

# Constructing Structural Profiles for Protein Torsion Angle Prediction

Zafer Aydin<sup>1</sup>, David Baker<sup>2</sup> and William Stafford Noble<sup>3</sup>

<sup>1</sup>Department of Computer Engineering, Abdullah Gul University, 38080, Kayseri, Turkey

<sup>2</sup>Department of Biochemistry, University of Washington, Seattle, WA 98195 U.S.A.

<sup>3</sup>Department of Genome Sciences, Department of Computer Science and Engineering, University of Washington, Seattle, WA 98195 U.S.A.

**Keywords:** Protein Torsion Angle Prediction, Structural Frequency Profiles, Template Scoring, Profile-Profile Alignment.

**Abstract:** Structural frequency profiles provide important constraints on structural aspects of a protein and is receiving a growing interest in the structure prediction community. In this paper, we introduce new techniques for scoring templates that are later combined to form structural profiles of 7-state torsion angles. By employing various parameters of target-template alignments we improve the quality and accuracy of structural profiles considerably. The most effective technique is the scaling of templates by integer powers of sequence identity score in which the power parameter is adjusted with respect to the similarity interval of the target. Incorporating other alignment scores as multiplicative factors further improves the accuracy of profiles. After analyzing the individual strengths of various structural profile methods, we combine them with ab-initio predictions of 7-state torsion angles by a linear committee approach. We show that incorporating template information improves the accuracy of ab-initio predictions significantly at all levels of target-template similarity even when templates are distant from the target. Template scaling methods developed in this work can be applied in many other prediction tasks and in more advanced methods designed for computing structural profiles.

## 1 INTRODUCTION

Protein 3D structure prediction benefits greatly from prediction of various 1D and 2D structural attributes such as secondary structure, backbone torsion (dihedral) angles, solvent accessibility, disordered regions, and contact maps (Cheng et al., 2008). Methods that predict structural properties of proteins typically employ sequence-based frequency profiles in their feature sets to utilize information in similar proteins. These profiles can be in the form of position specific scoring matrices (PSSM) or hidden Markov models (HMM) and can be derived by aligning the amino acid sequence of the query with sequences in a large protein database using an efficient algorithm such as PSI-BLAST (Altschul et al., 1997) or HHblits (Remmert et al., 2011). Despite the many efforts for improving the quality of sequence-based alignments and their profiles, the accuracy of 1D and 2D predictions has come to saturation due to the difficulty of eliminating false positives especially when the query sequence diverges from those in the protein database considerably. Recently, there has been a growing interest in using structural profiles as input features for predict-

ing various structural characteristics of proteins. A structural frequency profile is a position specific scoring matrix (PSSM) that is constructed from the structural labels of templates (*i.e.*, hit proteins) obtained by aligning the target (*i.e.*, query) against a set of proteins. To date structural profiles have been derived mainly for protein secondary structure (Li et al., 2012; Cong et al., 2013); backbone structural motifs, solvent accessibility, contact density (Mooney and Pollastri, 2009); and shape strings (Sun et al., 2012).

To construct a structural profile, the occurrence frequencies of template residues are accumulated followed by a normalization step. Methods that have been developed for this task mainly use Laplacian counts, which is a technique that gives equal weights to templates (Li et al., 2012). As an alternative to the Laplacian count method, a new scoring technique has been proposed which scale the templates by the third power of sequence identity score and the structural quality information (Pollastri et al., 2007; Walsh et al., 2009). In this paper, we propose new structural profile methods for 7-state torsion angles of proteins by incorporating various score terms of HH-search alignments (Soding, 2005) and by adjusting

the power parameter according to the target-template similarity.

Despite the variety of methods proposed for predicting backbone torsion angles of proteins (Singh et al., 2014; Song et al., 2012; Wu and Zhang, 2008a; Faraggi et al., 2012; Shen et al., 2009; Bjanskii et al., 2006), less effort has been made to systematically incorporate structurally related templates into torsion angle predictions (Mooney and Pollastri, 2009). To the best of our knowledge, there is no work in the literature that inspects the accuracy of torsion angle predictions at all levels of target-template similarity (*i.e.*, from easy to difficult targets). Therefore after deriving structural profiles, we combine them with ab-initio predictions of a two-stage classifier using a linear committee approach. Our method is able to generate specific and effective predictions for targets at all difficulty levels. We achieve this by adjusting the power parameter and the weight of the structural profiles with respect to the similarity interval of the target.

## 2 METHODS

### 2.1 Backbone Torsion Angles

Each residue (*i.e.*, amino acid) has three associated torsion angles:  $\phi$ ,  $\psi$ , and  $\omega$ . The angle  $\phi$  denotes rotation about the  $C_{\alpha}$ -N bond of the residue,  $\psi$  denotes rotation about the bond linking  $C_{\alpha}$  and the carbonyl carbon, and  $\omega$  denotes rotation about the bond between the carbonyl carbon of the current residue and the nitrogen of the next residue. We compute  $\phi$ ,  $\psi$ , and  $\omega$  from the 3-D coordinate information in Protein Data Bank (PDB), which is the database of solved protein structures. Each of these angles is constrained to the range  $[-180, 180]$ .

Following (Blum et al., 2008), we first subdivided residues into five torsion angle classes, which represent the major clusters observed in PDB. However, to reduce the imbalance in the sizes of these classes, we further subdivided the two most common labels (A and B) according to whether the secondary structure class is loop or not. The resulting seven labels are described in Table 1.

### 2.2 Torsion Angle Class Prediction

Based on the definition given in Table 1, the 7-state torsion angle prediction problem can be stated as follows. For a given protein, the goal is to assign to each amino acid a torsion angle label from the alphabet  $\{L, A, M, B, E, G, O\}$  as shown in Fig. 1.

```
LWGLVKQGLKCEDCGMNVHKKREKVANLC
MMELMGLBBBBLLLGMBBMAAAAALLMMLMO
```

Figure 1: **7-state torsion angle class prediction problem.** The first row shows the amino acid sequence of the target and the second row is the sequence of 7-state torsion angle labels, which are defined according to Table 1.

## 2.3 Alignment Methods

### 2.3.1 Deriving Templates for Structural Profiles

In this paper, we used the HHsearch method (Soding, 2005) to detect the templates that are similar to a given target. HHsearch first derives an HMM-profile for the target and aligns it against a database of HMM-profiles (Soding, 2005). At the end of the alignment, it ranks the templates (*i.e.* hits) according to a probability score ranging from 0% to 100% and reports the ones that score above a threshold. An example alignment is shown in Fig. 2. We used the following commandline to compute HMM-HMM alignments for each target: `./hhsearch -i protein.hhm -d hhm3 -o protein.hhr -cpu 2 -mact 0.05 -ssw 0.11 -atab protein.start.tab -realign -E 100 -cov 20 -b 20`. We then selected the HMM-HMM alignments that score above the given threshold as the templates. Note that, HHsearch uses predicted secondary structure to be able to compute sensitive HMM-HMM alignments. We used the PSIPRED version 2.61 (Jones, 1999) to predict secondary structures. All these alignments were generated in 2011. Further details on HHsearch and the HMM-HMM alignments can be found in the corresponding documentation (Soding, 2006; Soding et al., 2012).

### 2.3.2 Generating Position-Specific Scoring Matrices for the Ab-Initio Method

We employed PSSMs generated by the PSI-BLAST (Altschul et al., 1997) and HHMAKE (Soding, 2005) algorithms as input features. We used BLAST version 2.2.20 and the NCBI's non-redundant (NR) database dated June 2011 to generate PSI-BLAST PSSMs. We generated the HMM-profiles by HHsearch version 1.5.1 (Soding, 2005). Note that in deriving the HHMAKE PSSMs we did not perform any HMM-HMM alignments. After deriving PSSMs we scaled them to the interval  $[0, 1]$  by applying a sigmoidal transformation. Detailed descriptions of the PSSMs and the sigmoidal transformation can be found in (Aydin et al., 2011).



such that  $i \in \{L, A, M, B, E, G, O\}$  is the torsion angle class,  $j$  is the residue position of the target,  $A(j, k)$  is the residue of the  $k^{\text{th}}$  database protein aligned to the  $j^{\text{th}}$  position,  $T(j, k)$  is the corresponding torsion angle class of the template, and  $\delta(T(j, k), i)$  is the Kronecker delta function defined as

$$\delta(t, i) = \begin{cases} 1 & \text{if } t = i \\ 0 & \text{otherwise} \end{cases} \quad (2)$$

In other words, each torsion angle label that is aligned to the  $j^{\text{th}}$  position of the target contributes by a count of 1, which is also known as the Laplacian count method. Once the count matrix is obtained it is normalized so that each column sums to 1. This is formulated as

$$M_a(i, j) = \begin{cases} \frac{C(i, j)}{\sum_i C(i, j)} & \text{if } |A(j)| > 0 \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

where  $A(j)$  is the set of all residues aligned to the  $j^{\text{th}}$  residue of the target,  $|A(j)|$  is the number of residues in  $A(j)$ , and  $M_a$  is the normalized count matrix. In this formulation, the residues of the target are divided into two categories. In the first group, we have ‘‘aligned’’ positions (represented by the condition  $|A(j)| > 0$ ) where at least one residue is aligned from a database protein and in the second group there is the set of ‘‘un-aligned’’ positions (*i.e.*, the case where  $|A(j)| = 0$ ) for which no residues are aligned from any hits. The second condition is realized for positions that correspond to gapped regions and for positions that are left out of the aligned regions when a local alignment algorithm is employed.

After the normalization step, the structural profile matrix can be computed as

$$M(i, j) = \begin{cases} M_a(i, j) & \text{if } |A(j)| > 0 \\ M_b(i, j) & \text{otherwise} \end{cases} \quad (4)$$

where  $M_b(i, j)$  is the background probability of aligning a template residue with torsion class  $i$  to the  $j^{\text{th}}$  residue of the target. In this paper, we use predictions from the ab-initio classifier for the background distribution of torsion angle labels.

#### 2.4.2 Weighing Hits by Integer Powers of Sequence Identity Scores

A second method for computing structural profiles weights templates by integer powers of the sequence identity score. In HHsearch, this score is computed for each target template pair and is represented by the ‘‘Identities’’ field as shown in Fig. 2. We first divide this score by 100 and convert it to a weight value. We

then compute an integer power of this weight, which is used to scale templates that contribute to the structural profile. This is expressed in the equations below

$$C(i, j) = \sum_{A(j, k)} \theta(T(j, k), i) \quad (5)$$

$$\theta(t, i) = \begin{cases} I^a & \text{if } t = i \\ 0 & \text{otherwise} \end{cases} \quad (6)$$

where  $C$  is the count matrix,  $\theta$  is the new occurrence count function replacing the Kronecker delta in Eq. 2,  $I$  is the sequence identity score of the  $k^{\text{th}}$  template and  $a$  is an integer that represents the strength of the amplification one wishes to impose on the structurally similar templates. The remaining terms are the same as their counterparts in Eqs. 1 and 2. This type of template scaling has two benefits. The first one is related to scaling templates by sequence identity scores, which increases the contribution of structurally closer templates while reducing the votes of distant ones. The second benefit comes by taking integer powers of  $I$ , which manages the situation where a handful of structurally similar templates are followed by many less similar or distant templates. In such a scenario, if we use the Laplacian counts as in Section 2.4.1 or weigh templates by sequence identity scores only (*i.e.*,  $a = 1$ ) the contribution of the similar templates would be suppressed by many structurally less similar candidates. To further amplify the effect of structurally similar templates and to reduce the contribution of false positives (*i.e.*, noise) it is useful to take integer powers of the sequence identity scores as formulated in Eq. 6. Once we compute the count matrix, we normalize it as in Eq. 3. All the other steps in deriving the structural profile are the same as in Section 2.4.1.

#### 2.4.3 Incorporating Quality of Templates

In addition to taking integer powers of the sequence identity score, it is also possible to include other weight factors to the score function in Eq. 6. One such measure assesses the experimental quality of the templates and is proposed in (Pollastri et al., 2007; Walsh et al., 2009). When we employ this approach to score the templates, Eq. 6 takes the following form:

$$\theta(t, i) = \begin{cases} \frac{I^a}{q} & \text{if } t = i \\ 0 & \text{otherwise} \end{cases} \quad (7)$$

where  $q$  is the quality of the template computed as X-ray resolution + R-factor / 20 as proposed in (Hobohm and Sander, 1994). According to this measure, a template with a higher experimental quality has a lower  $q$  parameter. Since this measure requires the X-ray

resolution of the template, we apply it to those templates that have been solved by the X-ray method only ignoring the remaining templates for the target.

#### 2.4.4 Incorporating Other Alignment Scores

When two proteins are aligned to each other, typically several score terms are calculated for assessing the statistical significance including e-value, raw similarity score, and percentage of sequence identity. Employing these terms in scaling the templates could also be useful in constructing a structural profile. With this motivation, we first incorporated the e-value score into the occurrence count function by converting it to a multiplicative weight factor as in (Wu and Zhang, 2008b). This is formulated as

$$\theta(t, i) = \begin{cases} w_e I^a & \text{if } t = i \\ 0 & \text{otherwise} \end{cases} \quad (8)$$

where  $w_e$  is the e-value weight defined as

$$w_e = \begin{cases} 1 & \text{if } E < 10^{-10} \\ -0.05 \log_{10}(E) + 0.5 & \text{if } 10^{-10} \leq E < 10^{10} \\ 0 & \text{otherwise} \end{cases} \quad (9)$$

such that  $E$  is the E-value of the alignment. Note that we dropped the quality term as it did not bring any significant benefits, which is verified by our simulations. According to Eq. 9,  $w_e$  is set to 1 when the target-template similarity is above a certain threshold ( $E$ -value  $< 10^{-10}$ ) and decreases linearly as the E-value of the target-template alignment is greater than  $10^{-10}$  until it becomes considerably high (*i.e.*,  $10^{10}$ ), in which case it is set to zero.

In addition to the E-value, we also considered incorporating the overall raw similarity score of the alignment into our structural profiles. For this purpose, we modified the occurrence count function as

$$\theta(t, i) = \begin{cases} s w_e I^a & \text{if } t = i \\ 0 & \text{otherwise} \end{cases} \quad (10)$$

where  $s$  is the raw score of the alignment. For HH-search, this is the overall similarity score obtained at the end of the HMM-HMM alignment, which is computed as the sum of the similarities of the aligned profile columns minus the gap penalties (Soding, 2005). A slight variation of this approach normalizes the raw score with the length of the aligned region as

$$\theta(t, i) = \begin{cases} \frac{s}{L} w_e I^a & \text{if } t = i \\ 0 & \text{otherwise} \end{cases} \quad (11)$$

such that  $L$  is the length of the aligned region and is given as the field denoted as ‘‘Aligned\_columns’’ in HHsearch’s output (see Fig. 2).

#### 2.4.5 Scaling Columns of the Alignment

Up to this point, we scaled the templates uniformly throughout the aligned positions without discriminating the individual columns of the alignment. In this section, we explain an approach for amplifying local regions within an alignment that could potentially contribute more accurate torsion label information and suppressing those that could be locally more distant. For this purpose, we include the similarity score between the aligned residues from a BLOSUM matrix into the occurrence count function as formulated below

$$\theta(t, i) = \begin{cases} \frac{s}{L} w_e I^a e^b & \text{if } t = i \\ 0 & \text{otherwise} \end{cases} \quad (12)$$

where  $b$  is the similarity score such as BLOSUM matrix score between the  $j^{\text{th}}$  residue of the target and the residue in the template that is aligned to the target residue. When the two residues are biologically similar to each other we expect this score term to be larger than the term obtained from dissimilar pairings. This approach has the potential to amplify local matches between motifs that are common both in the target and the template. Note that normalizing the sequence alignment score with  $L$  is optional. For instance, a slightly modified version of Eq. 12 does not perform such type of normalization:

$$\theta(t, i) = \begin{cases} s w_e I^a e^b & \text{if } t = i \\ 0 & \text{otherwise} \end{cases} \quad (13)$$

## 2.5 Prediction Model

### 2.5.1 Ab-initio Predictor

Our *ab-initio* torsion class predictor is a hybrid architecture, in which four dynamic Bayesian network (DBN) models are combined with a neural network. Two of the DBNs use PSSMs derived by PSI-BLAST (Altschul et al., 1997) and the other two use PSSMs from the HHMAKE module of the HHsearch method (Soding, 2005). Details of the *ab-initio* predictor can be found in (Aydin et al., 2011; Aydin et al., 2012). For simplicity we treat the output signal of the neural network as a probability distribution due to the constraints it satisfies (*i.e.*, it sums to 1 and takes values from 0 to 1). This distribution is denoted as  $P_a(t_j|x)$ , which represents the *ab-initio* likelihood of the  $j^{\text{th}}$  residue to have  $t_j$  as the torsion angle label given  $x$ , the set of input features around position  $j$ . Hence our neural network predicts the torsion angle label of the residue at the center of the feature window by select-

ing the particular label with the maximum discriminant score at the output layer. This is formulated as

$$t_j^* = \arg \max_{t_j} P_a(t_j|x). \quad (14)$$

### 2.5.2 Committee Predictor

The committee predictor combines the ab-initio predictions of torsion angle classes with the structural frequency profile according to the following equation:

$$P_c(t_j|x,y) = \begin{cases} (1-\lambda)P_a(t_j|x) + \lambda P_s(t_j|y) & \text{if aligned} \\ P_a(t_j|x) & \text{if unaligned} \end{cases} \quad (15)$$

where  $P_c(t_j|x,y)$  is the combined likelihood of having torsion angle label  $t_j$  for the residue at position  $j$ ,  $\lambda$  is the weight assigned to the structural profile,  $x$  is the feature set of the ab-initio predictor described in Section 2.5.1,  $y$  is the amino acid sequence of the target,  $P_a(t_j|x)$  is the distribution of torsion angle classes obtained from the ab-initio predictor, and  $P_s(t_j|y)$  is the structural profile computed from the templates. According to this equation, the combined likelihood is the weighted average of the ab-initio predictions and the structural profile for positions that are aligned to at least one template residue and it becomes equal to the likelihood of the ab-initio predictor only for the remaining positions.

After computing  $P_c(t_j|x,y)$ , we predict the torsion angle class of the  $j^{\text{th}}$  residue as the particular label that maximizes  $P_c(t_j|x,y)$  as

$$t_j^* = \arg \max_{t_j} P_c(t_j|x,y), \quad (16)$$

where  $t_j^*$  is the final prediction of the committee method.

## 2.6 Datasets

### 2.6.1 PDB-PC90 Dataset

To obtain the PDB-PC90 dataset, we used the PISCES server (Wang and Dunbrack, Jr., 2003; Wang and Dunbrack, Jr., 2005) with the following set of criteria: percent identity threshold of 90%, resolution cut-off of 2.5 Å, and R-value cutoff of 1.0. We also used PISCES to filter out non-X-ray and  $C_\alpha$ -only structures and to remove short (< 30 amino acids) and long (> 10000 amino acids) chains. This dataset contained 17056 chains.

### 2.6.2 Training and Test Set

We randomly selected 5161 proteins from the PDB-PC90 dataset. Among those, we randomly selected

a set of 994 proteins for the test set. From the set of 5161 proteins, we then removed those proteins that are similar to the test set using a 10% sequence identity threshold. The remaining set contained 4205 chains, which is used to train our ab-initio prediction method. We computed the HHsearch alignments for the set of 994 proteins, which are used for computing the structural profiles of torsion angle classes and for predicting the torsion angle classes.

### 2.6.3 Similarity Intervals and Subsets of the Test Set

To distinguish easy targets from difficult ones, we defined similarity intervals using the HHsearch alignments from half of the proteins in our test set (see Section 2.6.2). For each target, we first selected the maximum sequence identity score from the set of target-template alignments. Then we ranked those scores and defined percentile intervals of sequence identity with increments of 5%. This initially produced a total of 20 intervals. We combined the eighth and ninth intervals as the maximum sequence identity scores for those targets were very close to each other. We also combined the tenth up to the twentieth intervals since the maximum sequence identity score was 100% for all the targets in those bins. This procedure resulted in a total of 9 sequence identity intervals. In the last step, we further reduced the number of intervals to 5 by combining the 2nd and 3rd, 4th and 5th, and 7th up to 9th intervals. The resulting intervals are tabulated below

According to Table 2, the first interval represents targets with the most distant templates (*i.e.*, those with the maximum sequence identity score of 26% or less) and the last interval represents targets that contain highly similar templates (*i.e.*, those with the maximum sequence identity score greater than 80%). Based on this binning, we further divided our test set of 994 proteins (see the previous section) into 5 subsets such that each contained those targets that fall into one and only one of the intervals defined in Table 2. The number of proteins and amino acids in each of these subsets are summarized in Table 3.

Note that the number of proteins in each subset is not uniformly the same (especially true for the last set that contains targets from the ‘‘High’’ category) mainly because datasets have been constructed by random sampling from PDB without enforcing specific constraints for having equal number of samples in each interval. Nonetheless we have enough samples in each subset mainly because the torsion angle class prediction is performed on each residue separately. Furthermore the proportion of target positions that are aligned to at least one template residue is con-

Table 2: **Intervals of sequence identity scores** The intervals are defined by selecting the target-template alignments with maximum sequence identity scores followed by sorting these scores in ascending order. Percentile increments of 5% results in a total of 20 bins, which are further reduced to 5 intervals.

Interval	Percentiles (%)	Max Identity (%)
Low	0-5	0.0-26.0
Medium-Low	5-15	26.0-35.0
Medium	15-25	35.0-50.0
Medium-High	25-30	50.0-80.0
High	30-100	80.0-100.0

Table 3: **The five subsets of the test set with 994 proteins.** The number of proteins, the total number of amino acids and the number of amino acids that are aligned to at least one template residue are shown for each subset. The subsets are derived based on the intervals defined in Table 2.

Subset	# proteins	# residues	# aligned res.
Low	56	12903	12665
Medium-Low	99	21792	21682
Medium	95	22596	22561
Medium-High	62	15326	15295
High	682	160037	159993
Total	994	232654	232196

siderably high. This shows that a structural profile column is computed using the aligned templates for most of the target residues.

### 3 RESULTS

We first compare the torsion angle label accuracy of the structural profiles on positions that are aligned to templates only. For this purpose, we implemented the profile methods summarized in Table 4.

In SP4, we modify the power of the sequence identity score term (*i.e.*,  $a$ ) according to the similarity interval the target belongs to. For this purpose, we use the following mapping to define  $a$ :

$$a = \begin{cases} 1 & \text{if target in Low Interval} \\ 3 & \text{if target in Medium-Low Interval} \\ 5 & \text{if target in Medium Interval} \\ 7 & \text{if target in Medium-High Interval} \\ 9 & \text{if target in High Interval} \end{cases} \quad (17)$$

where the interval of the target is defined according to Table 2. In SP9 and SP10 we use the BLOSUM62 matrix to scale individual columns of the alignment as formulated in Eqs. 12 and 13. In SP9, we employed the BLOSUM scores uniformly for all columns of the alignment whereas in SP10, we utilized the BLOSUM scores for targets in the ‘‘High’’ interval only. If the target belongs to one of the remaining four intervals then we turn off this score term in Eq. 13. In that case this equation takes the following form

$$\theta(t, i) = \begin{cases} sw_e I^a & \text{if } t = i \\ 0 & \text{otherwise} \end{cases} \quad (18)$$

Once we compute a structural profile, we predict the torsion angle class of the aligned target residues by selecting the particular label that yields the maximum value in the corresponding column of the profile.

Following these definitions, the torsion angle class prediction accuracy of the structural profiles listed in Table 4 is summarized in Table 5 below. In this table, Overall 1 is the number of correctly predicted amino acids divided by the total number of amino acids for which a structural profile column is computed (*i.e.*, those that are aligned to at least one template). The second up to the sixth columns show accuracies for the five similarity intervals and are computed on each subset of the test set. Overall 2 is the average of the five accuracies obtained for Low to High intervals. It estimates the accuracy we would obtain had we used equal number of amino acids for each of the five intervals. Based on these results, the most accurate structural profile method is SP10 though other methods such as SP7, SP4 and SP8 are also quite effective. SP10 outperforms SP1, (the Laplacian count method), by 14.53%, which is a statistically significant improvement. It is also better than SP5 by 1.2%, which was proposed in (Pollastri et al., 2007; Walsh et al., 2009). This improvement is also statistically significant (with a p-value < 0.0001 from a two-tailed Z-test at a significance level of 0.01).

After establishing that the new structural profile methods contain more accurate torsion angle information than the approaches proposed in the literature, we evaluated the accuracy when the structural profiles are combined with an ab-initio predictor. For this purpose, we trained our ab-initio method us-

Table 4: **The implemented structural profile methods and their descriptions.** Further details can be found in Section 2.4.

Structural Profile and Description
SP1: Eq. 1, Laplacian
SP2: Eq. 6, $a = 1$
SP3: Eq. 6, $a = 3$
SP4: Eq. 6, $a$ varies wrt similarity interval
SP5: Eq. 7, quality, $a = 3$ , (Pollastri et al.)
SP6: Eq. 8, E-value, $a = 3$
SP7: Eq. 10, E-value, align. score, $a = 3$
SP8: Eq. 11, E-value, norm. align. score, $a = 3$
SP9: Eq. 12, E-value, norm. align. score, BLOSUM62, $a = 3$
SP10: Eq. 13, E-value, align. score, $a$ varies, BLOSUM62 scores if target is in High interval only

Table 5: **7-state torsion angle prediction accuracy of structural profiles.** L: Low, ML: Medium-Low, M: Medium, MH: Medium-High, H: High. Only target residues that are aligned to at least one template are considered.

Profile	L	ML	M	MH	H	Overall 1	Overall 2
SP1	66.49	69.83	71.94	74.46	75.66	74.17	71.68
SP2	66.85	70.71	73.93	80.07	81.96	79.18	74.70
SP3	66.68	72.20	77.82	86.64	92.87	87.64	79.24
SP4	66.85	72.20	79.40	87.92	93.47	88.30	79.97
SP5	66.53	72.13	77.28	87.07	92.73	87.50	79.15
SP6	67.15	73.26	78.89	87.13	93.06	88.03	79.90
SP7	66.48	74.11	80.23	87.66	93.30	88.40	80.36
SP8	66.59	73.91	80.13	87.31	93.25	88.33	80.24
SP9	65.05	72.19	78.30	86.39	93.68	88.14	79.12
SP10	66.55	74.11	80.71	87.36	93.69	88.70	80.48

ing the training set described in Section 2.6.2 and computed 7-state torsion angle predictions on all the amino acids of the test set. We then combined those predictions with a structural profile as in Eq. 15. For the target residues that were not aligned to any template, we simply took predictions from the ab-initio method. Regarding the  $\lambda$  parameter (*i.e.*, the weight of the structural profile) we considered two possibilities. The first approach sets  $\lambda$  to 0.5 and the second one modifies it according to the similarity interval of the target according to the following function

$$\lambda = \begin{cases} 0.5 & \text{if target in Low Interval} \\ 0.6 & \text{if target in Medium-Low Interval} \\ 0.7 & \text{if target in Medium Interval} \\ 0.8 & \text{if target in Medium-High Interval} \\ 0.9 & \text{if target in High Interval} \end{cases} \quad (19)$$

In this equation,  $\lambda$  is gradually increased as the similarity interval of the target approaches to the ‘‘High’’ interval thereby giving more weight to the structural profile than the ab-initio predictor. Table 6 summarizes the accuracy of committee predictors that combine the ab-initio method with various structural profiles. In addition to the overall accuracy measure, we also included the segment overlap (SOV) measure that is used in 1D structure prediction to assess the accuracy at the segmental level (Zemla et al., 1999). The SOV measure depicts how well the predicted torsion label segments match the true segments and is biologically more meaningful than the residue level

accuracy.

According to this table, the ab-initio+SP10 method (with variable  $\lambda$  parameter) is better than the ab-initio+SP5 method in all categories. The improvements are 2.78% in Overall 1, 3.40% in Overall 2, 3.57% in SOV, 1.65% in Low interval, 2.74% in Medium-Low interval, 4.90% in Medium interval, 4.82% in Medium-High interval and 2.86% in High interval. When ab-initio+SP10 is compared with ab-initio+SP1 (*i.e.*, the Laplacian method) for  $\lambda = 0.5$ , the improvements are 12.20% in Overall 1, 6.68% in Overall 2, 13.62% in SOV, 0.12% in Low interval, 2.59% for Medium-Low interval, 4.59% in Medium interval, 10.43% in Medium-High interval and 15.66% in High interval. Adjusting the  $\lambda$  parameter with respect to the similarity interval was particularly useful for the ab-initio+SP5 method. In other words, when  $\lambda$  is set to 0.5 uniformly for all similarity intervals, the accuracy of ab-initio+SP5 dropped significantly higher than the ab-initio+SP10. This shows that the proposed structural profile SP10 is more useful than SP5 when combined with the ab-initio method. This is because torsion label errors of SP10 and the ab-initio method overlap less as compared to SP5 and therefore SP10 provides a better complement to the ab-initio predictor. Another observation one can make is the improvement over the ab-initio method when structural profiles are incor-

Table 6: **7-state torsion angle prediction accuracy of methods that incorporate structural profiles with ab-initio predictions.** L: Low, ML: Medium-Low, M: Medium, MH: Medium-High, H: High. All target residues in the test set are considered.

Method	L	ML	M	MH	H	Overall 1	Overall 2	SOV
Ab-initio	72.36	73.96	73.58	73.35	74.01	73.83	73.45	71.33
Ab-initio + SP1, $\lambda = 0.5$	74.47	75.42	76.10	76.94	77.96	77.28	76.18	74.49
Ab-initio + SP5, $\lambda = 0.5$	72.39	74.19	74.19	75.53	78.28	76.99	74.92	74.52
Ab-initio + SP10, $\lambda = 0.5$	74.59	78.01	80.69	87.37	93.62	89.10	82.86	88.11
Ab-initio + SP5, $\lambda$ as in Eq. 19	72.94	75.21	77.22	83.86	90.99	86.70	80.04	85.19
Ab-initio + SP10, $\lambda$ as in Eq. 19	74.59	77.95	82.12	88.68	93.85	89.48	83.44	88.76

porated. This is true even for the Low interval (an improvement of 2.23%) and is partly because of the sensitive nature of HMM-HMM profile alignments and also because HHsearch uses predicted secondary structure from PSIPRED (Jones, 1999) to align a pair of HMMs.

In addition to the structural profile methods described above, we also considered three other approaches. The first one incorporates the probability score of HMM-profile alignments into Eq. 13 as a multiplicative factor to globally scale the templates and the second method incorporates the column score of HHsearch alignments to amplify local regions that are well conserved (e.g. motifs). These two approaches did not bring any reasonable change in the accuracy measures (result not shown). As a third approach we considered employing Henikoff weights (Henikoff and Henikoff, 1994) to scale the count information of templates before constructing the structural profiles. We applied this weighting procedure for the following three scenarios: (1) weights based on matched residues only, (2) weights based on torsion angle labels only, (3) weights based on residue and torsion angle tuples. Unfortunately, in all three cases, torsion angle prediction accuracy was significantly lower than the level achieved by other scaling methods considered in this paper (result not shown). Note that we did not consider utilizing a background distribution in Eq. 4 for torsion angle labels mainly because we use predictions from our ab-initio method, which would eventually contain a more accurate torsion angle representation than a simple background distribution.

Finally, we would like to state that we are unable to compare our torsion angle class predictor with the literature mainly because there is no other work that performs 7-state torsion angle prediction on the same set of alphabet. However we had shown in an earlier paper that a 5-state version of our predictor provides results comparable to the state-of-the-art in the ab-initio setting (Aydin et al., 2012).

## 4 CONCLUSIONS

In this paper, we propose novel methods for scaling templates to construct structural profiles of torsion angle states. Though we use the score terms in HHsearch method, our approach is generic and most of the structural profile methods proposed in this work can also be implemented using other alignment methods including PSI-BLAST. Second, the scaling techniques can be applied in many other tasks such as secondary structure prediction, solvent accessibility prediction, contact map prediction, and 3D structure prediction. Third, they can easily be incorporated into other methods that have been developed for deriving structural profiles from templates.

The proposed methods can be improved in several ways. First of all, certain parameters of the method can be optimized such as the power of sequence identity score and the weight that is used to combine structural profiles with ab-initio predictions. Additionally, the templates can be scaled in a position-specific manner using the confidence scores, which are now available in HHblits (the new version of HHsearch). A third technique uses templates that score within a window only instead of taking all the templates that scored above a threshold. This approach can be combined with the existing techniques presented in this work to reduce the computational cost. Finally, instead of using a linear model, the ab-initio predictions can be combined with structural profiles using more advanced models such as neural networks. We believe that all these efforts will potentially improve the accuracy of structure prediction tasks further.

## ACKNOWLEDGEMENTS

This work is supported by grant 113E550 from 3501 TUBITAK National Young Researchers Career Award.

## REFERENCES

- Altschul, S. F., Madden, T. L., Schaffer, A. A., Zhang, J., Zhang, Z., Miller, W., and Lipman, D. J. (1997). Gapped BLAST and PSI-BLAST: A new generation of protein database search programs. *25*:3389–3402.
- Aydin, Z., Singh, A., Bilmes, J., and Noble, W. S. (2011). Learning sparse models for a dynamic Bayesian network classifier of protein secondary structure. *BMC Bioinformatics*, *12*:154.
- Aydin, Z., Thompson, J., Bilmes, J., Baker, D., and Noble, W. S. (2012). Protein torsion angle class prediction by a hybrid architecture of bayesian and neural networks. In *13th International Conference on Bioinformatics and Computational Biology*.
- Berjanskii, M. V., Neal, S., and Wishart, D. S. (2006). PREDITOR: a web server for predicting protein torsion angle restraints. *Nucleic Acids Research*, *34*(Web Server Issue):W:63–69.
- Blum, B., Jordan, M., Kim, D., Das, R., Bradley, P., and Baker, D. (2008). Feature selection methods for improving protein structure prediction with Rosetta. In Platt, J., Koller, D., Singer, Y., and Roweis, S., editors, *Advances in Neural Information Processing Systems 20*, pages 137–144. MIT Press, Cambridge, MA.
- Cheng, J., Tegge, A. N., and Baldi, P. (2008). Machine learning methods for protein structure prediction. *IEEE Reviews in Biomedical Engineering*, *1*:41–49.
- Cong, P., Li, D., Wang, Z., Tang, S., and Li, T. (2013). Spssm8: An accurate approach for predicting eight-state secondary structures of proteins. *Biochimie*, *95*(12):2460–2464.
- Faraggi, E., Zhang, T., Yang, Y., Kurgan, L., and Zhou, Y. (2012). SPINE X: improving protein secondary structure prediction by multistep learning coupled with prediction of solvent accessible surface area and backbone torsion angles. *PLoS One*, *7*(2):e30361.
- Henikoff, S. and Henikoff, J. G. (1994). Position-based sequence weights. *243*:574–578.
- Hobohm, U. and Sander, C. (1994). Enlarged representative set of protein structures. *Protein Science*, *3*:522–524.
- Jones, D. T. (1999). Protein secondary structure prediction based on position-specific scoring matrices. *292*:195–202.
- Li, D., Li, T., Cong, P., Xiong, W., and Sun, J. (2012). A novel structural position-specific scoring matrix for the prediction of protein secondary structures. *Bioinformatics*, *28*(1):32–39.
- Mooney, C. and Pollastri, G. (2009). Beyond the twilight zone: Automated prediction of structural properties of proteins by recursive neural networks and remote homology information. *Proteins: Structure, Function, and Bioinformatics*, *77*:181–190.
- Pollastri, G., Martin, A. J. M., Mooney, C., and Vullo, A. (2007). Accurate prediction of protein secondary structure and solvent accessibility by consensus combiners of sequence and structure information. *BMC Bioinformatics*, *8*(201).
- Remmert, M., Biegert, A., Hauser, A., and Soding, J. (2011). Hhblits: lightning-fast iterative protein sequence searching by hmm-hmm alignment. *Nature Methods*, *9*(2):173–175.
- Shen, Y., Delaglio, F., Cornilescu, G., and Bax, A. (2009). TALOS+: a hybrid method for predicting protein backbone torsion angles from nmr chemical shifts. *Journal of Biomolecular NMR*, *44*(4):213–223.
- Singh, H., Singh, S., and Raghava, G. P. S. (2014). Evaluation of protein dihedral angle prediction methods. *PLoS One*, *9*(8):e105667.
- Soding, J. (2005). Protein homology detection by HMM-HMM comparison. *Bioinformatics*, *21*:951–960.
- Soding, J. (2006). Quick guide to HHsearch. <ftp://toolkit.genzentrum.lmu.de/pub/HHsearch/old/HHsearch/HHsearch1.5.1/HHsearch-guide.pdf>.
- Soding, J., Remmert, M., and Hauser, A. (2012). HH-suite for sensitive sequence searching based on hmm-hmm alignment. <ftp://toolkit.genzentrum.lmu.de/pub/HH-suite/hhsuite-userguide.pdf>.
- Song, J., Tan, H., Wang, M., Webb, G. I., and Akutsu, T. (2012). TANGLE: two-level support vector regression approach for protein backbone torsion angle prediction from primary sequences. *PLoS One*, *7*(2):e30361.
- Sun, J., Tang, S., Xiong, W., Cong, P., and Li, T. (2012). Dsp: a protein shape string and its profile prediction server. *Nucleic Acids Research*, *40*(W1):W298–W302.
- Walsh, I., Bau, D., Martin, A. J. M., Mooney, C., Vullo, A., and Pollastri, G. (2009). Ab initio and template-based prediction of multi-class distance maps by two-dimensional recursive neural networks. *BMC Structural Biology*, *9*(5).
- Wang, G. and Dunbrack, Jr., R. L. (2003). PISCES: a protein sequence culling server. *Bioinformatics*, *19*:1589–1591. Web server at <http://dunbrack.fccc.edu/PISCES.php>.
- Wang, G. and Dunbrack, Jr., R. L. (2005). PISCES: recent improvements to a pdb sequence culling server. *Nucleic Acids Res.*, *33*:W94–W98. Web server at <http://dunbrack.fccc.edu/PISCES.php>.
- Wu, S. and Zhang, Y. (2008a). ANGLOR: A composite machine-learning algorithm for protein backbone torsion angle prediction. *PLoS One*, *3*(10):e3400.
- Wu, S. and Zhang, Y. (2008b). MUSTER: Improving protein sequence profile-profile alignments by using multiple sources of structure information. *Proteins: Structure, Function, and Bioinformatics*, *72*(2):547–556.
- Zemla, A., Venclovas, C., Fidelis, K., and Rost, B. (1999). A modified definition of Sov, a segment-based measure for protein secondary structure prediction assessment. *Proteins*, *34*:220–223.