

# Probabilistic Prediction of Protein Secondary Structure Using Causal Networks

(Extended Abstract)

**Arthur L. Delcher\***

Computer Science Dept.  
Loyola College  
Baltimore, MD 21210

**Simon Kasif\***

Dept. of Computer Science  
Johns Hopkins University  
Baltimore, MD 21218

**Harry R. Goldberg**

Mind-Brain Institute  
Johns Hopkins University  
Baltimore, MD 21218

**William H. Hsu**

Dept. of Computer Science  
Johns Hopkins University  
Baltimore, MD 21218

## Abstract

In this paper we present a probabilistic approach to analysis and prediction of protein structure. We argue that this approach provides a flexible and convenient mechanism to perform general scientific data analysis in molecular biology. We apply our approach to an important problem in molecular biology—predicting the secondary structure of proteins—and obtain experimental results comparable to several other methods. The causal networks that we use provide a very convenient medium for the scientist to experiment with different empirical models and obtain possibly important insights about the problem being studied.

## Introduction

Scientific analysis of data is an important potential application of Artificial Intelligence (AI) research. We believe that the ultimate data analysis system using AI techniques will have a wide range of tools at its disposal and will adaptively choose various methods. It will be able to generate simulations automatically and verify the model it constructed with the data generated during these simulations. When the model does not fit the observed results the system will try to explain the source of error, conduct additional experiments, and choose a different model by modifying system parameters. If it needs user assistance, it will produce a simple low-dimensional view of the constructed model and the data. This will allow the user to guide the system toward constructing a new model and/or generating the next set of experiments. We believe that flexibility, efficiency and direct representation of causality are key issues in the choice of representation in such a system.

---

\* Supported by NSF DARPA Grant CCR-8908092 and AFOSR Grant AFOSR-89-1151

As a first step, in this paper we present a probabilistic approach to analysis and prediction of protein structure. We argue that this approach provides a flexible and convenient mechanism to perform general scientific data analysis in molecular biology. We apply our approach to an important problem in molecular biology: predicting the secondary structure of proteins [?; ?]. A number of methods have been applied to this problem with various degree of success [?; ?; ?; ?; ?]. In addition to obtaining experimental results comparable to other methods, there are several theoretically and practically important observations that we have made in experimenting with our system.

- It has been claimed in several papers that probabilistic (statistical) approaches have been outperformed by neural network methods and memory-based methods by a wide margin. We show that probabilistic methods are comparable to other methods in prediction quality. In addition, the predictions generated by our methods have precise quantitative semantics which is not shared by other classification methods. Specifically, all the causal and statistical independence assumptions are made explicit in our networks thereby allowing biologists to study causal links in a convenient manner. This generalizes correlation studies that are normally used in statistical analysis of data.
- Our method provides a very flexible tool to experiment with a variety of modelling strategies. This flexibility allows a biologist to perform many practically important statistical queries which can yield important insight into a problem.
- From the theoretical point of view we found that different ways to model the domain produce practically different results. This is an experience that AI researchers encounter repeatedly in many knowledge-representation schemes: different coding of the prob-

lem in the architecture results in dramatic differences in performance. This has been observed in production systems, neural networks, constraint networks and other representations. Our experience reinforces the thesis that while knowledge representation is a key issue in AI, a knowledge-representation system typically provides merely the programming language in which a problem must be expressed. The coding, analogous to an algorithm in procedural languages, is perhaps of equally great importance. However, the importance of this issue is grossly underestimated and not studied as systematically and rigorously as knowledge representation languages.

- Previous methods for protein folding were based on the window approach. That is, the learning algorithm attempted to predict the structure of the central amino acid in a “window” of  $k$  amino acids residues. It is well recognized that in the context of protein folding, very minimal mutations (amino acid substitutions) often cause significant changes in the secondary structure located far from the mutation site. Our method is aimed at capturing this behavior.

## Protein Folding

Proteins have a central role in essentially all biological processes. They control cellular growth and development, they are responsible for cellular defense, they control reaction rates, they are responsible for propagating nerve impulses, and they serve as the conduit for cellular communication. The ability of proteins to perform these tasks, *i.e.*, the *function* of a protein, is directly related to its *structure*. The results of Christian Anfinsen’s work in the late 1950’s indicated that a protein’s unique structure is specified by its amino acid sequence. This work suggested that a protein’s conformation could be specified if its amino acid sequence was known, thus defining the protein folding problem. Unfortunately, nobody has been able to put this theory into practice.

The biomedical importance of solving the protein folding problem cannot be overstressed. Our ability to design genes—the molecular blueprints for specifying a protein’s amino acid sequence—has been refined. These genes can be implanted into a cell and this cell can serve as the vector for the production of large quantities of the protein. The protein, once isolated, potentially can be used in any one of a multitude of applications—uses ranging from supplementing the human defense system to serving as a biological switch for controlling abnormal cell growth and development. A critical aspect of this process is the ability to specify the amino acid sequence which defines the required

conformation of the protein.

Traditionally, protein structure has been described at three levels. The first level defines the protein’s amino acid sequence, the second considers local conformations of this sequence, *i.e.*, the formation of rod-like structures called  $\alpha$ -helices, planar structures called  $\beta$ -sheets, and intervening sequences often categorized as coil. The third level of protein structure specifies the global conformation of the protein. Due to limits on our understanding of solutions to the protein folding problem, most of the emphasis on structure prediction has been at the level of secondary structure prediction.

There are fundamentally two approaches that have been taken to predict the secondary structure of proteins. The first approach is based on theoretical methods and the second is based on data derived empirically. Theoretical methods rely on our understanding of the rules governing amino acid interactions, they are mathematically sophisticated and computationally time-intensive. Conversely, empirically based techniques combine a heuristic with a probabilistic schema in determining structure. Empirical approaches have reached prediction rates approaching 70%—the apparent limit given our current base of knowledge.

The most obvious weakness of empirically based prediction schemes is their reliance on exclusively local influences. Typically, a window that can be occupied by 9-13 amino acids is passed along the protein’s amino acid sequence. Based on the context of the central amino acid’s sequence neighbors, it is classified as belonging to a particular structure. The window is then shifted and the amino acid which now occupies the central position of the window is classified. This is an iterative process which continues until the end of the protein is reached. In reality, the structure of an amino acid is determined by its local environment. Due to the coiled nature of a protein, this environment may be influenced by amino acids which are far from the central amino acid in sequence but not in space. Thus, a prediction scheme which considers the influence of amino acids which are, in sequence, far removed from the central amino acid of the window may improve our ability to successfully predict a protein’s conformation.

## Notation

For the purpose of this paper, the set of proteins is assumed to be a set of sequences (strings) over an alphabet of twenty characters (different capital letters) that correspond to different amino acids. With each protein sequence of length  $n$  we associate a sequence of secondary structure descriptors of the same length. The structure descriptors take three values:  $h$ ,  $e$ ,  $c$  that correspond to  $\alpha$ -helix,  $\beta$ -sheet and coil. That

is, if we have a subsequence of  $hh\dots h$  in positions  $i, i+1, \dots, i+k$  it is assumed that the protein sequence in those positions folded as a helix. The classification problem is typically stated as follows. Given a protein sequence of length  $n$ , generate a sequence of structure predictions of length  $n$  which describes the secondary structure of the protein sequence. Almost without exception all previous approaches to the problem have used the following approach. The classifier receives a window of length  $2K + 1$  (typically  $K < 12$ ) of amino acids. The classifier then predicts the secondary structure of the central amino acid (*i.e.*, the amino acid in position  $K$ ) in the window.

## A Probabilistic Framework for Protein Analysis

When making decisions in the presence of uncertainty, it is well-known that Bayes rule provides an optimal decision procedure, assuming we are given all prior and conditional probabilities. There are two major difficulties with using the approach in practice. The problem of reasoning in general Bayes networks is  $\mathcal{NP}$ -complete, and we often do not have accurate estimates of the probabilities. However, it is known that when the structure of the network has a special form it is possible to perform a complete probabilistic analysis efficiently. In this section we show how to model probabilistic analysis of the structure of protein sequences as belief propagation in causal trees. In the full version of the paper we also describe how we dealt with problems such as undersampling and regularization. The general schema we advocate has the following form. The set of nodes in the networks are either protein-structure nodes (*PS*-nodes) or evidence nodes (*E*-nodes). Each *PS*-node in the network is a discrete random variable  $X_i$  that can take values which correspond to descriptors of secondary structure, *i.e.*, segments of  $h$ 's,  $e$ 's and  $c$ 's. With each such node we associate an evidence node that again can assume any of a set of discrete values. Typically, an evidence node would correspond to an occurrence of a particular subsequence of amino acids at a particular location in the protein. With each edge in the network we will associate a matrix of conditional probabilities. The simplest possible example of a network is given in Figure 1.

We assume that all conditional dependencies are represented by a causal tree. This assumption violates some of our knowledge of the real-world problem, but provides an approximation that allows us to perform an efficient computation. For an exact definition of a causal tree see Pearl [?].

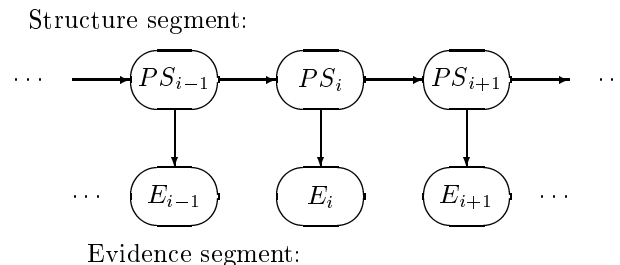


Figure 1: Causal tree model.

## Protein Modeling Using Causal Networks

As mentioned above, the network is comprised of a set of protein-structure nodes and a set of evidence nodes. Protein-structure nodes are finite strings over the alphabet  $\{h, e, c\}$ . For example the string  $hhhhhh$  is a string of six residues in an  $\alpha$ -helical conformation, while  $eecc$  is a string of two residues in a  $\beta$ -sheet conformation followed by two residues folded as a coil. Evidence nodes are nodes that contain information about a particular region of the protein. Thus, the main idea is to represent physical and statistical rules in the form of a probabilistic network. We note that the main point of this paper is advocating the framework of causal networks as an experimental tool for molecular biology applications rather than focusing on a particular network. The framework allows us flexibility to test causal theories by orienting edges in the causal network.

For our initial experiments we have chosen the simplest possible models. In this paper we describe two that we feel are particularly important: a classical Hidden Markov Model using the Viterbi algorithm and causal trees using Pearl's belief updating. We shall show that the second approach is better and matches in accuracy other methods that have a less explicitly quantitative semantics.

In our first set of experiments we converged on the following model that seems to match in performance many existing approaches. The network looks like a set of *PS*-nodes connected as a chain. To each such node we connect a single evidence node. In our experiments the *PS*-nodes are strings of length two or three over the alphabet  $\{h, e, c\}$  and the evidence nodes are strings of the same length over the set of amino acids. The following example clarifies our representation. Assume we have a string of amino acids  $GSAT$ . We model the string as a network comprised of three evidence nodes  $GS, SA, AT$  and three *PS*-nodes. The network is shown in Figure 2. A correct prediction will assign

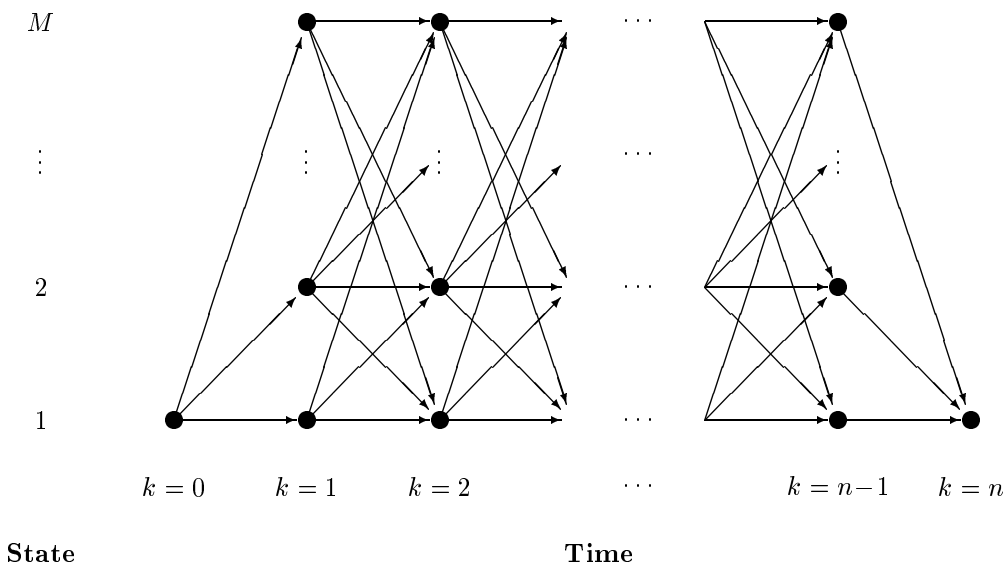


Figure 3: Modelling the Viterbi algorithm as a shortest path problem.

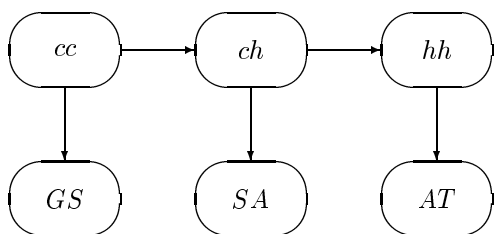


Figure 2: Example of causal tree model using pairs, showing protein segment *GSAT* with corresponding secondary structure *cch*

the values *cc*, *ch*, and *hh* to the *PS*-nodes as shown in the figure.

Let  $X_0, X_1, \dots, X_n$  be a set of *PS*-nodes connected as in Figure 1. Generally, speaking the distribution for the variable  $X_i$  in the causal network as below can be computed using the following formulae. Let  $e_{X_i}^- = e_i, e_{i+1}, \dots, e_n$  denote the set of evidence nodes to the right of  $X_i$ , and let  $e_{X_i}^+ = e_1, e_2, \dots, e_{i-1}$  be the set of evidence nodes to the left of  $X_i$ . By the assumption of independence explicit in the network we have

$$P(X_i | X_{i-1}, e_{X_i}^+) = P(X_i | X_{i-1})$$

Thus,

$$P(X_i | e_{X_i}^+, e_{X_i}^-) = \alpha P(e_{X_i}^- | X_i) P(X_i | e_{X_i}^+)$$

where  $\alpha$  is some normalizing constant. For length consideration we will not describe the algorithm to compute the probabilities. The reader is referred to Pearl for a detailed description [?]. Pearl gives an efficient procedure to compute the belief distribution of every node in such a tree. Most importantly, this procedure operates by a simple efficient propagation mechanism that operates in linear time.

### Protein Modeling Using the Viterbi Algorithm

In this section we describe an alternative model for prediction. This model has been heavily used in speech understanding systems, and indeed was suggested to us by Kai Foo Lee whose system using similar ideas achieves remarkable performance on speaker-independent continuous speech understanding.

We implemented the Viterbi algorithm and compare its performance to the method outlines above. We briefly describe the method here. We follow the discussion by Forney [?].

We assume a Markov process which is characterized by a finite set of state transitions. That is, we assume the process at time  $k$  can be described by a random variable  $X_k$  that assumes a discrete number of values (states)  $1, \dots, M$ . The process is Markov, *i.e.*, the probability  $P(X_{k+1} | X_0, \dots, X_k) = P(X_{k+1} | X_k)$ . We denote the process by the sequence  $X = X_0, \dots, X_k$ . We are given a set of observations  $Z = Z_0, \dots, Z_k$  such that  $Z_i$  depends only on the transition  $T_i = (X_{i+1}, X_i)$ . Specifically,  $P(Z | X) = \prod_{k=0}^n (Z_k | X_k)$ . The Viterbi al-

gorithm is a solution to the maximum a posteriori estimation of  $X$  given  $Z$ . In other words we are seeking a sequence of states  $X$  for which  $P(Z|X)$  is maximized.

An intuitive way to understand the problem is in graph theoretic terms. We build a  $n$ -level graph that contains  $nM$  nodes (see Figure 3). With each transition we associate an edge. Thus, any sequence of states has a corresponding path in the graph. Given the set of observations  $Z$  with any path in the graph we associate a length  $L = -\ln P(X, Z)$ . We are seeking a shortest path in the graph. However, since

$$\begin{aligned} P(X, Z) &= P(X)P(Z|X) \\ &= \prod_{k=0}^{n-1} P(X_{k+1}|X_k) \prod_{k=0}^{n-1} P(Z_k|X_{k+1}, X_k) \end{aligned}$$

if we define  $\lambda(T_k) = -\ln P(X_{K+1}|X_K) - \ln P(Z_k|T_k)$  we obtain that  $-\ln P(Z, X) = \sum_{k=0}^{n-1} \lambda_k$ .

Now we can compute the shortest path through this graph by a standard application of shortest path algorithms specialized to directed acyclic graphs. For each time step  $i$  we simply maintain  $M$  paths which are the shortest path to each of the possible states we could be in at time  $i$ . To extend the path to time step  $i + 1$  we simply compute the lengths of all the paths extended by one time unit and maintain the shortest path to each one of the  $M$  possible states at time  $i + 1$ .

Our experimentation with the Viterbi algorithm was completed in Spring 1992. We recently learned that David Haussler [?] and his group suggested the Viterbi algorithm framework for protein analysis as well. They experimented on a very different problem and also obtain interesting results. We document the performance of Viterbi on our problem even though, as described below, the causal-tree method outperformed Viterbi. The difference between the methods is that the Viterbi algorithm predicts the most likely complete sequence of structure elements, whereas the causal-tree method makes separate predictions about individual  $PS$ -nodes.

## Experiments

The experiments we conducted were performed to allow us to make a direct comparison with previous methods that have been applied to this problem. We followed the methodology described in [?; ?] which did a thorough cross-validated testing of various classifiers for this problem. Since it is known that two proteins that are homologous (similar in chemical structure) tend to fold similarly and therefore generate accuracies of predictions that are often overly optimistic, it is important to document the precise degree of homology between the training set and the testing

Trial	Positions	Correct Using:	
		Pairs	Triples
1	2339	1518 (64.9%)	1469 (62.8%)
2	2624	1567 (59.7%)	1518 (57.9%)
3	2488	1479 (59.5%)	1435 (57.7%)
4	2537	1666 (65.7%)	1604 (63.2%)
5	2352	1437 (61.1%)	1392 (59.2%)
6	2450	1510 (61.6%)	1470 (60.0%)
7	2392	1489 (62.3%)	1447 (60.5%)
8	2621	1656 (63.2%)	1601 (61.1%)
<b>All</b>	19803	12322 (62.2%)	11936 (60.3%)

Table 1: Causal tree results for 8-way cross-validation using segments of length 2 and length 3.

set. In our experiments the set of proteins was divided into eight subsets. We perform eight experiments in which we train the network on seven subsets and then predict on the remaining subset. The accuracies are averaged over all eight experiments. This methodology is referred to as  $k$ -way cross validation.

## Experimental Results

We report the accuracy of prediction on individual residues and also on predicting runs of helices and sheets. Table 1 shows the prediction accuracy of our methods using the causal network method for each one of the eight trials in our 8-way cross-validation study. In the pairs column we document the performance of the causal network described earlier using  $PS$ -nodes and  $E$ -nodes that represent protein segments of length 2. The triples column gives the results for the same network with segments of length 3. The decrease in accuracy for triples is a result of undersampling.

Table 2 shows the performance of our method in predicting the secondary structure at each amino acid position in comparison with other methods. In Table 3 we report the performance of our method on predicting runs of helices and sheets and compare those with other methods that were applied to this problem. To summarize, our method yields performance comparable to other methods on predicting runs of helices and sheets. It seems to have particularly high accuracy in predicting individual helices.

## Discussion

In this paper we have proposed causal networks as a general and efficient framework for data analysis in molecular biology. We have reported our initial ex-

Description	Chain-Pair	FSKBANN	ANN	Chou-Fasman
Average length of predicted helix run	9.4	8.52	7.79	8.00
Average length of actual helix run	10.3	–	–	–
Percentage of actual helix runs overlapped by predicted helix runs	66%	67%	70%	56%
Percentage of predicted helix runs that overlap actual helix runs	62%	66%	61%	64%
Average length of predicted sheet run	3.8	3.80	2.83	6.02
Average length of actual sheet run	5.0	–	–	–
Percentage of actual sheet runs overlapped by predicted sheet runs	56%	54%	35%	46%
Percentage of predicted sheet runs that overlap actual sheet runs	60%	63%	63%	56%

Table 3: Precision of run (segment) predictions. Comparative method results from [?].

Method	Total	Helix	Sheet	Coil
Chou-Fasman	57.3%	31.7%	36.9%	76.1%
ANN	61.8%	43.6%	18.6%	86.3%
w/ state	61.7%	39.2%	24.2%	86.0%
FSKBANN	63.4%	45.9%	35.1%	81.9%
w/o state	62.2%	42.4%	26.3%	84.6%
Viterbi	58.5%	48.3%	47.0%	69.3%
Chain-Pairs	62.2%	55.9%	51.7%	67.4%
Chain-Triples	60.3%	53.0%	45.5%	70.8%

Table 2: Overall prediction accuracies for various prediction methods. Comparative method results from [?].

periments applying this approach to the problem of protein secondary structure prediction. One of the main advantages of the probabilistic approach we described here is our ability to perform detailed experiments where we can experiment with different causal models. We can easily perform local substitutions (mutations) and measure (probabilistically) their effect on the global structure. Window-based methods do not support such experimentation as readily. Our method is efficient both during training and during prediction, which is important in order to be able to perform many experiments with different networks.

Our initial experiments have been done on the simplest possible models where we ignore many known dependencies. For example, it is known that in  $\alpha$ -

helices hydrogen bonds are formed between every  $i^{\text{th}}$  and  $(i + 4)^{\text{th}}$  residue in a chain. This can be incorporated in our model without losing efficiency. We also can improve our method by incorporating additional correlations among particular amino acids as in [?]. We achieve prediction accuracy similar to many other methods such as neural networks. We are confident that with sufficient fine tuning we can improve our results to equal the best methods. Typically, the current best prediction methods involve complex hybrid methods that compute a weighted vote among several methods using a combiner that learns the weights. *E.g.*, the hybrid method described by [?] combines neural networks, a statistical method and memory-based reasoning in a single system and achieves an overall accuracy of 66.4%.

Bayesian classification is a well-studied area and has been applied frequently to many domains such as pattern recognition, speech understanding and others. Statistical methods also have been used for protein structure prediction. What characterizes our approach is its simplicity and the explicit modeling of causal links. We believe that for scientific data analysis it is particularly important to develop tools that clearly display all the causal independence assumptions. Causal networks provide a very convenient medium for the scientist to experiment with different empirical models and obtain possibly important insights into a problem.