

# MATCHING OBSERVED ALPHA HELIX LENGTHS TO PREDICTED SECONDARY STRUCTURE\*

Brian D. Cloteaux\* and Nadezhda Serova\*\*

## Abstract

Because of the complexity in determining the 3D structure of a protein, the use of partial information determined from experimental techniques can greatly reduce the overall computational expense. We investigate the problem of matching experimentally observed lengths of helices to the predicted secondary structure of a protein. We give a simple and fast algorithm for producing a library of potential solutions. We test our algorithm by performing a series of computational experiments for predicting the alpha helix placement of proteins with an already known order. These tests seem to demonstrate that our method, if given a good prediction of the protein's secondary structure, can generate high quality lists of potential placements of the helix lengths onto the protein sequence.

## Key Words

Protein structure, alpha helix placement

## 1. Introduction

Understanding how specific proteins fold, or arrange themselves in three-dimensional space (3D) based on environmental and internal chemical constraints, is necessary to determine how these proteins function. But even when the amino acid sequence of the proteins (1D structure) is known, the prediction of their corresponding 3D structure is an extremely challenging problem.

This challenge is both from an experimental and computational viewpoint. Proteins require precise environments to fold properly. Because of the numerous complications in measuring protein under the correct environment, experimental methods for ascertaining the 3D arrangements are expensive, time consuming,

and of limited accuracy. X-ray crystallography, for example, is a powerful technique in the determination of these structures; however, it is ineffective in proteins that are not easily crystallized, such as membrane proteins. Many other methods provide only partial information about the protein's 3D structure.

At the same time, computationally determining the 3D structure of proteins is, in general, intractable. To reduce the difficulty of this problem, a recent approach has been to computationally match experimental observations to the 3D structure of the protein [2–4]. This paper extends an original investigation by He, Lu, and Pontelli [5] into the problem of matching observed lengths of the alpha helices from the electron cryomicroscopy technique to the predicted areas of secondary structure.

Electron cryomicroscopy can be used to produce a density map of some proteins. Although with current technology the resolutions of such maps are relatively low, certain secondary structures such as alpha helices can still be identified. Using electron cryomicroscopy, the lengths of the alpha helices can be observed, but the exact location of these helices on the protein sequence is not clear. To help overcome this limitation, He, Lu, and Pontelli suggested matching these observed lengths to the predicted probabilities of the protein's secondary structure. An example of the correspondence between protein sequences and the observed secondary structure is shown in Fig. 1. These probabilities on the placement of alpha helices onto the 1D sequence are generally based on the placement for similar sequences in other known proteins and have inherently limited accuracy. Thus, the result from matching observed length to the predicted secondary structure is to produce a set of probable arrangements of the observed lengths that can then be used as a starting point in determining 3D structure.

This article offers two contributions to the matching of observed lengths to their placement on the 1D protein structure. The first is an examination of the complexity and necessity of computing the optimal length placement. We give evidence that computing optimal solutions may not be worth the computational expense.

A second contribution is to introduce a new approach to computing possible arrangements. He, Lu, and Pontelli

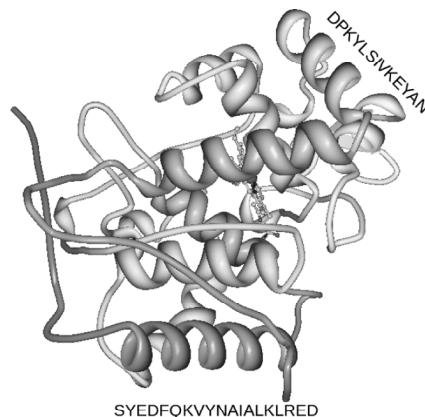
\* Applied and Computational Mathematics Division, National Institute of Standards and Technology, Gaithersburg, Maryland, USA; e-mail: brian.cloteaux@nist.gov

\*\* Department of Computer Science, University of Maryland, Baltimore County, Baltimore, Maryland, USA; e-mail: nserova1@umbc.edu

\* Official contribution of the National Institute of Standards and Technology; not subject to copyright in the United States. A conference version of this paper was published at the 2009 Computational Structural Bioinformatics Workshop [1].

1	30
TTPLVHVASVEKGRSYEDFQKVYNAIALKL	
REDDEYDNYIGYGPVLRVRLAWHTSGTWDKH	
DNTGGSYGGTYRFFKKEFNDP SNAGLQNGFK	
FLEPIHKEFPWISSGDLFSLGGVTAVQEMQ	
GPKIPWRCGRVDTPEDTTPDNGRLPDADKD	
ADYVRTFFQRLNMNDREVVALMGAHALGKT	
HLKNSGYEGPWGAANNVFTNEFYLNLLNED	
WKLEKNDANNEQWDSKSGYMLLP TDYSLIQ	
DPKYLSIVKEYANDQDKFFKDFSKAFEKLL	
ENGITFPKDAPSPFIFKTLEEQGL	
271	294

(a) Protein sequence of 2CYP with helices marked



(b) 3D structure of 2CYP with two of its alpha-helix sequence marked

Figure 1. Protein sequence and structure of the protein 2CYP. In this example, two of the alpha helices are highlighted on the structure to show their correspondence with the protein sequence.

introduced a method to produce a library of likely mappings to serve as starting points for a researcher. Our approach is similar to the He, Lu, and Pontelli method in the sense that it does not generate all of the possible mappings nor does it try to find an optimal solution. Instead, we give a simple heuristic algorithm that gives a good approximation of the placement of the lengths and then using this approximation as a starting point, we randomly modify it to look for other possible solutions. We collect the best arrangements to use as a library of possible mappings.

To test our approach, we compared the predicting length placement produced by our algorithm to actual ordering on several known proteins. These tests show that our method can be a fast and simple approach for producing high quality possible placements of observed alpha helix lengths onto a protein's sequence.

## 2. The Maximal Cover Sum Problem

Before we examine the problem of mapping the observed alpha helices to the predicted secondary structure, we first will consider a closely related problem that we call the *maximal cover sum problem*. This problem consists of a set  $\mathcal{C}$  of covers with an associated function  $\omega : \mathcal{C} \rightarrow \mathbb{N}$  that gives the length of each cover. There is also an  $n$ -length string  $P$  of positive real numbers, i.e.,  $P \in \mathbb{R}_+^n$ . The expression  $P_i$  is used to denote the  $i$ th value in the string  $P$ .

We define a placement of the covers on the string using an index function  $I : \mathcal{C} \rightarrow \{1 \dots |P|\}$ . The value  $I(c)$  gives the index in  $P$  of the first location to place the cover  $c$ . Any index function has the following three restrictions:

1.  $\forall c_1, c_2 \in \mathcal{C}, I(c_1) = I(c_2)$  if and only if  $c_1 = c_2$ .
2.  $\forall c_1, c_2 \in \mathcal{C}$ , if  $I(c_1) < I(c_2)$  then  $I(c_1) + \omega(c_1) < I(c_2)$ .
3.  $\forall c \in \mathcal{C}, I(c) + \omega(c) \leq |P|$ .

The first two restrictions say that covers are not allowed to overlap as they cover the string  $P$ . The last restriction prevents covers from extending beyond the length of the string  $P$ . These restrictions trivially imply that if for a problem instance the condition:

$$|P| \geq \sum_{c \in \mathcal{C}} \omega(c) \quad (1)$$

does not hold, then no index function can exist for that instance.

The maximal cover sum problem is then to find an index function that maximizes the expression:

$$\sum_{c \in \mathcal{C}} \sum_{j=0}^{\omega(c)-1} P_{j+I(c)} \quad (2)$$

In other words, find an arrangement of the covers in  $\mathcal{C}$  that covers the largest total of values on  $P$ .

We are interested in this problem because we can view matching the observed lengths of helices to the predicted secondary structure of a protein as a maximal cover sum problem. When examining possible arrangements of the helices onto the protein string, we would expect that the arrangements that cover the maximal value for the predicted probabilities of helices would be the most likely to occur in the actual protein. Thus, we are typically interested in examining these optimal arrangements first when trying to determine the 3D structure of the protein.

In considering how to obtain an optimal arrangement, we first notice that for any given ordering of the covers, we can compute an optimal covering using that ordering in polynomial time. To show this, we define  $\pi$  as an ordering of the cover set  $\mathcal{C} = \{c_1, c_2, \dots, c_n\}$ , i.e., the value of  $\pi(i)$  is the position of the element  $c_i \in \mathcal{C}$  in the order  $\pi$ . The inverse function  $\pi^{-1}$  then takes a position  $i$  in the ordering and returns the cover in that position. Using a given order  $\pi$ , we can define the following recurrence equation that determines the size of an optimal covering using  $\pi$ :

$$m(a, b) = \begin{cases} 0 & \text{if } a \text{ or } b = 0, \\ \max \left( m(a, b-1), m(a-1, b - \omega(\pi^{-1}(a))) \right. \\ \quad \left. + \sum_{k=b-\omega(\pi^{-1}(a))+1}^b P_k \right) \end{cases} \quad (3)$$

Under this definition, the term  $m(a, b)$  is the value of the maximal covering of the first  $a$  covers in the order  $\pi$  onto

the first  $b$  positions of the string  $P$ . This recursion is based on determining whether or not there exists an optimal covering of the  $b$  position of  $P$  with  $\pi^{-1}(a)$ . By using dynamic programming, the recurrence and its associated index function can be computed in  $O(|\mathcal{C}| \cdot |P|)$  time.

Thus, the complexity in the maximal cover set problem stems from finding an order that produces an optimal covering. In general, finding an optimal ordering is super-polynomial in the number of covers unless  $P = NP$ . This is a consequence of the fact that we can reduce the set partition problem, which is  $NP$ -complete [6], to the maximal cover sum problem. To see this, consider an instance of a set partition problem with a multiset of values  $S$ . Using the multiset  $S$ , we can create an instance of the maximal cover problem by creating a set of covers  $C$  where  $|C| = |S|$  and where the length of the covers are the values in  $S$ . We then create a string  $P$  where:

$$P = \left( \underbrace{1, 1, 1, \dots, 1, 1}_\ell, 0, \underbrace{1, 1, 1, \dots, 1, 1}_\ell \right) \quad (4)$$

and  $\ell = \frac{\sum_{c \in C} \omega(c)}{2}$ . The point of this construction is that the maximal cover sum is equal to  $\sum_{c \in C} \omega(c) = 2\ell$  if and only if the multiset  $S$  can be equally partitioned.

Although the set partition problem is  $NP$ -complete, our reduction does not necessarily prove that the maximal cover sum problem is  $NP$ -hard. This is because the number of bits needed in our created string  $P$  can potentially be exponential to the number of bits in  $S$  and so the size of the instance for the maximal cover sum problem can be exponentially larger than that of the set partition problem. But this reduction does tell us two facts about the maximal cover sum problem. The first is that if there is a reduction to an  $NP$ -complete problem where the maximum value in the cover set of these instances is bounded by a polynomial based on the size of the cover sets, then the problem itself would have to be  $NP$ -hard. But more importantly, even if this problem is in polynomial time, we would still not expect for there to be any algorithm to determine an optimal order that is polynomial in the size of the cover set. In other words, for some instances we probably cannot do much better than brute force checking of the permutations to find an optimal one.

### 3. Examining Optimal Solutions

As it appears that finding an optimal solution, or especially the top  $k$  optimal solutions, to the maximal cover sum problem is computationally difficult, it is necessary to ask whether we need to find these solutions to perform the matching. To test this assumption, we computed the optimal solutions for a number of proteins and then compared those results with the actual ordering.

The list of proteins we used were originally selected in the He, Lu, and Pontelli paper [5]. The set of covers for each protein were generated by taking all the alpha helices in the protein with length greater than 7. This was done to match the cover sets produced in the He, Lu, and Pontelli experiments. The predictions of the secondary

structure came from using the PHD algorithm [7] which has been implemented in the PredictProtein server [8]. This server returns a value for each amino acid within a given protein sequence corresponding to the likelihood of its participation in an alpha helix. These values range from 0 to 9, where a prediction value of 0 means that it is highly unlikely for that amino acid to be a part of an alpha helix and a value of 9 corresponds to a very high likelihood.

While there are a number of distance measures for list orders (see, e.g., Chapter 6 of Diaconis [9] and Fagin, Kumar, and Sivakumar [10]), we focused on two metrics. The first is the Hamming distance between the orders. This is a measure of the number of items that are in the same position between two lists. If  $\pi$  and  $\sigma$  are orderings of the cover set  $\mathcal{C} = \{c_1, c_2, \dots, c_n\}$ , then the Hamming distance is defined as:

$$d_{Ham}(\pi, \sigma) = n - |\{c_i \in \mathcal{C} : \pi(c_i) = \sigma(c_i)\}| \quad (5)$$

To compare the distance between orders of different lengths, we used a normalized version of the Hamming distance that is obtained by dividing the value  $d_{Ham}$  by  $n$ . This normalization maps all distances to interval  $[0, 1]$ , where a distance of 0 represents identical orders and a distance of 1 are orders that are maximally disordered to each other. For the Hamming distance, maximally disordered lists share no common item for any position in the list.

A second metric that we will consider is the Kendall-tau distance. This distance is defined as:

$$d_{K\tau}(\pi, \sigma) = |\{(c_i, c_j) \in \mathcal{C}^2 : \pi(c_i) < \pi(c_j) \wedge \sigma(c_i) > \sigma(c_j)\}| \quad (6)$$

The Kendall-tau distance is sometimes called the BubbleSort distance because it is equivalent to the number of flips needed in BubbleSort to transform one order to the second. For the Kendall-tau distance, two maximally disordered lists are in reverse order of each other. Again, we normalized the Kendall-tau distance to the interval  $[0, 1]$  by dividing the value  $d_{K\tau}$  by  $\binom{n}{2}$ .

Our experiment involved computing all the optimal orders for each of the test proteins. In most instances, there are multiple optimal solutions. We then computed the minimum and mean Kendall-tau and Hamming distances from the set of optimal solutions to the actual arrangement for the protein. These results are shown in Table 1.

A point to notice is that for all the given proteins, none of the optimal solutions were the actual ordering of the proteins. In fact, we can see examples, like 1L58, where a number of optimal solutions exist and also where all of these solutions are relatively distant from the actual order. These results call into question whether it is worth the computational expense to compute the optimal coverings of the predicted secondary structure.

### 4. A Greedy Heuristic for the Maximal Cover Sum

Because of the uncertainty inherent in predicting secondary structure strictly from 1D structure, we should not be

Table 1

Distances from the Optimal Covering Orders to the Actual Arrangement of the Alpha Helices for the given Proteins

Protein id	Number of Optimal Solutions	Minimum Kendall-Tau Distance	Average Kendall-Tau Distance	Minimum Hamming Distance	Average Hamming Distance
1CC5	1	0.667	0.667	0.750	0.750
6TMN_E	2	0.238	0.310	0.286	0.286
3TIM_A	1	0.100	0.100	0.400	0.400
2TSC_A	5	0.095	0.305	0.429	0.600
1ECA	6	0.190	0.294	0.429	0.714
1GD1_O	2	0.067	0.100	0.333	0.500
1L58	8	0.393	0.429	0.625	0.688
2PHH	1	0.476	0.476	1.000	1.000
2CYP	3	0.467	0.489	0.800	0.833

All distances have been normalized so that a distance of zero is an identical ordering and a distance of one is maximally dissimilar.

surprised that the optimal ordering is often a large distance from the actual order. This suggests that finding optimal orders may not be worth the computational expense, and that a simpler approximation method can be used. Towards this goal, we introduce a simple and fast heuristic (Algorithm 1) for producing a covering. The produced cover will not necessarily be optimal, but it will cover a large sum on the string.

**Algorithm**  $cov(\mathcal{C}, P)$

**Input:** a set of covers  $\mathcal{C}$  with an associated length function  $\omega$  and a string  $P$  of positive reals

**Output:** an array  $I$  that maps set of covers  $\mathcal{C}$  to indices in  $P$

Create empty index array  $I$

**If**  $\mathcal{C} \neq \emptyset$  **then**

$n \leftarrow |\mathcal{C}|$   
 $c \leftarrow$  a cover in  $\mathcal{C}$  with maximal length  
 $i \leftarrow$  smallest index of a maximal cover of  $c$  on  $P$   
 $I[n] \leftarrow i$

$I' \leftarrow cov(\mathcal{C} - \{c\}, P_{1,i} \oplus P_{i+C_n,|P|})$

**for**  $j \leftarrow 1$  **to**  $n - 1$  **do**

**If**  $I'[j] < i$  **then**

$I[j] \leftarrow I'[j]$

**if**  $I'[j] + \omega(c) > i$  **then**

$c_j \leftarrow$  cover in position  $j$

$I[n] \leftarrow I[n] + \omega(c_j)$

**end**

**else**

$I[j] \leftarrow I'[j] + \omega(p)$

**end**

**end**

**end**  
**return**  $I$

**Algorithm 1:** A heuristic algorithm  $cov$  for producing a covering of the string with a high sum. In this algorithm, the symbol  $\oplus$  denotes string concatenation.

The idea behind our algorithm is straightforward. For a set of covers  $\mathcal{C}$ , we select a cover  $c \in \mathcal{C}$  having maximal length and place it on the string  $P$  where it covers the greatest sum of values. We then create a new string  $P'$

which concatenates together the sections of  $P$  that are not covered by  $c$ . For the remaining covers  $\mathcal{C} - \{c\}$ , we obtain an index function  $I'$  for placing them on the string  $P'$  by calling our routine recursively. Using this new index function  $I'$ , we construct the index function  $I$  for the cover set  $\mathcal{C}$  and  $P$ . If  $i$  is the starting index of  $c$  on  $P$ , then this construction breaks into three cases. The first is for every cover  $c' \in \mathcal{C} - \{c\}$  that is completely placed on  $P'$  before  $i$ , in other words  $I'(c') + \omega(c') \leq i$ . For this case, we simply copy the index over, i.e.,  $I(c') = I'(c')$ . The second case is when the cover is completely after  $i$ , ( $I'(c') > i$ ), then we can insert the cover in  $P$  by offsetting its index by  $\omega(c)$  to make room for the cover  $c$  ( $I(c') = I'(c') + \omega(c)$ ). The final case is when a cover overlaps the index  $i$  ( $I'(c') < i$  and  $I'(c') + \omega(c') > i$ ). For this case, we notice that we can cover the same sum on  $P$  by keeping  $c'$  at the current index and sliding  $c$  over by the length of  $c'$  (i.e.,  $I(c') = I'(c')$  and  $I(c) = I'(c') + \omega(c')$ ). An example of this algorithm is shown in Fig. 2.

The advantage of this heuristic is that it gives a fast and reasonable covering of the string. As Fig. 2 demonstrates, this algorithm does not necessarily produce an optimal placement of the covers, but if there exists optimal solutions where the placement of covers are all separated on the string, this heuristic can often return one of those solutions. As the recursion depth of the algorithm is the number of covers, and for each cover we need to check the string  $P$  for its optimal placement, we can implement this algorithm to run in time  $O(|\mathcal{C}| \cdot |P|)$ .

In Table 2, we compare the distance of the cover order produced by our heuristic to the actual ordering of the helices on the proteins. Surprisingly, it seems that this heuristic often produces a result that is closer to the protein helix arrangement than using the set of optimal solutions.

## 5. Randomized Orderings based on BubbleSearch

As the complexity in the maximal cover sum problem is in determining the correct order, our approach is, instead

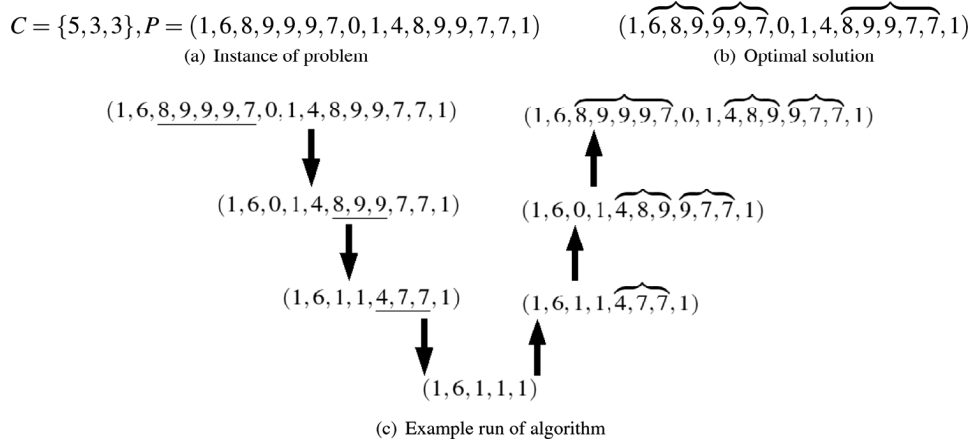


Figure 2. An example of how the heuristic algorithm constructs an ordering. Starting with the set  $C$  of cover lengths and a string  $P$  in 2(a), the algorithm (shown in 2(c)) greedily places the largest cover onto  $P$  and then removes that covered section from the string. When it has selected all cover lengths, it then inserts the covers while sliding over any covers that overlap. Figure 2(b) gives the optimal solution for this covering showing that the algorithm is a heuristic.

Table 2  
Distances to the Actual Ordering from the Ordering  
Produced by the Greedy Heuristic Algorithm *cov*

Protein id	Kendall-Tau Distance	Hamming Distance
1CC5	0.167	0.500
6TMN_E	0.571	1.000
3TIM_A	0.100	0.400
2TSC_A	0.095	0.571
1ECA	0.381	0.714
1GD1_O	0.133	0.667
1L58	0.250	0.750
2PHH	0.143	0.714
2CYP	0.578	0.900

of trying to compute every possible placement, to use a randomized process to find the high value orderings and then use only the optimal placement of those orderings. Our method is based on the randomized BubbleSearch algorithm of Lesh and Mitzenmacher [11] which is a type of priority algorithm [12]. Their approach is a heuristic used to search for optimal orderings especially in  $NP$ -hard applications. The basic algorithm starts with a simple approximation of the optimal order, called the *base order*, which is usually derived by using some type of greedy method. Using the base order, a series of random orders are then created. Starting from the first cover in the base order and then progressing through the order sequentially, each cover is added to the new order with some probability  $P$ . If the end of the order is reached without placing all the covers, then the selection process restarts from the first nonplaced cover. Denoting the base order as  $\pi$  and the new order as  $\sigma$ , the probability of producing  $\sigma$  is proportional to

$(1 - P)^{d_{K\tau}(\pi, \sigma)}$ . In other words, the closer the two orders are in Kendall-tau distance, the more likely it is to produce  $\sigma$  from  $\pi$ . The choice of  $P$  controls how large the radius of probable orders will be around the base order. Thus as  $P$  approaches 1, the Kendall-tau distance between the orders approaches 0.

For the observed helix matching problem, we conducted a series of computational experiments where we generated a base order using our heuristic and then generated 200 random orders each for a series of values for  $P$ . We then computed the Kendall-tau and Hamming distances between the randomly generated orders and the actual order of the protein. These results are shown in Table 3. We notice from these test results that in every instance the BubbleSearch method was able to find an ordering that was either the correct ordering or very close.

As the value of  $P$  tends to be domain specific, it is not specified for the randomized BubbleSearch algorithm. Instead, the user must tune this parameter for their individual problems. In examining the results of our experiments, it seems that  $P=0.6$  provides a good initial value for examining protein orders. This value seems to be a compromise between keeping most orders reasonably close to the base ordering, but still allowing the structure to be modified into a fundamentally different orderings to avoid being trapped in a local minima. Using this idea, we propose the following workflow to produce  $k$  potential matchings when investigating how to relate the observed helix lengths to the 3D structure:

1. Generate a base order using our greedy heuristic.
2. Run the randomized BubbleSearch algorithm with  $P=0.6$  for  $c \cdot k$  times for some  $c \geq 1$ .
3. For each order generated, find its optimal placement and compute the sum covered, saving the  $k$  top valued placements.

In the Lesh and Mitzenmacher paper, a second form of randomized BubbleSearch is suggested where the base

Table 3  
 Normalized Distances from the Randomized BubbleSearch Orders to the Actual Arrangement of the Alpha Helices on the given Proteins

Protein id	Val	<i>P</i> -values								
		0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9
1CC5	Min. $K\tau$	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.167
	Mean $K\tau$	0.379	0.332	0.338	0.334	0.307	0.300	0.261	0.243	0.209
	Std. $K\tau$	0.191	0.186	0.183	0.165	0.153	0.155	0.128	0.101	0.075
	Min. Ham.	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.500
	Mean Ham.	0.640	0.645	0.636	0.675	0.660	0.659	0.650	0.649	0.601
	Std. Ham.	0.248	0.264	0.260	0.254	0.239	0.243	0.236	0.200	0.186
6TMN_E	Min. $K\tau$	0.000	0.143	0.095	0.143	0.143	0.143	0.143	0.190	0.286
	Mean $K\tau$	0.474	0.468	0.466	0.501	0.497	0.511	0.518	0.535	0.553
	Std. $K\tau$	0.152	0.139	0.138	0.147	0.145	0.127	0.114	0.100	0.074
	Min. Ham.	0.000	0.286	0.429	0.286	0.571	0.429	0.429	0.429	0.714
	Mean Ham.	0.824	0.826	0.843	0.856	0.857	0.852	0.858	0.897	0.949
	Std. Ham.	0.156	0.153	0.139	0.140	0.120	0.140	0.146	0.116	0.092
3TIM_A	Min. $K\tau$	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.100
	Mean $K\tau$	0.439	0.412	0.380	0.339	0.328	0.295	0.284	0.220	0.164
	Std. $K\tau$	0.206	0.180	0.186	0.172	0.163	0.174	0.161	0.144	0.100
	Min. Ham.	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.400
	Mean Ham.	0.705	0.692	0.697	0.664	0.671	0.634	0.657	0.591	0.526
	Std. Ham.	0.230	0.212	0.224	0.215	0.188	0.212	0.206	0.199	0.176
2TSC_A	Min. $K\tau$	0.048	0.048	0.000	0.000	0.000	0.000	0.048	0.048	0.048
	Mean $K\tau$	0.348	0.325	0.307	0.260	0.252	0.228	0.215	0.174	0.142
	Std. $K\tau$	0.131	0.131	0.154	0.127	0.131	0.138	0.132	0.116	0.106
	Min. Ham.	0.286	0.286	0.000	0.000	0.000	0.000	0.286	0.286	0.286
	Mean Ham.	0.678	0.647	0.629	0.596	0.576	0.551	0.528	0.496	0.519
	Std. Ham.	0.178	0.157	0.186	0.185	0.185	0.196	0.160	0.150	0.124
1ECA	Min. $K\tau$	0.095	0.048	0.095	0.143	0.143	0.143	0.143	0.143	0.190
	Mean $K\tau$	0.436	0.442	0.435	0.424	0.416	0.419	0.418	0.400	0.391
	Std. $K\tau$	0.151	0.148	0.144	0.151	0.124	0.110	0.112	0.092	0.072
	Min. Ham.	0.286	0.286	0.429	0.429	0.286	0.286	0.286	0.286	0.571
	Mean Ham.	0.824	0.794	0.824	0.834	0.841	0.854	0.864	0.816	0.801
	Std. Ham.	0.160	0.172	0.148	0.146	0.153	0.138	0.141	0.145	0.132
1GD1_O	Min. $K\tau$	0.067	0.000	0.000	0.067	0.067	0.067	0.067	0.067	0.067
	Mean $K\tau$	0.438	0.423	0.361	0.365	0.338	0.294	0.250	0.238	0.186
	Std. $K\tau$	0.170	0.170	0.167	0.165	0.153	0.148	0.129	0.128	0.097
	Min. Ham.	0.333	0.000	0.000	0.333	0.333	0.333	0.333	0.333	0.333
	Mean Ham.	0.772	0.769	0.713	0.712	0.717	0.688	0.673	0.663	0.663
	Std. Ham.	0.191	0.174	0.189	0.175	0.167	0.164	0.140	0.136	0.093

(continued.)

Table 3  
(Continued)

Protein id	Val	<i>P</i> -values								
		0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9
1L58	Min. $K\tau$	0.071	0.036	0.071	0.071	0.036	0.036	0.071	0.107	0.143
	Mean $K\tau$	0.322	0.298	0.301	0.306	0.308	0.296	0.280	0.272	0.262
	Std. $K\tau$	0.113	0.117	0.104	0.092	0.102	0.089	0.083	0.063	0.044
	Min. Ham.	0.250	0.250	0.250	0.375	0.250	0.250	0.250	0.375	0.375
	Mean Ham.	0.676	0.674	0.672	0.686	0.666	0.664	0.649	0.651	0.676
	Std. Ham.	0.144	0.134	0.130	0.121	0.124	0.123	0.125	0.105	0.095
2PHH	Min. $K\tau$	0.095	0.048	0.048	0.000	0.048	0.000	0.048	0.048	0.048
	Mean $K\tau$	0.437	0.407	0.394	0.381	0.350	0.309	0.277	0.254	0.189
	Std. $K\tau$	0.159	0.153	0.156	0.153	0.146	0.158	0.134	0.134	0.095
	Min. Ham.	0.286	0.286	0.286	0.000	0.286	0.000	0.286	0.286	0.286
	Mean Ham.	0.759	0.739	0.745	0.725	0.715	0.692	0.633	0.628	0.625
	Std. Ham.	0.166	0.175	0.178	0.188	0.174	0.192	0.194	0.184	0.163
2CYP	Min. $K\tau$	0.089	0.044	0.111	0.200	0.178	0.200	0.244	0.244	0.289
	Mean $K\tau$	0.410	0.410	0.433	0.448	0.463	0.463	0.493	0.512	0.549
	Std. $K\tau$	0.108	0.111	0.103	0.102	0.099	0.102	0.095	0.083	0.061
	Min. Ham.	0.400	0.400	0.500	0.400	0.600	0.500	0.500	0.600	0.500
	Mean Ham.	0.826	0.830	0.837	0.838	0.841	0.823	0.826	0.843	0.879
	Std. Ham.	0.120	0.116	0.111	0.113	0.107	0.113	0.107	0.100	0.080

For each protein and *P*-value, 200 random orders were produced, using the greedy heuristic as the base order. This table shows the minimum, mean, and standard deviation of the Kendall-tau and Hamming distances for these random orders.

ordering is modified whenever an ordering is found that produces a better value. Intuitively, it would seem that if the optimal solutions are not good indicators of the actual ordering, then progressing toward an optimal solution during the computation would not necessarily produce better solutions. We examined this notion by computing the differences in solutions produced by the version of BubbleSearch which only uses the greedy heuristic as the base order versus the version which updates base order during the computation. The results of this experiment is shown in Table 4. One point we notice is that, other than the values for the protein 6TMN\_E, using the greedy heuristic as the base order tended to give better solutions on average. It is also worth noting that in several instances, while the minimal Kendall-tau difference to the actual order was almost negligible, the minimal Hamming distance was much lower when the base value was allowed to be modified. It is an open question on whether this observed behaviour is strictly a product of the protein that were selected or not.

One final point to mention about this algorithm is that it can typically produce a number of potential solutions very quickly. But if additional speed is needed, this entire process can be trivially parallelized.

Table 4  
The Difference in Distances to the Actual Proteins Orders between the BubbleSearch Algorithm with a Fixed Base Order (in this case, Produced by the Greedy Algorithm) and the BubbleSearch Algorithm while Modifying the Base Order

	Minimum Kendall-Tau Difference	Average Kendall-Tau Difference	Minimum Hamming Difference	Average Hamming Difference
1CC5	0.000	-0.073	0.000	-0.176
6TMN_E	0.095	0.018	0.286	0.129
3TIM_A	0.000	-0.037	0.000	-0.105
2TSC_A	0.000	-0.060	0.000	-0.089
1ECA	0.000	-0.007	0.000	0.060
1GD1_O	0.067	-0.003	0.333	0.013
1L58	0.071	-0.046	0.250	-0.042
2PHH	0.000	-0.025	0.000	-0.047
2CYP	0.067	-0.012	0.400	0.050

Each value was the result of 200 random samples using a *P*-value of 0.6 for all calculations.

## 6. Conclusions

In this article, we have examined the problem of matching observed lengths of alpha helices to their predicted location on a protein's amino acid sequence. This potentially is a first step towards determining the 3D structure of the protein. Our first conclusion is that finding an optimal solution does not seem to be worth the computational effort. In particular, we showed evidence that, because of the uncertainty of the helix prediction, the optimal coverings can be relatively distance from the actual ordering on the protein.

Instead, we introduced a simple greedy heuristic for estimating the order. Using this heuristic as a starting point, we chose random orders around it using the BubbleSearch method of Lesh and Mitzenmacher. When compared to the actual orderings, this method was able to either find the correct ordering or an ordering that was very close. Thus, we believe that our method is a fast and efficient algorithm for determining a set of potential placements of helix lengths onto a protein sequence.

A future direction in this line of research would be to develop methods to incorporate additional constraints in the order building process. Specifically, a problem that has been brought to our attention is the building of a library of potential solutions in which all solutions obey some given partial ordering of the covers. But the addition of any type of constraint provides interesting new research directions with these algorithms.

## Acknowledgements

This research was performed while the second author was at NIST supported by the Summer Undergraduate Research Fellowship (SURF) program. We would like to thank the anonymous referees for their helpful comments on our paper.

## References

- [1] B. Cloteaux & N. Serova, Matching observed alpha helix lengths to predicted secondary structure, *BIBMW 2009. IEEE International Conference on Bioinformatics and Biomedicine Workshop*, Washington, DC, USA, 2009, 113–119.
- [2] Y. Wu, M. Chen, M. Lu, Q. Wang, & J. Ma, Determining protein topology from skeletons of secondary structures, *Journal of Molecular Biology*, 350(3), 2005, 571–586.
- [3] W. Sun, S. Al-Haj, & J. He, Parallel computing in protein structure topology determination, *Proceedings of the 26th Army Science Conference*, Orlando, Florida, USA, 2008, cp8.
- [4] Y. Lu, J. He, & C.E.M. Strauss, Deriving topology and sequence alignment for the helix skeleton in lowresolution protein density maps, *Journal of Bioinformatics and Computational Biology*, 6(1), 2008, 183–201.
- [5] J. He, Y. Lu, & E. Pontelli, A parallel algorithm for helix mapping between 3d and 1d protein structure using the length constraints, *ISPA '04: Proceedings of Second International Symposium on Parallel and Distributed Processing and Applications*, Hong Kong, China, 2004, 746–756.
- [6] M.R. Garey & D.S. Johnson, *Computers and intractability: A guide to the theory of NP-completeness* (New York, NY, USA: W. H. Freeman & Co., 1990).
- [7] B. Rost & C. Sander, Prediction of protein secondary structure at better than 70% accuracy, *Journal of Molecular Biology*, 232(2), 1993, 584–599.

- [8] B. Rost, G. Yachdav, & J. Liu, The PredictProtein server, *Nucleic Acids Research*, 32, 2004, W321–W326.
- [9] P. Diaconis, *Group representations in probability and statistics*, Lecture notes – Monograph series, vol. 11 (Hayward, California, USA: Institute of Mathematical Statistics, 1988).
- [10] R. Fagin, R. Kumar, & D. Sivakumar, Comparing top k lists, *SIAM Journal on Discrete Mathematics*, 17(1), 2003, 134–160.
- [11] N. Lesh & M. Mitzenmacher, Bubblesearch: A simple heuristic for improving priority-based greedy algorithms, *Information Processing Letters*, 97(4), 2006, 161–169.
- [12] A. Borodin, M.N. Nielsen, & C. Rackoff, (Incremental) priority algorithms, *Proceedings of the Thirteenth Annual ACM-SIAM Symposium on Discrete Algorithms*, San Francisco, California: Society for Industrial and Applied Mathematics, 2002, 752–761.

## Biographies



*Brian D. Cloteaux* is a computer scientist in the Discrete Mathematical Analysis Group of the Information Technology Laboratory at the National Institute of Standards and Technology in Gaithersburg, Maryland. He holds a Ph.D. degree in Computer Science from New Mexico State University. Previously, he has worked at ExxonMobil in both computing and research

capacities.



*Nadezhda Serova* is currently in her third year of undergraduate studies at University of Maryland, Baltimore County with a major in Computer Science and a minor in Physics. After graduating, she is planning to pursue a M.S. degree in Computer Science. Her current areas of interest include mobile and wearable computing, computational photography, and robotics. She has taken part in

the Summer Undergraduate Research Fellowship at NIST for two summers and will be participating again in summer of 2010.