

Predicting backbone C# angles and dihedrals from protein sequences by stacked sparse auto-encoder deep neural network

Author

Lyons, James, Dehzangi, Abdollah, Heffernan, Rhys, Sharma, Alok, Paliwal, Kuldip, Sattar, Abdul, Zhou, Yaoqi, Yang, Yuedong

Published

2014

Journal Title

Journal of Computational Chemistry

DOI

<https://doi.org/10.1002/jcc.23718>

Copyright Statement

© 2014 Wiley Periodicals, Inc.. This is the accepted version of the following article: Predicting backbone Ca angles and dihedrals from protein sequences by stacked sparse auto-encoder deep neural network Journal of Computational Chemistry, Vol. 35(28), 2014, pp. 2040-2046, which has been published in final form at [dx.doi.org/10.1002/jcc.23718](https://doi.org/10.1002/jcc.23718).

Downloaded from

<http://hdl.handle.net/10072/64247>

Griffith Research Online

<https://research-repository.griffith.edu.au>

Predicting backbone C α angles and dihedrals from protein sequences by stacked sparse auto-encoder deep neural network

James Lyons^a, Abdollah Dehzangi^{a,b}, Rhys Heffernan^a, Alok Sharma^{a,c}, Kuldip Paliwal^a, Abdul Sattar^{a,b}, Yaoqi Zhou^{d*}, Yuedong Yang^{d*}

^aInstitute for Integrated and Intelligent Systems, Griffith University, Brisbane, Australia

^bNational ICT Australia (NICTA), Brisbane, Australia

^cSchool of Engineering and Physics, University of the South Pacific, Private Mail Bag, Laucala Campus, Suva, Fiji

^dInstitute for Glycomics and School of Information and Communication Technique, Griffith University, Parklands Dr. Southport, QLD 4222, Australia

*Corresponding authors (yaoqi.zhou@griffith.edu.au, yuedong.yang@griffith.edu.au)

ABSTRACT

Because a nearly constant distance between two neighbouring C α atoms, local backbone structure of proteins can be represented accurately by the angle between C α_{i-1} -C α_i -C α_{i+1} (θ) and a dihedral angle rotated about the C α_i -C α_{i+1} bond (τ). θ and τ angles, as the representative of structural properties of 3 to 4 amino-acid residues, offer a description of backbone conformations that is complementary to ϕ and ψ angles (single residue) and secondary structures (>3 residues). Here, we report **the first** machine-learning technique for sequence-based prediction of θ and τ angles. Predicted angles based on an independent test have a mean absolute error of 9° for θ and 34° for τ with a distribution on the θ - τ plane close to that of native values. The average root-mean-square distance of 10-residue fragment structures constructed from predicted θ and τ angles is only 1.9Å from **their corresponding** native structures. Predicted θ and τ angles are expected to be complementary to predicted ϕ and ψ angles and secondary structures for using in model validation and template-based as well as template-free structure prediction. The deep neural network learning technique is available as an on-line server called SPIDER (Structural Property prediction with Integrated DEep neuRal network) at <http://sparks-lab.org>.

INTRODUCTION

Template-based and template-free protein-structure prediction relies strongly on prediction of local backbone structures [1,2]. Protein local structure prediction is dominated by secondary structure prediction with its accuracy stagnant around 80% for more than a decade [3,4]. However, secondary structures are only a coarse-grained description of protein local structures in three states (helices, sheets and coils) that are somewhat arbitrarily defined because helices and sheets are often not in their ideal shapes in protein structures. This arbitrariness has limited the theoretically achievable accuracy of three-state prediction to 88-90% [4,5]. Moreover, predicted coil residues do not have a well-defined structure.

An alternative approach to characterize the local backbone structure of a protein is to employ three dihedral or rotational angles about the N-C α bond (ϕ), the C α -C bond (ψ), and the C-N bond (ω). A schematic illustration is shown in Figure 1. Because ω angles are restricted to 180° (the majority) or 0° due to rigid planar peptide bonds, two dihedral angles (ϕ and ψ) essentially determine the overall backbone structure. Unlike secondary structures, these dihedral angles (ϕ and ψ) can be predicted as continuous variables and their predicted accuracy has been improved over the years [6-8] so that it is closer to dihedral angles estimated according to NMR chemical shifts [9]. Predicted backbone dihedral angles were found to be more useful than predicted secondary structure as restrains for *ab initio* structure prediction [9,10]. It has also been utilized for **improving sequence alignment** [11], **secondary structure prediction** [3,12,13] and template-based structure prediction and fold

recognition [14-16]. However, unlike the secondary structure of proteins, ϕ and ψ are limited to the conformation of a single residue.

Two different angles can also be employed for representing protein backbones. As shown in Figure 1, they are the angle between $C\alpha_{i-1}-C\alpha_i-C\alpha_{i+1}$ (θ_i) and a dihedral angle rotated about the $C\alpha_i-C\alpha_{i+1}$ bond (τ_i). This **two-angle representation** is possible because neighbouring $C\alpha$ atoms mostly have a fixed distance (3.8Å) due to the fixed plane in $C\alpha_{i-1}-C-N-C\alpha_i$. These two inter-residue angles (θ and τ) reflect the conformation of four connected, neighbouring residues **that is** longer than **a** single-residue conformation represented by ϕ and ψ angles. By comparison, **a** conformation represented by helical **or** sheet residues **involves in an** undefined **number of residues** (4 for 3_{10} helix, 5 for α -helix, and **an** undefined number of residues for sheet residues). Thus, secondary structure, ϕ/ψ , and θ/τ provide complimentary local structural information along the backbone. Indeed, both predicted ϕ/ψ and secondary structure are useful for template-based structure prediction [14].

In this paper, we will develop the **first machine-learning** technique to predict θ and τ from protein sequences. **This** tool is needed not only because these two angles yield local **structural information** complementary to secondary structure and ϕ/ψ angles, but also because they have been widely employed in coarse-grained models for protein dynamics [17], folding [18], structure prediction [19,20], conformational analysis [21], and model validation [22]. **That is, accurate prediction of θ and τ** will be useful for template or template-free structure prediction as well as validation of predicted models. Using 4590 proteins for training and cross validation and 1199 proteins for an independent test, we have developed **a deep-learning neural-network-based method** that **achieved θ and τ angles within 9 and 34 degrees, in average, of their native values.**

METHOD

Data sets: In this study, we obtained a dataset of 5840 proteins with less than 25% sequence identity and X-ray resolution better than 2Å from the protein sequence culling server PISCES [23]. After removing 51 proteins with obsolete IDs or missing data, the final data set consists of 5789 proteins with 1,246,420 residues. We randomly selected 4590 proteins from this data set for training and cross-validation (TR4590) and employed the remaining 1199 proteins for an independent test (TS1199).

Deep neural-network learning. An Artificial Neural Network (ANN) consists of highly interconnected, multi-layer processing units called neurons. Each neuron combines its inputs with a non-linear sigmoid activation function to produce an output. Deep neural networks refer to feed-forward ANNs with three or more hidden layers. Multi-layer networks were not widely used because of the difficulty to train neural-network weights. This has changed due to recent advances through unsupervised weight initialization, followed by fine-tuned supervised training [24,25]. In this study, unsupervised weight initialization was done by stacked sparse auto-encoder. A stacked auto-encoder treats each layer as an auto-encoder that maps the layer's inputs back to themselves. During training auto-encoders a sparsity penalty **was** utilized to prevent learning of the identity function [26]. Initialised weights **were** then refined by standard back-propagation. The stacked sparse auto-encoder used in this study consists of three hidden layers with 150 hidden nodes in each layer (Figure 2). The input data **was** normalised so that each feature is in the range of 0 to 1. For residues near the ends of a protein, the features of the amino acid residue at the other end of the protein were duplicated so that a full window could be used. The learning rate was initialised to start at 0.5, and was then decreased as training progressed. In this study, we **used** the deep neural network Matlab toolbox implemented by Palm [27].

Input features. Each amino acid **was** described by a vector of input features that include 20 values from the Position Specific Scoring Matrix (PSSM) generated by PSI-BLAST [28] with three iterations of searching against non-redundant (NR) sequence database with an E-value cut off of 0.001. We also used seven representative amino-acid properties: a steric parameter (graph shape index), hydrophobicity, volume, polarizability, isoelectric point, helix probability, and sheet probability [29]. In addition, we employed predicted secondary structures (three probability values for helix, sheet and coils) and predicted solvent accessible surface area (one value) from SPINE-X [3]. That is, this is a vector of 31 dimensions per amino acid residue. As before, we also employed a window size of 21 amino acids (10 amino acids at each side of the target amino acid). This led to a total of 651 input features (21×31) for a given amino acid residue.

Output. Here we attempt to predict two angles. One is θ , the angle between three consecutive C α atoms of a protein backbone. The other one is τ , the dihedral angle between four consecutive C α atoms of protein backbone. Two angles are predicted at the same time. To remove the effect of periodicity, we employed four output nodes that correspond to $\text{Sin}(\theta)$, $\text{Cos}(\theta)$, $\text{Sin}(\tau)$, and $\text{Cos}(\tau)$, respectively. Predicted sine and cosine values were converted back to angles by using $\theta = \tan^{-1}[\text{sin}(\theta)/\text{cos}(\theta)]$ and $\tau = \tan^{-1}[\text{sin}(\tau)/\text{cos}(\tau)]$. Such transformation is widely used in signal processing and speech recognition [30].

Evaluation Methods: We investigated the effectiveness of our proposed method using 10-fold cross validation (TR4590) and independent test sets (TS1199). In 10-fold cross validation, TR4590 was divided into 10 groups. Nine groups were used as a training data set while the remaining group was used for test. This process was repeated 10 times until all the 10 groups were used once as the test data set. In addition to 10-fold cross validation, TR4590 was used as the training set and TS1199 was employed as an independent test set. Comparison between 10 fold cross validation and the test gives an indicator for the generality of the prediction tool. We evaluated the accuracy of our prediction by mean absolute error (MAE), the average absolute difference between predicted and experimentally determined angles. The periodicity of τ angles was taken care of by utilizing the smaller value of the absolute difference $d_i (= |\tau_i^{\text{Pred}} - \tau_i^{\text{Expt}}|)$ and $360 - d_i$ for average.

RESULT

Table 1 compares the results of ten-fold cross validation based on TR4590 and the independent test (TS1199). θ angles with a range of 0 to 180° were predicted significantly more accurate than τ angles with a range of -180° to 180°. The MAE is <9° for θ but 33-34° for τ . This level of accuracy can be compared to the baseline MAE values of 18.8° for θ and 86.2° for τ if θ and τ are assigned randomly according to their respective distributions. Accuracy for angles differs significantly in secondary structure types. The angles for helical residues have the highest accuracy (MAE<5° for θ and 17° for τ). The MAE for sheet residues is about twice larger than that for helical residues. Angles for coil residues have the largest error (τ in particular). Different levels of accuracy in different secondary structural types reflect the fact that helical structures are more locally stabilized than sheet structures while coil residues do not have a well-defined conformation. Similar trends were observed for prediction of backbone ϕ and ψ angles [6-9]. We also noted that MAEs from ten-fold cross validation and from the independent test are essentially the same. This indicates the robustness of the method trained. Thus, here and hereafter, we will present the result based on the independent test only.

Actual and predicted distributions of θ and τ angles for TS1199 are shown in Figure 3. Predicted and actual distributions agree with each other very well. Both predicted and actual peaks for θ angles are located at 92° and 119°, respectively. Actual peaks for τ angles are also in good agreement with those predicted ones at 50° and -164°, respectively. Predicted peaks, however, are slightly narrower than native peaks for all cases. Predicted and actual angle distributions also agree in a two-dimensional plane of θ and τ . As shown in Figure 4, the locations of three major populations were well captured by predicted distributions.

Table 2 lists the MAEs for 20 individual residue types along with their frequencies of occurrence in the TS1199 dataset. Glycine (G) has the largest MAE, corresponding to the fact that it is the most flexible residue due to lack of a side chain. Leucine (L), on the other hand, has the smallest MAE and interestingly also the most frequently occurred residue (9.2%). The angles for several other small hydrophobic residues [isoleucine (I), valine (V), and alanine (A)] are also in the pack of residues with smallest errors. There is no strong correlation between the MAE of an amino acid residue type and its frequency of occurrence.

In Figure 5, MAEs for predicted angles are shown as a function of relative solvent accessible surface area. MAEs for θ and τ have similar trend: two peaks separated by a valley (although in a smaller magnitude for θ). Both angles have the highest accuracy (the smallest error) at an intermediate range of solvent accessibility and the lowest accuracy (the largest error) at 90-100% solvent accessibility. The lowest accuracy at 90-100% solvent accessibility is likely due to the smallest number of residues at 90-100% solvent-accessible and 20% more coil residues in fully exposed residues [3].

Figure 6 displays the fraction of proteins with more than a given fraction of correctly predicted angles (θ and τ). Here, a correct prediction is defined as 36° or less from the actual angle. We use 36° as a cut off value because it is relatively easy for a conformational sampling technique to sample conformational changes within 36° . θ angles are always predicted within 36° for all residues in all proteins. 70% or more τ angles are predicted correctly for nearly 90% proteins. However, less than 10% proteins have 100% correctly predicted θ and τ .

θ and τ can also be calculated from backbone torsion angles ϕ and ψ by assuming $\omega = 180^\circ$. Thus, it is of interest to compare the accuracy of θ and τ predicted in this work with those calculated from predicted ϕ and ψ . For the TS1199 dataset, we found that the MAE values for θ and τ derived from ϕ and ψ predicted by SPINE X [3] are 9.6° and 37.7° , respectively. Thus, the angles predicted in this work (MAE = 8.6° and 33.6° , respectively) are about 10% more accurate in θ or τ than those calculated from ϕ and ψ predicted by SPINE X. The largest improvement by direct prediction of θ or τ as shown in Table 1 is in coil residues. The MAE for a coil residue is reduced from 13.8° to 11.4° for θ and from 56.4° to 50.2° for τ .

One application of predicted θ and τ angles is to utilize them for direct construction of local structures whose accuracies can be measured by the root-mean-square distance (RMSD) from their corresponding native conformations. Fragment structures of a length L are derived from predicted angles using the TS1199 dataset with a sliding window (1 to L , 2 to $L+1$, 3 to $L+2$, and etc.). For $L=15$, a total of 229681 fragments are constructed. Each fragment structure was built by using the standard $C\alpha$ - $C\alpha$ distance of 3.8\AA , and predicted θ and τ angles. We compared the accuracy of local structures from predicted θ and τ angles to those from ϕ and ψ angles predicted by SPINE X in Figure 7A. The RMSD between a native local structure (15 residue fragment) and its corresponding local structure constructed from predicted θ and τ angles (X-axis) is plotted against the RMSD between a native local structure and its corresponding structure constructed from predicted ϕ and ψ angles (Y-axis) in a density plot. The majority of RMSD values are less than 6\AA . The average RMSD values of local structures from predicted θ and τ angles are 1.9\AA for 10mer, 3.1\AA for 15mer, 4.3\AA for 20mer and 7.0\AA for 30mer. By comparison, the average RMSD values from predicted ϕ and ψ angles are 2.2\AA for 10mer, 3.4\AA for 15mer, 4.8\AA for 20mer and 7.7\AA for 30mer. The improvement of θ/τ derived structures over ϕ/ψ derived structures is more than 10%. More local structures from predicted θ and τ angles are more accurately predicted than those from predicted ϕ and ψ angles as demonstrated by the size of the triangle at the bottom-right corner. The spread from the diagonal line confirms the complementary role of these four predicted angles.

The difference (RMSD) between local structures generated by predicted θ and τ angles and those by predicted ϕ and ψ angles can serve as an effective measure of how accurate a predicted local structure is. Figure 7B shows the density plot of the RMSD from the native (Y-axis) versus the RMSD from the ϕ and ψ -derived structure (X-axis) for 15-residue fragments. There is a trend that the larger the structural difference from different types of angles is, the less accurate the predicted local structure (larger RMSD) will be. For example, if the RMSD between θ/τ -derived and ϕ/ψ -derived local structures is less than 2\AA , the RMSD of a θ/τ -derived structure from its native structure is most likely less than 4\AA based on the most populated region in red.

DISCUSSION

This study developed the first machine-learning technique for prediction of the angle between $C\alpha_{i-1}$ - $C\alpha_i$ - $C\alpha_{i+1}$ (θ) and a dihedral angle rotated about the $C\alpha_i$ - $C\alpha_{i+1}$ bond (τ). These angles reflect a local structure of 3 to 4 amino acid residues. By comparison, ϕ and ψ angles are the property of a single residue while secondary helical and sheet structures involve more than 3 residues. Thus, direct prediction of θ and τ angles is complementary to sequence-based prediction of ϕ and ψ angles and secondary structures. Predicting θ and τ angles also has one advantage over ϕ and ψ angles because θ has a narrow range of 0 to 180° while ϕ and ψ , similar to τ are both dihedral angles ranging from -180° to $+180^\circ$. Indeed, by using the stacked sparse auto-encoder deep neural network, we achieved MSE values of 9° for θ and 34° for τ . By comparison, MAE is 22° for ϕ and 33° for ψ by SPINE-X. As a result, θ and τ calculated from predicted ϕ and ψ angles are less accurate with an MAE of 10° for θ and 38° for τ , 10% higher than direct prediction of θ and τ .

Complementarity between predicted θ/τ angles and predicted ϕ/ψ angles is demonstrated from the accuracy of local structures constructed based on **these** predicted angles. As shown in Figure 7A, some local structures are more accurately constructed by θ and τ angles while others are more accurately constructed by ϕ and ψ angles. Moreover, RMSD values between θ/τ -derived and ϕ/ψ -derived structures can be utilized as a measure for the accuracy of **a** predicted local structure (Figure 7B). Usefulness of predicted angles for fragment structure prediction is illustrated by the fact that the average RMSD of 15-residue fragments is only 3Å from the corresponding native fragment structures. Currently, the most successful techniques in structure prediction (e.g. ROSETTA [31] and TASSER [32]) are based on mixing and matching of known native structures either in whole (template-based modelling) or in part (fragment assembly) [33,34]. Fragment **structures** based on predicted **θ and τ** angles provide an alternative but complementary approach to the homolog-based approach for generating fragment structures. In addition to fragment-based structure prediction, predicted θ and τ angles can also be employed directly as a constraint for fragment-free or *ab initio* structure prediction [1,2] as predicted ϕ and ψ angles **did** [9].

How to handle the periodicity of torsion angles is an issue facing angle prediction (-180° is same as 180°). In our previous work for predicting ϕ and ψ angles, we employed a simple angle shift [7], and prediction of peaks (two-state classification), followed by prediction of deviation from the peaks [9]. Here we introduced a sine and cosine transformation of θ and τ angles, a technique commonly employed in signal processing and speech recognition [30]. We have compared the sine and cosine transformation with angle shifting and its combination of two-state classification because the distributions of θ and τ angles also have two peaks (Figure 3). We found that the MAE of τ is 54° by direct prediction, 41° by angle shifting and 36° by combining two-peak prediction with angle shifting. Thus a MAE of 34° by sine and cosine transformation has the highest accuracy. We also examined the use of arcsine or arccosine, rather than arctangent. We found that using arccosine (with sine for phase determination) yields similar prediction accuracy as using arctangent but using arcsine leads to significantly worse prediction. We expect that such sine and cosine transformation of ϕ and ψ angles will also likely improve over existing SPINE-X prediction. For SPINE-X, MAE values are 33° for ψ angles and 22° for ϕ angles, respectively.

We also examined how much improvement in angle prediction is due to the use of deep learning neural networks. We found that when only one hidden layer (150 nodes) is utilised, MAE values are 8.8° for θ angles and 34.1° for τ angles, respectively. Thus, using deep 3-layer neural networks yields minor but statistically significant improvement over simple neural networks.

The most difficult angles to predict are the angles of coil residues (Table 1). This is true for θ and τ angles as well as for ϕ and ψ angles. Angles in coil regions have a mean absolute error of 11° for θ and 50° for τ , compared to 32° for ϕ and 56° for ψ . This is likely because coil regions are structurally least defined. Despite of large errors, predicted ϕ and ψ angles in coil regions have been proved to significantly improve the accuracy of predicted structures [9]. Thus, we expect that predicted θ and τ angles in coil regions will also be useful as restraints for *ab initio* structure prediction [9] or template-based structure prediction [14].

Acknowledgement. This work was supported in part by National Health and Medical Research Council (1059775) of Australia to Y.Z. We also gratefully acknowledge the support of the Griffith University eResearch Services Team and the use of the High Performance Computing Cluster "Gowonda" to complete this research. This research/project has also been undertaken with the aid of the research cloud resources provided by the Queensland Cyber Infrastructure Foundation (QCIF).

REFERENCES

1. Zhou YQ, Duan Y, Yang YD, Faraggi E, Lei HX (2011) Trends in template/fragment-free protein structure prediction. Theoretical Chemistry Accounts 128: 3-16.

2. Guo JT, Ellrott K, Xu Y (2008) A historical perspective of template-based protein structure prediction. *Methods Mol Biol* 413: 3-42.
3. Faraggi E, Zhang T, Yang Y, Kurgan L, Zhou Y (2011) SPINE X: Improving protein secondary structure prediction by multi-step learning coupled with prediction of solvent accessible surface area and backbone torsion angles. *J Computational Chemistry* 33: 259-263.
4. Rost B (2001) Review: Protein secondary structure prediction continues to rise. *Journal of Structural Biology* 134: 204-218.
5. Kihara D (2005) The effect of long-range interactions on the secondary structure formation of proteins. *Protein Science* 14: 1955-1963.
6. Dor O, Zhou Y (2007) Real-SPINE: an integrated system of neural networks for real-value prediction of protein structural properties. *Proteins-Structure Function and Bioinformatics* 68: 76-81.
7. Xue B, Dor O, Faraggi E, Zhou Y (2008) Real-value prediction of backbone torsion angles. *Proteins-Structure Function and Bioinformatics* 72: 427-433.
8. Faraggi E, Xue B, Zhou Y (2009) Improving the prediction accuracy of residue solvent accessibility and real-value backbone torsion angles of proteins by guided-learning through a two-layer neural network. *Proteins-Structure Function and Bioinformatics* 74: 847-856.
9. Faraggi E, Yang YD, Zhang SS, Zhou Y (2009) Predicting Continuous Local Structure and the Effect of Its Substitution for Secondary Structure in Fragment-Free Protein Structure Prediction. *Structure* 17: 1515-1527.
10. Simons KT, Bonneau R, Ruczinski I, Baker D (1999) Ab initio protein structure prediction of CASP III targets using ROSETTA. *Proteins-Structure Function and Genetics*: 171-176.
11. Huang YM, Bystroff C (2006) Improved pairwise alignments of proteins in the Twilight Zone using local structure predictions. *Bioinformatics* 22: 413-422.
12. Mooney C, Vullo A, Pollastri G (2006) Protein structural motif prediction in multidimensional phi-psi space leads to improved secondary structure prediction. *Journal of Computational Biology* 13: 1489-1502.
13. Wood MJ, Hirst JD (2005) Protein secondary structure prediction with dihedral angles. *Proteins-Structure Function and Bioinformatics* 59: 476-481.
14. Yang Y, Faraggi E, Zhao H, Zhou Y (2011) Improving protein fold recognition and template-based modeling by employing probabilistic-based matching between predicted one-dimensional structural properties of the query and corresponding native properties of templates. *Bioinformatics* 27: 2076-2082.
15. Karchin R, Cline M, Mandel-Gutfreund Y, Karplus K (2003) Hidden Markov models that use predicted local structure for fold recognition: Alphabets of backbone geometry. *Proteins-Structure Function and Bioinformatics* 51: 504-514.
16. Zhang W, Liu S, Zhou Y (2008) SP5: Improving protein fold recognition by using predicted torsion angles and profile-based gap penalty. *PLoS ONE* 6: e2325.
17. Korkut A, Hendrickson WA (2009) A force field for virtual atom molecular mechanics of proteins. *Proc Natl Acad Sci U S A* 106: 15667-15672.
18. Zhou Y, Karplus M (1997) Folding thermodynamics of a model three-helix-bundle protein. *Proc Natl Acad Sci U S A* 94: 14429-14432.
19. Kihara D, Lu H, Kolinski A, Skolnick J (2001) TOUCHSTONE: An ab initio protein structure prediction method that uses threading-based tertiary restraints. *Proceedings of the National Academy of Sciences of the United States of America* 98: 10125-10130.
20. Liwo A, He Y, Scheraga HA (2011) Coarse-grained force field: general folding theory. *Phys Chem Chem Phys* 13: 16890-16901.
21. Flocco MM, Mowbray SL (1995) C alpha-based torsion angles: a simple tool to analyze protein conformational changes. *Protein Sci* 4: 2118-2122.
22. Kleywegt GJ (1997) Validation of protein models from C α coordinates alone. *J Mol Biol* 273: 371-376.
23. Wang G, Dunbrack RL, Jr. (2005) PISCES: recent improvements to a PDB sequence culling server. *Nucleic Acids Res* 33: W94-98.
24. Hinton GE (2007) Learning multiple layers of representation. *Trends in Cognitive Sciences* 11: 428-434.
25. Bengio Y (2009) Learning deep architectures for AI. *Foundations and trends[®] in Machine Learning* 2 2: 1-127.
26. Bengio Y, Lamblin P, Popovici D, Larochelle H (2007) Greedy layer-wise training of deep networks. *Advances in neural information processing systems* 19: 153.
27. Palm RB (2012) Prediction as a candidate for learning deep hierarchical models of data: Technical University of Denmark.
28. Altschul SF, Madden TL, Schaffer AA, Zhang JH, Zhang Z, et al. (1997) Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Research* 25: 3389-3402.
29. Meiler J, Müller M, Zeidler A, Schmäsckhe F (2001) Generation and evaluation of dimension-reduced amino acid parameter representations by artificial neural networks. *Journal of Molecular Modeling* 7: 360-369.
30. Bozkurt B, Couvreur L, Dutoit T (2007) Chirp group delay analysis of speech signals. *Speech Communication* 49: 159-176.
31. Simons KT, Kooperberg C, Huang E, Baker D (1997) Assembly of protein tertiary structures from fragments with similar local sequences using simulated annealing and Bayesian scoring functions. *Journal of Molecular Biology* 268: 209-225.
32. Zhang Y, Skolnick J (2004) Automated structure prediction of weakly homologous proteins on a genomic scale'. *Proceedings of the National Academy of Sciences of the United States of America* 101: 7594-7599.
33. Bujnicki JM (2006) Protein-structure prediction by recombination of fragments. *Chembiochem* 7: 19-27.
34. Zhang Y (2009) Protein structure prediction: when is it useful? *Current Opinion in Structural Biology* 19: 145-155.

Table 1. Performance of θ and τ angle prediction based on the mean absolute error (MAE) as compared to θ and τ angles calculated from ϕ and ψ angles predicted by SPINE-X for two datasets (ten-fold cross validation for TR4590 and independent test for TS1199).

MAE	TR4590(°)	TS1199(°)	TS1199(°) from predicted ϕ and ψ.
θ -All	8.57±0.01	8.6	9.6
θ -Helix	4.50±0.02	4.5	4.5
θ -Sheet	10.45±0.02	10.6	11.3
θ -Coil	11.437 ±0.01	11.4	13.8
τ	33.4±0.3	33.6	37.7
τ -Helix	17.1±0.9	16.9	17.8
τ -Sheet	32.4±0.1	33.1	39.1
τ -Coil	50.1±0.3	50.2	56.4

Table 2: The mean absolute errors (MAEs) of θ and τ prediction for 20 amino acid residue types along with their frequency of occurrence in the TS1199 dataset.

Amino acids	Frequency	Theta	Tau
A	8.3	7.5	28.5
C	1.4	10.1	38.1
D	5.9	8.4	38.9
E	6.7	7.1	29.7
F	4.0	9.3	34.2
G	7.2	12.3	51.5
H	2.3	9.6	37.9
I	5.6	7.2	26.2
K	5.8	7.8	30.9
L	9.2	6.9	25.9
M	2.1	7.9	29.2
N	4.4	9.0	41.0
P	4.6	8.5	33.5
Q	3.8	7.6	30.6
R	5.1	8.0	31.2
S	5.9	10.7	40.4
T	5.6	9.9	35.6
V	7.1	7.7	27.7
W	1.5	9.2	35.3
Y	3.6	9.3	34.5
Average		8.6	33.6

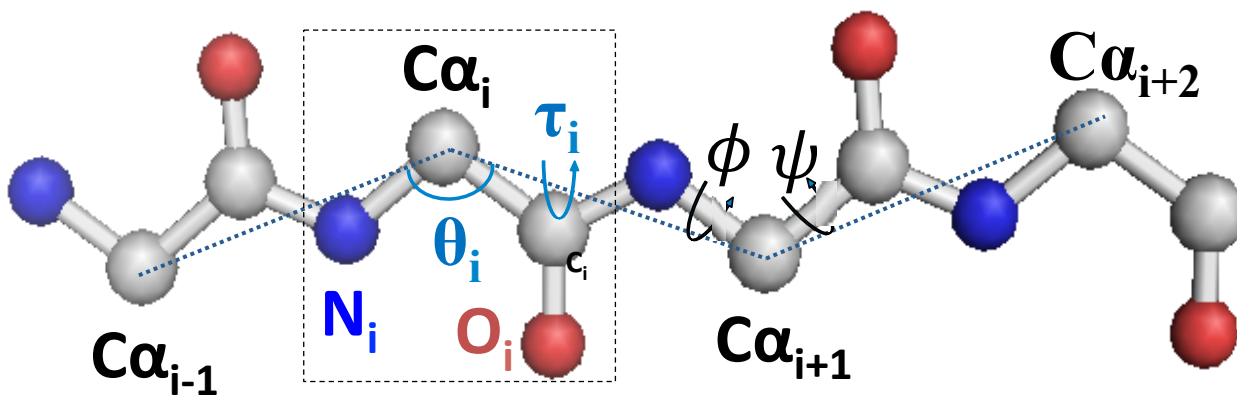


Figure 1. The schematic illustration of the protein backbone and associated angles.

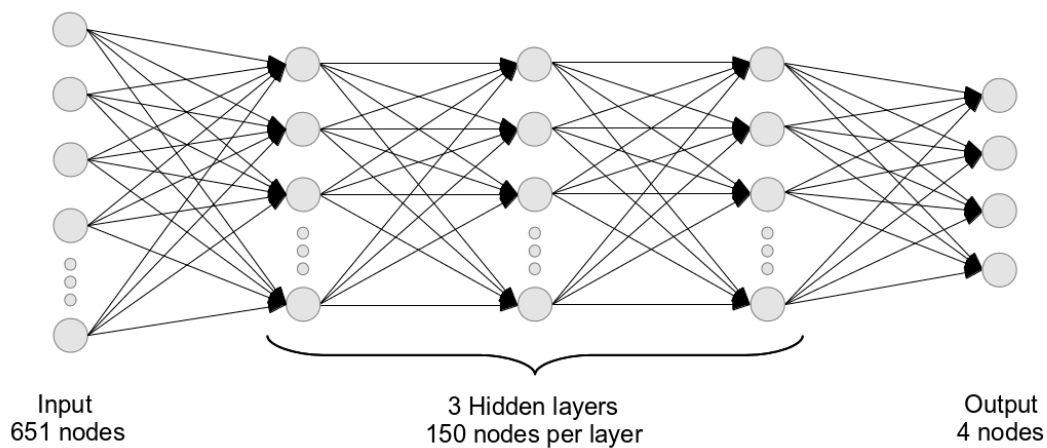


Figure 2: The general architecture of the stacked sparse auto-encoder deep neural network. Four output nodes are $\text{Sin}(\theta)$, $\text{Cos}(\theta)$, $\text{Sin}(\tau)$, and $\text{Cos}(\tau)$, respectively.

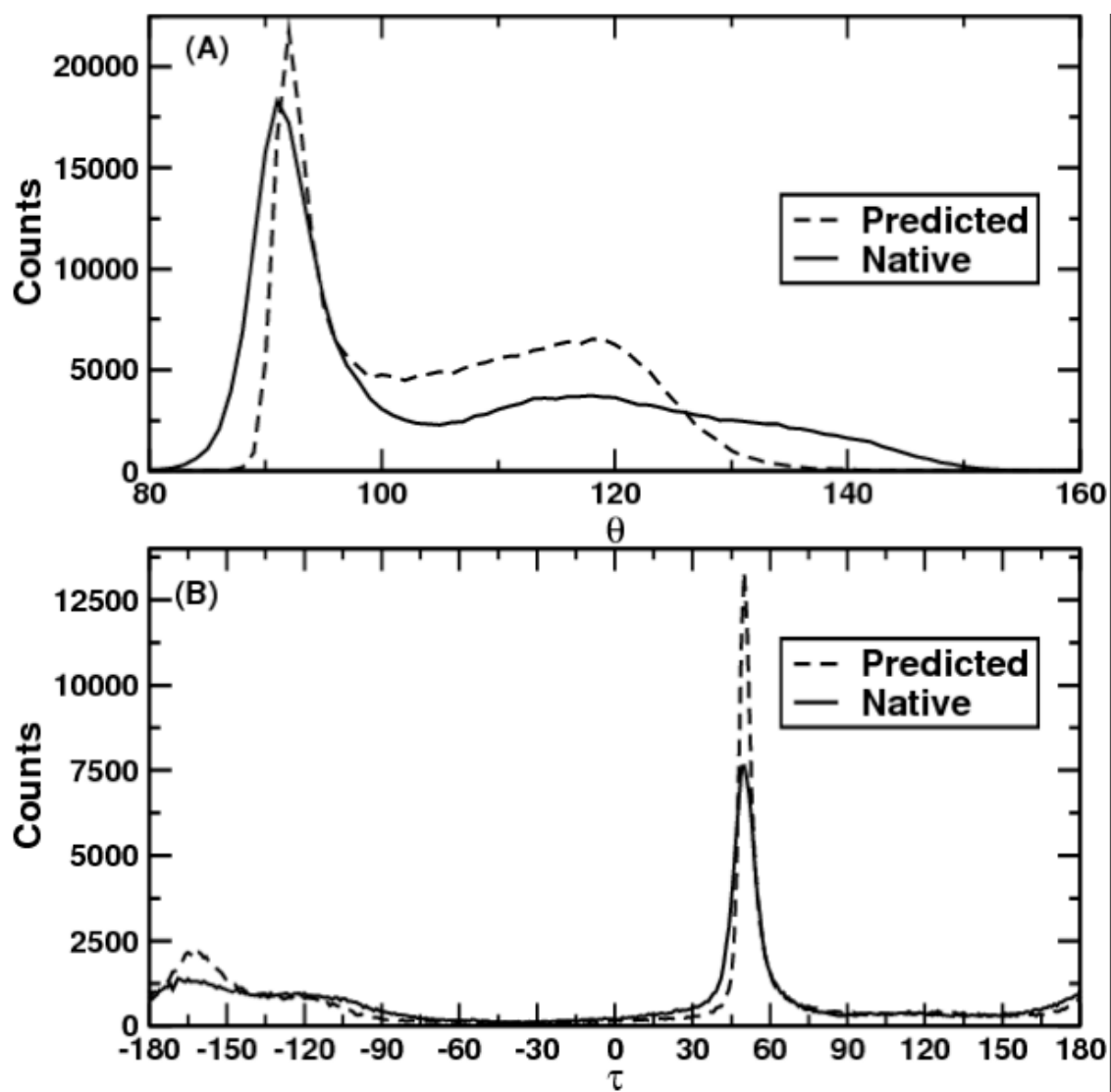


Figure 3 Predicted and actual distributions of θ (A) and τ (B) angles for the TS1199 dataset.

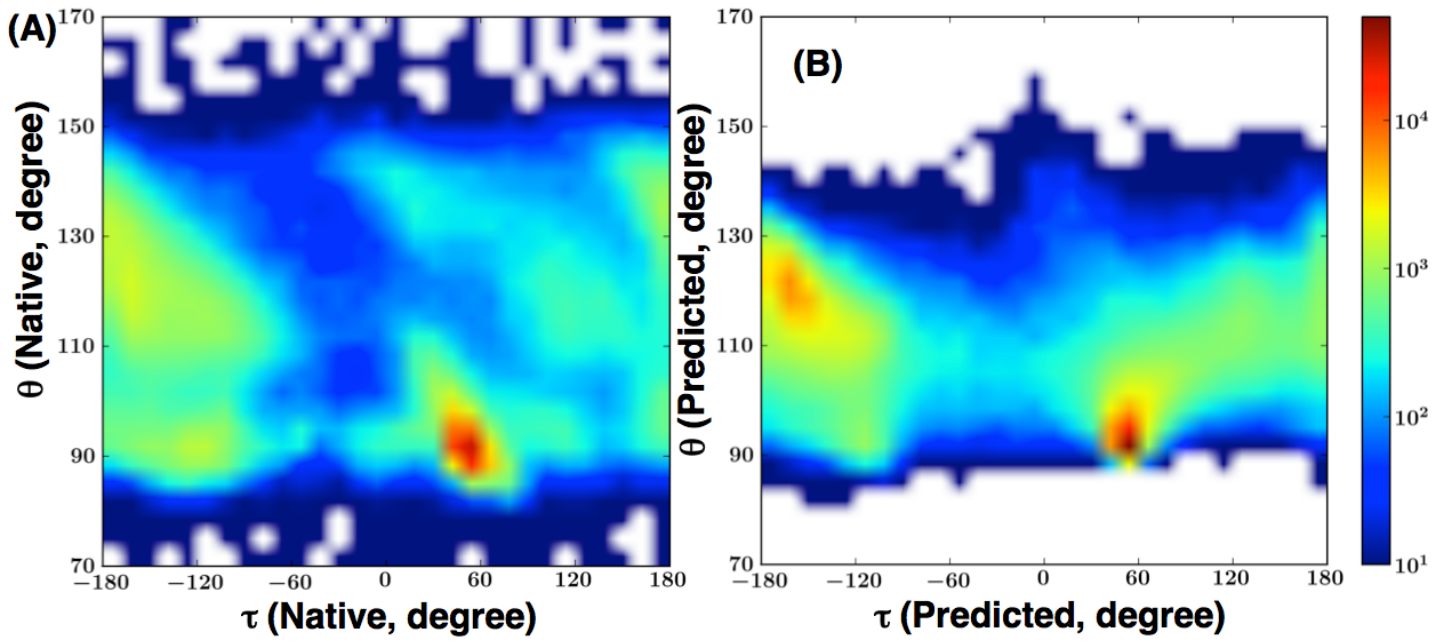


Figure 4 Actual (A) and predicted (B) distributions in the θ - τ plane for the TS1199 dataset.

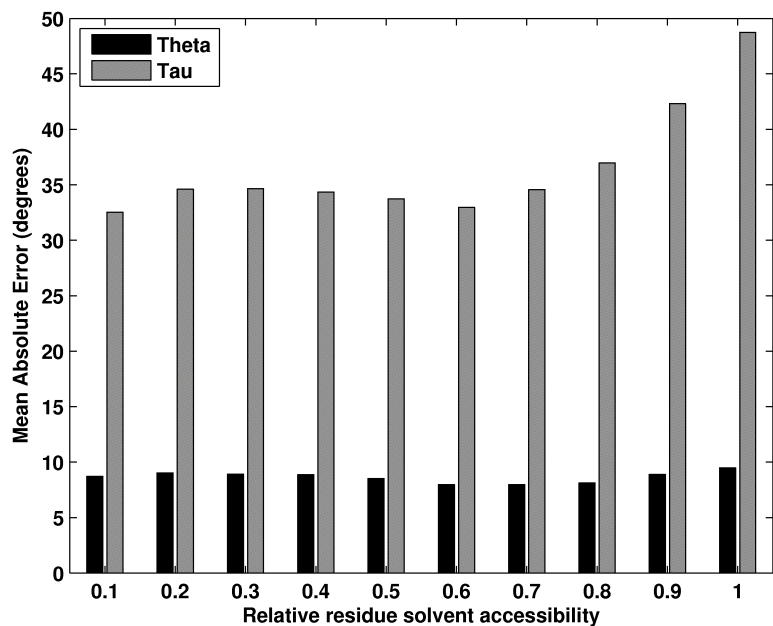


Figure 5 Mean absolute errors as a function of relative solvent accessibility for the TS1199 dataset.

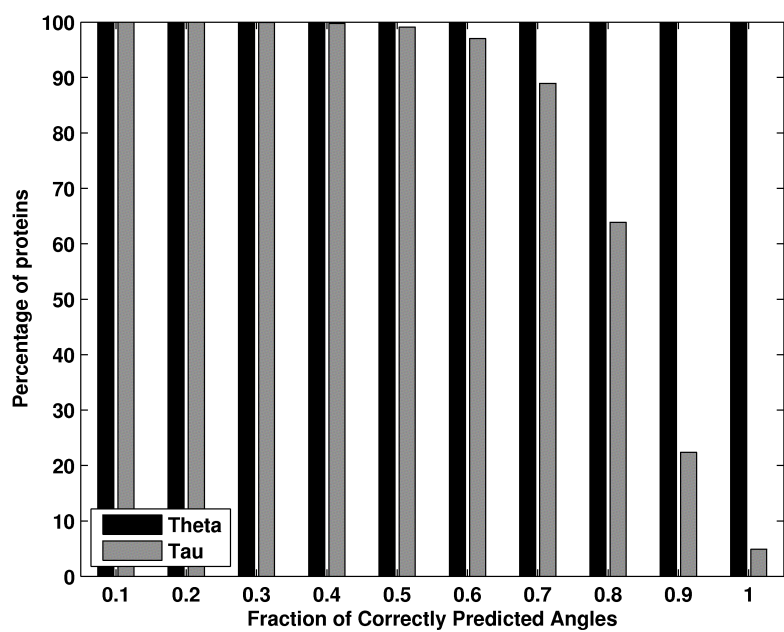


Figure 6 Percentage of proteins with more than a fraction of correctly predicted angles (θ and τ angles are less than 36° from native values, respectively) for the TS1199 dataset.

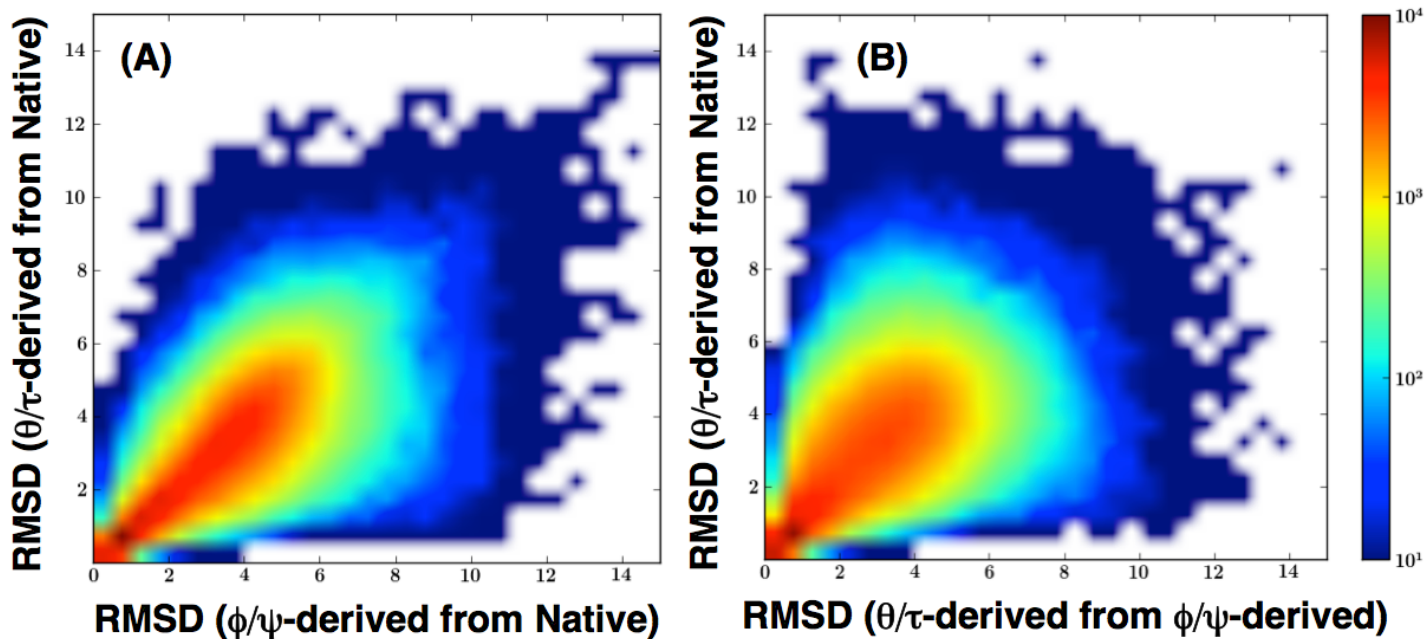


Figure 7 (A) Consistency between 15-residue local fragment structures derived from predicted ϕ/ψ (X-axis) and those from predicted θ/τ angles (Y-axis) in term of their root-mean-square distance (RMSD in Å) from the native structure for the TS1199 dataset. (B) RMSD values between two angle-derived local structures (X-axis) are compared to RMSD of a θ/τ -derived structure from its native structure.