

# Improving the Prediction of Protein Secondary Structure in Three and Eight Classes Using Recurrent Neural Networks and Profiles

Gianluca Pollastri,<sup>1</sup> Darisz Przybylski,<sup>2</sup> Burkhard Rost,<sup>2</sup> and Pierre Baldi<sup>3\*</sup>

<sup>1</sup>*Department of Information and Computer Science, Institute for Genomics and Bioinformatics, University of California, Irvine, Irvine, California*

<sup>2</sup>*Department of Biochemistry and Molecular Biophysics, CUBIC, Columbia University, New York, New York*

<sup>3</sup>*Department of Information and Computer Science, Institute for Genomics and Bioinformatics, University of California, Irvine, Irvine, California*

**ABSTRACT** Secondary structure predictions are increasingly becoming the workhorse for several methods aiming at predicting protein structure and function. Here we use ensembles of bidirectional recurrent neural network architectures, PSI-BLAST-derived profiles, and a large nonredundant training set to derive two new predictors: (a) the second version of the SSpro program for secondary structure classification into three categories and (b) the first version of the SSpro8 program for secondary structure classification into the eight classes produced by the DSSP program. We describe the results of three different test sets on which SSpro achieved a sustained performance of about 78% correct prediction. We report confusion matrices, compare PSI-BLAST to BLAST-derived profiles, and assess the corresponding performance improvements. SSpro and SSpro8 are implemented as web servers, available together with other structural feature predictors at: <http://promoter.ics.uci.edu/BRNN-PRED/>. *Proteins* 2002;47:228–235.

© 2002 Wiley-Liss, Inc.

**Key words:** recurrent neural networks; profiles; evolutionary information; PSI-BLAST

## INTRODUCTION

Secondary structure predictions are increasingly becoming the workhorse for several methods aiming at predicting protein structure and function, especially on a genomic scale.<sup>1–5</sup> Several threading techniques aiming at the identification of structural similarities between proteins with different sequences are based on predictions of secondary structure.<sup>6,7</sup> Predicting contact maps from primary sequence, secondary structure, and other structural features has also emerged as a key possible strategy for predicting protein structure.<sup>8–10</sup>

Methods predicting protein secondary structure have improved substantially in the 1990s through the use of machine learning methods and evolutionary information taken from the divergence of proteins in the same structural family.<sup>11–17</sup> In recent years, increases in the available training data and progress in algorithmic approaches have boosted prediction accuracy to about 76% of all

residues predicted correctly in one of the three states: helix, strand, and “other.”<sup>10,18</sup> On the algorithmic front, improvements of various kinds have resulted mostly from combining predictors, from using more sensitive methods for deriving evolutionary profiles, and from developing more flexible machine learning architectures.

It is well known that combining predictors usually improves prediction accuracy. Current methods for predicting secondary structure typically combine multiple neural networks, sometimes several hundreds of them,<sup>19</sup> trained more or less independently. Combination of different systems rather than networks has also been used.<sup>20</sup> At the alignment level, the ability to produce profiles that include increasingly remote homologs using PSI-BLAST<sup>15</sup> has also contributed to performance improvement. Divergent evolutionary profiles contain not only enough information to substantially improve prediction accuracy but even to predict long stretches of identical residues observed in alternative secondary structure states depending on nonlocal conditions.<sup>21–23</sup> An example is a method automatically identifying structural switches, and thus finding a remarkable connection between predicted secondary structure and aspects of function.<sup>22,23</sup> Finally, at the algorithmic level, new bidirectional recurrent neural network architectures were introduced in Ref. 14 in combination with BLAST profiles to produce a first-generation secondary structure predictor SSpro 1.0.

Here we develop the second version, SSpro 2.0, by using an ensemble of bidirectional recurrent neural networks and PSI-BLAST profiles. For the purpose of comparison, we use the same training set as SSpro 1.0, but with larger validation sets that have become available since the first version. We show improved performance, to about 78% correct prediction under stringent conditions. Secondary

Grant sponsor: National Institutes of Health; Grant number: R01-GM63029-01.

\*Correspondence to: Pierre Baldi, Department of Information and Computer Science, Institute for Genomic and Bioinformatics, University of California, Irvine, Irvine, CA 92697-3425. E-mail: pfbaldi@ics.uci.edu

Received 12 July 2001; Accepted 30 November 2001

**TABLE I. Eight- and Three-Class Assignment Statistics for the Four Sets Adopted<sup>†</sup>**

3class	8class	TRAIN		R126		EVA		CASP4	
H	H	91911	32.56	6573	28.13	16421	34.67	—	—
	G	11173	3.96	862	3.69	1751	3.70	—	—
	I	67	0.02	5	0.02	8	0.02	—	—
		103151	36.54	7440	31.85	18180	38.38	3600	39.79
E	E	59302	21.01	5068	21.69	8940	18.87	—	—
	B	3634	1.29	353	1.51	489	1.03	—	—
		62936	22.29	5421	23.20	9429	19.91	2048	22.64
	S	25036	8.87	2539	10.87	4536	9.58	—	—
C	T	33452	11.85	2718	11.63	5581	11.78	—	—
	.	57728	20.45	5245	22.45	9644	20.36	—	—
		116216	41.17	10502	44.95	19769	41.73	3399	37.57
	All	282303	100.00	23363	100.00	47370	100.00	9047	100.00

<sup>†</sup>For each set, the first column gives the total number of residues and the second the corresponding percentages. Three-class: H = helix; E = strand; and C = coil. Eight-class: H = alpha helix; B = residue in isolated  $\beta$ -bridge; E = extended strand, participates in  $\beta$ -ladder; G = 3-helix (3/10 helix); I = 5 helix (pi helix); T = hydrogen bonded turn; S = bend; and ".".

structure classification for resolved structures is typically obtained by collapsing the eight-class output of the the DSSP program<sup>24</sup> into the standard three classes. Because useful information may be present in the eight classes, we also develop SSpro8, a secondary structure predictor into eight classes.

## MATERIALS AND METHODS

### Data: Training and Testing

The assignment of the SS categories to the experimentally determined three-dimensional (3D) structure is non-trivial and typically performed by the widely used DSSP program.<sup>24</sup> DSSP works by assigning potential backbone hydrogen bonds (based on the 3D coordinates of the backbone atoms) and subsequently by identifying repetitive bonding patterns.

The DSSP program classifies each residue into eight classes (H = alpha helix; B = residue in isolated beta-bridge; E = extended strand, participates in beta ladder; G = 3-helix [3/10 helix]; I = 5 helix [pi helix]; T = hydrogen bonded turn; S = bend; and "."). These are typically collapsed into the three standard classes associated with helices,  $\beta$ -strand, and coils. In the CASP experiments,<sup>1,10,25</sup>  $\alpha$  contains H and G,  $\beta$  contains E and B, and  $\gamma$  contains everything else. This assignment is known to be somewhat "harder" to predict than the other ones used in the literature where, for instance,  $\alpha$  is formed by DSSP class H,  $\beta$  by E, and  $\gamma$  by everything else (including DSSP classes G, S, T, B, I, and "."). Other assignments used in the literature include,<sup>26</sup> where  $\alpha$  contains DSSP classes H, G, and I. A study of the effect of various assignments on prediction performance can be found in Ref. 20. It is clear, however, that it may also be of interest to build a finer-grained predictor for the eight classes produced by DSSP. Because some classes are fairly rare, lack of data may have been an obstacle in the past. But with the steady stream of new structures deposited in the PDB<sup>27</sup> every week, the time may have come to revisit this issue.

Four main data sets are used to develop and test the approach: a training set (TRAIN) and three test sets (R126, EVA, and CASP4) to assess algorithmic perfor-

mance in the most objective way. The distribution of the eight and three classes in these sets are summarized in Table I.

### TRAIN

To ensure fair comparison with SSpro 1.0, we use the same training test, originally developed at the end of 1999. We constructed a data set containing all proteins in PDB which are (a) at least 30 amino acids long, (b) have no chain breaks (defined as neighboring amino acids in the sequence having  $C^\alpha$ -distances exceeding 4.0 Å), (c) produce a DSSP output, and (d) are obtained by X-ray diffraction methods with a resolution of at least 2.5 Å. Internal homology is reduced by using an all-against-all alignment approach,<sup>28</sup> keeping the PDB sequences with the best resolution. A 50% threshold curve is used for homology reduction. Furthermore, the proteins in the set have <25% identity with the sequences in the set R126. The resulting training set consists of 1180 sequences corresponding to 282,303 amino acids.

### R126

As a first independent test set, we use the original set of 126 sequences of Rost and Sander, currently corresponding to a total of 23,363 amino acid positions (this number has varied slightly over the years because of changes and corrections in the PDB).

### EVA

A novel test set is provided by all the sequences available from the real-time evaluation experiment EVA (<http://cubic.bioc.columbia.edu/eva/>), which compares a number of prediction servers on a regular basis by using the sequences deposited in the PDB every week. In particular, we use the set labeled "common3" published on 3/3/2001, the largest EVA data set on which SSpro 1.0 had been tested together with the other main servers at that date. The set consists of 223 proteins with a total of 47,370 residues and contains sequences with no homology to proteins previously stored in PDB. This set was uploaded in PDB between 3/2000 and 3/2001; thus, it has no

homology to the TRAIN set extracted from PDB in September 1999.

### CASP4

The last novel test set is provided by the 40 CASP4 sequences available at <http://predictioncenter.llnl.gov/casp4/> and corresponding to 9047 residues. Some of the sequences show homology to PDB structures. In this case we could not obtain the PDB file of three of the proteins (T91, T92, T93), hence were unable to run DSSP on them and did not use this set for testing classification into eight classes. The secondary structure assignment into three classes for this set was downloaded directly from the CASP4 web site.

### Profiles

To improve prediction, we use both BLAST and PSI-BLAST input profiles. Using profiles at the input level generally has been shown to yield better results than using profiles at the output level.<sup>14</sup>

### BLAST

Input profiles for SSpro 1.0 were constructed primarily by running the BLAST program<sup>29</sup> against the NR (nonredundant) database,<sup>27,35</sup> with standard default parameters ( $E = 10.0$ , BLOSUM62 matrix). The version used was available online in October 1999 and contained approximately 420,000 protein sequences. For redundancy reduction, instead of applying a hard threshold that requires an arbitrary cutoff choice, we used a continuous weighting scheme. In this scheme, the weight of a sequence measures how different the sequence is from the profile in terms of the Kullback-Leibler divergence.<sup>17</sup> More specifically, for any given sequence, the weight is the sum over all columns in the profile of the Kullback-Leibler divergence between the delta distribution associated with the composition of the sequence in the column and the corresponding profile distribution. This is also a measure of the Shannon information in the sequence, given the current profile. Formally, the weight of sequence  $s$  is

$$W(s) = - \sum_c \log P[s(c)] = - \log P[s] \quad (1)$$

where  $P[s(c)]$  is the probability according to the profile of the letter  $s$  has in column  $c$ . Highly redundant sequences have higher probabilities and, therefore, are assigned a lower weight. In summary, every sequence in a given alignment is assigned a weight proportional to the information the sequence carries with respect to the unweighted profile. A weighted profile matrix is then recompiled and used as input to the prediction algorithm (see also Ref. 14).

### PSI-BLAST

Here we derive new profiles by aligning all proteins against the NR database using PSI-BLAST<sup>31</sup> with the following four-step protocol.<sup>32</sup> First, we filter and remove all database sequences with COILS to mark coiled-coil regions<sup>33</sup> and SEG to mark regions of low complexity.<sup>34</sup> Second, we align the query protein against this filtered

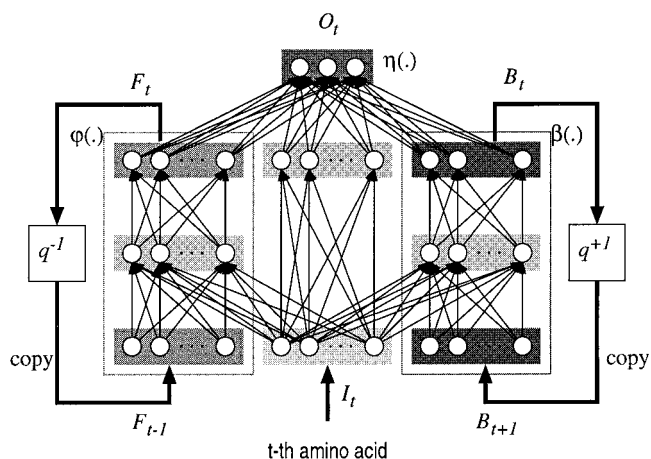


Fig. 1. A BRNN architecture with a left (forward) and right (backward) context,  $F_t$  and  $B_t$ , associated with two recurrent networks (wheels). The output layer  $O_t$  has three normalized exponential units associated with membership in each one of the three secondary structure classes for the current residue at position  $t$ . The functions  $\phi$ ,  $\beta$ , and  $\eta$  are implemented by feed-forward neural networks.

database with an E-value threshold for the iteration of  $10^{-10}$  (PSI-BLAST “h” threshold) and a final threshold of  $E \leq 10^{-3}$  to accept hits. The number of iterations is restricted to three to avoid drift.<sup>15,32</sup> Third, we align the query against the unfiltered NR database by using the previously found, position-specific profile. Finally, we use the same weighting scheme as in the case of BLAST profiles to balance the profile and remove redundancy.

### Recurrent Neural Network Architectures

In Ref. 14, BRNNs (Bidirectional Recurrent Neural Networks) were proposed as a class of recurrent neural network architectures that can address some of the limitations of simple feed-forward networks associated with small fixed-length input windows. A typical BRNN architecture is represented in Figure 1. In this architecture, the output decision or classification is determined by three components. First, there is a central component associated with the local window at the location  $t$  of the current prediction, as in standard feed-forward neural networks for secondary structure prediction. The main difference between the BRNN and the standard approach is the contribution of the left and right “contexts.” These are produced by two similar recurrent networks which, intuitively, can be thought of in terms of two “wheels” being rolled along the protein chain, from the N- and the C-terminus all the way to the point of prediction.

Architectural variants can be obtained by changing the size of the input windows, the size of the window of hidden states used to determine the output, the number of hidden layers, the number of hidden units in each layer, and so forth. As in standard secondary structure and other protein prediction architectures, we use sparse encoding for the 20 amino acids.

In what follows, we use the following notation:

$C_t$  = size of semiwindow of context states considered by the output network

**TABLE II. Parameters and Total Weights of the 11 BRNN Models<sup>†</sup>**

Model	Ct	NFB	NHO	NHC	Weights
0	3	8	11	9	2181
1	2	9	11	8	1899
2	3	12	11	9	2949
3	3	12	9	9	2565
4	3	15	12	13	4167
5	3	17	12	15	4831
6	3	17	14	15	5355
7	3	15	14	15	4839
8	4	15	14	15	5679
9	4	25	30	27	18107
10	3	25	30	27	15107

<sup>†</sup>Ct = size of semiwindow of context states; NFB = number of output units in the forward and backward context networks; NHO = number of hidden units in the output network; NHC = number of hidden units in the context networks.

NFB = number of output units in the left (forward) and right (backward) context networks (wheels)  
 NHO = number of hidden units in the output network  
 NHC = number of hidden units in the context networks

In the three (resp. 8) class-prediction applications considered here, there are three (resp. 8) normalized exponential output units that enable us to estimate the class membership probability for the residue being considered. The output error is the relative entropy between the output and target probability distributions.<sup>17</sup> All the weights of the BRNN architecture, including the weights in the recurrent wheels, can be trained from examples in a supervised fashion by using a generalized form of gradient descent or backpropagation through time, or rather space, because of the forward and backward nature of the chains.

BRNNs were used to develop the first version of the SSpro predictor.<sup>14</sup> They have also been used for the prediction of amino acid partners in  $\beta$ -sheets,<sup>35</sup> contact numbers,<sup>9</sup> and relative solvent accessibility. In Ref. 14, evidence is provided that, in the case of secondary structure prediction, these architectures extend the range over which information can be effectively integrated with respect to a feed-forward neural network, up to an effective window size of about 30 amino acids.

## Experiments

In terms of architectures, both the three- and eight-class predictors use the same ensemble of 11 BRNNs, as in the online version of SSpro 1.0. Parameters of the 11 networks are given in Table II. The total number of parameters varies over one order of magnitude from 1,899 to 18,107. Here we train SSpro 2.0 using PSI-BLAST profiles, SSpro8 1.0 using BLAST profiles, and SSpro8 2.0 using PSI-BLAST profiles. SSpro predicts the CASP three-class assignment, whereas SSpro8 predicts the eight classes that are produced by the DSSP program.

The training strategy adopted is the same for all the systems and derives from the preliminary studies reported

**TABLE III. Performances of All the 11 Models and the Ensemble, for SSpro 1.0 and SSpro 2.0, on the R126 Test Set and on the Training Set Using the Percentage Q3 of Correctly Classified Residues**

Model	SSpro 1.0		SSpro 2.0	
	Q3	Q3train	Q3	Q3train
0	74.91	77.24	76.85	78.27
1	74.82	76.93	76.53	78.36
2	75.20	78.03	76.62	79.02
3	74.98	77.62	76.35	78.94
4	75.15	79.02	76.30	80.39
5	74.54	79.80	75.90	80.79
6	74.50	78.93	76.08	80.87
7	74.80	78.66	76.61	80.91
8	74.57	78.55	76.32	80.91
9	74.40	81.00	74.84	85.21
10	73.37	83.26	75.16	84.21
Ensemble	76.62	81.01	78.13	83.02

in Ref. 14. In a typical case, we use a hybrid learning approach that combines online and batch training, with 430 batch blocks (two or three proteins each) per training set. Thus, weights are updated 430 times per epoch after each block. The learning rate per block is initially set at about  $1.5 \times 10^{-4}$ , corresponding to the number of blocks divided by 10 times the number of residues ( $0.1 \times 430/280000$ ) and is progressively decreased. The training set is also shuffled at each epoch. When the error does not decrease for 50 consecutive epochs, the learning rate is divided by 2, and training is restarted from the lowest error model. Training stops after eight or more reductions, corresponding to a learning rate that is 256 times smaller than the initial one, which usually happens after 1500–2500 epochs.

Several indices can be used to score the efficiency of the algorithm.<sup>36</sup> Here we use Q3 or Q8, the number of correctly predicted residues divided by the total number of residues, as well as the corresponding per-class version Qclass (percentage of residues in a given structural class that are correctly predicted) and the Matthew's correlation coefficient.

## RESULTS AND DISCUSSION

### SSpro

Performance for each individual model and for the ensemble average is given in Table III for the TRAIN and R126 sets. In all cases, PSI-BLAST profiles provide a Q3 improvement of at least 1.5%, and often more. On the EVA set, SSpro 2.0 achieves 77.7%, better than all the other evaluated systems. Incidentally, training on BLAST profiles and testing on PSI-BLAST profiles also leads to some performance improvement, although not as much.

In Table IV, we give the performance of SSpro 1.0 and SSpro 2.0 on the three test sets measured by Q3 and the percentage per class (Qclass). Again, the PSI-BLAST profiles lead to significant improvement in all categories. In some cases (E in EVA and CASP4), the improvement



**TABLE IV. Performances of SSpro 1.0 and SSpro 2.0 on the R126, EVA, and CASP4 Test Sets<sup>†</sup>**

		SSpro 1.0	SSpro 2.0
126	H	80.79	82.38
	E	63.23	66.19
	C	80.56	81.26
EVA	Q3	76.62	78.13
	H	80.76	82.48
	E	62.50	65.56
CASP4	C	78.05	79.03
	Q3	76.00	77.67
	H	83.86	86.08
	E	61.87	68.51
	C	80.99	82.20
	Q3	77.80	80.65

<sup>†</sup>Q3 and Qclass percentages.

**TABLE V. Performances of SSpro 1.0 and SSpro 2.0 on the R126, EVA, and CASP4 Test Sets<sup>†</sup>**

		SSpro 1.0	SSpro 2.0
R126	H	0.732	0.752
	E	0.598	0.634
	C	0.571	0.59
EVA	H	0.695	0.722
	E	0.594	0.632
	C	0.568	0.588
CASP4	H	0.749	0.788
	E	0.609	0.674
	C	0.599	0.634

<sup>†</sup>Matthews' correlation coefficients.

exceeds 3%. Similar results using Matthew's correlation coefficients are given in Table V.

SSpro 2.0 achieves Q3 of 78.13% on R126 and 77.67% on EVA, which, to the best of our knowledge, at the time of this writing is second to none. (The results reported on EVA as of 3/3/01 at [http://cubic.bioc.columbia.edu/eva/sec/bup\\_common/2001\\_03\\_03/common1.html](http://cubic.bioc.columbia.edu/eva/sec/bup_common/2001_03_03/common1.html) for the other tested predictors are PROF1 76.8%, PSIPred 76.5%, JPRED 74.7%, PHDpsi 74.7%, and PHD 71.5%). The error margin on the performance for single residues range from 0.19% for the EVA set to 0.42% for the CASP4 set, so the second decimal for Q3 is not particularly significant. The results currently reported on the EVA web site correspond to SSpro 1.0 only, because SSpro 2.0 came online in April 2001 and EVA does not offer an automated procedure to evaluate a newly entered predictor using all sequences up to the time of entry. Although SSpro 2.0 is derived after the CASP4 experiment, it is trained by using the same training set as SSpro 1.0 before CASP4. SSpro achieves a Q3 performance of 80.65% on the CASP4 set, which is known to contain sequences with a wide difficulty range.<sup>10</sup> For comparison, the predictor PSIPRED achieves 79.9% correct prediction per residue, computed by using the official predictions reported at CASP4 for 39 sequences and submitting the remaining sequence (T0106) directly to the

**TABLE VI. Confusion Matrices for SSpro 1.0 and SSpro 2.0 on the Set R126<sup>†</sup>**

	SSpro 1.0			SSpro 2.0		
	Hpred	Epred	Cpred	Hpred	Epred	Cpred
Hobs	80.79	2.64	16.57	82.38	1.83	15.79
Eobs	4.37	63.23	32.40	3.30	66.19	30.51
Cobs	9.85	9.60	80.55	9.64	9.10	81.26
	Hobs	Eobs	Cobs	Hobs	Eobs	Cobs
Hpred	82.50	3.25	14.25	83.69	2.44	13.87
Epred	4.24	73.90	21.86	2.91	76.57	20.52
Cpred	10.74	15.27	73.99	10.31	14.49	75.20

<sup>†</sup>Xpred = structure X is predicted. Yobs = structure Y is observed. Rows sum to 100%. The number in row Xpred and column Yobs represents the percentage of times structure Y is observed, given that structure X has been predicted.

**TABLE VII. Performances of SSpro8 1.0 and SSpro8 2.0 on the R126 and EVA Test Sets<sup>†</sup>**

		SSpro8 1.0	SSpro8 2.0
R126	H	89.21	89.93
	G	6.38	8.70
	I	0.00	0.00
	E	76.85	78.77
	B	0.00	0.00
	S	6.58	7.48
	T	43.34	45.44
	.	57.79	61.33
	Q8	60.74	62.58
	H	88.38	89.49
EVA	G	4.63	6.85
	I	0.00	0.00
	E	74.49	76.10
	B	0.00	0.00
	S	4.43	5.78
	T	39.67	40.75
	.	58.59	60.51
	Q8	61.89	63.31

<sup>†</sup>Q8 percentage and Qclass percentage of observed residues.

PSIPRED server. An overall performance of 80% also has been reported in Ref. 19, but an easier mapping of the eight DSSP classes into three is used rather than the one used for the CASP experiments. With this easier assignment, SSpro 2.0 achieves performances >80%.

Table VI provides the confusion matrices of SSpro 1.0 and SSpro 2.0 measured on the R126 test set. Perhaps not surprisingly,  $\beta$ -strands continue to remain the most difficult class probably because of a number of factors including the involvement of long-ranged interactions with respect to the primary sequence and the fact that they are underrepresented in the data (roughly 20% strands and 35% helices).

### SSpro8

For the eight-class prediction, it is first worth noting that class I is almost irrelevant because it represents

**TABLE VIII. Performances of SSpro8 1.0 and SSpro8 2.0 on the R126 and EVA Test Sets<sup>†</sup>**

		SSpro8 1.0	SSpro8 2.0
R126	H	75.55	78.38
	G	36.91	41.44
	I	—	—
	E	62.39	66.02
	B	—	—
	S	49.41	48.47
	T	42.04	43.89
	.	49.94	50.36
EVA	Q8	60.74	62.58
	H	76.06	77.11
	G	33.61	38.46
	I	—	—
	E	61.67	64.70
	B	—	—
	S	38.58	44.33
	T	42.92	43.70
	.	48.83	49.91
	Q8	61.89	63.31

<sup>†</sup>Q8 percentage and Qclass percentage of predicted residues (percentage of correctly classified X, given that X is predicted). Note that no numbers are available for classes I and B because they are never predicted.

0.02% of cases. Class B is small (1–1.5%), and the number of training examples is not yet large enough to yield any reliable prediction. Class G is also relatively small (~4%), but some generalization appears to be possible. The same is true for class S, which represents ~9% of cases. Each one of the other four classes contains at least 8% of the total residues.

The performances of SSpro8 1.0 and 2.0 are reported in Tables VII and VIII. For both versions of SSpro8, no residue is predicted as being in either class I or class B. Residues observed in class I are predicted as being in class H, and residues in class B are predicted somewhat evenly as being in classes in “.” or E, as shown by the confusion matrix in Table IX. Classes G and S are underpredicted. Roughly 30% of the residues classified in G by DSSP are classified as H by SSpro8, and 21% are classified as being in a turn (T). When the system predicts G or S, however, there is about a 40–50% chance it is correct, which is considerably better than a random prediction. Although class T represents only 11% of the residues, it tends to be slightly overpredicted. Approximately 43–45% of the observed turns are predicted correctly, and 40% of the predicted turns are actual turns. More precisely, bends (S) tend to be confused with turns (T). Of the observed bends, >20% are predicted as turns (<10% as bends). If a T is predicted, the probability of the actual residue being in either a T or an S structure is approximately 65%. A class obtained by merging turns and bends would be classified by SSpro8 2.0 with a 65% correct percentage of predicted residues. The overall performance shows a gain of 1.4–1.8% with PSI-BLAST profiles, reaching the 62.6–63.3% range. The error margins on the performance for single residues range from 0.22% for the EVA set to 0.32% for the R126 set.

Table VIII compares the two versions of SSpro8 on the other two test sets by using Q8 and Qclass. The confusion matrix of SSpro8 2.0 is reported in Table IX.

If we combine the predictions of SSpro8 using the CASP assignment, the performance obtained is inferior to SSpro by more than a percentage point. This is perhaps not too surprising, because SSpro8 is trained for a different task. Performance very close to SSpro (within 0.2%) can be achieved by cascading SSpro8 with a small neural network trained with threefold cross-validation on the test set. Overall, prediction into eight classes does not seem to improve prediction of secondary structure into three classes. However, current results are encouraging, especially for turns, and are bound to improve as more data becomes available.

## CONCLUSION

We have developed two state-of-the-art predictors for secondary structure in three and eight categories, using ensembles of bidirectional recurrent neural networks and PSI-BLAST profiles. We have assessed the gains attributable to the use of PSI-BLAST profiles over BLAST profiles and have implemented the programs in the form of two web servers, SSpro and SSpro8, available over the Internet at <http://promoter.ics.uci.edu/BRNN-PRED/>. Users can enter primary amino acid sequences, and predictions are e-mailed back to them after a short period of time, depending on server load. SSpro and SSpro8 are also part of a broader suite of programs aimed at predicting protein 3D structure via contact map prediction, and contact map prediction via prediction of structural features, such as secondary structure, relative solvent accessibility (ACCpro), and contact numbers (CONpro).<sup>37</sup>

Improvements of a few percentage points are significant, especially in the context of the Human Genome Project and other genome sequencing projects as well as high-throughput structural proteomics projects. It is encouraging to witness gradual performance improvements year after year, as a result of algorithmic improvements and data growth. Perfect prediction cannot be expected for a number of reasons, including (a) dynamic properties of protein chains; (b) quaternary structures; (c) existence of proteins that do not fold spontaneously; (d) errors and variability in databases, as well as in DSSP program output; and (e) effects of external variables, such as pH, that currently are not taken into consideration. Thus, some degree of prediction saturation is likely to emerge in the coming years, although the exact level at which it will occur remains unclear. For the time being, efforts toward exhaustive prediction and exhaustive taxonomy of protein folds should continue to advance in synergy.

## ACKNOWLEDGMENTS

The work of PB and GP is supported by a Laurel Wilkening Faculty Innovation award and a Sun Microsystems award to PB at UCI. The work of DP and BR is supported by NIH grant R01-GM63029-01.

**TABLE IX. Confusion Matrices for SSpro8 2.0 on the Set R126<sup>†</sup>**

	SSpro8 2.0							
	Hpred	Gpred	Ipred	Epred	Bpred	Spred	Tpred	.pred
Hobs	89.93	0.36	0.00	2.58	0.00	0.14	3.06	3.93
Gobs	28.19	8.70	0.00	11.49	0.00	0.58	20.88	30.16
Iobs	100.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
Eobs	3.99	0.08	0.00	78.77	0.00	0.45	2.82	13.89
Bobs	7.37	0.00	0.00	30.59	0.00	0.85	7.37	53.82
Sobs	9.89	0.94	0.00	13.59	0.00	7.48	23.12	44.98
Tobs	20.46	1.36	0.00	7.69	0.00	2.35	45.44	22.70
.obs	6.62	0.32	0.00	21.43	0.00	1.87	8.43	61.33

	Hobs	Gobs	Iobs	Eobs	Bobs	Sobs	Tobs	.obs
Hpred	78.38	3.23	0.07	2.68	0.34	3.33	7.37	4.60
Gpred	13.26	41.44	0.00	2.21	0.00	13.26	20.44	9.39
Ipred	—	—	—	—	—	—	—	—
Epred	2.81	1.63	0.00	66.02	1.79	5.71	3.46	18.58
Bpred	—	—	—	—	—	—	—	—
Spred	2.29	1.28	0.00	5.87	0.76	48.47	16.33	25.00
Tpred	7.14	6.40	0.00	5.08	0.92	20.86	43.89	15.71
.pred	4.04	4.07	0.00	11.02	2.97	17.88	9.66	50.36

<sup>†</sup>Xpred = structure X is predicted. Yobs = structure Y is observed. Rows sum to 100%. For instance, the number on Row Xpred, Column Yobs represents the percentage of times structure Y is observed, given that structure X has been predicted. Note that no numbers are available in the rows for classes I and B predicted, because they are never predicted.

## REFERENCES

- Moult J, Hubbard T, Bryant SH, Fidelis K, Pedersen JT. Critical assessment of methods of protein structure prediction (CASP): Round II. *Proteins* 1997;Suppl 1:29:2–6.
- Sanchez R, Sali A. Large-scale protein structure modeling of the *Saccharomyces cerevisiae* genome. *Proc Natl Acad Sci USA* 1998; 95:13597–13602.
- Moult J, Hubbard T, Bryant SH, Fidelis K, Pedersen JT. Critical assessment of methods of protein structure prediction (CASP): round III. *Proteins* 1999;Suppl 3:22–29.
- Jones DT. Protein structure prediction in the postgenomic era. *Curr Opin Struct Biol* 2000;10:371–379.
- Rost B. Review: protein secondary structure prediction continue to rise. *J Struct Biol* 2001;134:204–218.
- Fischer D, Barret C, Bryson K, Elofsson A, Godzik A, Jones D, Karplus KJ, Kelley LA, MacCallum RM, Pawowski K, Rost B, Rychlewski L, Sternberg M. CAFASP-1: critical assessment of fully automated structure prediction methods. *Proteins* 1999; Suppl 3:209–217.
- Fischer D, Elofsson A, Rychlewski L, Pazos F, Valencia A, Godzik A, Rost B, Ortiz AR, Dunbrack RL. CAFASP-2: the second critical assessment of fully automated structure prediction methods. *Proteins* 2001;Suppl 5:171–183.
- Baldi P, Pollastri G. Machine learning structural and functional proteomics. *IEEE Intelligent Systems. Special Issue on Intelligent Systems in Biology* 2001. Forthcoming.
- Pollastri G, Baldi P, Fariselli P, Casadio R. Improved prediction of the number of residue contacts in proteins by recurrent neural networks. *Proceedings of the ISMB 2001 Conference. Bioinformatics* 2001;17:S234–S242.
- Lesk AM, Lo Conte L, Hubbard TJP. Assessment of novel fold targets in CASP4: predictions of three-dimensional structures, secondary structures and interresidue contacts. *Proteins* 2001; Suppl 5:98–118.
- Rost B, Sander C. Combining evolutionary information and neural networks to predict protein secondary structure. *Proteins* 1994;19: 55–72.
- Barton GJ. Protein secondary structure prediction. *Curr Opin Struct Biol* 1995;5:372–376.
- Rost B, Sander C. Bridging the protein sequence-structure gap by structure predictions. *Annu Rev Biophys Biomol Struct* 1996;25: 113–136.
- Baldi P, Brunak S, Frasconi P, Pollastri G, Soda G. Exploiting the past and the future in protein secondary structure prediction. *Bioinformatics* 1999;15:937–946.
- Jones DT. Protein secondary structure prediction based on position-specific scoring matrices. *J Mol Biol* 1999;292:195–202.
- Cuff JA, Barton GJ. Application of multiple sequence alignments profiles to improve protein secondary structure prediction. *Proteins* 2000;40:502–511.
- Baldi P, Brunak S. *Bioinformatics: the machine learning approach*. 2nd ed. Cambridge, MA: MIT Press; 2001.
- Rost B, Eyrich V. EVA: large-scale analysis of secondary structure prediction. *Proteins* 2001;Suppl 5:192–199.
- Pedersen TN, Lundegaard C, Nielsen M, Bohr H, Bohr J, Brunak S, Gippert GP, Lund O. Prediction of protein secondary structure at 80% accuracy. *Proteins* 2000;41:17–20.
- Cuff JA, Barton GJ. Evaluation and improvement of multiple sequence methods for protein secondary structure prediction. *Proteins* 1999;34:508–519.
- Jacoboni I, Martelli PL, Fariselli P, Compiani M, Casadio R. Predictions of protein segments with the same aminoacid sequence and different secondary structure: a benchmark for predictive methods. *Proteins* 2000;41:535–544.
- Kirshenbaum K, Young M, Highsmith S. Predicting allosteric switches in myosins. *Protein Sci* 1999;8:1806–1815.
- Young M, Kirshenbaum K, Dill KA, Highsmith S. Predicting conformational switches in proteins. *Protein Sci* 1999;8:1752–1764.
- Kabsch W, Sander C. Dictionary of protein secondary structure: pattern recognition of hydrogen-bonded and geometrical features. *Biopolymers* 1983;22:2577–2637.
- CASP3. Third community wide experiment on the critical assessment of techniques for protein structure prediction. Unpublished results available in <http://predictioncenter.llnl.gov/casp3>, December 1998.
- Riis SK, Krogh A. Improving prediction of protein secondary structure using structured neural networks and multiple sequence alignments. *J Comput Biol* 1996;3:163–183.

27. Berman HM, Westbrook J, Feng Z, Gilliland G, Bhat TN, Weissig H, Shindyalov IN, Bourne PE. The Protein Data Bank. *Nucleic Acids Res* 2000;28:235–242.
28. Hobohm U, Scharf M, Schneider R, Sander C. Selection of representative data sets. *Protein Sci* 1992;1:409–417.
29. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool. *J Mol Biol* 1990;215:403–410.
30. Bairoch A, Apweiler R. The SWISS-PROT protein sequence database and its supplement TrEMBL in 2000. *Nucleic Acids Res* 2000;28:45–48.
31. Altschul SF, Madden TL, Schaffer AA. Gapped blast and psi-blast: a new generation of protein database search programs. *Nucleic Acids Res* 1997;25:3389–3402.
32. Przybylski D, Rost B. Alignments grow, secondary structure prediction improves. *Proteins* 2002;46:195–205.
33. Lupas A, Van Dyke M, Stock J. Predicting coiled coils from protein sequences. *Science* 1991;252:1162–1164.
34. Wootton JC, Federhen S. Analysis of compositionally biased regions in sequence databases. *Methods Enzymol* 1996;266:554–571.
35. Baldi P, Pollastri G, Andersen CAF, Brunak S. Matching protein  $\beta$ -sheet partners by feedforward and recurrent neural networks. In *Proceedings of the 2000 Conference on Intelligent Systems for Molecular Biology (ISMB00)*, La Jolla, CA, Menlo Park, CA, AAAI Press; 2000. p 25–36.
36. Baldi P, Brunak S, Chauvin Y, Andersen CAF, Nielsen H. Assessing the accuracy of prediction algorithms for classification: an overview. *Bioinformatics* 2000;16:412–424.
37. Pollastri G, Baldi P, Fariselli P, Casadio R. Prediction of coordination number and relative solvent accessibility in proteins. *Proteins* 2002;47:142–153.