinomial Distribution:* We assume that M_k comes from a multinomial distribution with 20 possible outcomes and outcome probabilities θ_k , a 20×1 vector. The outcomes refer to the types of amino acids occurring at the current residue position, while M_k is a 20×1 vector counting the occurrence of these outcomes. Thus, the probability of getting M_k can be evaluated by

$$\mathcal{P}(M_k|\theta_k, T_i) = \frac{\Gamma(\sum_a M_k^a + 1)}{\prod_a \Gamma(M_k^a + 1)} \prod_{a \in \mathcal{A}} (\theta_k^a)^{M_k^a} \quad (5)$$

where \mathcal{A} is the set of 20 amino acids, M_k^a is the element in M_k for the amino acid a , and θ_k^a denotes the probability of the outcome a with the constraint $\sum_a \theta_k^a = 1$. $\Gamma(\cdot)$ is the Gamma function defined as $\Gamma(x) = \int_0^\infty t^{x-1} \exp(-t) dt$.²

b) Dirichlet Prior: As shown in the dependency window of Figure 2, the multinomial distribution at the k -th residue is dependent upon preceding observations within the dependency window, the segment type, and the current capping position within the segment. The underlying causal impact on the current multinomial distribution, where we observed M_k , can be captured by a prior distribution over the latent variables θ_k . A natural choice for the prior distribution over θ_k is a Dirichlet, which has also been used to define priors for protein family HMMs [23]. In our case, this can be explicitly parameterized by weight matrices with *positive elements* as follows:

$$\mathcal{P}(\theta_k | M_{[1:k-1]}, T_i) = \frac{\Gamma(\sum_a \gamma_k^a)}{\prod_a \Gamma(\gamma_k^a)} \prod_{a \in \mathcal{A}} (\theta_k^a)^{\gamma_k^a - 1} \quad (6)$$

where γ_k is a 20×1 vector defined as

$$\gamma_k = W^0 + \sum_{j=1}^{\ell_k} W_{intra}^j M_{k-j} + \sum_{j=\ell_k+1}^{\ell} W_{inter}^j M_{k-j} \quad (7)$$

with ℓ is the length of dependency window,³ $\ell_k = \min(k - e_{i-1} - 1, \ell)$,⁴ and a 20×1 weight vector W^0 is used for local contributions. Weight matrices W_{intra} and W_{inter} of size 20×20 are used to capture both intra-segmental and inter-segmental dependency respectively, where the superscript denotes the residue interval. The constraint $\gamma_k^a > 0 \forall a$ is guaranteed by constraining the weight variables to have positive values. In total we have three sets of weights for $\tau \in \mathcal{T}$ individually. For a segment type τ , we get the set of weight parameters, $\mathbf{W}_\tau = \{W^0, W_{intra}^1, \dots, W_{intra}^\ell, W_{inter}^1, \dots, W_{inter}^\ell\}$.

c) Dirichlet-Multinomial Distribution: The quantity of interest, $\mathcal{P}(M_k | M_{[1:k-1]}, T_i)$ in (3), can be finally obtained as an integral over the space of the latent variables θ_k , which is given by

$$\begin{aligned} \mathcal{P}(M_k | M_{[1:k-1]}, T_i) &= \int_{\theta_k} \mathcal{P}(M_k | \theta_k, T_i) \mathcal{P}(\theta_k | M_{[1:k-1]}, T_i) d\theta_k \\ &= \frac{\Gamma(\sum_a \gamma_k^a) \cdot \prod_a \Gamma(\gamma_k^a + M_k^a)}{\Gamma(\sum_a (\gamma_k^a + M_k^a)) \cdot \prod_a \Gamma(\gamma_k^a)} \cdot \frac{\Gamma(\sum_a M_k^a + 1)}{\prod_a \Gamma(M_k^a + 1)} \end{aligned} \quad (8)$$

where $\Gamma(\cdot)$ denotes the Gamma function, and γ_k is defined as in (7).

²Note that $\Gamma(x+1) = x!$ for positive integers x .

³The window length may be specified individually for segment types.

⁴ $\min(a, b)$ means a if $a \leq b$, otherwise b .

2) *An Expert for Helical Capping Signals*: Helical capping signals [17] refer to the preference for particular amino acids at the N- and C-terminal ends which terminate helices through side chain-backbone hydrogen bonds or hydrophobic interactions. Thus amino acid distributions around the segment ends differ significantly from those of the internal positions, which provide important information for identifying α -helix segments in protein sequences.

The component $\mathcal{P}(R_k|R_{[1:k-1]}, T_i)$ in (3) can be simply modelled as $\mathcal{P}(R_k|T_i)$, which represents the probability of observing R_k at the particular capping position in segments with type T_i . The capping position of each residue within a segment can be determined uniquely (see Figure 1 for an illustration).⁵ The probability distribution of amino acids on a specific capping position c in segments with type τ , denoted as $\mathcal{P}_c(R|\tau)$, can be directly estimated from the training data set, where $c \in \{N1, N2, \dots, \text{Internal}, \dots, C2, C1\}$, $R \in \mathcal{A}$ and $\tau \in \mathcal{T}$.

In summary, the segmental likelihood function we proposed can be explicitly written as

$$\mathcal{P}(O|m, e, T) = \prod_{i=1}^m \mathcal{P}(S_i|T_i, S_{-i}) = \prod_{i=1}^m \prod_{k=e_{i-1}+1}^{e_i} \mathcal{P}(M_k|M_{[1:k-1]}, T_i) \mathcal{P}_c(R_k|T_i) \quad (9)$$

where $\mathcal{P}(M_k|M_{[1:k-1]}, T_i)$ is defined as in (8), and $\mathcal{P}_c(R_k|T_i)$ is the position-specific distribution of capping signals.

Winther and Krogh [24] have demonstrated that optimized potential functions learned from training data can provide very strong restrictions on the spatial arrangement of protein folding. As a very promising direction for future work, the introduction of an additional “steric expert” into our likelihood function could provide global restrictions on secondary structure and fulfill the potential of the Bayesian segmental model for tertiary structure prediction.

D. Posterior Distribution

All inferences about the segmental variables (m, e, T) defining secondary structure are derived from the posterior probability $\mathcal{P}(m, e, T|O)$. Using Bayes’ theorem,

$$\mathcal{P}(m, e, T|O) = \frac{\mathcal{P}(O|m, e, T)\mathcal{P}(m, e, T)}{\mathcal{P}(O)} \quad (10)$$

where $\mathcal{P}(O) = \sum_{\{m, e, T\}} \mathcal{P}(O|m, e, T)\mathcal{P}(m, e, T)$ as the normalizing factor. From the posterior distribution over segmental variables $\mathcal{P}(m, e, T|O)$, we can obtain two different ways of estimating the secondary structure of a given sequence:

⁵Note that we have used two sets of positioning indices for each residue: a sequential number k where $1 \leq k \leq n$, and a capping position cap where $cap \in \{N1, N2, \dots, \text{Internal}, \dots, C2, C1\}$.

- The most probable segmental variables in the posterior distribution: $\arg \max_{m,e,T} \mathcal{P}(m, e, T|O)$, known as the MAP estimate;
- The posterior distribution of the segment type at each residue: $\mathcal{P}(T_{O_i}|O)$ where we denote T_{O_i} as the segment type at the i -th observation. The marginal posterior mode estimate is defined as $\arg \max_T \mathcal{P}(T_{O_i}|O)$.

The Viterbi and forward-backward algorithms for SSMM [25] can be employed for the MAP and marginal posterior mode estimate respectively (refer to Appendix A for a summary).

III. INCORPORATING LONG RANGE INTERACTIONS IN β -SHEETS

We have set up a Bayesian framework to predict the secondary structure. However, the secondary structure might be affected not only by local sequence information, but also by long range interactions between distal regions of the amino acid sequence. This is particularly important in the case of β sheets, which are built up from several interacting regions of β -strands. The strands align so that the NH groups on one strand can form hydrogen bonds with the CO groups on the distal strand and vice versa. The alignment can happen in two ways: either the direction of the polypeptide chain of β -strands is identical, a *parallel* β -sheet, or the strand alignment is in the alternative direction, an *anti-parallel* β -sheet. In Figure 3, we present the two cases for a pair of interacting segments, S_i and S_j with $i < j$. A binary variable is used to indicate alignment direction; $d_{ij} = +1$ for parallel and $d_{ij} = -1$ for anti-parallel. A integer variable a_{ij} is used to indicate the alignment position. The endpoint of S_i , known as e_i , is used as the origin, and then a_{ij} is defined as the shift between e_i and e_j for parallel cases, while for anti-parallel cases it is the shift between e_i and the beginning point of S_j , i.e. $e_{j-1} + 1$.⁶ The challenge for a predictive approach is how to introduce these long range interactions into the model. In this section, we extend the parametric model to incorporate information on long range interactions in β -sheets.

A. Prior Specification for Distal Interactions

A set of random variables is introduced to describe the long range interactions, collected as $\mathcal{I} = \{\{S_j \leftrightarrow S_{j'}, d_{jj'}, a_{jj'}\}_{j=1}^r\}$, where r is the number of interacting pairs and $\{S_j \leftrightarrow S_{j'}, d_{jj'}, a_{jj'}\}$ is a pair of interacting segments together with their alignment information. We can expand the prior probability as $\mathcal{P}(m, e, T, \mathcal{I}) = \mathcal{P}(\mathcal{I}|m, e, T)\mathcal{P}(m, e, T)$, where $\mathcal{P}(m, e, T)$ is defined as in (1) and the conditional

⁶We assume interaction parts to be contiguous, e.g. excluding the case of β -bulges.

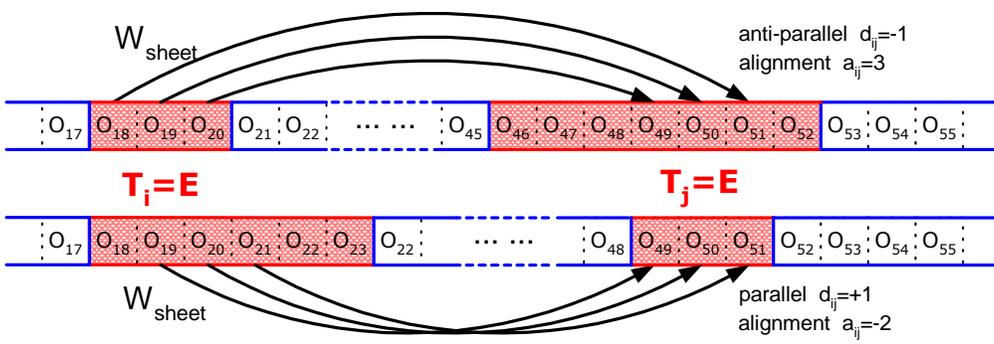


Fig. 3. Anti-parallel (top), and parallel (bottom), pairs of interacting segments, S_i and S_j . d_{ij} is the binary variable for alignment direction, and a_{ij} is the integer variable for alignment position. A weight matrix W_{sheet} is introduced to capture the distal residue interactions.

probability $\mathcal{P}(\mathcal{I}|m, e, T)$ can be further factored as

$$\mathcal{P}(\mathcal{I}|m, e, T) = \mathcal{P}(r|k)\mathcal{P}(\{S_j \leftrightarrow S_{j'}\}_{j=1}^r) \prod_{j=1}^r \mathcal{P}(d_{jj'}|S_j \leftrightarrow S_{j'})\mathcal{P}(a_{jj'}|S_j \leftrightarrow S_{j'}, d_{jj'}) \quad (11)$$

where r is the number of interacting pairs, k is the number of β -strands, and $\{S_j \leftrightarrow S_{j'}\}_{j=1}^r$ denotes a combination for β -strands to form r interacting pairs. Various specifications for these distributions in (11) are applicable provided that they satisfy $\sum_{\mathcal{I}} \mathcal{P}(\mathcal{I}|m, e, T) = 1$. In the present work, we assumed a uniform distribution, $\mathcal{P}(\{S_j \leftrightarrow S_{j'}\}_{j=1}^r) = \frac{1}{c(r, k)}$ if the combination is valid, where $c(r, k)$ is the total number of *valid combinations*,⁷ otherwise $\mathcal{P}(\{S_j \leftrightarrow S_{j'}\}_{j=1}^r) = 0$. $\mathcal{P}(r|k)$, $\mathcal{P}(d_{jj'}|S_j \leftrightarrow S_{j'})$ and $\mathcal{P}(a_{jj'}|S_j \leftrightarrow S_{j'}, d_{jj'})$ are discrete distributions depending on the distance between the two β -strands and their lengths, which were learned from training data by counting the relative occurrence frequencies.

B. Joint Segmental Likelihood

It is straightforward to extend the parametric model (8) to include long range interactions in β -sheets, which can be regarded as an extension of the dependency window to include the distal pairing partners. We introduce another 20×20 weight matrix W_{sheet} to capture the correlation between distal interacting pairs. The segmental likelihood function (3) for the β -strands can be enhanced as

$$\mathcal{P}(S_i|T_i = E, S_{-i}, \mathcal{I}) = \prod_{k=e_{i-1}+1}^{e_i} \frac{\Gamma(\sum_a \tilde{\gamma}_k^a) \prod_a \Gamma(\tilde{\gamma}_k^a + M_k^a)}{\Gamma(\sum_a (\tilde{\gamma}_k^a + M_k^a)) \prod_a \Gamma(\tilde{\gamma}_k^a)} \frac{\Gamma(\sum_a M_k^a + 1)}{\prod_a \Gamma(M_k^a + 1)} \mathcal{P}_c(R_k|T_i = E) \quad (12)$$

with $\tilde{\gamma}_k = \gamma_k + \sum_{\{k^*\}} W_{sheet} M_{k^*}$ where γ_k is defined as in (7) and $\{k^*\}$ denotes the set of interacting residues of O_k that can be determined by \mathcal{I} .

⁷A valid combination requires that each β -strand interacts with at least one and at most two other strands. This constraint comes from the chemical structure of amino acids, i.e. the CO and NH groups.

C. β -Sheet Contact Maps

Contact maps represent the pairwise, inter-residue contacts, as a symmetrical, square, boolean matrix. Pollastri and Baldi [26] have previously applied ensembles of bidirectional recurrent neural network architectures to the prediction of such contact maps. In this section, we describe the capability of this parametric SSMM model to carry out inference on contact maps. This capability is one of the advantages of the probabilistic modelling approach over the traditional discriminative approach (e.g. neural networks) to protein secondary structure prediction. β -sheets are built up from pairs of β -strands with hydrogen bonds, which are prominent features in contact maps. The set of β -sheet interactions is associated with a β -sheet contact map defined by a $n \times n$ matrix \mathcal{C} whose ij -th entry \mathcal{C}^{ij} defined as

$$\mathcal{C}^{ij}(\mathcal{I}) = \begin{cases} 1 & \text{if } O_i \text{ and } O_j \text{ are paired in } \mathcal{I}; \\ 0 & \text{otherwise} \end{cases} \quad (13)$$

We may estimate the marginal predicted \mathcal{C} from the posterior distribution of $\mathcal{P}(m, e, T, \mathcal{I}|O)$, given by

$$\mathcal{P}(\mathcal{C}^{ij} = 1|O) = \sum_{m,e,T,\mathcal{I}} \mathcal{C}^{ij}(\mathcal{I}) \mathcal{P}(m, e, T, \mathcal{I}|O) \quad (14)$$

where the indicator function $\mathcal{C}^{ij}(\mathcal{I})$ is defined as in (13). Using the samples we have collected in the distributions $\mathcal{P}(m, e, T|O)$ and $\mathcal{P}(\mathcal{I}|m, e, T)$ (see Section V and Appendix B for details), (14) can be estimated by

$$\mathcal{P}(\mathcal{C}^{ij} = 1|O) = \sum_{m,e,T} \sum_{\mathcal{I}} \mathcal{C}^{ij}(\mathcal{I}) \mathcal{P}(m, e, T, \mathcal{I}|O) \approx \frac{1}{\mathcal{N}} \sum_{\{m,e,T\}} \sum_{\{\mathcal{I}\}} \mathcal{C}^{ij}(\mathcal{I}) \frac{\mathcal{P}(O|m, e, T, \mathcal{I})}{\sum_{\{\mathcal{I}\}} \mathcal{P}(O|m, e, T, \mathcal{I})} \quad (15)$$

where the samples $\{\mathcal{I}\}$ are collected from $\mathcal{P}(\mathcal{I}|m, e, T)$, and \mathcal{N} samples of $\{m, e, T\}$ are from $\mathcal{P}(m, e, T|O)$.

IV. PARAMETER ESTIMATION

The probabilistic model we describe above has five classes of latent variables, parameters and hyperparameters, which are inferred or specified in different ways:

- Latent variables related to the location and length of secondary structure elements $\{m, e, T\}$
 - number of segments: m
 - the end points of each segment that specify the segment lengths: e
 - secondary structure classes of each segment: T

We infer these latent variables by sampling in the posterior distribution $\mathcal{P}(m, e, T|O)$ (see Section V for details).

- Latent variables related to distal interactions in β -sheets $\{\mathcal{I}\}$
 - number of interacting pairs: r
 - the interacting pairs of β -strands: $\{S_j \leftrightarrow S_{j'}\}$
 - orientation indicators: $\{d_{jj'}\}$
 - the indicators of alignment positions: $\{a_{jj'}\}$

These interacting variables can be sampled in the conditional distribution $\mathcal{P}(\mathcal{I}|m, e, T)$ (see Section V and Appendix B for details).

- Parameters that specify discrete distributions
 - State transition probabilities for $\mathcal{P}(T_i|T_{i-1})$ as defined in (1)⁸
 - Segmental length distributions $\mathcal{P}(e_i|e_{i-1}, T_i)$ as defined in (1)
 - Position-specific distributions of amino acids $\mathcal{P}_c(R|T_i)$ as defined in (9) for capping signals
 - The conditional distribution of the number of interacting pairs $\mathcal{P}(r|k)$ as defined in (11)
 - The conditional distribution of the orientation indicators $\mathcal{P}(d_{jj'}|S_j, S_{j'})$ as defined in (11)⁹
 - The conditional distribution of the alignment position $\mathcal{P}(a_{jj'}|S_j, S_{j'}, d_{jj'})$ as defined in (11)

These parameters specifying discrete distributions can be directly estimated by their relative frequency of occurrence in the training data set.¹⁰ We present the results of state transition probabilities and segmental length distributions, estimated from our training data, in Figure 4 and Figure 5 respectively as an illustration. $\mathcal{P}(\{S_j, S_{j'}\}_{j=1}^r)$ defined as in (11) is uniformly distributed.

- Weight parameters in the likelihood function for segmental dependency (8) and (12) were estimated by penalized maximum likelihood, which is presented with details in Section IV-A below.
- Model parameters
 - N-capping and C-capping length for capping signals.

The capping components result in amino acid distributions at the end-segment positions which differ significantly from the overall distribution. In Table I, we presented the Kullback-Leibler divergence from the amino acid distribution at capping positions to their overall distribution. Based on these divergences, the N-capping and C-capping length were both determined as 4.

- The length of dependency window in (7).

Crooks and Brenner [27] have examined the entropy densities of protein primary and secondary structure sequences, and the local inter-sequence mutual information density. They found that the

⁸The initial state probabilities $\mathcal{P}(T_0)$ can simply be set to be equal.

⁹The distribution is actually only conditional on the distance between the β -strand pair and the segment lengths of the two β -strands.

¹⁰An appropriate prior might be used for smoothing.

TABLE I

KULLBACK-LEIBLER DIVERGENCE FROM THE AMINO ACID DISTRIBUTION AT CAPPING POSITIONS TO THEIR OVERALL DISTRIBUTION. BOLD FACE WAS USED TO INDICATE DIFFERENCE ABOVE THE CUTOFF 0.01. THE DIVERGENCE FROM TWO DISTRIBUTIONS \mathcal{P} AND \mathcal{Q} IS EVALUATED BY $\sum_R \mathcal{Q}(R) \log\left(\frac{\mathcal{Q}(R)}{\mathcal{P}(R)}\right)$. HERE, $R \in \mathcal{A}$, $\mathcal{P}(R)$ IS THE AMINO ACID DISTRIBUTION AT CAPPING POSITIONS, AND $\mathcal{Q}(R)$ IS THE OVERALL SEGMENTAL DISTRIBUTION.

Capping Position	α -helix	β -strand	Coil
N1	0.0290	0.0059	0.0064
N2	0.0535	0.0093	0.0043
N3	0.0136	0.0069	0.0008
N4	0.0478	0.0038	0.0018
N5	0.0030	0.0076	0.0032
\vdots	\vdots	\vdots	\vdots
C5	0.0091	0.0050	0.0037
C4	0.0181	0.0018	0.0025
C3	0.0086	0.0089	0.0019
C2	0.0044	0.0028	0.0015
C1	0.0066	0.0155	0.0046

inter-sequence interactions important to secondary structure prediction are short-ranged. Based on their results, we decided to fix the window length at 5 in the present work.

A. Estimates on Weight Parameters

The weight parameters consist of three sets for different segmental types, i.e. $\{\mathbf{W}_\tau\}$ for $\tau \in \mathcal{T}$. For each segment type τ , there are $|\mathcal{A}|^2\ell$ parameters of W_{intra} 's, $|\mathcal{A}|^2\ell$ parameters of W_{inter} 's and $|\mathcal{A}|$ parameters in the vector W^0 , where the types of amino acid residues $|\mathcal{A}| = 20$ and the length of dependency window $\ell = 5$ in the present work. Thus the total number of weight parameters is 4020. β -strands have an additional $|\mathcal{A}|^2$ parameters in W_{sheet} if the long-range interactions are incorporated.

The maximum a posteriori (MAP) estimate of its associated weights \mathbf{W}_τ can be obtained as

$$\arg \max_{\mathbf{W}_\tau} \mathcal{P}(\{O, m, e, T\} | \mathbf{W}_\tau) \mathcal{P}(\mathbf{W}_\tau) \quad (16)$$

under the condition of positive elements, where $\mathcal{P}(\mathbf{W}_\tau)$ is the prior probability usually specified by $\mathcal{P}(\mathbf{W}_\tau) \propto \exp(-\frac{C}{2} \|\mathbf{W}_\tau\|_2^2)$ with $C \geq 0$, and $\mathcal{P}(\{O, m, e, T\} | \mathbf{W}_\tau)$ is the product of the joint probabilities over all protein chains in training data set. The optimal \mathbf{W}_τ is therefore the minimizer of the negative

logarithm of (16), which can be obtained by

$$\min_{\mathbf{W}_\tau} \mathcal{L}(\mathbf{W}_\tau) = - \sum_{\{O\}} \sum_{\{\tau\}} \ln \mathcal{P}(S_i|S_{-i}, \tau) + \frac{C}{2} \|\mathbf{W}_\tau\|_2^2 \quad (17)$$

subject to $w > 0, \forall w \in \mathbf{W}_\tau$, where $\sum_{\{O\}}$ means the sum over all the protein chains, $\sum_{\{\tau\}}$ denotes the sum over all the segments of type τ , and $\mathcal{P}(S_i|S_{-i}, \tau)$ is defined as in (3). A set of auxiliary variables $\mu = \ln w$ can be introduced to convert the constrained optimization problem into an unconstrained problem. The derivatives of $\mathcal{L}(\mathbf{W}_\tau)$ with respect to μ are given as follows:

$$\frac{\partial \mathcal{L}(\mathbf{W}_\tau)}{\partial \mu} = w \left(\sum_{\{O\}} \sum_{\{\tau\}} \sum_{k=e_{i-1}+1}^{e_i} \psi_k^T \frac{\partial \gamma_k}{\partial w} + Cw \right) \quad (18)$$

where γ_k is defined as in (7), and $\psi_k = \frac{\partial -\ln \mathcal{P}(M_k|M_{[1:k-1]}, \tau)}{\partial \gamma_k}$ is a 20×1 vector whose a -th element is

$$\psi_k^a = \Psi(\gamma_k^a) - \Psi(\gamma_k^a + M_k^a) + \Psi(\sum_a (\gamma_k^a + M_k^a)) - \Psi(\sum_a (\gamma_k^a))$$

where $\Psi(x) = \frac{d}{dx} \ln(\Gamma(x))$ is known as the digamma function. Then standard gradient-based optimization methods are employed to minimize (17).

The optimal value of the regularization factor C in the regularized functional $\mathcal{L}(\mathbf{W}_\tau)$ was determined by standard k-fold cross validation [28] [29]. We carried out 7-fold cross validation as follows. The original training data were randomly partitioned into 7 almost equal folds with each fold having an almost equal percentage of different segments and amino acid residues. Given a particular value of C , one fold was left out as a validation set in turn, the weight parameters were estimated by minimizing $\mathcal{L}(\mathbf{W}_\tau) \forall \tau$ over the protein chains in the other 6 folds, and the resulting model was tested on the left-out fold to obtain the validation error. The average of the validation errors on the 7 left-out folds indicates the predictive performance of the regularization factor C . We tried the set of C values: $C = \{10^{-3}, 10^{-2}, \dots, 10^{+2}\}$, and found the best validation performance was achieved when $C = 0.01$. The optimal weight parameters in the model were finally obtained by optimizing on the whole training data set with the best C value.

It is possible to specify different values of C for the segment types, but it increases the computational cost of cross validation massively. Approximate Bayesian techniques could also be used to further specify different C values on weight matrices individually, while the computational difficulty lies in evaluating the integral over the high-dimensional weight space. This is an interesting and worthwhile issue for further investigation.

V. SAMPLING SCHEME FOR PREDICTION

Without the incorporation of long range interactions, the quantities of the segmentation variables can be inferred exactly by the Viterbi and forward-backward algorithms in the segmental semi-Markov framework

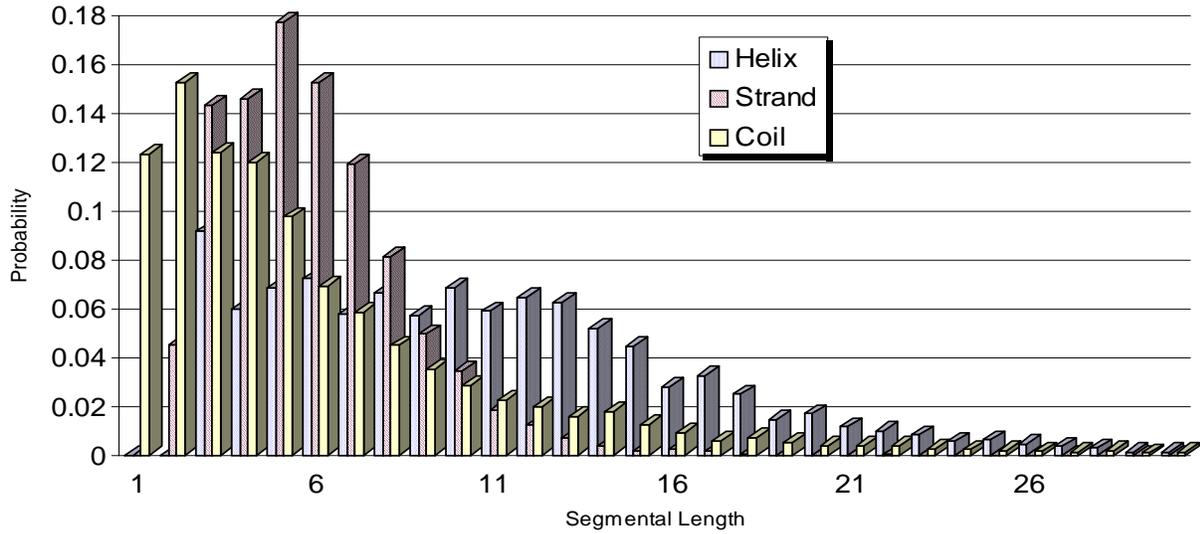


Fig. 4. The distributions of segmental length for the three segment types, $\mathcal{P}(e_i|e_{i-1}, T_i)$ defined as in (1). Note that the three distributions are quite different.

(see Appendix A for the details). Generally, the introduction of long range interactions into the segmental model makes exact calculation of posterior probabilities intractable. Markov Chain Monte Carlo (MCMC) algorithms can be applied here to obtain approximate inference. The latent variables of segmentation $\{m, e, T\}$ are sampled from the posterior distribution $\mathcal{P}(m, e, T|O)$ with MCMC, keeping the weight parameters and the model parameters fixed. In our model, the dimension of the variable vectors e and T is the latent variable m . The dimensionality of the variable space, indexed by m , could be changed in the Markov chain simulation. In this case, the Metropolis-Hasting scheme can be applied with a reversible-jump approach [30], which ensures that jumps between variable spaces of differing dimension are reversible.¹¹

What we are interested in here is the posterior distribution $\mathcal{P}(m, e, T|O)$ which is proportional to the joint distribution $\mathcal{P}(m, e, T, O)$. The joint distribution can be evaluated as

$$\mathcal{P}(m, e, T, O) = \mathcal{P}(m, e, T) \prod_{S_i \notin \mathcal{I}} \mathcal{P}(S_i|S_{-i}, T_i) \sum_{\mathcal{I}} \mathcal{P}(\mathcal{I}|m, e, T) \prod_{S_i \in \mathcal{I}} \mathcal{P}(S_i|S_{-i}, T_i) \quad (19)$$

where $\mathcal{P}(m, e, T)$ is defined as in (1), and only the segments of β -strands are in the interaction set \mathcal{I} . The following set of Metropolis proposals are defined for the construction of a Markov chain on the space of segmentations, denoted as $\mathcal{V} = (m, e, T)$:

¹¹Schmidler [11] attempted to collect samples in the joint posterior distribution $\mathcal{P}(m, e, T, \mathcal{I}|O)$, while the dependency between (m, e, T) and \mathcal{I} makes it complicated to design Metropolis proposals jointly.

Segment Type Transition Probabilities

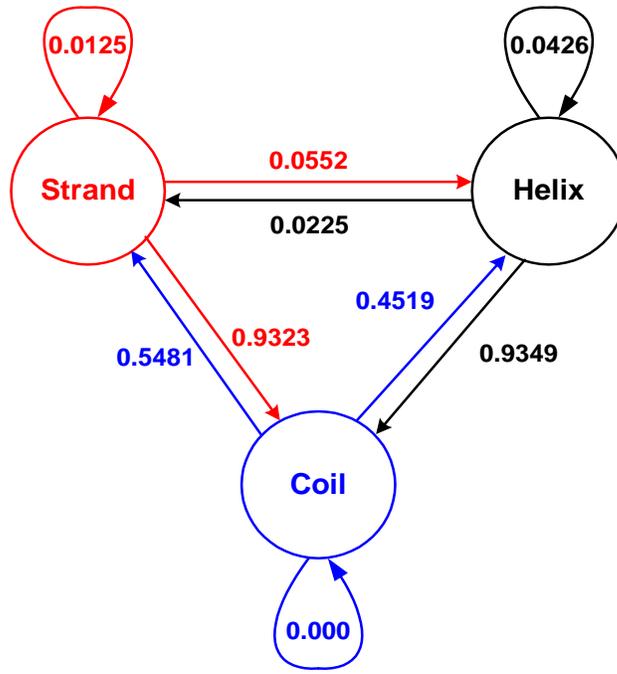


Fig. 5. The segment type transition probabilities, $\mathcal{P}(T_i|T_{i-1})$ defined as in (1). The self-transitions are obtained from the annotations in the training database.

- *Segment split*: propose $\mathcal{V}^* = (m^*, e^*, T^*)$ with $m^* = m + 1$ by splitting segment S_k into two new segments (S_{k^*}, S_{k^*+1}) with $k \sim \text{Uniform}[1 : m]$, $e_{k^*} \sim \text{Uniform}[e_{k-1} + 1 : e_k - 1]$, $e_{k^*+1} = e_k$, $T_{k^*} \sim \text{Uniform}[H, E, L]$, and $T_{k^*+1} \sim \text{Uniform}[H, E, L]$.¹²
- *Segment merge*: propose $\mathcal{V}^* = (m^*, e^*, T^*)$ with $m^* = m - 1$ by merging the two segments S_k and S_{k+1} into one new segment S_{k^*} with $k \sim \text{Uniform}[1 : m - 1]$, $e_{k^*} = e_{k+1}$, and $T_{k^*} \sim \text{Uniform}[H, E, L]$.
- *Type change*: propose $\mathcal{V}^* = (m, e, T^*)$ with $T^* = [T_1, \dots, T_{k-1}, T_k^*, T_{k+1}, \dots, T_m]$ where $T_k^* \sim \text{Uniform}[H, E, L]$.
- *Endpoint change*: propose $\mathcal{V}^* = (m, e^*, T)$ with $e^* = [e_1, \dots, e_{k-1}, e_k^*, e_{k+1}, \dots, e_m]$ where $e_k^* \sim \text{Uniform}[e_{k-1} + 1 : e_{k+1} - 1]$.

The acceptance probability for *Type change* and *Endpoint change* depends on the ratio of likelihood $\frac{\mathcal{P}(\mathcal{V}^*, O)}{\mathcal{P}(\mathcal{V}, O)}$, where the likelihood is defined as in (19). *Segment split* and *Segment merge* jumps between segmentations of different dimension are accepted or rejected according to a reversible-jump Metropolis criteria. According to the requirement of detailed balance, the acceptance probability for a new proposal

¹²Here $\sim \text{Uniform}[H, E, L]$ denotes uniformly sampling in the set $\{H, E, L\}$.

\mathcal{V}^* should be $\rho(\mathcal{V}, \mathcal{V}^*) = \frac{\mathcal{P}(\mathcal{V}^*, O)}{\mathcal{P}(\mathcal{V}, O)} \times \frac{\mathcal{P}(\mathcal{V} \leftarrow \mathcal{V}^*)}{\mathcal{P}(\mathcal{V}^* \leftarrow \mathcal{V})}$. Therefore, the acceptance probability for *Segment split* and *Segment merge* should respectively be

$$\begin{aligned} \rho_{split(k)}(\mathcal{V}, \mathcal{V}^*) &= \frac{\mathcal{P}(\mathcal{V}^*, O)}{\mathcal{P}(\mathcal{V}, O)} \times |\mathcal{I}| \cdot (e_k - e_{k-1} - 1) \\ \rho_{merge(k)}(\mathcal{V}, \mathcal{V}^*) &= \frac{\mathcal{P}(\mathcal{V}^*, O)}{\mathcal{P}(\mathcal{V}, O)} \times \frac{1}{|\mathcal{I}| \cdot (e_{k+1} - e_{k-1} - 1)} \end{aligned} \quad (20)$$

where $\mathcal{P}(\mathcal{V}, O)$ is defined as in (19) and $|\mathcal{I}| = 3$ denotes the number of segment types.

Due to the factorizations in (19), only the changed segments require evaluation in computing the acceptance probability for the new proposal \mathcal{V}^* . Once the β -strands are changed in the new proposal, the interacting set \mathcal{I} is changed too. The joint segmental likelihood of the β -strands has to be calculated again, which is a sum $\sum_{\mathcal{I}} \mathcal{P}(\mathcal{I}|m, e, T) \prod_{S_i \in \mathcal{I}} \mathcal{P}(S_i|T_i = E, S_{-i})$. Although the set \mathcal{I} is composed of finite elements, it might be too expensive to enumerate all of them for the marginalization. We again apply sampling methods here to approximate the sum by randomly walking in the distribution $\mathcal{P}(\mathcal{I}|m, e, T)$ that is defined as in (11). A sampling scheme is described in Appendix B for this purpose. These samples can be reused to estimate the β -sheet contact map as in (15).

In the model training, we need to solve three optimization problems to estimate the weight parameters by gradient-descent methods. This is required tens of times in cross validation.¹³ Once the optimal regularization parameter is found, we solve the minimization problems once more to get the final weight parameters. The cost on counting the occurrence frequencies of these discrete distributions is relatively negligible. In the inference without long-range interactions, the computational complexity is presented in Appendix A. With the incorporation of long-range interactions, we employed the sampling scheme described above to collect 10000 samples in the posterior distribution.¹⁴ To approximate the marginalization over the interacting set, we randomly collected 40 samples in the β -sheet space $\mathcal{P}(\mathcal{I}|m, e, T)$.

Ideally, the inference problem should be formulated as a Bayesian hierarchical model and all quantities, which includes the latent variables, parameters and model parameters, could be sampled from the joint posterior distribution by MCMC methods [31]. However, the computational cost could be prohibitively expensive, since it would involve sampling in several high-dimensional spaces jointly. For example, the three spaces of the weight parameters contain over ten thousand variables. Nevertheless, we believe that complete Bayesian inference could achieve a genuine improvement, which is well worth further investigation.

¹³In the present work, it is required 6×7 times for 7-fold cross validation on 6 different C values.

¹⁴The first 1000 samples in the Markov chain were discarded for "burn in".

We implemented the proposed algorithm in ANSI C.¹⁵ In this implementation, the length of dependency window was fixed at 5, and the length of N- and C-capping was fixed at 4, and the regularization factor C was fixed at 0.01. We normalized the M_i vectors so that $\sum_a M_i^a = 1$ for both the multiple sequence alignment profile and the PSSM based profile. We used the following quantities as performance measures:

- Overall 3-state Accuracy Q_3 ,
- Sensitivity $Q^{obs} = \frac{TruePositive}{TruePositive+FalseNegative}$,
- Positive Predictive Value $Q^{pred} = \frac{TruePositive}{TruePositive+FalsePositive}$,
- Matthew's correlation $C = \frac{TP \times TN - FP \times FN}{\sqrt{(TP+FN)(TP+FP)(TN+FP)(TN+FN)}}$ defined by Matthews [32],
- Segment Overlap Measure (SOV) as defined by Zemla et al. [33].

A. Validation on CB513

The data set we used is CB513, a non-redundant set of 513 non-homologous protein chains with structures determined to a resolution of $\leq 2.5\text{\AA}$ generated by Cuff and Barton [34].¹⁶ This data set has been used as a common benchmark for a number of different secondary structure prediction algorithms. We used 3-state DSSP definitions of secondary structure [35], calculated from the PDB files.¹⁷ We removed the proteins that are shorter than 30 residues, or longer than 550 residues, following [5], to leave 480 proteins for 7-fold cross validation. Seven folds were created randomly, and validation outputs were carried out on the left-out fold while the weight parameters were optimized on the other 6 folds with the regularization factor $C = 0.01$ in turn. We used two kinds of alignment profiles: the multiple sequence alignment profiles (MSAP) used by Cuff and Barton [5], and position-specific score matrices (PSSM) as in [4]. For comparison purposes, we also implemented the algorithm proposed by Schmidler [16], which uses the single sequence information only.¹⁸ The validation results are recorded in Table II. We also cite the results reported by Cuff and Barton [5] in Table III for reference. The results obtained from our model show a substantial improvement over those of Schmidler [16] on all evaluation criteria. Compared with the performance of the neural network methods with various alignment profiles as shown in Table III, the

¹⁵The web server of our algorithm is available at <http://public.kgi.edu/~chuwei/eva/submiteva.html>.

¹⁶The data set and the multiple sequence alignments profiles generated by Cuff and Barton [5], can be accessed at <http://www.compbio.dundee.ac.uk/~www-jpred/data/>.

¹⁷In DSSP definitions, H and G were assigned as α -helix segments, E and B were assigned as β -strands, and the others were assigned as coil. The segments with only one residue were also labelled as coil.

¹⁸The source code in ANSI C can be accessed at <http://www.gatsby.ucl.ac.uk/~chuwei/code/bspss.tar.gz>.

prediction accuracy of our model is also competitive.¹⁹ Due to small sample errors and the variation due to changes in secondary structure assignment by different methods, reported accuracies separated by less than about two percentage points are unlikely to be statistically significant [36], [37] and our results are comparable to many other prediction methods which have been tested on this benchmark data. Crooks and Brenner [27] point out that this is probably due to the fact that most contemporary methods for secondary structure prediction all utilize local sequence correlations, which contain only about one quarter of the total information necessary to determine secondary structure .

We did observe that the marginal posterior mode is more accurate than the MAP estimate, which shows that averaging over all the possible segmentations helps. According to the class definitions of the Structural Classification of Proteins database (SCOP) [38], we divided the 480 chains of CB513 into four groups: α , β , α/β and $\alpha + \beta$. The validation results of marginal posterior mode estimate on these groups are recorded separately in Table IV. We note that the performance on α/β and α proteins is relatively better than that on $\alpha + \beta$ and β .

B. Blind Test on CASP Targets

The meetings of Critical Assessment of Techniques for Protein Structure Prediction (CASP) facilitate large-scale experiments to assess protein structure prediction methods. To perform a blind test experiment, we extracted protein chains from the latest three meetings from the public web page of the Protein Structure Prediction Center.²⁰ With the model parameters specified in Section IV, we optimized the weight parameters of our model using all the 480 chains from CB513 and their PSSM profiles, and then carried out prediction on these CASP target proteins. We also prepared a larger training dataset using the CullerPDB list with the percentage identity cutoff 25%, the resolution cutoff 1.8 angstroms, and the R-factor cutoff 0.25.²¹ There are 2147 chains in this expanded list. We used the same model parameters, and optimized the weights parameters on the subset of 1814 chains.²² The predictive results of marginal posterior mode estimate of our two models are reported in Table V, indexed by meeting, along with the marginal posterior mode estimate of the Schmidler's algorithm [16]. The predictive results of CASP 5 are presented in Table VI in more detail. We cite the average performance of the participants from the CASP5 website for comparative purposes.

¹⁹It is also possible to further improve performance by constructing smoothers over current predictive outputs as Cuff and Barton [5] did in their Jury networks.

²⁰<http://predictioncenter.llnl.gov/>

²¹The protein list is accessible at http://dunbrack.fccc.edu/Guoli/pisces_download.php.

²²The reduction was caused by removing the protein chains that are shorter than 30 residues, or longer than 550 residues, following [5].

TABLE II

VALIDATION RESULTS FOR SECONDARY STRUCTURE PREDICTION ON 480 PROTEIN SEQUENCES FROM CB513. ‘SEQUENCE ONLY’ DENOTES THE ALGORITHM OF SCHMIDLER [16]; MSAP DENOTES OUR APPROACH USING MULTIPLE SEQUENCE ALIGNMENT PROFILES; PSSM DENOTES OUR APPROACH USING POSITION SPECIFIC SCORE MATRICES. Q_3 DENOTES THE OVERALL ACCURACY. $Q^{obs} = \frac{TruePositive}{TruePositive+FalseNegative}$ AND $Q^{pred} = \frac{TruePositive}{TruePositive+FalsePositive}$. \mathcal{C} DENOTES MATTHEWS’ CORRELATION COEFFICIENT DEFINED BY MATTHEWS [32]. SOV DENOTES THE SEGMENT OVERLAP MEASURE [33]. THE SUBSCRIPTS DENOTE THE SECONDARY STRUCTURE TYPE. MAP DENOTES THE MOST PROBABLE POSTERIOR ESTIMATE, WHILE MARG DENOTES MARGINAL POSTERIOR MODE ESTIMATE.

	Sequence Only		with MSAP		with PSSM	
	MAP	MARG	MAP	MARG	MAP	MARG
Q_3	59.23%	65.08%	68.11%	71.31%	62.54%	72.23%
Q_H^{obs}	66.34%	66.73%	78.17%	78.71%	66.74%	71.56%
Q_E^{obs}	20.74%	46.32%	41.40%	57.61%	26.18%	54.93%
Q_C^{obs}	72.80%	73.19%	73.28%	72.11%	77.45%	81.54%
Q_H^{pred}	61.87%	68.64%	69.91%	73.51%	66.82%	79.90%
Q_E^{pred}	56.45%	58.88%	70.15%	67.67%	72.16%	70.96%
Q_C^{pred}	57.77%	64.72%	66.06%	70.94%	58.54%	67.91%
\mathcal{C}_H	0.3709	0.4621	0.5457	0.5927	0.4350	0.6085
\mathcal{C}_E	0.2194	0.3809	0.4306	0.5120	0.3359	0.5203
\mathcal{C}_C	0.2821	0.3945	0.4253	0.4820	0.3338	0.5064
SOV_H	48.10%	61.18%	61.56%	67.83%	50.03%	68.58%
SOV_E	58.55%	64.41%	68.78%	74.26%	57.39%	70.71%
SOV_C	33.13%	60.90%	53.08%	69.42%	39.50%	68.92%
SOV	37.50%	60.97%	58.42%	63.79%	48.54%	68.31%

The results of these blind tests indicate that our algorithm based on generative modelling gives comparable results to other contemporary methods.²³ The performance of Q_3 and SOV on the target proteins of CASP 5 are shown in Figure 6. We found that the model trained on the larger dataset can achieve better generalization performance, especially on SOV .

C. Prediction of Contact Maps

In the inference with long range interactions, we approximated the marginalization over the β -sheet space by randomly collecting 40 samples in $\mathcal{P}(\mathcal{I}|m, e, T)$ as described in Appendix B. We present the trace plot of ten test proteins in the MCMC sampling to show the convergence of the Markov chains, and compare the results to those without long range interactions in Figure 7. We found that the Markov

²³The predictive results produced by other contemporary methods, indexed by CASP meeting, are available at <http://predictioncenter.llnl.gov/>

THE RESULTS OF 7-FOLD CROSS VALIDATION ON 480 PROTEINS OF CB513 REPORTED BY [5], ALONG WITH OUR RESULTS. Q_3

DENOTES THE OVERALL ACCURACY.

METHOD DESCRIPTION	Q_3
NETWORKS USING FREQUENCY PROFILE FROM CLUSTALW	71.6%
NETWORKS USING BLOSUM62 PROFILE FROM CLUSTALW	70.8%
NETWORKS USING PSIBLAST ALIGNMENT PROFILES	72.1%
ARITHMETIC SUM BASED ON THE ABOVE THREE NETWORKS	73.4%
NETWORKS USING PSIBLAST PSSM	75.2%
OUR ALGORITHM WITH MSAP OF [5]	71.3%
OUR ALGORITHM WITH PSIBLAST PSSM	72.2%

chains converge well after 6000 samples in all cases.

We prepared a dataset with long range interaction information specified by the Protein Data Bank (PDB) files. The dataset, a subset of CB513, is composed of 198 protein chains along with β -sheet definitions.²⁴ This reduction in size was caused by the incompleteness in the long range interaction information in many of the original PDB files. In MCMC sampling we collected 9000 samples. 30-fold cross validation was carried out on this subset. Surprisingly, we have not yet observed significant improvement on secondary structure prediction accuracy in the sampling results over exact inference without long range interactions. This indicates either some limitations in our current implementation or sampling scheme, or the small size of the training data set used in this set of experiments, which we will investigate further by re-training the model on a larger data set. However, this observation is consistent with the findings of Cline et al. [39] and Crooks et al. [40], who examined the mutual information content of interacting amino acid residues distantly separated by sequence but proximate in three-dimensional structure, and concluded that, for the purposes of tertiary structure prediction, these interactions were essentially uninformative. The analysis of Cline et al. [39] and Crooks et al. [40] also suggests that a modification to our method, which captures distal interactions between *secondary structure elements* rather than amino acid residues, should provide a distinct improvement.

However, it is interesting that we can infer β -sheet contacts based on the predicted secondary structure. We present predicted contact maps in Figure 8 as an example, where the colour scale indicates the probability $\mathcal{P}(C^{ij} = 1|O)$. It can be seen that, in the case of 1PGA (Protein G), which contains 2 parallel and 2 anti-parallel β -strands, and 1DTX (α -dendrotoxin), which contains 2 anti-parallel β -strands, the

²⁴The list of these proteins can be found at http://www.gatsby.ucl.ac.uk/~chuwei/biopss/jcbmc_list.txt.

TABLE IV

VALIDATION RESULTS OF MARGINAL POSTERIOR MODE ESTIMATE FOR SECONDARY STRUCTURE PREDICTION ON 480 PROTEIN SEQUENCES FROM CB513, CATEGORIZED BY STRUCTURAL CLASSES OF PROTEINS (SCOP). MSAP DENOTES OUR APPROACH USING MULTIPLE SEQUENCE ALIGNMENT PROFILES; PSSM DENOTES OUR APPROACH USING POSITION SPECIFIC SCORE MATRICES. Q_3 DENOTES THE OVERALL ACCURACY. $Q^{obs} = \frac{TruePositive}{TruePositive+FalseNegative}$ AND $Q^{pred} = \frac{TruePositive}{TruePositive+FalsePositive}$. C DENOTES MATTHEWS' CORRELATION COEFFICIENT DEFINED BY MATTHEWS [32]. SOV DENOTES THE SEGMENT OVERLAP MEASURE [33]. THE SUBSCRIPTS DENOTE THE SECONDARY STRUCTURES.

	α proteins (65 chains)		β proteins (92 chains)		α/β proteins (129 chains)		$\alpha + \beta$ proteins (87 chains)	
	MSAP	PSSM	MSAP	PSSM	MSAP	PSSM	MSAP	PSSM
Q_3	73.97%	74.78%	68.35%	69.95%	72.97%	74.47%	69.81%	70.72%
Q_H^{obs}	79.53%	76.02%	64.45%	55.82%	81.02%	73.36%	76.91%	69.72%
Q_E^{obs}	47.45%	45.60%	57.06%	52.71%	61.90%	63.44%	55.92%	52.39%
Q_C^{obs}	67.50%	76.48%	78.00%	86.82%	69.48%	79.94%	72.66%	82.44%
Q_H^{pred}	85.36%	89.13%	46.01%	56.58%	75.72%	82.14%	71.19%	77.49%
Q_E^{pred}	25.00%	25.45%	77.43%	81.60%	65.67%	69.97%	68.85%	74.84%
Q_C^{pred}	67.35%	66.00%	71.22%	68.17%	72.82%	70.00%	69.17%	65.63%
C_H	0.5328	0.5694	0.4424	0.4792	0.5945	0.6131	0.5770	0.5910
C_E	0.2917	0.2894	0.4842	0.4965	0.5517	0.5894	0.4880	0.5100
C_C	0.4956	0.5339	0.4416	0.4654	0.4985	0.5340	0.4652	0.4858
SOV_H	67.04%	71.23%	65.44%	65.74%	71.77%	72.94%	66.84%	69.02%
SOV_E	74.36%	75.91%	60.25%	52.11%	82.95%	77.97%	77.04%	75.54%
SOV_C	75.98%	76.85%	66.54%	65.57%	73.32%	76.03%	65.17%	62.88%
SOC	59.98%	67.62%	66.82%	69.03%	64.96%	70.71%	62.55%	67.69%

position and direction of the β -strands are predicted correctly, but have a shorter range than in the true contact maps. The false positive predictions in the case of 1DTX (α -dendrotoxin) are due to errors in the prediction of which residues are in the β -strands.

To assess the overall prediction accuracy, we have also computed the area under the ROC curve (AUC) [41] for β -sheet contact prediction. The average AUC over these protein chains is 0.90 ± 0.10 . The average ROC curves categorized by SCOP classes are presented in Figure 11. The averaged AUC of 44 β proteins is 0.87 ± 0.07 , the averaged AUC of 64 α/β proteins is 0.93 ± 0.06 and the averaged AUC of 37 $\alpha + \beta$ proteins is 0.90 ± 0.10 . Based on these ROC curves, we find that this algorithm performs better on the α/β class.

TABLE V

PREDICTIVE RESULTS OF MARGINAL POSTERIOR MODE ESTIMATE OF OUR ALGORITHM USING PSSM ON THE PROTEIN DATA OF CASP.

CASP 3 HAS 36 CHAINS, CASP 4 HAS 40 CHAINS, AND CASP 5 HAS 56 CHAINS. "SEQUENCE ONLY" DENOTES THE ALGORITHM OF SCHMIDLER [16]; "CB513 WITH PSSM" DENOTES OUR MODEL TRAINED ON THE 480 CHAINS FROM CB513 WITH PSSM PROFILES; "CULLEDPDB WITH PSSM" DENOTES OUR MODEL TRAINED ON THE 1814 CHAINS OF CULLEDPDB DATA. Q_3 DENOTES THE OVERALL

ACCURACY. $Q^{obs} = \frac{TruePositive}{TruePositive+FalseNegative}$ AND $Q^{pred} = \frac{TruePositive}{TruePositive+FalsePositive}$. C DENOTES MATTHEWS' CORRELATION

COEFFICIENT DEFINED BY MATTHEWS [32]. SOV DENOTES THE SEGMENT OVERLAP MEASURE [33].

	SEQUENCE ONLY			CB513 WITH PSSM			CULLEDPDB WITH PSSM		
	CASP3	CASP4	CASP5	CASP3	CASP4	CASP5	CASP3	CASP4	CASP5
Q_3	63.92%	66.29%	67.13%	72.03%	74.87%	74.49%	72.28%	74.68%	74.86%
Q_H^{obs}	61.55%	67.64%	67.47%	78.61%	85.17%	85.85%	71.13%	77.42%	77.53%
Q_E^{obs}	42.45%	44.92%	48.15%	58.30%	58.74%	60.74%	56.19%	58.42%	60.46%
Q_C^{obs}	78.57%	77.13%	77.46%	75.63%	73.69%	72.40%	82.78%	81.13%	80.54%
Q_H^{pred}	64.27%	72.49%	71.25%	71.17%	77.88%	73.73%	78.30%	82.82%	79.74%
Q_E^{pred}	69.81%	65.06%	67.96%	80.12%	77.07%	79.41%	78.67%	76.54%	78.18%
Q_C^{pred}	62.03%	61.79%	63.97%	69.29%	70.58%	73.06%	66.73%	67.35%	69.92%
C_H	0.4200	0.4721	0.4904	0.5961	0.6573	0.6442	0.6111	0.6509	0.6451
C_E	0.4072	0.4157	0.4478	0.5766	0.5804	0.6022	0.5589	0.5784	0.5966
C_C	0.3881	0.4344	0.4378	0.4956	0.5362	0.5330	0.5050	0.5445	0.5472
SOV_H	62.58%	59.63%	65.60%	70.31%	69.24%	73.35%	70.52%	70.49%	74.31%
SOV_E	56.08%	60.39%	65.61%	67.55%	73.55%	77.23%	66.81%	70.69%	75.61%
SOV_C	61.55%	60.92%	60.87%	66.13%	67.73%	71.51%	71.90%	74.52%	73.73%
SOV_C	65.26%	65.39%	66.55%	64.91%	64.85%	68.18%	74.52%	68.95%	73.08%

VII. CONCLUSION

In this paper, we have described a novel parametric Bayesian segmental semi-Markov model for proteins which incorporates the information in multiple sequence alignment profiles. Long range interaction information in β -sheets can be directly incorporated. The numerical results show that the generalization performance of this generative model is similar to other contemporary methods. However, contact map prediction can also be carried out in the Bayesian segmental framework, which represents a considerable advantage over discriminative methods. Moreover, with the inclusion of potential functions with dihedral angle information in the joint sequence-structure probability distribution, this probabilistic model also has the potential for tertiary structure prediction, and this is the focus of our current work.

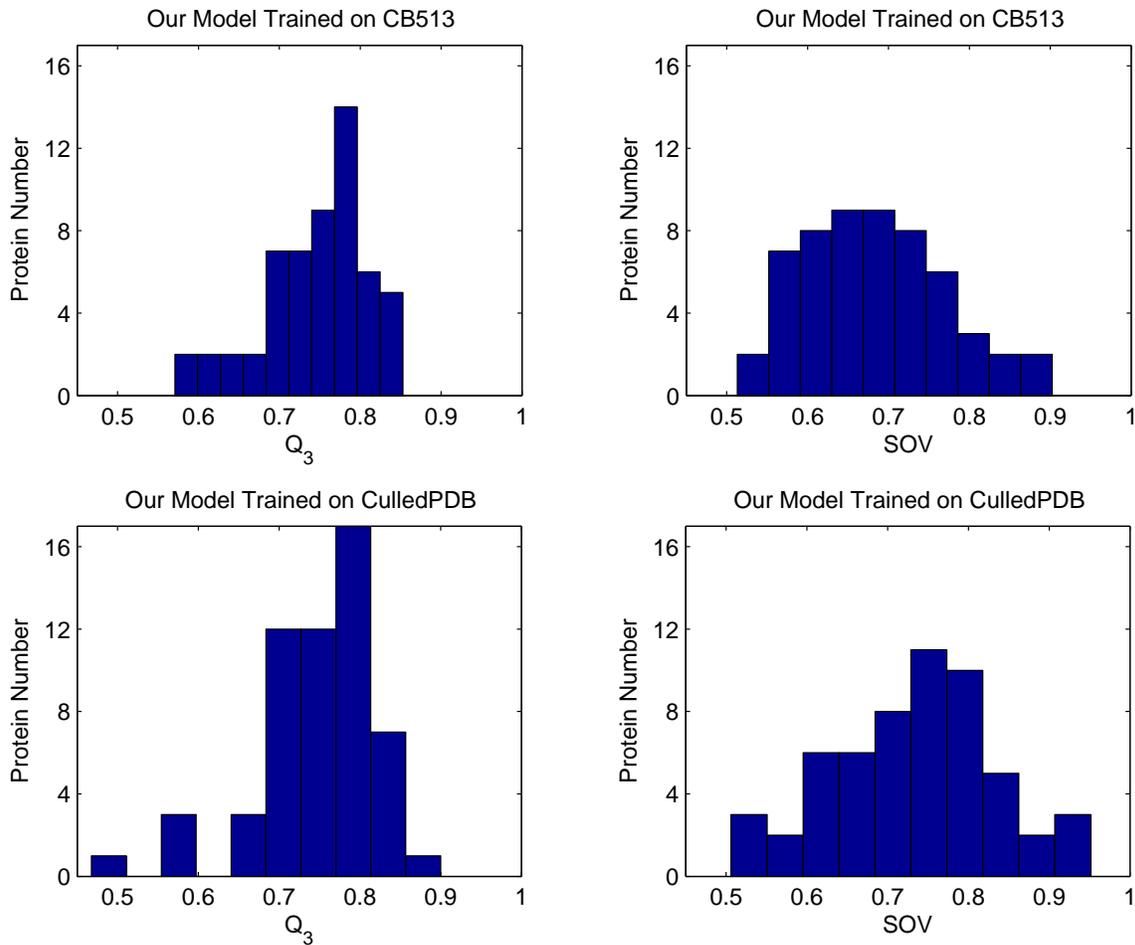


Fig. 6. The histogram of Q_3 and SOV to visualize the marginal posterior mode estimate of our models on CASP5 data. One model was trained on CB513 dataset and another was trained on CullerPDB dataset. The vertical axes are indexed by the number of proteins falling in the bins.

ACKNOWLEDGMENT

We would like to thank the Institute of Applied Mathematics (IPAM) at UCLA, where part of this work was carried out. This work was supported by the National Institutes of Health Grant Number 1 P01 GM63208.

APPENDIX

A. Exact Inference

We give the outlines of the forward-backward algorithm for the marginal posterior mode and the Viterbi algorithm for the MAP estimate respectively, see [25] for more details.

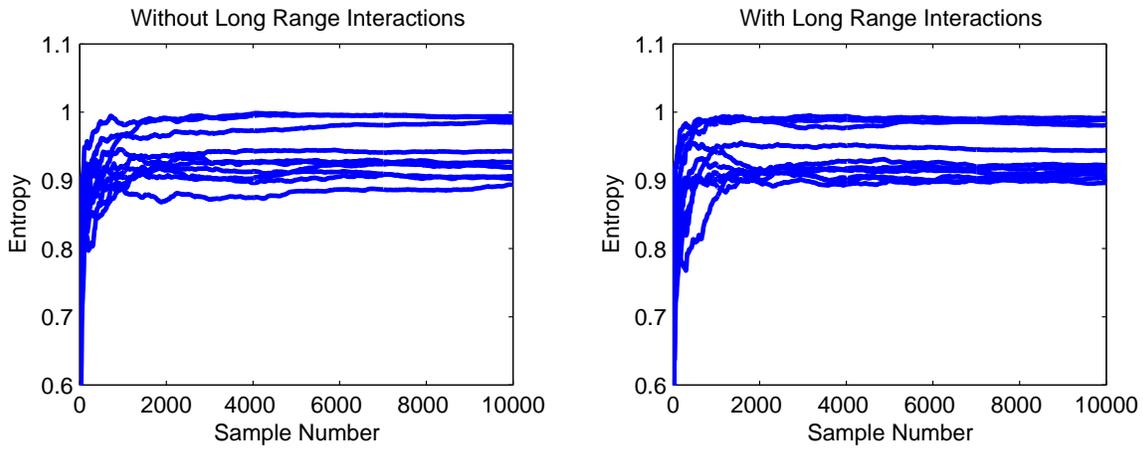


Fig. 7. The trace plot of ten protein chains in MCMC sampling. ‘Entropy’ denoted the average entropy of the posterior distribution on the chain, i.e. $-\frac{1}{n} \sum_{i=1}^n \mathcal{P}(T_{O_i}|O) \log \mathcal{P}(T_{O_i}|O)$ where $\mathcal{P}(T_{O_i}|O)$ is the predictive probability of the segment type on the i -th amino acid of the protein chain. The sampling results of the cases with/without long range interaction are presented respectively.

1) *Forward-Backward*: Let us suppose that the parameters θ have been given. The forward-backward algorithm computes the following quantities:

$$\alpha(j, t) = \mathcal{P}(e_c = j, T_c = t, R_{[1:j]}) \quad (21)$$

$$\beta(j, t) = \mathcal{P}(R_{[j+1:n]} | e_c = j, T_c = t) \quad (22)$$

where $0 \leq j \leq n$, $t \in \{H, E, L\}$ and c denotes the index of the current segment containing R_j . We may assign the starting point $\alpha(0, t) = \mathcal{P}(T_0 = t) \forall t$,²⁵ and then we compute other $\alpha(j, t)$ in a forward pass as

$$\alpha(j, t) = \sum_{\nu=0}^{j-1} \sum_{\tau} \alpha(\nu, \tau) \mathcal{P}(R_{[\nu+1:j]} | e_{c-1} = \nu, e_c = j, T_c = t) \times \mathcal{P}(e_c = j | e_{c-1} = \nu, T_c = t) \mathcal{P}(T_c = t | T_{c-1} = \tau) \quad (23)$$

where j increments from 1 to n sequentially.²⁶ Afterwards, we initialize $\beta(n, t) = 1 \forall t$, and compute β 's in a backward pass as

$$\beta(j, t) = \sum_{\nu=j+1}^n \sum_{\tau} \beta(\nu, \tau) \mathcal{P}(R_{[j+1:\nu]} | e_c = j, e_{c+1} = \nu, T_{c+1} = \tau) \times \mathcal{P}(e_{c+1} = \nu | e_c = j, T_{c+1} = \tau) \mathcal{P}(T_{c+1} = \tau | T_c = t) \quad (24)$$

²⁵Simply we set $\mathcal{P}(T_0 = t) = 1/|T|$ with $|T| = 3$ in this paper.

²⁶In practice, we always normalize $\alpha(j, t)$ to sum to one at each step j of the recursions, since α might become vanishingly small for long sequences. The normalized $\alpha(j, t)$ and its normalizer $\sum_t \alpha(j, t)$ will be saved at each step j . The normalizers $\sum_t \alpha(j, t)$ can be reused in the backward recursions of β . $\beta(j, t)$ might be divided by the normalizer $\sum_t \alpha(j, t)$, and then be saved at each step j . This is known as the issue of scaling in [25].

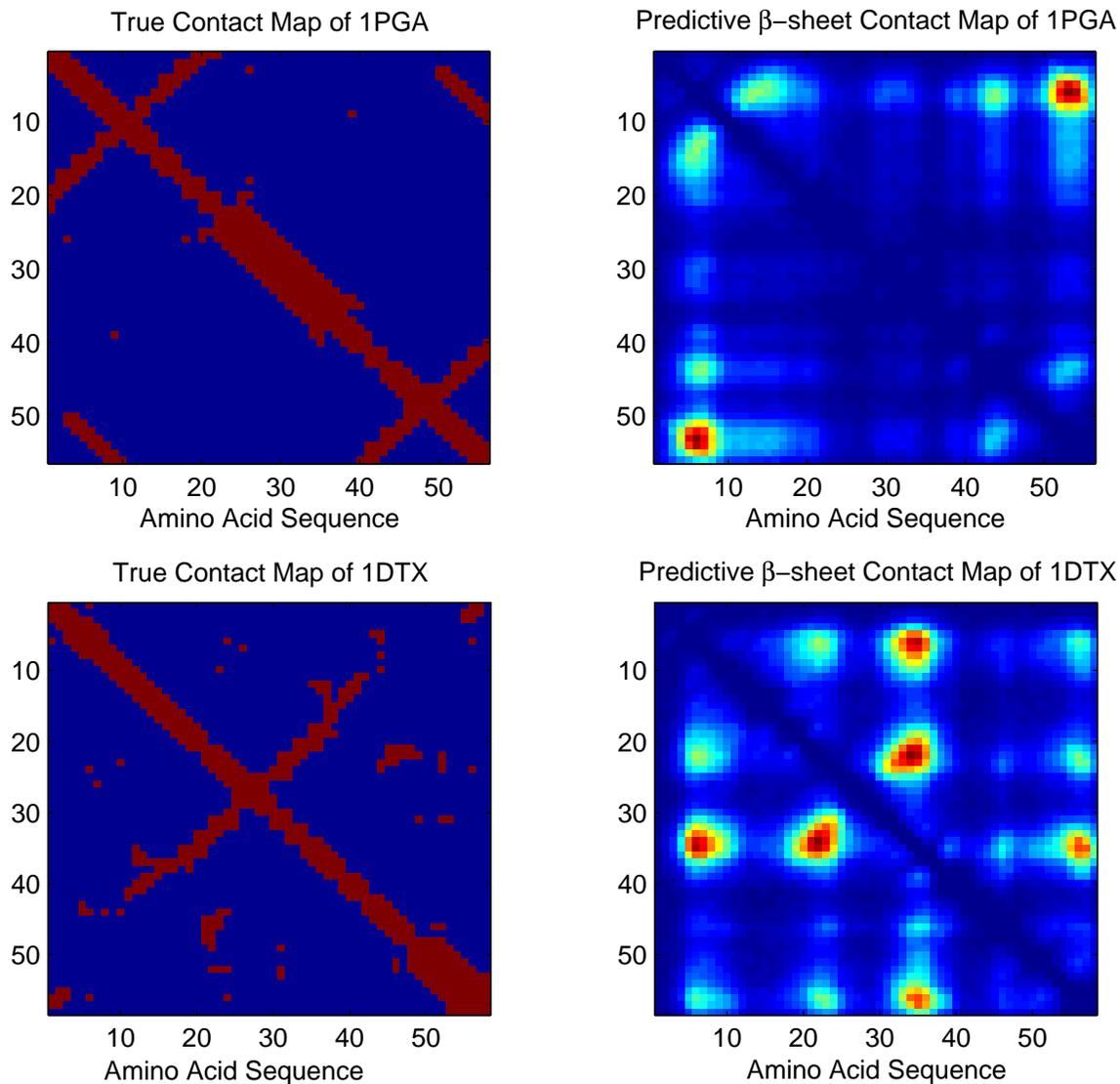


Fig. 8. True β -sheet contact maps and predicted maps for protein chains 1PGA and 1DTX. The true contact maps were produced with a threshold of 7 Å. The colour scale indicates the probability $\mathcal{P}(C^{ij} = 1|O)$.

where j decrements from $n-1$ to 0 sequentially. This algorithm has complexity $\mathcal{O}(|\mathcal{T}|^2 n^3)$, but in practice we may limit the maximum size of any one segment to some length L .²⁷ Thus, the first summation in (23) begins at $\max(j-L, 0)$ and the first summation in (24) ends at $\min(j+L, n)$, which reduces the complexity to $\mathcal{O}(|\mathcal{T}|^2 L^2 n)$. The complexity of computational time can be further reduced to $\mathcal{O}(|\mathcal{T}|^2 Ln)$ if we cache more intermediate values. Using the outputs of the Forward-Backward algorithm, we can

²⁷In this work, we set L at 30.

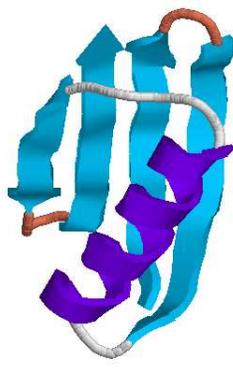


Fig. 9. The structure of protein 1PGA (Protein G)

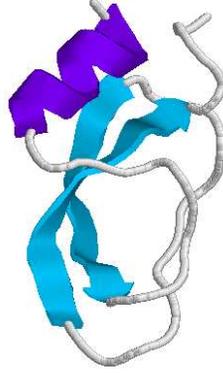


Fig. 10. The structure of protein 1DTX (α -dendrotoxin)

compute the probability of the segment type at each residue:

$$\begin{aligned} \mathcal{P}(T_{R_i} = t | R, \theta) &= \sum_{j=0}^{i-1} \sum_{k=i}^n \sum_{\tau} \alpha(j, \tau) \beta(k, \tau) \mathcal{P}(R_{[j+1, k]} | e_{c-1} = j, e_c = k, T_c = t) \\ &\quad \times \mathcal{P}(e_c = k | e_{c-1} = j, T_c = t) \mathcal{P}(T_c = t | T_{c-1} = \tau) / \mathcal{P}(R | \theta) \end{aligned} \quad (25)$$

where $\mathcal{P}(R | \theta) = \sum_{\tau} \mathcal{P}(T_{R_i} = \tau, R | \theta)$ is a normalizer that can be evaluated by $\sum_{\tau} \alpha(0, \tau) \beta(0, \tau)$.

2) *Viterbi Algorithm*: A procedure analogous to the Viterbi algorithm [42] can be used to find the optimal state sequence associated with the given observation sequence. Let us define the quantity

$$\delta(j, t) = \max_{e_1, \dots, e_{c-1}, T_1, \dots, T_{c-1}} \mathcal{P}(e_1, \dots, e_c = j, T_1, \dots, T_c = t, R_{[1:j]}) \quad (26)$$

which is the highest probability along a single path to account for the first j observation residues and ends in state t . To actually retrieve the state sequence, we need to keep track of the argument which

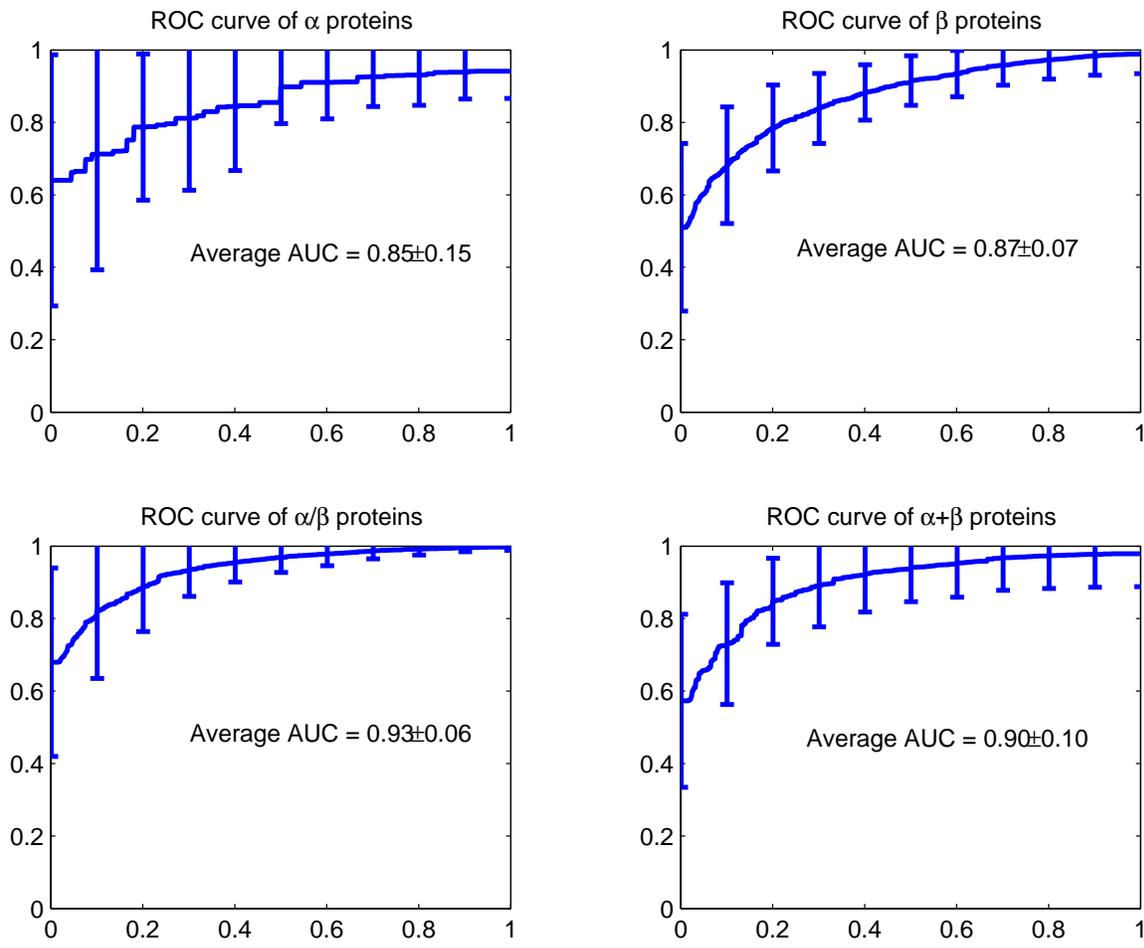


Fig. 11. The ROC curves of the proteins categorized by SCOP. The vertical lines bounded by bars in these graphs indicate the standard deviation at those positions. For these four graphs, the horizontal axes are indexed by $1.0 - \textit{specificity}$ evaluated by $\frac{\text{Number of False Positive}}{\text{Number of Negative Samples}}$, and the vertical axes are of $\textit{sensitivity}$ evaluated by $\frac{\text{Number of True Positive}}{\text{Number of Positive Samples}}$. For each of the structural classes, the average value and its standard deviation of AUC (the area under the ROC curve) are given by text in the corresponding graph.

maximized (26), for each j and t . The record array ψ might be used for this, with entries $\psi(j, t)$ which contain two elements $\{\arg \max_{e_{c-1}, T_{c-1}} \delta(j, t)\}$. The whole procedure can now be described as follows:

1) Initialization: $\delta(0, t) = \mathcal{P}(T_0 = t) \forall t$

2) Recursion with j from 1 to n , $\forall t$:

$$\delta(j, t) = \max_{0 \leq \nu \leq j-1, \tau \in \{H, E, C\}} \delta(\nu, \tau) \mathcal{P}(R_{[\nu+1:j]} | e_{c-1} = \nu, e_c = j, T_c = t) \times \mathcal{P}(e_c = j | e_{c-1} = \nu, T_c = t) \mathcal{P}(T_c = t | T_{c-1} = \tau) \quad (27)$$

$$\psi(j, t) = \arg \max_{e_{c-1}, T_{c-1}} \delta(\nu, \tau) \mathcal{P}(R_{[\nu+1:j]} | e_{c-1} = \nu, e_c = j, T_c = t) \times \mathcal{P}(e_c = j | e_{c-1} = \nu, T_c = t) \mathcal{P}(T_c = t | T_{c-1} = \tau) \quad (28)$$

3) Termination: $\mathcal{P}^* = \max_t \delta(t, n)$, $i = m$, $e_m^* = n$ and $T_m^* = \arg \max_t \delta(t, n)$

4) State Back-tracing with $i = i - 1$: $\{e_i^*, T_i^*\} = \psi(e_{i+1}^*, T_{i+1}^*)$ till $e_{i=0}^* = 0$.²⁸

Noted that the Viterbi algorithm is similar to the forward calculation except for the maximization in (27) over previous states in place of the summation in (23).

B. Sampling in β -sheet Space

Given a specific segmentation, i.e. a set of $\{m, S, T\}$, there is a corresponding β -sheet space defined by a set of interaction variables \mathcal{I} that specifies the interactions within these β -strands. The total number of β -strands is known, denoted as k . The distribution of the β -sheet space, $\mathcal{P}(\mathcal{I}|m, e, T)$, is defined as in

(11). There are four steps to collect a sample in $\mathcal{P}(\mathcal{I}|m, e, T)$:

- 1) generate a sample of r in $\mathcal{P}(r|k)$.
- 2) collect a valid combination of r pairs from the k β -strands. The valid combination requires that each β -strand should be used at least once and at most twice.
- 3) for each pair $\{S_j, S_{j'}\}$, generate the alignment direction by $\mathcal{P}(d_{jj'}|S_j, S_{j'})$.
- 4) for each pair $\{S_j, S_{j'}\}$, generate the alignment position by $\mathcal{P}(a_{jj'}|S_j, S_{j'})$.

REFERENCES

- [1] N. Qian and T. J. Sejnowski, "Predicting the secondary structure of globular proteins using neural network models," *Journal of Molecular Biology*, vol. 202, pp. 865–884, 1988.
- [2] B. Rost and C. Sander, "Prediction of protein secondary structure at better than 70% accuracy," *Journal of Molecular Biology*, vol. 232, pp. 584–599, 1993.
- [3] S. F. Altschul, T. L. Madden, A. A. Schaeffer, J. Zhang, Z. Zhang, W. Miller, and D. J. Lipman, "Gapped BLAST and PSI-BLAST: a new generation of protein database search programs," *Nucleic Acids Research*, vol. 25, pp. 3389–3402, 1997.
- [4] D. Jones, "Protein secondary structure prediction based on position-specific scoring matrices," *Journal of Molecular Biology*, vol. 292, pp. 195–202, 1999.
- [5] J. A. Cuff and G. J. Barton, "Application of multiple sequence alignment profiles to improve protein secondary structure prediction," *Proteins: Structure, Function and Genetics*, vol. 40, pp. 502–511, 2000.
- [6] A. L. Delcher, S. Kasif, H. R. Goldberg, and W. H. Hsu, "Protein secondary structure modelling with probabilistic networks," in *Proc. of Int. Conf. on Intelligent Systems and Molecular Biology*, 1993, pp. 109–117.
- [7] C. Burge and S. Karlin, "Prediction of complete gene structures in human genomic DNA," *Journal of Molecular Biology*, vol. 268, no. 1, pp. 78–94, 1997.
- [8] R. F. Yel, L. P. Lim, and C. B. Burge, "Computational inference of homologous gene structures in the human genome," *Genome Res.*, vol. 11, no. 5, pp. 803–816, 2001.
- [9] L. Zhang, V. Pavlovic, C. R. Cantor, and S. Kasif, "Human-mouse gene identification by comparative evidence integration and evolutionary analysis," *Genome Res.*, vol. 13, pp. 1190–1202, 2003.

²⁸ m is unknown in Termination, but can be retrieved after we reach $e_{i=0}^* = 0$.

- [10] I. Korf, P. Flicek, D. Duan, and M. R. Brent, "Integrating genomic homology into gene structure prediction," *Bioinformatics*, vol. 17 Suppl 1, pp. S140–S148, 2001.
- [11] C. S. Schmidler, "Statistical models and monte carlo methods for protein structure prediction," Ph.D. Thesis, Stanford University, May 2002.
- [12] M. Ostendorf, V. Digalakis, and O. Kimball, "From HMM to segment models: a unified view of stochastic modelling for speech recognition," *IEEE Trans. on Speech and Audio Processing*, vol. 4, no. 5, pp. 360–378, 1996.
- [13] K. Simons, C. Kooperberg, E. Huang, and D. Baker, "Assembly of protein tertiary structures from fragments with similar local sequences using simulated annealing and bayesian scoring functions," *Journal of Molecular Biology*, vol. 268, pp. 209–225, 1997.
- [14] Y. Ye, L. Jaroszewski, W. Li, and A. Godzik, "A segment alignment approach to protein comparison," *Bioinformatics*, vol. 19, pp. 742–749, 2003.
- [15] J. D. Thompson, D. G. Higgins, and T. J. Gibson, "CLUSTAL W: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice," *Nucleic Acids Res.*, vol. 22, pp. 4673–4680, 1994.
- [16] C. S. Schmidler, J. S. Liu, and D. L. Brutlag, "Bayesian segmentation of protein secondary structure," *Journal of Computational Biology*, vol. 7, no. 1/2, pp. 233–248, 2000.
- [17] R. Aurora and G. D. Rose, "Helix capping," *Protein Science*, vol. 7, pp. 21–38, 1998.
- [18] D. Eisenberg, R. M. Weiss, and T. C. Terwilliger, "The hydrophobic moment detects periodicity in protein hydrophobicity," *Proceedings of the National Academy of Sciences, USA*, vol. 81, pp. 140–144, 1984.
- [19] W. Chu, Z. Ghahramani, and D. Wild, "Protein secondary structure prediction using sigmoid belief networks to parameterize segmental semi-markov models," in *the proc. of 12th European Symposium on Artificial Neural Networks*, 2004.
- [20] N. C. Fitzkee and G. D. Rose, "Steric restrictions in protein folding: An α -helix cannot be followed by a contiguous β -strand," *Protein Science*, Feb. 2004.
- [21] G. E. Hinton, "Products of experts," in *Proceedings of the Ninth International Conference on Artificial Neural Networks*, 1999, pp. 1–6.
- [22] T. M. Klingler and D. L. Brutlag, "Protein science," *Discovering structural correlations in α -helices*, vol. 3, pp. 1847–1857, 1994.
- [23] K. Sjölander, K. Karplus, M. Brown, R. Hughey, A. Krogh, I. S. Mian, and D. Haussler, "Dirichlet mixtures: A method for improving detection of weak but significant protein sequence homology," *Computing Applications in the Biosciences*, vol. 12, no. 4, pp. 327–345, 1996.
- [24] O. Winther and A. Krogh, "Teaching computers to fold proteins," *Phys. Rev. E*, 2004, to appear.
- [25] R. L. Rabiner, "A tutorial on hidden Markov models and selected applications in speech recognition," *Proceedings of The IEEE*, vol. 77, no. 2, pp. 257–286, 1989.
- [26] G. Pollastri and P. Baldi, "Prediction of contact maps by GIOHMMs and recurrent neural networks using lateral propagation from all four cardinal corners," *Bioinformatics*, vol. 18 Suppl 1, pp. S62–S70, 2002.
- [27] G. E. Crooks and S. E. Brenner, "Protein secondary structure: Entropy, correlations and prediction," *Bioinformatics*, vol. 20, pp. 1603–1611, 2004.
- [28] P. Burman, "A comparative study of ordinary cross validation, v-fold cross validation and the repeated learning-testing methods," *Biomatrika*, vol. 76, no. 3, pp. 503–514, 1989.
- [29] M. Stone, "Cross-validatory choice and assessment of statistical predictions (with discussion)," *Journal of the Royal Statistical Society B*, vol. 36, pp. 111–147, 1974.
- [30] P. J. Green, "Reversible jump Markov chain Monte Carlo computation and Bayesian model determination," *Biometrika*, vol. 82, no. 4, pp. 711–732, 1995.

- [31] J. S. Liu, *Monte Carlo Strategies in Scientific Computing*. Springer, 2001.
- [32] B. W. Matthews, "Comparison of the predicted and observed secondary structure of t4 phage lysozyme," *Biochem. Biophys.*, vol. 405, pp. 442–451, 1975.
- [33] A. Zemla, C. Venclovas, K. Fidelis, and B. Rost, "A modified definition of sov, a segment-based measure for protein secondary prediction assessment," *Proteins: Structure, Function, and Genetics*, vol. 34, pp. 220–223, 1999.
- [34] J. A. Cuff and G. J. Barton, "Evaluation and improvement of multiple sequence methods for protein secondary structure prediction," *Proteins: Structure, Function and Genetics*, vol. 34, pp. 508–519, 1999.
- [35] W. Kabsch and C. Sander, "A dictionary of protein secondary structure: Pattern recognition of hydrogen-bonded and geometrical features," *Biopolymers*, vol. 22, pp. 2577–2637, 1983.
- [36] B. Rost and V. Eyrich, "Eva: large-scale analysis of secondary structure prediction," *Proteins*, vol. 45, Suppl. 5, pp. 192–199, 2001.
- [37] D. Przybylski and B. Rost, "Alignments grow, secondary structure prediction improves," *Proteins*, vol. 46, pp. 197–205, 2002.
- [38] A. Murzin, S. Brenner, T. Hubbard, and C. Chothia, "SCOP: a structural classification of proteins database for the investigation of sequences and structures," *J. Mol. Biol.*, vol. 247, pp. 536–540, 1995.
- [39] M. S. Cline, K. Karplus, R. Lathrop, T. Smith, R. Rogers Jr., and D. Haussler, "Information-theoretic dissection of pairwise contact potentials," *Proteins: Structure, Function, and Bioinformatics*, vol. 49, pp. 7–14, 2002.
- [40] G. E. Crooks, J. Wolfe, and S. E. Brenner, "Measurements of protein sequence-structure correlations," *Proteins: Structure, Function, and Bioinformatics*, vol. 57, pp. 804–810, 2004.
- [41] A. P. Bradley, "The use of the area under the roc curve in the evaluation of machine learning algorithms," *Pattern Recognition*, vol. 30, no. 7, pp. 1145–1159, 1997.
- [42] G. D. Forney, "The Viterbi algorithm," *Proceedings of The IEEE*, vol. 61, pp. 268–278, March 1973.

TABLE VI

THE DETAILED PREDICTIVE RESULTS OF MARGINAL POSTERIOR MODE ESTIMATE OF OUR ALGORITHM USING PSSM ON THE 56 CASP5

CHAINS. "SCOP" DENOTES THE STRUCTURAL CLASSES IN SCOP; THE SUPERSCRIP "CASP5" DENOTES THE AVERAGE RESULTS OF

CASP 5 PARTICIPANTS; "CB513" AND "CULLED" DENOTE OUR MODEL TRAINED ON CB513 AND CULLEDPDB DATA RESPECTIVELY.

Target Index	PDB ID	SCOP	Chain Length	Q_3^{casp5}	SOV^{casp5}	Q_3^{cb513}	SOV^{cb513}	Q_3^{culled}	SOV^{culled}
T0129	1IZM		170	78.4±10.8%	79.3±12.0%	82.4%	67.6%	81.8%	69.3%
T0130			100	76.2±8.4%	78.7±11.1%	84.0%	90.2%	90.0%	92.2%
T0132			146	81.9±11.1%	75.7±10.5%	83.6%	68.8%	84.2%	77.4%
T0133			293	73.3±11.0%	70.8±11.9%	73.4%	60.0%	70.0%	64.2%
T0134			233	59.8±8.7%	62.4±9.4%	57.1%	57.5%	57.9%	62.0%
T0135			106	70.8±17.6%	66.6±20.0%	77.4%	73.0%	78.3%	78.3%
T0136	1ON3	α/β	520	75.5±10.2%	74.6±10.9%	77.9%	78.0%	77.9%	73.1%
T0137	1O8V	β	133	86.2±12.7%	88.4±14.1%	68.4%	58.8%	59.4%	54.4%
T0138	1M2E	α/β	135	69.3±9.0%	70.9±10.8%	75.6%	64.3%	84.4%	77.6%
T0139	1KOY	α	62	64.8±14.8%	67.1±16.4%	61.3%	79.9%	46.8%	76.7%
T0140	1IYA	$\alpha + \beta$	86	66.5±8.8%	67.6±14.0%	61.6%	59.8%	65.1%	69.6%
T0141			187	72.7±7.9%	66.1±11.6%	73.8%	67.5%	77.5%	72.8%
T0142			280	73.6±8.8%	69.6±11.8%	74.3%	62.9%	75.7%	79.4%
T0143			215	72.8±9.0%	68.7±11.1%	64.2%	58.2%	67.4%	70.7%
T0146			299	65.6±8.7%	51.5±7.5%	75.6%	54.6%	77.6%	51.3%
T0147	1M65	α/β	235	76.7±10.9%	76.3±12.4%	78.3%	72.8%	79.6%	83.8%
T0148	1INO	β	162	74.0±11.2%	74.3±12.3%	71.0%	68.5%	70.4%	73.0%
T0149	1NU	$\alpha + \beta \alpha/\beta$	317	71.9±10.1%	71.6±12.3%	72.9%	77.6%	72.6%	83.4%
T0150	1H7M	$\alpha + \beta$	97	79.7±11.3%	85.0±16.0%	81.4%	86.3%	78.4%	84.0%
T0151			107	77.9±12.9%	81.0±14.0%	74.8%	59.2%	73.8%	76.6%
T0152			197	69.2±11.3%	67.6±10.8%	77.2%	70.8%	78.7%	73.5%
T0153	1MQ7	β	134	79.8±12.4%	76.6±12.7%	81.3%	72.9%	85.1%	81.6%
T0154	1MOP	α/β	288	76.9±9.0%	73.7±10.2%	77.4%	64.0%	76.7%	70.4%
T0155	1NBU		117	79.1±11.2%	80.3±14.1%	75.2%	63.9%	76.1%	66.8%
T0156	1NXJ	α/β	156	69.8±8.9%	66.9±9.9%	72.4%	67.1%	76.9%	83.1%
T0157			121	73.9±10.5%	76.4±12.2%	82.6%	73.4%	83.5%	94.5%
T0159			309	72.7±10.0%	68.2±12.5%	77.0%	64.6%	74.4%	81.7%
T0160			126	79.5±10.6%	80.0±15.0%	80.2%	88.5%	80.2%	87.8%
T0161	1MW5		154	66.0±8.5%	65.0±10.1%	64.9%	51.3%	59.7%	62.1%
T0162	1IZN	$\alpha + \beta(e)$	275	70.9±8.1%	67.9±8.7%	67.3%	63.6%	70.2%	78.9%
T0165	1L7A	α/β	318	76.8±11.9%	78.5±15.6%	78.9%	74.9%	79.2%	72.9%
T0167	1M3S	α/β	181	82.5±9.3%	79.3±9.3%	79.6%	70.0%	81.2%	65.1%
T0168	1MKI	$\alpha + \beta$	313	76.0±10.4%	75.1±11.3%	75.7%	74.2%	72.2%	77.1%
T0169	1MK4	$\alpha + \beta$	156	71.3±10.2%	71.3±11.5%	71.8%	65.2%	71.8%	70.1%
T0170	1H40	α	68	85.2±10.8%	83.4±14.0%	85.3%	81.5%	83.8%	90.3%
T0172	1M6Y	α/β	293	76.2±8.1%	73.1±8.4%	70.0%	55.9%	71.0%	56.7%
T0173	1Q74		289	72.7±9.0%	65.5±8.9%	72.3%	61.3%	75.8%	74.7%
T0174	1MG7	$\alpha + \beta$	354	59.7±6.0%	60.1±6.4%	59.9%	55.7%	64.7%	70.4%
T0176	1N91	$\alpha + \beta$	100	74.9±9.5%	81.2±13.3%	69.0%	69.0%	70.0%	64.2%
T0177	1MW7	$\alpha + \beta$	220	71.9±10.4%	72.4±11.4%	70.5%	58.4%	70.5%	67.1%
T0178	1MZH	α/β	219	80.7±10.8%	83.3±14.2%	84.5%	83.2%	81.3%	84.2%
T0179	1IY9	α/β	274	75.5±11.5%	75.7±15.2%	67.9%	64.4%	71.2%	73.7%
T0181	1NYN	$\alpha + \beta$	110	75.3±10.3%	72.6±14.7%	79.1%	76.2%	84.5%	95.1%
T0182	1O0X	$\alpha + \beta$	249	81.1±12.4%	72.9±15.6%	78.3%	55.8%	77.9%	63.9%
T0183	1O0Y	α/β	247	76.9±10.7%	78.2±14.3%	81.8%	76.8%	79.8%	80.5%
T0184	1O0W	$\alpha + \beta$	237	77.4±10.3%	75.7±12.6%	75.9%	60.4%	73.0%	58.2%
T0185	1J6U	α/β	431	78.4±9.9%	79.5±13.7%	77.5%	72.2%	78.9%	79.1%
T0186	1O12	β	363	73.2±9.1%	67.8±9.6%	72.2%	65.7%	69.4%	69.2%
T0187	1O0U	α/β	413	77.8±7.9%	80.7±10.2%	78.5%	78.6%	80.4%	77.1%
T0188	1O13	α/β	107	77.8±11.6%	74.3±14.3%	80.4%	73.6%	77.6%	73.3%
T0189	1O14	α/β	319	77.9±8.9%	80.5±11.1%	76.2%	67.6%	75.2%	61.0%
T0190			111	81.1±14.1%	77.7±17.6%	79.3%	78.3%	73.0%	50.6%
T0191	1NVT	α/β	282	77.7±10.8%	76.1±13.4%	78.0%	61.3%	75.9%	61.4%
T0192			170	69.8±8.9%	69.3±12.6%	71.2%	59.8%	72.4%	67.5%
T0193	1R72		206	71.5±8.7%	71.4±11.4%	69.9%	63.2%	74.3%	60.1%
T0195			292	73.3±9.2%	68.7±10.6%	74.7%	69.5%	77.7%	78.9%
Average			215.75	74.6±10.3 %	73.4±12.3 %	74.7±6.4%	68.2±9.0%	74.9±7.5%	73.1±10.3%