# Bidirectional IOHMMs and Recurrent Neural Networks for Protein Secondary Structure Prediction

Pierre Baldi[1], Soren Brunak[2], Paolo Frasconi[3], and Gianluca Pollastri[1]

[1] Dept. of Information and Computer Science, University of California at Irvine,
Irvine, CA 92697-3425
[pfbaldi|gpollast]@ics.uci.edu

[2] Center for Biological Sequence Analysis, The Technical University of Denmark
DK-2800 Lyngby, Denmark
brunak@cbs.dtu.dk

[3] Dept. of Systems and Computer Science, University of Florence
Via di Santa Marta 3, 50139 Firenze, Italy
paolo@dsi.unifi.it

**Abstract.** Prediction of protein secondary structure (SS) is one of the classical problems in bioinformatics that are best solved using computational prediction methods based on machine learning. Current state-of-the-art predictors are based on feedforward artificial neural networks fed by a *fixed-width* window of amino acids, centered on the predicted residue. Using a fixed-width small window offers the advantage of architectural simplicity and allows controlling parameter overfitting. On the other hand, relevant information is also contained in distant portions of the proteins and current methods cannot exploit this information. In this chapter, we describe two alternative architectures based on noncausal (bidirectional) dynamics. These architectures can be seen as generalizations of input-output hidden Markov models or recurrent neural networks. Unlike their conventional counterparts, their outputs depend on both upstream and downstream information. This novel algorithmic idea is a first step towards architectures capable of making predictions based on *variable* ranges of dependencies.

*Introduction*

During the last few years, computational techniques for prediction of the structure and the function of proteins have received increasingly attention. One of the classical problems is the prediction of protein secondary structure (SS), a useful approximation to the full knowledge of the three-dimensional folding that can help to elucidate the function of a protein. Human genome and other sequencing projects are enormously increasing the availability of biological sequence data, whilst first principle analyses capable of predicting SS with high accuracy are still lacking. In this situation one can expect that machine learning techniques be particularly adequate. Indeed, the present state-of-the-art predictors are based on artificial neural networks (ANNs) [3].

The influential early work of Qian and Sejnowski [29] proposes a fully connected feedforward ANN, with a local input window of typical length 13 amino acids, orthogonal encoding, and a single hidden layer. The output layer consists of three sigmoidal units with orthogonal encoding of the SS classes for the residue located at the center of the input window. A significant improvement is obtained by cascading the previous architecture with a second network to clean up the output of the lower network. The cascaded architecture reaches a performance of $Q_3 = 64.3\%$, with the correlations $C_\alpha = 0.41$ for helices, $C_\beta = 0.31$ for sheets, and $C_\gamma = 0.41$ for coils (see [4] for a review of the standard performance measures used in this chapter). Unless we specify otherwise, $Q_3$ percentages are measured on a per residue basis. Prediction of SS based on single sequences and local windows seem to be limited to $< 65-69\%$ accuracy. Increasing the size of the window, however, does not lead to improvements because of the overfitting problem associated with large networks. Building upon the work of Qian and Sejnowski [29], for a long time the best SS prediction performance has been achieved by the PHD scheme by Rost and Sander [32,33]. Rost and Sander use a number of machine learning techniques including early stopping, ensemble averages of different networks, and a weighting scheme to compensate for the well known composition biases of large low-similarity databases (roughly 30% helices, 20% sheets, and 50% coils). Most of the improvements, however, seem to result from the use of input profiles, derived from multiple alignments, that can leverage evolutionary information in the SS being considerably more conserved than the primary structure. In the 1996 Asilomar blind prediction competition CASP2 (Critical Assessment of Protein Structure Prediction), this method outperformed all others, reaching a performance level of 74%. While input profiles contain information not present in individual sequences, it is worth noting that they also discard information by losing intra-sequence corre-

lations. In [31], further NN architectural and machine learning refinements are used. One of these is the adaptive encoding of the input amino acids by the NN weight sharing technique to reduce the number of free parameters. Different networks are designed for each SS class by leveraging biological knowledge, such as the periodicity of alpha helices, with output filtering and ensemble averaging. Finally, predictions made from individual sequences are combined at the output level, using both multiple alignments and a maximum entropy weighting scheme [22]. In spite of a considerable amount of architectural design, the final performance with multiple alignments is practically identical to the one attained by PHD. The overall accuracy is $Q_3 = 71,3\%$, and correlation coefficients correlations $C_\alpha = 0.59$, $C_\beta = 0.50$, and $C_\gamma = 0.41$. More recently, Cuff and Barton [12] have compared and combined the main existing predictors. On the particular data sets used in their study, the best isolated predictor is still PHD with $Q_3 = 71.9\%$. At the 1998 CASP3 competition, the best results were obtained by one of the two programs entered by D. Jones, using a relatively simple NN architecture [20]. Out of the 35 blind sequences, the program selected 23 and achieved a performance of $Q_3 = 77.6\%$ per protein, or $Q_3 = 75.5\%$ per residue. Thus it appears today that to further improve SS prediction one should use distant information, in sequences and alignments, that is not contained in local input windows. This is particularly clear in the case of beta sheets where stabilizing bonds can be formed between amino acids far apart. This, however, poses two related challenges: (1) avoiding the overfitting associated with large-input-windows; (2) detecting sparse and weak long-ranged signals to modulate the significant local information, while ignoring the additional noise found over larger distances.

In this chapter, we approach the prediction problem in a new way, introducing an algorithm that uses the whole protein sequence rather than a short substring. To begin with, protein SS prediction can be formulated as the problem of learning a synchronous sequential translation from strings in the amino acid alphabet to strings in the SS alphabet. This task is a special form of grammatical inference. Although several symbolic algorithms exist for learning grammars [1], to the best of our knowledge they have not led to successful protein SS predictors presumably because of their scarce robustness in the presence of noisy data. Connectionist approaches, on the other hand, are based on statistical learning and therefore tend to exhibit greater robustness. The main connectionist architectures that have been investigated for grammatical inference are recurrent neural networks (RNN), with both first [11] and second-order [17] connections, as well as and input-output hidden Markov models (IOHMM) [6,7].

Both RNNs and IOHMMs are sensible alternatives to methods based on a fixed-width input window. The expressive power of these models enables them to capture distant information in the form of contextual knowledge stored into

hidden state variables. In this way, they can overcome the main disadvantage of feedforward networks, namely the linear growth of the number of parameters with the window size. Intuitively, these models are parsimonious because of the implicit weight sharing resulting from their stationarity, i.e. parameters do not vary over time. Thus, it would make sense to tackle the SS prediction problem using RNNs or IOHMMs. A more careful analysis, however, reveals a basic limitation of standard RNNs and IOHMMs in computational biology. In fact, both classes of models are causal in the sense that the output at time t does not depend on future inputs. Causality is easy to justify in dynamics that attempt to model the behavior of physical systems, or that need to operate in real time. Clearly, in these cases the response at time t cannot depend on stimulae that the system has not yet encountered. But biological sequences are not really temporal: the conformation and function of a region in a sequence may strongly depend on events located both upstream and downstream. Thus, to tackle the SS prediction problem, we develop a connectionist architecture that provides a non-causal generalization of RNNs. Our proposal is motivated by the assumption that both adaptive dynamics and non-causal processing are needed to overcome the drawbacks of local fixed-window approaches. Furthermore, we leverage evolutionary information, both at the input and output levels, using a mixture-of-estimators approach. While our current system achieves an overall performance exceeding 75% correct prediction (at least comparable to the best existing systems) the main emphasis here is on the development of new algorithmic ideas.

*Datasets*

The assignment of the SS categories to the experimentally determined 3D structure is nontrivial and is usually performed by the widely used DSSP program [21]. DSSP works by assigning potential backbone hydrogen bonds (based on the 3D coordinates of the backbone atoms) and subsequently by identifying repetitive bonding patterns. Two alternatives to this assignment scheme are the programs STRIDE and DEFINE. In addition to hydrogen bonds, STRIDE uses also dihedral angles [16]. DEFINE uses difference distance matrices for evaluating the match of interatomic distances in the protein to those from idealized SS [30]. While assignment methods impact prediction performance to some extent [12], here we concentrate exclusively on the DSSP assignments. A number of data sets were used to develop and test our algorithms. We will refer to each set using the number of sequences contained in it. The first high quality data used in this study was extracted from the Brookhaven Protein Data Bank (PDB) [9] release 77 and subsequently updated. We excluded entries if:

- They were not determined by X-ray diffraction, since no commonly used measure of quality is available for NMR or theoretical model structures.
- The program DSSP could not produce an output, since we wanted to use the DSSP assignment of protein secondary structure [21].
- The protein had physical chain breaks (defined as neighboring amino acids in the sequence having $C^{\alpha}$–distances exceeding 4.0Å).
- They had a resolution worse than 1.9Å, since resolutions better than this enables the crystallographer to remove most errors from their models.
- Chains with a length of less than 30 amino acids were also discarded.
- From the remaining chains, a representative subset with low pairwise sequence similarities was selected by running the algorithm #1 of Hobohm et al. [18], using the local alignment procedure search (rigorous Smith-Waterman algorithm) [27,28] using the pam120 matrix, with gap penalties -12, -4.

Thus we obtained a data set consisting of 464 distinct protein chains, corresponding to 123,752 amino acids, roughly 10 times more than what was available in [29]. Another set we used is the EMBL non-redundant PDB subsets that can be accessed by ftp at the site `ftp.embl-heidelberg.de`. Data details are in the file `/pub/databases/pdb_select/README`. The extraction is based on the file `/pub/databases/pdb_select/1998_june.25.gz` containing a set of non-redundant (25%) PDB chains. After removing 74 chains on which the DSSP program crashes, we obtained another set of 824 sequences, overlapping in part with the former ones. In addition, we also used the original set of 126 sequences of Rost and Sander (corresponding to 23,348 amino acid positions) as well as the complementary set of 396 non-homologue sequences (62,189 amino acids) prepared by Cuff and Barton [12]. Both sets can be downloaded at `http://circinus.ebi.ac.uk:8081/pred_res/`. Finally, we also constructed two more data sets, containing all proteins in PDB which are at least 30 amino acids long, produce DSSP output without chain breaks, and have a resolution of at least 2.5 Å. Furthermore the proteins in both sets have less than 25% identity to any of the 126 sequences of Rost and Sander. In both sets, internal homology is reduced again by Hobohm's #1 algorithm, keeping the PDB sequences with the best resolution. For one set, we use the standard 25% threshold curve for homology reduction. For the other set, however, we raise, the threshold curve by 25%. The set with 25% homology threshold contains 826 sequences, corresponding to a total of 193,249 amino acid positions, while the set with 50% homology threshold contains 1180 sequences (282,303 amino acids). Thus, to the best of our knowledge, our experiments are based on the currently largest available corpora of non-redundant data. In all but one experiment (see below), profiles were obtained from the HSSP database [35] available at `http://www.sander.embl-heidelberg.de/hssp/`.

*Bidirectional architectures*

Let us assume a probabilistic framework. A SS prediction algorithm esti-
mate, for each sequence position $t$, the posterior probabilities of secondary
structure classes (1,2, and 3, corresponding to alpha-helices, beta-sheets, and
coils), given the protein's sequence of amino acids. Formally, the output at resi-
due t is a vector of conditional probabilities $O_t=[o_{1,t}, o_{2,t}, o_{3,t}]$. In most neural
networks based approaches, $O_t$ is computed as a function of a substring of
amino acids centered on the $t$-th residue. In our method, the output is computed
as

$$O_t = \eta(F_t, B_t, I_t) \tag{1}$$

and depend on the forward (upstream) context $F_t$, the backward (downstream)
context $B_t$ and the input $I_t$. The vector $I_t \in \Re^k$ encodes a protein substring cen-
tered on $t$. In the most simple case, where the input is limited to a single amino
acid, $k = 20$ by using one-hot encoding. Larger input windows extending over
several amino acids are of course also possible. The function $\eta()$ is realized by a
neural network $N_\eta$ (see center and top connections in Figure 2). Thus to guar-
antee a consistent probabilistic interpretation, the three output units of network
$N_\eta$ are obtained as normalized exponentials (or softmax ):

$$O_{i,t} = \exp(net_{i,t}) / \Sigma_l \exp(net_{l,t}) \tag{2}$$

where $net_{i,t}$ is the activation of the $i$-th output unit at position $t$. The perform-
ance of the model can be assessed using the usual relative entropy between the
estimated and the target distribution.

The novelty of the model is in the contextual information contained in the
vectors $F_t \in \Re^n$ and especially in $B_t \in \Re^m$. These vectors are defined by the re-
current bidirectional equations:

$$F_t = \varphi(F_{t-1}, I_t) \tag{3}$$
$$B_t = \beta(B_{t+1}, I_t).$$

Here $\varphi()$ and $\beta()$ are two adaptive nonlinear state transition functions. They
can be implemented in different forms but here we assume that they are realized
by two NNs, $N_\varphi$ and $N_\beta$ (left and right subnetworks in Figure 2), with n and m
logistic output units, respectively. Thus, $N_\varphi$ and $N_\beta$ are fed by $n+k$ and $m+k$ in-
puts, respectively. Here also larger input windows are possible, especially in

combination with the weight sharing approach described in [31], and different inputs could be used for the computation of $F_t$, $B_t$, and $O_t$. The forward chain $F_t$ stores contextual information contained at the left of index $t$ and plays the same role as the internal state in standard RNNs. The novel part of the model is the presence of an additional backward chain $B_t$, in charge of storing contextual information contained at the right of index $t$, i.e. in the future. The actual form of the bidirectional dynamics is controlled by the connection weights in the sub-networks $N_\varphi$ and $N_\beta$. As we shall see, these weights can be adjusted using a maximum-likelihood approach. Since eq. 2 involves two recurrences, two corresponding boundary conditions must be specified, at the beginning and the end of the sequence. For simplicity, in our experiments we always used $F_t = 0$ and $B_{t+1} = 0$, but it is also possible to adapt the boundaries to the data, extending the technique suggested in [17] for standard RNNs.

The discrete time index $t$ ranges from 1 to $T$, the total length of the protein sequence being examined. Hence the probabilistic output $O_t$ is parameterized by a RNN and depends on the input $I_t$ and on the contextual information, from the entire protein, summarized into the pair $(F_t, B_t)$. In contrast, in a conventional NN approach this probability distribution depends only on a relatively short subsequence of amino acids. Intuitively, we can think of $F_t$ and $B_t$ as "wheels" that can be "rolled" along the protein. To predict the class at position $t$, we roll the wheels in opposite directions from the N and C terminus up to position $t$ and then combine what is read on the wheels with $I_t$ to calculate the proper output using function $\eta()$.
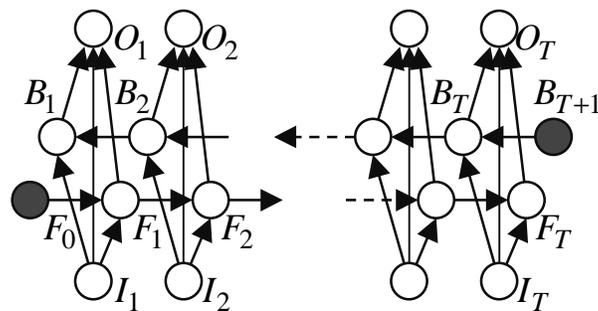


**Figure 1.** Direct dependencies among the variables involved in a bidirectional BRNN. The boundary conditions are provided by $F_0=0$, $B_{T+1}=0$ (dark nodes). To form a prediction, evidence is entered into input nodes, which are associated with the current protein sequence.

The graphical model shown in Figure 1 describes the global mapping from the input amino acid sequence to the output SS sequence. The network represents the direct dependencies among the variables $I_t$, $F_t$, $B_t$ and $O_t$, unrolled over time for $t=1,\ldots,T$. Each node is labeled by one of the variables and arcs represent direct functional dependencies. Interestingly, the same graph would represent a Bayesian network if the relationships amongst $I_t$, $F_t$, $B_t$, $O_t$ were probabilistic, rather than deterministic as in Eqs. 1, 2. In fact, such probabilistic version of the architecture would yield a bidirectional generalization of IOHMMs. Unlike RNNs, however, propagation of information in bidirectional IOHMMs is computationally expensive. The underlying Bayesian network contains undirected loops that require the use of the junction tree algorithm [19]. While inference in this network can be shown to be tractable, the corresponding time complexity of $O(n^3)$ for each time step (here $n$ is the typical number of states in the chains) limits their practical applicability to the SS prediction task [5]. An architecture resulting from Eqs. 2 and 1 is shown in Figure 2 where, for simplicity, all the NNs have a single hidden layer. The hidden state F(t) is copied back to the input. This is graphically represented in Figure 2.
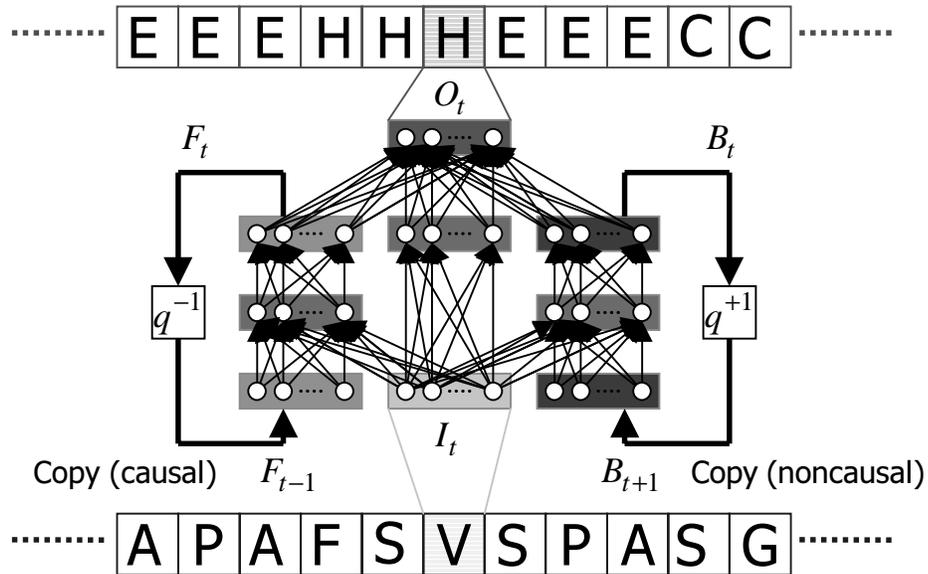


**Figure 2.** A BRNN architecture.

Removal of the backward chain would result in a standard causal RNN. The number of degrees of freedom of the model depends on the following factors: (1) the dimensions n and m of the forward and backward state vectors; (2) the number of hidden units in the three feedforward networks realizing the state transition and the output functions (see Fig. 2). It is important to remark that the BRNN has been defined as a *stationary* model, i.e. the connection weights in the networks realizing $N_\varphi$, $N_\beta$, and $N_\eta$ do not change with respect to position along the protein. This is a form of weight sharing that reduces the number of free parameter and the risk of overfitting, without necessarily sacrificing the capability to capture distant information. Since the graph shown in Fig. 1 is acyclic, nodes can be topologically sorted, defining unambiguously the global processing scheme. Using the network unrolled through time, the BRNN prediction algorithm updates all the states $F_t$ from left to right, starting from $F_0=0$. Similarly, states $B_t$ are updated from right to left. After forward and backward propagations have taken place, the predictions $O_t$ can be computed. The forward and backward propagations need to be computed from end to end only once per protein sequence. As a result, the time complexity of the algorithm is $O(TW)$, where $W$ is the number of weights and $T$ the protein length. This is the same complexity as feedforward networks fed by a fixed-size window. In the case of BRNNs, $W$ typically grows as $O(n^2)$ and the actual number of weights can be reduced by limiting the number of hidden units in the subnetworks $N_\varphi$ and $N_\beta$. Thus, inference in BRNNs is more efficient than in bidirectional IOHMMs, where complexity is $O(Tn^3)$ [5]. Learning can be formulated as a maximum likelihood estimation problem, where the log likelihood is essentially the relative entropy function between the predicted and the true conditional distribution of the secondary structure sequence given the input amino acid sequence:

$$L=\Sigma_{\text{sequences}} \ \Sigma_t z_{i,t} \log o_{i,t} \qquad (4)$$

with $z_{i,t}=1$ if the SS at position $t$ is $i$, and $z_{i,t}=0$ otherwise. The optimization problem can be solved by gradient ascent. The only difference with respect to standard RNNs is that gradients must be computed by taking into account non-causal temporal dependencies. Because the unrolled network is acyclic, the generalized backpropagation algorithm can be derived as a special case of the backpropagation through structure algorithm [15]. Intuitively, the error signal, is first injected into the leaf nodes, corresponding to the output variables $O_t$. The error is then propagated through time in both directions, by following any reverse topological sort of the unrolled network (see Figure 1). Obviously, this step also involves backpropagation through the hidden layers of the NNs. Since the model is stationary, weights are shared among the different replicas of the

NNs at different time steps. Hence, the total gradient is simply obtained by summing all the contributions associated with different time steps.

To speed-up convergence, we found it convenient to adopt a hybrid between batch and on-line weight updating strategy. Once gradients relative to a set of five protein have been computed, weights are immediately updated. This scheme was enriched also with a heuristic adaptive learning rate algorithm that progressively reduces the learning rate if the average error reduction within a fixed number of epochs falls below a given threshold.

*Long ranged dependencies*

One of the principal difficulties when training standard RNNs is the problem of vanishing gradients. This problem is related to the shape of the error surface associated with a RNN when sequences in the training set contain long time lags between inputs and corresponding supervision signals. A theoretical analysis reveals that a necessary condition for the network to robustly store long-term information is also a sufficient condition for gradient decay [8]. The consequence is that short-term information dominates long-term information, making the task of capturing distant information very hard. In the case of BRNNs, error propagation in both the forward and the backward chains is also subject to exponential decay. Thus, although the model has in principle the capability of storing remote information, such information cannot be learnt effectively.

Clearly, this is a theoretical argument and its practical impact needs to be evaluated on a per case basis. In practice, in the case of proteins, the BRNN can reliably utilize input information located within about ±15 amino acids (i.e., the total effective window size is about 31). This was empirically evaluated by feeding the model with increasingly long protein fragments. We observed that the average predictions at the central residues did not significantly change if fragments were extended beyond 41 amino acids. This is an improvement over standard NNs with input window sizes ranging from 11 to 17 amino acids [33,31]. Yet, there is presumably relevant information located at longer distances that our model may have not been able to discover so far. To limit this problem, we propose a remedy motivated by recent studies [24] suggesting that the vanishing gradients problem can be mitigated by the use of an explicit delay line applied to the output, which provides shorter paths for the effective propagation of error signals. This idea cannot be applied directly to BRNNs since output feedback, combined with bidirectional propagation, would generate cycles in the unrolled network. A similar mechanism, however, can be implemented using the following modified dynamics:

$$F_t = \varphi(F_{t-1}, F_{t-2}, \dots, F_{t-s}, I_t)$$
$$B_t = \beta(B_{t+1}, B_{t+2}, \dots, B_{t+s}, I_t). \tag{5}$$

The explicit dependence on forward or backward states introduces *shortcut* connections in the graphical model, forming shorter paths along which gradients can be propagated. This is akin to introducing higher order Markov chains in the probabilistic version. However, unlike Markov chains where the number of parameters would grow exponentially with $s$, in the present case the number of parameters grows only linearly with $s$. To reduce the number of parameters, a simplified version of the above dynamics limits the dependencies to state vectors located s residues apart from $t$:

$$F_t = \varphi(F_{t-1}, F_{t-s}, I_t)$$
$$B_t = \beta(B_{t+1}, B_{t+s}, I_t). \tag{6}$$

Another variant of the basic architecture which also allows to increase the effective window size consists in feeding the output networks with a window in the forward and backward state chains. In this case, the prediction is computed as

$$O_t = \eta(F_{t-k}, \dots, F_{t+k}, B_{t-k}, \dots B_{t+k}, I_t) \tag{7}$$

*Multiple Alignments and Mixture of Estimators*

Multiple alignments and mixture of estimators are two algorithmic ideas which have been shown to be very effective in the protein SS prediction task. Both of them have been incorporated in our BRNN-based system. Multiple predictors can be obtained by varying the size of the BRNN, as controlled by the dimensions of the state vectors and the number of hidden units. Moreover, different predictors can be obtained using profiles, or multiple alignments, in input mode [33], or in output mode [31]. In the first case, instead of a code for the current amino acid, the input $I_t$ contains the relative frequencies of amino acids at position t in the protein family. In the second case, a separate sequence of SS predictions is obtained for each aligned protein, and then all the predictions are averaged in each column. Since the two methods give different prediction errors − the input mode, for instance, yields slightly more accurate beta-sheet predictions − it is reasonable to build ensembles containing both types of predictors.

The setup developed in [25] has also been used to select the profiles. The setup contains methods similar to the ones applied earlier by Sander and Schneider ([34]), where the key parameter is the similarity threshold (in terms of identical residues in a particular pairwise alignment). In [25], the similarity threshold (dividing sequences with structural homology from those without) had the form $I<290/\sqrt{L}$, where $L$ is the length of the alignment. The threshold was then used to build a profile for all the relevant PDB entries from the matches found in the SWISS-PROT database.

*Implementation and Results*

We carried out several preliminary experiments to tune up and evaluate the prediction system. DSSP classes were assigned to three secondary structure classes $\alpha,\beta$, and $\gamma$ as follows: $\alpha$ is formed by DSSP class H, $\beta$ by E, and $\gamma$ by everything else (including DSSP classes F, S, T, B, and I). This assignment is slightly different from other assignments reported in the literature. For example, in [31], $\alpha$ contains DSSP classes H, G, and I. In the CASP competition [26,10], $\alpha$ contains H, and G, while $\beta$ contains E, and B.

In a first set of experiments, we used the 824 sequences dataset and reserved 2/3 of the available data for training, and 1/3 for testing. We trained several BRNNs of different sizes and different architectural details. In all experiments, we set $n=m$ and we tried different values for $n$ and $k$ (see Eq. 7). The number of free parameters varied from about 1400 to 2600. Qualitatively we observed that using $k>0$ can improve accuracy, but increasing $n$ beyond 12 does not help because of overfitting. Results for this method, without using profiles, are summarized in the first rows of Table 1. By comparison, we also trained several feedforward NNs on the same data. The best feedforward NN achieved $Q_3=67.2\%$ accuracy using a window of 13 amino acids. By enriching the feedforward architecture with adaptive input encoding and output filtering, as in [31], 68.5% accuracy was achieved (output filtering actually increases the length of the input window). Hence, the best BRNN outperforms our best feedforward network, even when additional architectural design is included.

Subsequent experiments included the use of profiles. Table 1 reports the best results obtained by using multiple alignments, both at the input and output levels. Profiles at the input level consistently yielded better results. The best feedforward networks trained in the same conditions achieve $Q_3= 73.0\%$ and 72.3%, respectively.

**Table 1.** Experimental results using a single BRNN and 1/3 of the data as test set. $h_\varphi$, $h_\beta$ and $h_\eta$ are the number of hidden units for the transition networks $N_\varphi$, $N_\beta$ and the output network $N_\eta$ respectively. We always set $h_\varphi = h_\beta$.

| Profiles | $n$ | $k$ | $h_\varphi$ | $h_\eta$ | $W$ | Accuracy ($Q_3$) |
|----------|-----|-----|-------------|----------|------|------------------|
| No | 7 | 2 | 8 | 11 | 1611 | 68.7% |
| No | 9 | 2 | 8 | 11 | 1899 | 68.8% |
| No | 7 | 3 | 8 | 11 | 1919 | 68.6% |
| No | 8 | 3 | 9 | 11 | 2181 | 68.8% |
| No | 20 | 0 | 17 | 11 | 2601 | 67.6% |
| Output | 9 | 2 | 8 | 11 | 1899 | 72.6% |
| Output | 8 | 3 | 9 | 11 | 2181 | 72.7% |
| Input | 9 | 2 | 8 | 11 | 1899 | 73.3% |
| Input | 8 | 3 | 9 | 11 | 2181 | 73.4% |
| Input | 12 | 3 | 9 | 11 | 2565 | 73.6% |

In a second set of experiments (also based on the 824 sequences), we combined several BRNNs to form an ensemble, as in [23], using a simple averaging scheme. Different networks were obtained by varying architectural details such as n, k, and the number of hidden units. Combining 6 networks using profiles at the input level we obtained the best accuracy $Q_3=75.1\%$, measured in this case using 7-fold cross validation. We also tried to include in the ensemble a set of 4 BRNNs using profiles at the output level but performance in this way slightly decreased to 75.0%.

A study for assessing the capabilities of the model in capturing long ranged information was also performed. Results indicate that the model is sensitive to information located within about ±15 amino acids. Although this value is not *very* high, it should be remarked that typical feedforward nets reported in the literature do not exploit information beyond $\tau = 8$.

To further explore the long-range information problem we conducted another set of experiments using BRNNs with simplified shortcuts (see eq. 6). In this case, as for the results reported in Table 1, we used a single model (rather than a mixture) and the test set method (1/3 of the available data) for measuring accuracy. We tried all values of s from 1 to 10, but in no case we could observe a significant performance improvement on the test set. Interestingly, our experiments showed that using shortcuts reduces the convergence difficulties associated with vanishing gradients: accuracy on the training set increased from 75.7% using no shortcuts to 76.9% with s=3. On the other hand, the gap between training set and test set performance also increased. Thus overfitting offset the convergence improvement, probably because long-range information is too sparse and noisy.

Another experiment was conducted by training on all the 824 sequences and using the official test sequences used at the 1998 CASP3 competition. In this case, we adopted a slightly different class assignment for training (DSSP classes H, G, and I were merged together). The CASP3 competition was won by one of the two programs entered by D. Jones, which selected 23 out of 35 proteins obtaining a performance of $Q_3$=77.6% per protein, or $Q_3$=75.5% per residue. We evaluated that system on the whole set of 35 proteins by using Jones' prediction server at `http://137.205.156.147/psiform.html`. It achieved $Q_3$=74.3% per residue and 76.2% per protein. On the same 35 sequences our system achieved $Q_3$=73.0% per residue and 74.6 per protein. A test set of 35 proteins is relatively small for drawing general conclusions. Still, we believe that this result confirms the effectiveness of the proposed model, especially in consideration of the fact that Jones' system builds upon more recent profiles from TrEMBL database [2]. These profiles contain many more sequences than our profiles, which are based on the older HSSP database, leaving room for further improvements of our system.

**Table 2.** First confusion matrix derived with an ensemble of 8 BRNNs with 2/3-1/3 data splitting. First row provides percentages of predicted helices, sheets, and coils within (DSSP-assigned) helices.

|   | pred $\alpha$ | pred $\beta$ | pred $\gamma$ |
|---|---|---|---|
| $\alpha$ | 78.61% | 3.13% | 18.26% |
| $\beta$ | 5.00% | 61.49% | 33.51% |
| $\gamma$ | 10.64% | 9.37% | 79.99% |

**Table 3.** Same as above. First row provides percentages of (DSSP-assigned) helices, sheets, and coils within the predicted helices.

|   | $\alpha$ | $\beta$ | $\gamma$ |
|---|---|---|---|
| pred $\alpha$ | 88.77% | 3.74% | 15.49% |
| pred $\beta$ | 5.11% | 73.17% | 21.72% |
| pred $\gamma$ | 11.71% | 15.63% | 72.66% |

To further compare our system with other predictors, as in [12], we also trained an ensemble of BRNNs using the 126 sequences in the Rost and Sander data set. The performance on the 396 test sequences prepared by Cuff and Barton is $Q_3 = 72.0\%$. This is slightly better than the 71.9% score for the single best predictor (PHD) amongst (DSC, PHD, NNSSP, and PREDATOR) reported in [12]. This result is also achieved with the CASP class assignment.

Finally, we also trained an ensemble of 6 BRNNs using the set containing 826 sequences with less than 25% identity to the 126 sequences of Rost and Sander. When tested on the 126 sequences, the system achieves $Q_3 = 74.7\%$ per residue, with correlation coefficients $C_\alpha = 0.692$, $C_\beta = 0.571$, and $C_\gamma = 0.544$. This is again achieved with the harder CASP assignment. In contrast, the $Q_3 = 75.1\%$ described above was obtained by 7 fold cross-validation on 824 sequences and with the easier class assignment (H→α, E→β, the rest →γ). The same experiment was performed using the larger training set of 1,180 sequences having also less than 25% identity with the 126 sequences of Rost and Sander, but with a less stringent redundancy reduction requirement. In this case, and with the same hard assignment, the results are $Q_3 = 75.3\%$ with correlation coefficients $C_\alpha = 0.704$, $C_\beta = 0.583$, and $C_\gamma = 0.550$. The corresponding confusion matrices are given in Tables 2 and 3. Table 4 provides a summary of the main results with different datasets.

**Table 4.** Summary of main performance results.

| Training sequences | Test sequences | Class assignment | Performance per residue |
|---|---|---|---|
| 824 (2/3) | 824 (1/3) | Default | $Q_3 = 75.1\%$ |
| 824 | 35 | CASP | $Q_3 = 73.0\%$ |
| 126 | 396 | CASP | $Q_3 = 72.0\%$ |
| 826 | 126 | CASP | $Q_3 = 74.7\%$ |
| 1180 | 126 | CASP | $Q_3 = 75.3\%$ |

*Conclusions*

Given the large number of protein sequences available through genome and other sequencing projects, even small percentage improvements in SS prediction can be significant. The system presented here achieves an overall performance of over 75% correct classification, at least comparable to the best existing predictors, but using a different NN approach based on recurrent networks and bidirectional dynamics.

Interestingly, we have circumstantial evidence that the two methods behave in significantly different ways: there exist sequences for which our method achieves over 80% correct prediction, while Jones method is below 70%, and vice versa. Such differences require further study, and suggest that both methods could be combined to further improve the results. In particular, if the ad-

vantage of the method of Jones resides in the type of alignments used, similar alignments could be incorporated in the BRNN approach. While there is room for performance improvement, one should also not forget that 100% correct prediction, from the primary sequence alone, is probably unachievable if nothing else because a minority of proteins may not fold spontaneously, or because beta-sheet partner strands may be located on a different chain.

Most importantly, perhaps, we have developed here new algorithmic ideas that begin to address the problem of long-range dependencies. Unlike feedforward networks, BRNNs can possibly prove advantageous from this point of view and our preliminary experiments encourage further investigations. This work could be extended in additional directions. These include architectural variations (such as the use of larger input windows in the BRNN architecture), non-symmetrical chains for the past and the future, and the use of priors on the parameters and/or the architecture together with a maximum a posteriori learning approach. It is also worth noting that using multi-layered perceptrons for implementing $\beta()$ and $\varphi()$ is just one of the available options. For example, a generalization of second-order RNN [17] is an easily conceivable alternative parameterization.

Finally, it is clear that the ideas introduced can be applied to other problems in bioinformatics, as well as other domains, where non-causal dynamical approaches are suitable. Obvious candidates for further tests of the general method include the prediction of DNA exon/intron boundaries and promoter regions, of RNA secondary structure, and of protein functional domains, such as signal peptides.

*References*

1. Angluin, D. & Smith, C. (1987). Inductive inference. In *Encyclopedia of Artificial Intelligence*. J. Wiley and Sons, New York, pp. 409–418.
2. Bairoch, A. & Apweiler, R. (1999). The SWISS-PROT protein sequence data bank and its supplement TrEMBL in 1999. *Nucleic Acids Res*, **27**, 49–54.
3. Baldi, P. & Brunak, S. (1998). *Bioinformatics: The Machine Learning Approach*. MIT Press, Cambridge, MA.
4. Baldi, P., Brunak, S., Chauvin, Y. & Nielsen, H. (1999a). Assessing the accuracy of prediction algorithms for classification: an overview. Submitted for publication.
5. Baldi, P., Brunak, S., Frasconi, P., Pollastri, G. & Soda, G. (1999b). Bidirectional dynamics for protein secondary structure prediction. IJCAI-99 workshop on neural, symbolic, and reinforcement methods for sequence learning (unpublished).
6. Baldi, P. & Chauvin, Y. (1996). Hybrid modeling, HMM/NN architectures, and protein applications. *Neural Computation*, **8**, 1541–1565.
7. Bengio, Y. & Frasconi, P. (1996). Input-output HMM's for sequence processing. *IEEE Trans. on Neural Networks*, **7**, 1231– 1249.

8. Bengio, Y., Simard, P. & Frasconi, P. (1994). Learning long-term dependencies with gradient descent is difficult. *IEEE Transactions on Neural Networks*, **5**, 157–166.

9. Bernstein, F. C. & et al. (1977). The protein data bank: A computer based archival file for macromolecular structures. *J. Mol. Biol.*, **112**, 535–542.

10. CASP3 (1998). Third community wide experiment on the critical assessment of techniques for protein structure prediction. Unpublished results available in `http://predictioncenter.llnl.gov/casp3`.

11. Cleeremans, A. (1993). *Mechanisms of Implicit Learning. Connectionist Models of Sequence Processing*. MIT Press.

12. Cuff, J. A. & Barton, G. J. (1999). Evaluation and improvement of multiple sequence methods for protein secondary structure prediction. *Proteins*, **34**, 508–519.

13. Eisenberg, D. (1997). Into the black night. *Nat. Struct. Biol.*, **4**, 95–97.

14. Forcada, M. & Carrasco, R. (1995). Learning the initial state of a second-order recurrent neural network during regular-language inference. *Neural Computation*, **7**, 923–930.

15. Frasconi, P., Gori, M. & Sperduti, A. (1998). A general framework for adaptive processing of data structures. *IEEE Trans. on Neural Networks*, **9**, 768–786.

16. Frishman, D. & Argos, P. (1995). Knowledge-based secondary structure assignment. *Proteins*, **23**, 566–579.

17. Giles, C. L., Miller, C. B., Chen, D., Chen, H. H., Sun, G. Z. & Lee, Y. C. (1992). Learning and extracting finite state automata with second-order recurrent neural networks. *Neural Computation*, **4**, 393–405.

18. Hobohm, U., Scharf, M., Schneider, R. & Sander, C. (1992). Selection of representative data sets. *Prot. Sci.*, **1**, 409–417.

19. Jensen, F. V., Lauritzen, S. L. & Olosen, K. G. (1990). Bayesian updating in recursive graphical models by local computations. *Comput. Stat. Quarterly*, **4**, 269–282.

20. Jones, D. T. (1999). Protein Secondary Structure Prediction Based on Position-specific Scoring Matrices. *J.Mol.Biol.*, **292**, 195–202.

21. Kabsch, W. & Sander, C. (1983). Dictionary of protein secondary structure: pattern recognition of hydrogen-bonded and geometrical features. *Biopolymers*, **22**, 2577–2637.

22. Krogh, A. & Mitchinson, G. (1995). Maximum entropy weighting of aligned sequences of proteins of DNA. In et al., C. R., (ed.) *Proc. 3rd Int. Conf. on Intelligent Systems for Molecular Biology*. Menlo Park, CA, pp. 215–221.

23. Krogh, A. & Vedelsby, J. (1995). Neural network ensembles, cross validation, and active learning. In Tesauro, G., Touretzky, D. & Leen, T., (eds.) NIPS 7. The MIT Press, pp. 231–238.

24. Lin, T., Horne, B. G., Tino, P. & Giles, C. L. (1996). Learning long-term dependencies in NARX recurrent neural networks. *IEEE Transactions on Neural Networks*, **7**, 1329–1338.

25. Lund, O., Frimand, K., Gorodkin, J., Bohr, H., Bohr, J., Hansen, J. & Brunak, S. (1997). Protein distance constraints predicted by neural networks and probability density functions. *Prot. Eng.*, **10**, 1241–1248.

26. Moult, J. & et al. (1997). Critical assessment of methods of protein structure prediction (CASP): Round II. *Proteins*, **29**, 2–6. Supplement 1.

27. Myers, E. W. & Miller, W. (1988). Optimal alignments in linear space. *Comput. Appl. Biosci.*, **4**, 11–7.

28. Pearson, W. R. (1990). Rapid and sensitive sequence comparison with FASTP and FASTA. *Meth. Enzymol.*, **183**, 63–98.

29. Qian, N. & Sejnowski, T. J. (1988). Predicting the secondary structure of glubular proteins using neural network models. *J. Mol. Biol.*, **202**, 865–884.

30. Richards, F. M. & Kundrot, C. E. (1988). Identification of structural motifs from protein coordinate data: secondary structure and first-level supersecondary structure. *Proteins*, **3**, 71–84.

31. Riis, S. K. & Krogh, A. (1996). Improving prediction of protein secondary structure using structured neural networks and multiple sequence alignments. *J. Comput. Biol.*, **3**, 163–183.

32. Rost, B. & Sander, C. (1993). Improved prediction of protein secondary structure by use of sequence profiles and neural networks. *Proc. Natl. Acad. Sci.* USA, **90**, 7558–7562.

33. Rost, B. & Sander, C. (1994). Combining evolutionary information and neural networks to predict protein secondary structure. *Proteins*, **19**, 55–72.

34. Sander, C & Schneider, R. (1991). Database of homology-derived protein structures and the structural meaning of sequence alignment. *Proteins*, **9**, 56–68.

35. Schneider, R., de Daruvar, A. & Sander, C. (1997). The HSSP database of protein structure-sequence alignments. *Nucleic Acids Res.*, **25**, 226–230.