

# 1 PSIST: A Scalable Approach to Indexing 2 Protein Structures using Suffix Trees

3 Feng Gao, Mohammed J. Zaki\*

4 *Department of Computer Science, Rensselaer Polytechnic Institute, 110 8th Street,*  
5 *Troy, NY, 12180*

---

## 6 Abstract

7 Approaches for indexing proteins, and for fast and scalable searching for struc-  
8 tures similar to a query structure have important applications such as protein struc-  
9 ture and function prediction, protein classification and drug discovery. In this paper,  
10 we develop a new method for extracting local structural (or geometric) features from  
11 protein structures. These feature vectors are in turn converted into a set of symbols,  
12 which are then indexed using a suffix tree. For a given query, the suffix tree index  
13 can be used effectively to retrieve the maximal matches, which are then chained to  
14 obtain the local alignments. Finally, similar proteins are retrieved by their align-  
15 ment score against the query. Our results show classification accuracy up to 50%  
16 and 92.9% at the topology and class level according to the CATH classification.  
17 These results outperform the best previous methods. We also show that PSIST is  
18 highly scalable due to the external suffix tree indexing approach it uses; it is able  
19 to index about 70,500 domains from SCOP in under an hour.

20 *Key words:*

21 Protein Structure Indexing, External Suffix Trees, Bioinformatics

---

## 22 1 Introduction

23 Proteins are composed of chains of basic building blocks called amino acids.  
24 Traditionally the problem of determining similar proteins was approached by  
25 finding the amount of similarity in their amino acid sequences. However bi-  
26 ologists have determined that even proteins which are remotely homologous  
27 in their sequence similarities can perform surprisingly very similar functions  
28 in living organisms [37]. This fact has been attributed to the dependency of  
29 the functional role of proteins on their actual three-dimensional (3D) struc-  
30 ture. That is, two proteins with remote sequence homology can be functionally  
31 classified as being similar if they exhibit structural homology.

32 Searching the growing database of protein structures for structural homologues  
33 is a difficult and time-consuming task. For example, we may want to retrieve  
34 all structures that contain sub-structures similar to the query, a specific 3D  
35 arrangement of surface residues, etc. Searches such as these are the first step  
36 towards building a systems level model for protein interactions. In fact, high  
37 throughput proteomics methods are already accumulating the protein inter-  
38 action data that we would wish to model, but fast computational methods for  
39 structural database searching lag far behind; biologists are in need of a means  
40 to search the protein structure databases rapidly, similar to the way BLAST  
41 [1] rapidly searches the sequence databases.

42 In this paper, we present a fast, novel protein structure indexing method called  
43 PSIST (which stands for **P**rotein **S**tructure **I**ndexing using **S**uffix **T**rees; note

---

\* Corresponding author.

*Email addresses:* `gaof@cs.rpi.edu` (Feng Gao), `zaki@cs.rpi.edu` (Mohammed J. Zaki).

44 that a preliminary version of this paper has appeared in [15]). As the name  
45 implies, our new approach transforms the local structural information of a  
46 protein into a “sequence” on which a suffix tree is built for fast matches. We  
47 first extract local structural feature vectors using a sliding window along the  
48 protein backbone. For a pair of residues, the distance between their  $C_\alpha$  atoms  
49 and the angle between the planes formed by the  $C_\alpha$ ,  $N$  and  $C$  atoms of each  
50 residue are calculated. The feature vectors for a given window include all the  
51 distances and angles between the first residue and the rest of the residues  
52 within the window. Compared with the local features from a single residue,  
53 our feature vectors contain both translational and rotational information. Af-  
54 ter normalizing the feature vectors, the protein structure is converted to a  
55 sequence (called the *structure-feature sequence* or *SF-sequence*) of discretized  
56 symbols.

57 We use external suffix trees to index the protein SF-sequences. For a given  
58 query, all the maximal matches are retrieved from the suffix tree and chained  
59 into alignments using dynamic programming. The top proteins with the high-  
60 est alignment scores are finally selected. Our results shows classification ac-  
61 curacy up to 50.0% and 92.9% at the topology and class level according to  
62 the CATH [33,34] classification. These results outperform the best previous  
63 method. We also show that PSIST is highly scalable due to the external suffix  
64 tree indexing approach it uses; it is able to index about 70,500 domains from  
65 SCOP [30] in under an hour!

## 66 2 Related Work

### 67 2.1 Structural Similarity Search

68 Protein structural similarity determination can be classified into three ap-  
69 proaches: pairwise alignment, multiple structure alignment, and database in-  
70 dexing.

71 Pair-wise structure alignment methods can be classified into three classes [13].  
72 The first class works at the residue level [20,40]. The second class focuses on  
73 using secondary structure elements (SSEs) such as  $\alpha$ -helices and  $\beta$ -strands  
74 to align two proteins approximately [26,29,32]. The third approach is to use  
75 geometric hashing, which can be applied at both the residue [25] and SSE  
76 level [21].

77 Previous work has also looked at multiple structure alignment. These meth-  
78 ods are based on geometric hashing [31], or SSE information [12]. A recent  
79 method [39] aims to solve the multiple structural alignment problem with  
80 detection of partial solutions; it computes the best scoring structural align-  
81 ments, which can be either sequential or sequence-order independent [44], if  
82 one seeks geometric patterns which do not follow the sequence order. Due  
83 to their time complexity, the pairwise and multiple structure alignment ap-  
84 proaches are not suitable for searching for similarity over thousands of protein  
85 structures. Database indexing and scalable searching approaches satisfy this  
86 requirement.

87 There are two classes of protein structure indexing approaches according to  
88 the representation of the local features. The first class focuses on indexing the

89 local features at the residue level directly, and the other class uses SSEs to  
90 approximate the local features of the proteins.

91 CTSS [6] approximates the protein  $C_\alpha$  backbone with a smooth spline with  
92 minimum curvature. The method then stores the curvature, torsion angle and  
93 the secondary structure that each  $C_\alpha$  atom in the backbone belongs to, in a  
94 hash-based index. ProGreSS [3] is a recent method, which extracts the fea-  
95 tures for both the structure and sequence, within a sliding window over the  
96 backbone. Its structure features are the same as the CTSS features (curvature,  
97 torsion angles, and SSE information); its sequence features are derived using  
98 scoring matrices like PAM [9] or BLOSUM [19].

99 The LFF profile algorithm [7] first extracts some representative local features  
100 from the distance matrix of all the proteins, and then each protein's distance  
101 matrix is encoded by the indices of the nearest representative features. Each  
102 structure is represented by a vector of the frequency of the representative  
103 local features. The structural similarity between two proteins is the Euclidean  
104 distance between their LFF profile vectors. This method is more suitable for  
105 global rather than local similarity between the query and database proteins.

106 There are also some methods that index the protein structures using SSEs.  
107 For each protein, PSI [5] uses a  $R^*$ -tree to index a nine-dimensional feature  
108 vector, a representation of all the triplet SSEs within a range. After retriev-  
109 ing the matching triplet pairs, a graph-based algorithm is used to compute  
110 the alignment of the matching SSE pairs. Another SSE-based method, Prot-  
111 Dex [2], obtains the sub-matrices of the SSE contact patterns from the dis-  
112 tance matrix of a protein structure. The grand sum of the sub-matrices and the  
113 contact-pattern type are indexed by an inverted file index. By their nature,

114 SSEs model the protein only approximately, and therefore these SSE-based  
115 approaches suffer in retrieval accuracy and as such are not very useful for  
116 small query proteins with few SSEs.

## 117 2.2 Suffix Trees

118 A suffix tree is a versatile data structure for substring problems [18], and it  
119 has been used for various problems such as protein sequence indexing [23,28],  
120 genome alignment [10,11] and structural motif detection [42]. Suffix trees can  
121 be constructed in  $O(n)$  time and space [27,43], and thus are an effective choice  
122 for indexing sequences.

123 Most suffix tree algorithms are not designed to scale as efficiently when the  
124 input sequence is extremely large. Also, the use of suffix links, which are a key  
125 feature in obtaining the linear construction time, can result in poor locality  
126 of reference [14,22,41]. To address these issues, several disk-based suffix tree  
127 algorithms have been proposed in the last few years. Some of the approaches  
128 [22,24,38,41] completely abandon the use of suffix links and sacrifice the theo-  
129 retically superior linear construction time in exchange for a better locality of  
130 reference. A linear algorithm to construct distributed suffix trees (DST) was  
131 proposed in [8]. DST introduces a new notion of sparse suffix links and uses  
132 different rules to follow a sparse suffix link to the tree root. The suffix tree  
133 can be distributed over a number of computing nodes. It can handle larger  
134 data than existing suffix trees, but it does assume that the input data cannot  
135 exceed the size of real memory.

## 136 3 The PSIST Approach to Protein Structure Indexing

### 137 3.1 Indexing Proteins

#### 138 3.1.1 Local Feature Extraction

139 A protein is composed of an ordered sequence of residues linked by peptide  
140 bonds. Each residue has  $C_\alpha$ ,  $N$  and  $C$  atoms, which constitute the backbone of  
141 the protein. Although the backbone is linear topologically, it is very complex  
142 geometrically. The bond lengths, bond angles and torsion angles completely  
143 define the conformation and geometry of the protein.

144 The bond length is the distance between the bonded atoms, and the bond  
145 angle is the angle between any two covalent bonds that include a common  
146 atom (see Figure 1). For instance, the bond length of  $N-C$  is  $1.32\text{\AA}$  ( $\text{\AA}$  denotes  
147 distance in angstroms), and the bond angle between the bonds  $C_\alpha-N$  and  $N-C$   
148 is  $123^\circ$ . Torsion angles are used to describe conformations around rotatable  
149 bonds (see Figure 2). Assume four consecutive atoms are connected by the  
150 three bonds  $b_{i-1}$ ,  $b_i$  and  $b_{i+1}$ . The torsion angle of  $b_i$  is defined as the smallest  
151 angle between the projections of  $b_{i-1}$  and  $b_{i+1}$  on the plane perpendicular to  
152 bond  $b_i$ . In Figure 2,  $\phi$ ,  $\psi$  and  $\omega$  are the torsion angles on the bonds  $N-C_\alpha$ ,  
153  $C_\alpha-C$  and  $C-N$ , respectively.

154 To capture the local structural features more accurately, we need to extract the  
155 features from a set of local residues. To obtain the local feature vector, we first  
156 represent each residue individually, and then consider the relationship between  
157 a pair of residues and a set of residues. For each residue, the length of  $C_\alpha-N$   
158 bond is  $1.47\text{\AA}$  and that of the  $C_\alpha-C$  bond is  $1.53\text{\AA}$ , and the angle between  $C_\alpha-$

159  $N$  and  $C_\alpha$ - $C$  bonds is  $110^\circ$ . Thus all the triangles formed by  $N$ - $C_\alpha$ - $C$  atoms in  
160 each residue are equivalent, and each residue can be represented by a triangle  
161 of the same size. The relationship between a pair of residues in 3D (three-  
162 dimensional) space can be fully described by the rigid transformation between  
163 two residues, which is a vector of 6 dimensions, containing 3 translational and  
164 3 rotational degrees of freedoms. To reduce the dimension of the vector, we  
165 use a distance and an angle to describe the transformation features between  
166 two residues.

We define the distance  $d$  between a pair of residues as the Euclidean distance between their  $C_\alpha$  atoms. The angle  $\theta$  between a pair of residues is defined as the angle between the planes that contain  $N$ - $C_\alpha$ - $C$  triangles representing each residue (see Figure 3). The distance and angle are invariant to displacement and rotation of the protein. The Euclidean distance between two  $C_\alpha$  atoms is calculated by their 3D coordinates directly. The angle between the two planes defined by the  $N$ - $C_\alpha$ - $C$  triangles, is calculated between their normals having  $C_\alpha$  as the origin. The normal of the plane defined by the triangle  $N$ - $C_\alpha$ - $C$  is given as

$$\vec{n} = \frac{\overrightarrow{NC_\alpha} \times \overrightarrow{C_\alpha C}}{\|\overrightarrow{NC_\alpha} \times \overrightarrow{C_\alpha C}\|}$$

The angle between the two normals  $\vec{n}_1$  and  $\vec{n}_2$  is then calculated as

$$\cos \theta = \frac{\|\vec{n}_1\|^2 + \|\vec{n}_2\|^2 - \|\vec{n}_2 - \vec{n}_1\|^2}{2\|\vec{n}_1\|\|\vec{n}_2\|}$$

167 To describe the local structural features between a set of residues, we slide  
168 a window of length  $w$  along the backbone of the protein. The distances and  
169 angles between the first residue  $i$  and all the other residues  $j$  (with  $j \in [i +$   
170  $1, i + w - 1]$ ) within the window are computed and added to a feature vector.



171 Each window is associated with one feature vector.

Let  $P = \{p_1, p_2, \dots, p_n\}$  represent a protein, where  $p_i$  is the  $i$ th-residue along the backbone. The set of feature vectors of the protein is given as  $P^v = \{p_1^v, p_2^v, \dots, p_{n-w+1}^v\}$ , where  $w$  is the sliding window size, and  $p_i^v$  is a feature vector

$$(d(p_i, p_{i+1}), \cos \theta(p_i, p_{i+1}), \dots, d(p_i, p_{i+w-1}), \cos \theta(p_i, p_{i+w-1}))$$

172 where  $d(p_i, p_j)$  is the distance between the residues  $p_i$  and  $p_j$ , and  $\cos \theta(p_i, p_j)$   
173 gives the cosine of the angle between the residues  $p_i$  and  $p_j$ . With window size  
174  $w$ , the dimension of each feature vector  $p_i^v$  is  $2(w - 1)$ .

### 175 3.2 Feature Normalization

176 Each structural feature vector is a combination of distances and angles, which  
177 have different measures. A normalization procedure is performed after the  
178 feature vectors are extracted. The angle  $\theta$  is in the range  $[0, \pi]$ , so  $\cos \theta \in$   
179  $[-1, 1]$ .

180 For normalizing the distances, we need to know the upper-bound on the dis-  
181 tance between the  $i$ -th and  $(i + w - 1)$ -th residue in the protein. From Fig-  
182 ure 1, the average distance between  $C_{\alpha 1}$ - $N$  atoms is  $d_1 = 1.47\text{\AA}$ , the average  
183 distance between  $N$ - $C$  atoms is  $d_2 = 1.32\text{\AA}$ , and the angle  $\alpha$  between  $C_{\alpha 1}$ -  
184  $N$  and  $N$ - $C$  bonds is  $123^\circ$ . The distance between  $C_{\alpha 1}$ - $C$  atoms is therefore  
185  $d(C_{\alpha 1}, C) = \sqrt{d_1^2 + d_2^2 - 2d_1d_2 \cos \alpha} = 2.453$ . The distance between  $C$ - $C_{\alpha 2}$   
186 atoms is  $d(C, C_{\alpha 2}) = 1.53$ , so the average distance between two  $C_\alpha$  atoms is:  
187  $d(C_{\alpha 1}, C_{\alpha 2}) \leq d(C_{\alpha 1}, C) + d(C, C_{\alpha 2}) = 2.453 + 1.57 = 4.023$ . If the distance

188 between two atoms is greater than 4.023, it is trimmed to 4.023. For a sliding  
 189 window of size  $w$ , the lower bound of the distance between any two atoms is  
 190 0, and the upper bound is  $4.023(w - 1)$ , so the distance between any pair of  
 191 residues within a  $w$  length window is in the range  $[0, 4.023(w - 1)]$ .

192 All the distances and angles are normalized and binned into an integer within  
 193 the range  $[0, b - 1]$ . We use the equation  $\lfloor \frac{d \times b}{4.023(w-1)} \rfloor$  to normalize and bin the  
 194 distances and  $\lfloor \frac{(\cos \theta + 1)b}{2} \rfloor$  to normalize and bin the angles. Table 1 shows 3  
 195 examples of normalized and binned feature vectors for  $w = 3$  and  $b = 10$ . The  
 196 size of each feature vector is  $2(w - 1) = 4$ , and the normalized value is within  
 197  $[0, 9]$ .

198 After normalization and binning, each feature vector is defined as  $p^s = \{p_0^s, p_1^s,$   
 199  $\dots, p_{2(w-1)-1}^s\}$ , where  $p_i^s$  is an integer within the range  $[0, b - 1]$ . Thus, the  
 200 structure of each protein  $P$  is converted into a structure-feature sequence  $P^s =$   
 201  $\{P_0^s, P_1^s \dots P_{n-w+1}^s\}$ , called the *SF-sequence*, where  $P_i^s$  is the  $i$ -th normalized  
 202 feature vector ( $p^s$ ) along the backbone. Note that each symbol within an SF-  
 203 sequence is a vector of length  $2(w - 1)$ , to which we assign a unique integer  
 204 identifier as its label. Thus the SF-sequences are over an alphabet of size  
 205  $b^{2(w-1)}$ .

### 206 3.3 Generalized Suffix Tree Index Construction

207 After obtaining the SF-sequences for all proteins in the database, we use a  
 208 generalized suffix tree (GST) as the indexing structure. A GST is a compact  
 209 representation of the suffixes of multiple sequences in a single tree, and can be  
 210 constructed in linear time [43]. A suffix can be located by following an unique

211 path from the root to a leaf.

212 To save the storage space of the suffix tree, we map each structure feature  
213 vector  $p^s$  to an unique key or symbol for the suffix tree construction, and  
214 map it back to the normalized vector when we need to compute the distance  
215 between two feature vectors. For instance, the three feature vectors in Table  
216 1 could be mapped to the symbols  $a$ ,  $b$  and  $x$  respectively.

217 Let  $GST$  be a generalized suffix tree, we use the following notation in the rest  
218 of the paper. We use  $N$  for a node in the suffix tree,  $E$  for an edge,  $C(E)$   
219 for a child node of the edge  $E$ ,  $L(E)$  for the label on edge  $E$ ,  $L(E[i])$  for the  
220  $i^{th}$  symbol of the edge label  $L(E)$ ,  $P(N)$  for the path-label of the node  $N$   
221 (formed by concatenating all the edge labels from the root node to  $N$ ), and  
222  $P(E[i])$  for the path-label of  $L(E[i])$ . Further, each leaf node in  $GST$  contains  
223 a sequence-position pair  $(x, p)$ , where  $x$  is a sequence identifier, and  $p$  is the  
224 start position of the suffix within sequence  $x$ . For any node  $N$ , we use the  
225 notation  $sp-list(N)$  for the collection of the sequence-position pairs for all  
226 the leaves under  $N$ .

227 Figure 4 shows an example of GST for two SF-sequences  $S_1 = xabxa$  and  
228  $S_2 = babxba$ , over the alphabet  $\{a, b, x\}$ , obtained by mapping each normalized  
229 feature vectors in Table 1 to a unique letter symbol. Node 0 is the root node,  
230 node 1 to 7 are internal nodes, and the rest are leaves. ‘\$’ is the unique  
231 termination character. The path label of node 7 is  $xa$ . The edge label  $L(E)$   
232 of the edge out of node 7 is  $bx$ , so its second character  $L(E[2])$  is  $x$ , and  
233 its path-label  $P(E[2])$  is  $xabx$ . The sequence-position identifier  $(1, 0)$  of the  
234 node 7 stands for  $xabxa$ , the suffix of sequence  $S_1$  that starts at position  
235 0. Thus  $sp-list(7) = \{(1, 3), (1, 0)\}$ , and the sp-list for node 6 is  $sp-list(6) =$

236  $\{(2, 3), (1, 3), (1, 0)\}$ .

### 237 3.4 *Parallel/Distributed External-Memory Suffix Tree Construction*

238 During construction of a typical in-memory suffix tree, a large amount of  
239 memory would be required to store the input strings and possibly some other  
240 bookkeeping information for large databases. This amount is normally too  
241 large for a typical computer; hence a disk-based suffix tree was selected as the  
242 method of indexing instead of an in-memory suffix tree.

243 In our experiments, TRELIS [36] was applied to create the disk-based suf-  
244 fix tree from all of the sequences in the database. TRELIS is an effective  
245 algorithm that builds the disk-based suffix tree based on a partitioning and  
246 merging method. It creates suffix trees for smaller substrings of the input se-  
247 quence(s), and stores the suffix trees according to their common prefixes. Then,  
248 it merges the subtrees of the same prefix together, and stores the subtrees sep-  
249 arately on disk. Let  $S$  denote the input sequence obtained by concatenating  
250 all sequences in the database. Our external-memory suffix indexing approach  
251 has three main steps:

- 252 (1) *Prefix Creation Phase*: The first step creates a list of variable-length pre-  
253 fixes  $\{P_0, P_1, \dots, P_{m-1}\}$ . Each prefix  $P_i$  is chosen so that its frequency in  
254 the input string  $S$  does not exceed a maximum frequency threshold (de-  
255 termined by the main-memory limit). This also means that the number  
256 of suffixes beginning with  $P_i$  as a prefix will fit in the main-memory.
- 257 (2) *Partitioning Phase*: In the second phase, the input string  $S$  is partitioned  
258 into segments  $R_i$  (Figure 5, step a). The segment size is chosen such that

259 the resulting suffix tree  $T_{R_i}$  from each segment (Figure 5, step b) fits  
 260 in main-memory. Each resulting suffix tree is further split into smaller  
 261 subtrees  $T_{R_i, P_j}$  (Figure 5, step c), that share a common prefix  $P_j$ , which  
 262 are then stored on the disk.

263 (3) *Merging Phase:* After processing all segments  $R_i$ , we merge all the sub-  
 264 trees  $T_{R_i, P_j}$  for each prefix  $P_j$  from the different partitions  $R_i$  into a  
 265 merged suffix subtree  $T_{P_j}$  (Figure 5, step d). The prefixes  $P_j$  were chosen  
 266 so that their suffix subtrees also fit entirely in memory. As each merged  
 267 subtree  $T_{P_j}$  is constructed, it is written to disk. The complete suffix tree  
 268 is simply a forest of these prefix-based subtrees ( $T_{P_j}$ ).

269 **Parallel/Distributed Suffix Tree Indexing:** The idea of prefix partition-  
 270 ing and merging is very suitable for parallel or distributed suffix tree con-  
 271 struction. For the prefix creation phase, let's assume initially that the set of  
 272 variable-length prefixes is known. In this case, the concatenated input sequence  
 273  $S$  can be partitioned among the available processors, and each processor can  
 274 obtain the local frequency of each prefix in its assigned segment (note that  
 275 some overlap has to be allowed among the sequence segments to take care  
 276 of boundary conditions). A summation over the processors yields the global  
 277 frequencies for the set of prefixes. Since the prefix set is, in fact, not known  
 278 *a priori*, the parallel prefix frequency computations can be done in multiple  
 279 count-reduce iterations. In each iteration, prefixes up to a given length are  
 280 counted (only those that exceed the frequency threshold in the last iteration),  
 281 and a reduction is done to obtain the global frequencies.

282 The partitioning phase is straightforward to parallelize, since each partition  
 283 is independent. Essentially, each processor builds the complete suffix tree  $T_{R_i}$   
 284 for partition  $R_i$  and splits them into the prefix-based suffix subtrees  $T_{R_i, P_j}$ ,

285 and stores them on disk. Since the partitions are all of the same size (with  
286 the exception of the last partition), a simple round-robin partition assignment  
287 scheme is sufficient to ensure good load balancing among the processors.

288 For the merging phase, we assign the variable-length prefixes among the pro-  
289 cessors. Each processor is responsible for merging the subtrees  $T_{R_i, P_j}$  from all  
290 the partitions  $R_0, R_1, \dots, R_{r-1}$ , for a given prefix  $P_j$ . The main complication  
291 here is that prefix-based suffix subtrees for partitions assigned to other proces-  
292 sors in the second phase, may not be available locally. Thus before the merge  
293 phase, each processor communicates its prefixed-based suffix subtrees for pre-  
294 fix  $P_j$  to the processor responsible for constructing the merged suffix tree  $T_{P_j}$ .  
295 Note that for the merging phase also a simple round-robin prefix assignment  
296 scheme suffices to achieve good a load balance, since each prefix yields suffix  
297 subtrees of approximately the same size.

## 298 4 Querying

299 Given a query  $Q$ , we first extract its feature vectors and convert it into a SF-  
300 sequence  $Q^s$  as described above. Then two phases are performed: searching and  
301 ranking. The searching phase retrieves all the matching segments/subsequences  
302 from the database within a distance threshold  $\epsilon$  (on a per symbol basis), and  
303 the ranking phase ranks all the proteins by chaining the matching segments.

### 304 4.1 Searching

305 For a given query SF-sequence  $Q^s = \{Q_1^s Q_2^s \dots Q_n^s\}$ , maximum feature dis-  
306 tance threshold  $\epsilon$ , and a minimum match length threshold  $l$ , the search algo-

307 rithm finds all maximal matching SF-subsequences  $P^s = \{P_1^s, P_2^s \dots P_m^s\}$  that  
308 occur in both the query SF-sequence and a database protein SF-sequence. A  
309 maximal match has the following properties:

- 310 (1) There exists a matching SF-subsequence  $Q_{i+1}^s \dots Q_{i+m}^s$  of  $Q^s$ , such that  
311  $dist(Q_{i+j}^s, P_j^s) < \epsilon$ , where  $j = 1, 2 \dots m$ ,  $Q_{i+j}^s$  and  $P_j^s$  are the normalized  
312 and binned feature vectors of length  $2(w-1)$ . The distance function used  
313 in our algorithm is Euclidean distance.
- 314 (2) The length of the match is at least as long as the length threshold, i.e.,  
315  $m \geq l$ .
- 316 (3) For any SF-subsequence  $P^s$  of protein  $R^s$  neither  $P^s v$  nor  $v P^s$  is a match-  
317 ing SF-subsequence of  $Q^s$  and  $R^s$  for any feature vector  $v$  (this ensures  
318 maximality).

319 For instance, *abx* is a maximal match between the SF-sequences *xabxa* and  
320 *babxba* in Figure 4. Note that our approach differs from MUMmer genome  
321 alignment method presented in [10] which finds *exact* maximal *unique* matches  
322 between two genomes.

323 To find all maximal matches within  $\epsilon$  between the query  $Q^s$  and suffix tree  
324  $GST_d$  built from the database proteins, one solution is to trace every SF-  
325 subsequence of  $Q^s$  from the root of  $GST_d$ . However in this approach, if there  
326 are common prefixes among the suffixes, they will be searched multiple times,  
327 leading to more comparisons than necessary. To reduce the number of compar-  
328 isons, we build another suffix tree  $GST_q$  for  $Q^s$ , and then traverse two suffix  
329 trees simultaneously to retrieve all the maximal matches. This way, each com-  
330 mon prefix is searched only once. In the discussion below, we use the subscript  
331  $q$  for the query, and  $d$  for the database. For instance,  $N_q$  stands for a query

332 suffix tree node, while  $N_d$  stands for a database suffix tree node.

333 The matching algorithm starts with the *MMS* procedure as shown in Figure 6,  
334 and its inputs are the root node ( $N_q$ ) of the query suffix tree  $GST_q$ , the root  
335 node ( $N_d$ ) of the database suffix tree  $GST_d$ , the distance tolerance  $\epsilon$  and the  
336 minimum length of the maximal match  $l$ . For every edge out of the query node  
337 and the database node, *MMS* calls the NodeSearch procedure (see Figure 7)  
338 to match their labels and follow the path to find all the matching nodes.

339 In the NodeSearch procedure (Figure 7), for two edges from different suffix  
340 trees, the distance between the corresponding pair of label symbols ( $L(E[i]_q)$   
341 and  $L(E[j]_d)$ ) is computed in step 2. If the distance is larger than  $\epsilon$ , which  
342 implies a mismatch, the procedure updates the *MMS*Set and proceeds to the  
343 next branch. If there is no mismatch, the short edge will reach the end first. If  
344 the child node of the short edge is a leaf, we need to update the *MMS*Set. If the  
345 child node is an internal node, two different procedures are called recursively.  
346 1) If the lengths of two edge labels are the same, then *MMS* procedure is called  
347 for two child nodes in step 3. 2) If one of the edge has a shorter label, the  
348 algorithm NodeSearch will be called recursively with the new input composed  
349 of all the edges out of the child node of the short edge (see steps 4 and 5).

350 Each matching SF-subsequence  $s$  is defined by two triplets  $(x, p, l)$  and  $(y, q, l)$ ,  
351 where  $p$  and  $q$  are the start positions of  $s$  in the query sequence  $Q_x$  and the  
352 protein sequence  $P_y$  respectively, and  $l$  is the length. If  $s$  is a maximal match,  
353 it will be added to the *MMS*Set in the *updateMMS* procedure. To identify  
354 a maximal match, we need to compare whether any extension of the match  
355 will result in a mismatch. In our algorithm, each common subsequence  $s$  is  
356 obtained either from a character mismatch or a leaf node, so we just need



357 to compare the characters before the common subsequence ( $Q_x[p - 1]$  and  
358  $P_y[q - 1]$ ) to identify the maximal matches.

359 We can also process multiple query SF-sequences at the same time by inserting  
360 them to the query suffix tree  $GST_q$ , so the nodes with the same path-label are  
361 visited only once and the performance will be improved.

## 362 4.2 Ranking

363 The maximal matches are obtained for the query sequence and reference se-  
364 quences in the database. Every maximal match is a diagonal run in the matrix  
365 formed by a query and reference sequence. We use the best diagonal runs de-  
366 scribed in the FASTA algorithm [35] as our ranking scheme. We calculate the  
367 alignment as a combination of the maximal matches with the highest score.  
368 The score of the alignment is the sum of the scores of the maximal matches  
369 minus the gap penalty. The length of a maximal match and a gap are used  
370 as the match score and gap penalty, respectively. Two maximal matches can  
371 be chained together if there is no overlap between them. We use a fast greedy  
372 algorithm to find the chains of maximal alignments. At first, the maximal  
373 matches are sorted by their length. The longest maximal match is chosen  
374 first, and we remove all other overlapping matches. Then we choose the next  
375 longest maximal match, remove its overlapping matches and repeat the above  
376 steps until no maximal matches are left. This way we find the longest chained  
377 maximal matches between the query and each retrieved database SF-sequence.  
378 Finally all the candidates with small alignment scores are screened out and  
379 only the top similar proteins are selected.

## 380 5 Results and Discussion

381 To evaluate the performance of our algorithm we conduct an extensive set of  
382 experiments. The first test compares the performance of PSIST with ProGreSS [3],  
383 a state-of-the-art protein indexing method. The second test compares the re-  
384 sults of suffix tree indexing using different pieces of information: sequence or  
385 structure. The third test shows the performance of indexing the whole set  
386 of proteins in SCOP [30], a database of proteins classified according to their  
387 structure. Our algorithm was implemented in C++ and all experiments re-  
388 ported below were done on a Power Mac G5 with 2.7GHz CPU, and 4GB  
389 Memory, running Darwin Kernel Version 8.0.0.

### 390 5.1 Comparison with ProGreSS

391 The CATH [33,34] database gives a hierarchical classification of protein do-  
392 main structures based on sequence and structure similarity. It operates on do-  
393 mains because domains are likely to be the fundamental evolutionary building  
394 blocks. CATH has four major levels of classification, namely Class, Archi-  
395 tecture, Topology and Homologous family. Homologous family is the lowest  
396 level; it contains either proteins having significant sequence similarity (35%)  
397 or high structural similarity and some sequence identity (20%). The sequences  
398 are aligned using dynamic programming and the structures are aligned using  
399 SSAP [32]. Protein domains that share a significant structure similarity but  
400 low sequence similarity are grouped into the same Topology. Architecture is  
401 assigned manually according to the gross arrangement of secondary structures  
402 in 3D space. At the top of the hierarchy, domains are clustered into four classes

403 automatically by the percentage of  $\alpha$ -helices or  $\beta$ -strands. The latest version  
404 (2.6.0) of the CATH database contains more than 67,000 domains classified  
405 into 6,003 homologous families.

406 We compare our approach with one of the best previous indexing approaches,  
407 ProGreSS [3], using the Java-based code provided by its authors. We choose  
408 the 35% representative dataset, consisting of 6003 domains from CATH, where  
409 sequence pairs have at most 35% amino acid sequence identity. Since ProGreSS  
410 can not handle a large dataset, we selected 2000 domains randomly out of the  
411 35% representative set of CATH as our dataset  $D$ . From topologies with least  
412 8 proteins, one protein is chosen randomly as the query, resulting in a query  
413 set  $D_q$ , having 42 proteins.

414 To evaluate our algorithm we perform three different tests: The *retrieval test*  
415 finds the number of correct matching structures from the same topology as  
416 the query among the top  $k$  scoring proteins, and the *classification test* tries  
417 to classify the query at the topology and class levels. The *performance test*  
418 compares the algorithms in terms of the total running time.

#### 419 5.1.1 Retrieval Test

420 We ran the experiments using PSIST and ProGreSS to obtain the number of  
421 proteins found from the same topology for each of the 42 queries. There are  
422 five parameters used in our approach:  $w$  is the size of the window used to  
423 index the local features,  $b$  is the range used to normalize the feature vectors,  
424  $\epsilon$  is the distance threshold based on the normalized feature vectors,  $l$  is the  
425 minimum length of the maximal matches, and  $k$  is the number of top scoring  
426 proteins reported. Based on our tuning experiments for PSIST we set  $w = 3$ ,

427  $b = 2$ ,  $\epsilon = 0$  and  $l = 10$ . For fair comparison, we tuned the parameter settings  
428 for ProGreSS to report its best results (we use sequence distance threshold  
429  $\epsilon_t = 0.05$ , the structure distance threshold  $\epsilon_q = 0.01$  and window size  $w = 3$ ).

430 Figure 9 and Table 2 show the number of proteins found from the same topol-  
431 ogy for different top- $k$  cutoffs. Note that the number of correct matches is an  
432 average over all 42 CATH topologies used in our test. We find that on average  
433 PSIST returns more correct matches; for example in the top 20 results, PSIST  
434 has 4 correct matches, whereas ProGreSS returns only 2 correct matches. For  
435 the top  $k = 100$ , PSIST returns around 10 matches, whereas ProGreSS returns  
436 only 7.3 correct matches.

### 437 5.1.2 Classification Test

438 In the classification test, we assume we do not know the topology or the class  
439 to which a query protein belongs. For each query we then classify it into one of  
440 the 42 CATH topologies and one of the four CATH classes (all  $\alpha$ , all  $\beta$ ,  $\alpha + \beta$   
441 and  $\alpha/\beta$ ) as follows. For each query, the top  $k$  similar proteins are selected  
442 from the database. The query itself is not counted in the top  $k$  matches.  
443 Each protein among the top  $k$  matches is assigned a score, a topology id,  
444 and a class id. The scores of the top  $k$  proteins from the same topology or  
445 class are accumulated. The query is assigned to the topology or class with  
446 the highest score. This classification approach can thus be thought of as  $k$   
447 nearest neighbor classification. Below we tabulate the results separately for  
448 the topology-level and class-level classification, and we report the percentage  
449 of correctly classified query proteins (out of the 42 queries). For PSIST and  
450 ProGreSS we use the best parameter settings reported in the last section.

451 Table 3 shows the CATH classification comparison at the topology and class  
452 level respectively. ProGreSS uses both the structure and sequence features to  
453 classify the proteins, and its accuracy is 7.14% and 57.1% at the topology  
454 and class levels. Without considering the sequence features, PSIST has much  
455 better performance than ProGreSS; its accuracy is 50.0% and 92.9% at the  
456 topology and class levels.

### 457 5.1.3 Performance Test

458 We compare the running time of different approaches in this section. Suppose  
459 a protein has  $n$  residues, the window size is  $w$ , then the number of feature  
460 vectors is  $n - w + 1$ , so the complexity of our approach is  $O(n - w - 1) = O(n)$   
461 per protein.

462 Both ProGreSS and PSIST provide a trade-off between the running time and  
463 the accuracy performance by adjusting the parameters such as window size  
464 and distance. Table 4 shows the running time for ProGreSS and PSIST. For  
465 ProGreSS, we choose the best sequence and structure distance thresholds and  
466 set window size  $w = 3$ . For PSIST, we set the same parameters  $w = 3$ ,  $b = 2$ ,  
467  $\epsilon = 0$  for all three cases, but different length of maximal matches:  $l = 18$   
468 for the first case,  $l = 14$  for the second case and  $l = 10$  for the third case.  
469 All three cases have better retrieval and classification performances compared  
470 to ProGreSS. The first case is 2.75 times faster than ProGreSS, the second  
471 case is 1.57 times faster, and the third case is the slowest, but it has the best  
472 performance.

## 473 5.2 Sequence and Structure Comparison

474 In this test, we choose the same 35% representative set as our database, which  
475 has 6003 domains. However, we select the queries from CATH using a different  
476 method. We choose all of the *singletons* from the 35% representative set of  
477 CATH domains. If one domain is the only member in a homology family, it is  
478 called a singleton. If a topology has only one homology family, it is impossible  
479 to obtain homologous proteins of the singleton in that homology family, so  
480 we need to prune out these impossible singletons. After pruning, there are  
481 370 singletons out of 6003 domains in the 35% representative set. These 370  
482 singletons comprise the query set. For any of these singletons it is very hard  
483 to obtain similar proteins from other homologue families.

484 To evaluate the performance of our algorithm, we use Receiver Operating  
485 Characteristic (ROC) [17] score as our measurement. The ROC score is the  
486 area under the ROC curve, which plots true positives versus true negatives in  
487 the retrieved set of proteins. It combines measures of sensitivity and specificity.  
488 A score of 1 indicates perfect separation of positives from negatives, while a  
489 score of 0 means that none of the selected proteins are in the same topology  
490 as the query.

491 Two approaches are considered, one using the amino acid sequences, the other  
492 using the structures. The average ROC score for sequences was 26%, while  
493 the average the ROC score for structures was 30%. Figure 10 shows the total  
494 number of queries whose ROC score exceeds a given ROC score threshold (on  
495 the x-axis). Not surprisingly, using the structural information leads to better  
496 retrieval quality.

498 The SCOP database [30] classifies proteins according to a four level hierarchi-  
499 cal classification, namely, family, super-family, fold and class. SCOP release  
500 1.69 (from July 2005) contains a total of 25973 proteins and 70859 domains,  
501 spanning 2845 families, 1539 super-families, and 945 folds. Since the SCOP  
502 database is curated by visual inspection it is considered to be extremely accu-  
503 rate. For our tests the target has all the proteins from four classes of SCOP:  
504 all  $\alpha$ , all  $\beta$ ,  $\alpha + \beta$  and  $\alpha/\beta$ . Our dataset  $D$  contains a total of 70,500 ASTRAL  
505 SCOP 1.69 genetic domains [4]. ProGreSS ran out of memory when building  
506 the index. For PSIST, the indexing time was 3184.4 seconds and the average  
507 running time for each query was about 104.7 seconds with the parameters  
508  $w = 3$ ,  $b = 2$ ,  $\epsilon = 0$  and  $l = 15$ .

## 509 6 Conclusion

510 In this paper, we presented a new local feature representation for protein struc-  
511 tures. We transform the structure indexing problem into a sequence indexing  
512 problem by directly indexing the structure-feature sequences using suffix trees.  
513 The suffix trees enable rapid retrieval of maximal matching segments, which  
514 are chained into longer local structural alignments. Finally the matches are  
515 ranked according to their alignment scores. Compared to ProGreSS, our ap-  
516 proach can index much larger databases, and at the same time it obtains  
517 higher retrieval accuracy. We also show that PSIST is highly scalable due to  
518 the distributed, and external suffix tree indexing approach it uses; it is able  
519 to index about 70,500 domains from SCOP [30] in under an hour!

## 520 **Acknowledgment**

521 This work was supported in part by NSF CAREER Award IIS-0092978, DOE  
522 Career Award DE-FG02-02ER25538, NSF grant EIA-0103708, and NSF grant  
523 EMT-0432098. We thank Benjarath Phoophakdee for her help in integrating  
524 the TRELIS external suffix tree method with PSIST. We also thank Tolga  
525 Can, Arnab Bhattacharya and Ambuj Singh for providing us the ProGreSS  
526 code and other assistance. Finally we thank Chris Bystroff for the many helpful  
527 suggestions.

## 528 **References**

- 529 [1] S. Altschul, T. Madden, A. Schaffer, J. Zhang, Z. Zhang, W. Miller, D. Lipman,  
530 Gapped BLAST and PSI-BLAST: a new generation of protein database search  
531 programs, *Nucleic Acids Research*, 25(17), 1997, pp. 3389-3402.
- 532 [2] Z. Aung, W. Fu, K.L. Tan, An efficient index-based protein structure database  
533 searching method, *Int'l Conf. on Database Systems for Advanced Applications*  
534 (DASFAA), 2003, pp. 311-318.
- 535 [3] A. Bhattacharya, T. Can, T. Kahveci, A.K. Singh, Y.F. Wang, ProGreSS:  
536 Simultaneous Searching of Protein Databases by Sequence and Structure, *Pacific*  
537 *Symp. on Bioinformatics*, 2004, pp. 264-275.
- 538 [4] S.E. Brenner, P. Koehl, M. Levitt, The ASTRAL compendium for sequence and  
539 structure analysis, *Nucleic Acids Research*, 28, 2000, pp. 254-256.
- 540 [5] O. Çamoğlu, T. Kahveci, A.K. Singh, Towards index-based similarity search for  
541 protein structure databases, *IEEE Computer Society Bioinformatics Conference*  
542 (CSB), 2003, pp. 148-158.
- 543 [6] T. Can, Y.F. Wang, CTSS: a robust and efficient method for protein structure  
544 alignment based on local geometrical and biological features, *IEEE Computer*



- 545 Society Bioinformatics Conference (CSB), 2003, pp. 169-179.
- 546 [7] I. Choi, J. Kwon, S. Kim, Local feature frequency profile: A method to measure  
547 structural similarity in proteins, *Proc. Nat'l Acad. Sci. USA*, 101(11), 2004, pp.  
548 3797-3802.
- 549 [8] R. Clifford, Distributed Suffix Trees, *Journal of Discrete Algorithms*, 3(2-4), June  
550 2005, pp. 176-197.
- 551 [9] M.O. Dayhoff, R.M. Schwartz, B.C. Orcutt, A model of evolutionary change in  
552 proteins. In *Atlas of Protein Sequence and Structure*, M.O. Dayhoff (ed.), pp.  
553 345-352, 1978.
- 554 [10] A.L. Delcher, S. Kasif, R.D. Fleischmann, J. Peterson, O. White, S.L. Salzberg,  
555 Alignment of whole genomes, *Nucleic Acid Research*, 27(11), 1999, pp. 2369-2376.
- 556 [11] A.L. Delcher, A. Phillippy, J. Carlton, S.L. Salzberg, Fast algorithms for large-  
557 scale genome alignment and comparison, *Nucleic Acid Research*, 30(11), 2002, pp.  
558 2478-2483.
- 559 [12] O. Dror, H. Benyamini, R. Nussinov, H. Wolfson, MASS: Multiple structural  
560 alignment by secondary structures, *Bioinformatics*, 19(12), 2003, pp. 95-104.
- 561 [13] I. Eidhammer, I. Jonassen, W.R. Taylor, Structure comparison and structure  
562 patterns, *Journal of Computational Biology*, 7(5), 2000, pp. 685-716.
- 563 [14] M. Farach-Colton, P. Ferragina, S. Muthukrishnan, On the sorting-complexity  
564 of suffix tree construction, *Journal of the ACM*, 47(6), 2000, pp. 987-1011.
- 565 [15] F. Gao, M.J. Zaki, PSIST: Indexing Protein Structures using Suffix Trees, *IEEE*  
566 *Computational Systems Bioinformatics Conference*, 2005.
- 567 [16] A. Godzik, The structural alignment between two proteins: is there a unique  
568 answer? *Protein Science*, 5, 1996, pp. 1325-1338.
- 569 [17] M. Gribskov, N. Robinson, Use of receiver operating characteristic (ROC)  
570 analysis to evaluate sequence matching, *Comput. Chem.*, 20, 1996, pp. 25-33.
- 571 [18] D. Gusfield, *Algorithms on strings, trees, and sequences: Computer science and*  
572 *computational biology*, Cambridge University Press, New York, 1997.

- 573 [19] S. Henikoff, J.G. Henikoff, Amino acid substitution matrices from protein  
574 blocks, Proc. Nat'l. Acad. Sci. USA, 89(22), 1992, pp. 10915-9.
- 575 [20] L. Holm, C. Sander, Protein structure comparison by alignment of distance  
576 matrices, J. Mol. Biol, 233, 1993, pp. 123-138.
- 577 [21] L. Holm, C. Sander, 3-D lookup: fast protein structure database searches at  
578 90% reliability, Int'l Conf. on Intelligent Systems for Molecular Biology (ISMB),  
579 1995, pp. 179-187.
- 580 [22] E. Hunt, M. Atkinson, R. Irving. A database index to large biological sequences,  
581 Int'l Conf. on Very Large Data Bases (VLDB), 2001.
- 582 [23] E. Hunt, M. Atkinson, R. Irving, Database indexing for large DNA and protein  
583 sequence collections, Int'l Conf. on Very Large Data Bases (VLDB), 2002, pp.  
584 256-271.
- 585 [24] R. Japp, The top-compressed suffix tree, 21st Annual British Nat'l Conf. on  
586 Databases, 2004.
- 587 [25] Y. Lamdan, H.J. Wolfson, Geometric hashing: a general and efficient model-  
588 based recognition scheme, Int'l Conf. on Computer Vision (ICCV), 1988, pp. 238-  
589 249.
- 590 [26] T. Madej, J.F. Gibrat, S.H. Bryant, Threading a database of protein cores,  
591 Proteins, 23, 1995, pp. 356-369.
- 592 [27] E.M. McCreight, A space-economic suffix tree construction algorithm, Journal  
593 of the Association for Computing machinery, 23(2), 1976, pp. 262-272.
- 594 [28] C. Meek, J.M. Patel, S. Kasetty, OASIS: An online and accurate Technique for  
595 local-alignment searches on biological sequences, Int'l Conf. on Very Large Data  
596 Bases (VLDB), pp. 910-923.
- 597 [29] K. Mizoguchi, N. Go, Comparison of spatial arrangements of secondary  
598 structural elements in proteins, Protein Engineering, 8, 1995, pp. 353-362.
- 599 [30] A.G. Murzin, S.E. Brenner, T. Hubbard, C. Chothia, SCOP: a structural  
600 classification of proteins database for the investigation of sequences and structures,  
601 J. Mol. Biol., 247, 1995, pp. 536-540.

- 602 [31] R. Nussinov, N. Leibowitz, H.J. Wolfson, MUSTA: a general, efficient, automated  
603 method for multiple structure alignment and detection of common motifs:  
604 Application to proteins, *J. Comp. Bio.*, 8(2), 2001, pp. 93-121.
- 605 [32] C.A. Orengo, W.R. Taylor, SSAP: Sequential structure alignment program for  
606 protein structure comparisons. *Methods in Enzymology*, 266, 1996, pp. 617-634.
- 607 [33] C. Orengo, A. Michie, S. Jones, D. Jones, M. Swindells, J. Thornton, CATH -  
608 a hierarchic classification of protein domain structures. *Structure*, 5(8), 1997, pp.  
609 1093-1108.
- 610 [34] F. Pearl, I. Todd, A. Sillitoe, M. Dibley, O. Redfern, T. Lewis, C. Bennett, R.  
611 Marsden, A. Grant, D. Lee, A. Akpor, M. Maibaum, A. Harrison, T. Dallman, G.  
612 Reeves, I. Diboun, S. Addou, S. Lise, C. Johnston, A. Sillero, J. Thornton, and  
613 C. Orengo, The CATH domain structure database and related resources gene3d  
614 and DHS provide comprehensive domain family information for genome analysis,  
615 *Nucleic Acids Research*, 33, 2005, pp. 247-251.
- 616 [35] W.R. Pearson, D.J. Lipman, Improved tools for biological sequence comparison,  
617 *Proc. Nat'l Acad. Sci. USA*, 85, 1988, pp. 2444-2448.
- 618 [36] B. Phoophakdee, M. J. Zaki, TRELLIS: Genome-scale Disk-based Suffix Tree  
619 Indexing, *ACM SIGMOD International Conference on Management of Data*, June  
620 2007.
- 621 [37] B. Rost, Twilight zone of protein sequence alignments, *Protein Engineering*,  
622 12(2), 1999, pp. 85-94.
- 623 [38] K.B. Schürmann, J. Stoye, Suffix tree construction and storage with limited  
624 main memory, *Tech. Report 0946-7831*, Universität Bielefeld, 2003.
- 625 [39] M. Shatsky, R. Nussinov, H.J. Wolfson, Multiprot - a multiple protein structural  
626 alignment algorithm, *Proteins*, 56, 2004, pp. 143-156.
- 627 [40] I.N. Shindyalov, P.E. Bourne, Protein structure alignment by incremental  
628 combinatorial extension (CE) of the optimal path, *Protein Engineering*, 11(9),  
629 1998, pp. 739-747.

- 630 [41] S. Tata, R. Hankins, J. Patel, Practical suffix tree construction, Int'l Conf. on  
631 Very Large Data Bases (VLDB), 2004, pp. 36-47.
- 632 [42] H. Taubig, A. Buchner, J. Griebisch, A method for fast approximate searching of  
633 polypeptide structures in the PDB, German Conference on Bioinformatics (GCB),  
634 2004.
- 635 [43] E. Ukkonen, On-line construction of suffix trees, *Algorithmica* 14(3), 1995, pp.  
636 249-260.
- 637 [44] X. Yuan, C. Bystroff, Non-sequential Structure-based Alignments Reveal  
638 Topology-independent Core Packing Arrangements in Proteins, *Bioinformatics*,  
639 21(7), 2005, pp. 1010-1019.

#### 640 **Author Biographies**

641 **Feng Gao** received his B.E. degree from Shandong University of Technology  
642 in 1995 and the M.E. degree from Harbin Institute of Technology in 1997.  
643 He received his Ph.D. degree in computer science from Rensselaer Polytechnic  
644 Institute in December, 2006. His research interests are in data mining and  
645 bioinformatics.

646 **Mohammed J. Zaki** is an Associate Professor of Computer Science at RPI.  
647 He received his Ph.D. degree in computer science from the University of  
648 Rochester in 1998. His research interests focus on developing novel data mining  
649 techniques and their applications, especially for bioinformatics. He has pub-  
650 lished over 160 papers on data mining, and co-edited several books (including  
651 "Data Mining in Bioinformatics, Springer-London, 2005). He is currently an  
652 associate editor for *IEEE Transactions on Knowledge and Data Engineering*,  
653 action editor for *Data Mining and Knowledge Discovery*, and is on the edito-

654 rial boards of Statistical Analysis and Data Mining, Scientific Programming,  
655 International Journal of Data Mining and Bioinformatics, International Jour-  
656 nal of Data warehousing and Mining, International Journal of Business and  
657 Systems Research, and The Open Artificial Intelligence Journal. He is a re-  
658 cipient of the NSF CAREER Award (2001) and DOE Early Career Award  
659 (2002). He received the ACM Recognition of Service Award in 2003, and an  
660 IEEE Certificate of Appreciation in 2005.

Table 1

Examples of normalized feature vectors for  $w = 3$  and  $b = 10$ 

	<b>Feature vector</b>			
	$d$	$\cos \theta$	$d$	$\cos \theta$
original	3.55	0.29	5.4	-0.23
normalized	4	6	6	3
original	4.04	0.11	5.75	-0.25
normalized	5	5	7	3
original	3.60	0.45	5.29	0.21
normalized	4	7	6	6

Table 2

Overall comparison of the number of proteins found from the same topology among the top k candidates

<b>Algorithm</b>	<b>top4</b>	<b>top10</b>	<b>top50</b>	<b>top100</b>
ProGreSS	1.17	1.52	3.81	7.33
PSIST	2.02	3.17	6.29	10.02

Table 3

CATH classification accuracy comparison at the topology (TO) and class (CL) level

<b>Algorithm</b>	<b>Topology</b>	<b>Class</b>
ProGreSS	7.14 %	57.1%
PSIST	50.0%	92.9 %



Table 4  
Running time comparison

<b>Algorithm</b>	<b>TO%</b>	<b>CL%</b>	<b>top10</b>	<b>time(s)</b>
ProGreSS	7.14%	57.1%	1.52	1.57
PSIST-1	33.3%	64.3%	2.57	0.57
PSIST-2	47.6%	88.0%	2.93	0.95
PSIST-3	50.0%	92.9 %	3.17	2.08

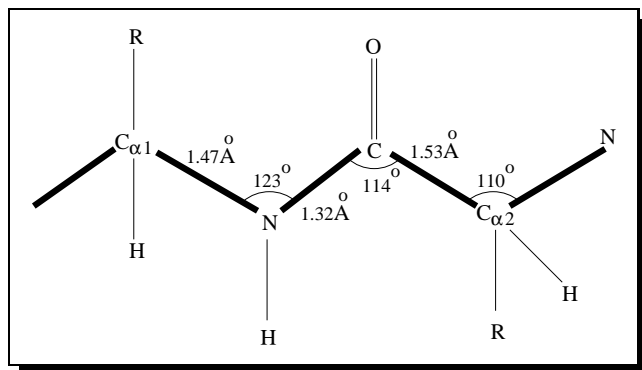


Fig. 1. Bond length and bond angles

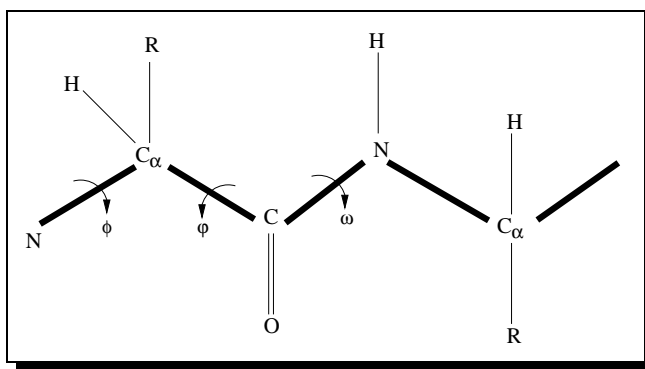


Fig. 2. Torsion angles

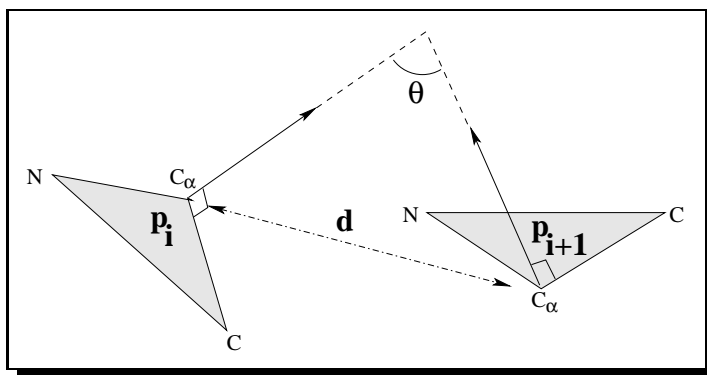


Fig. 3. The distance and angle between two residues

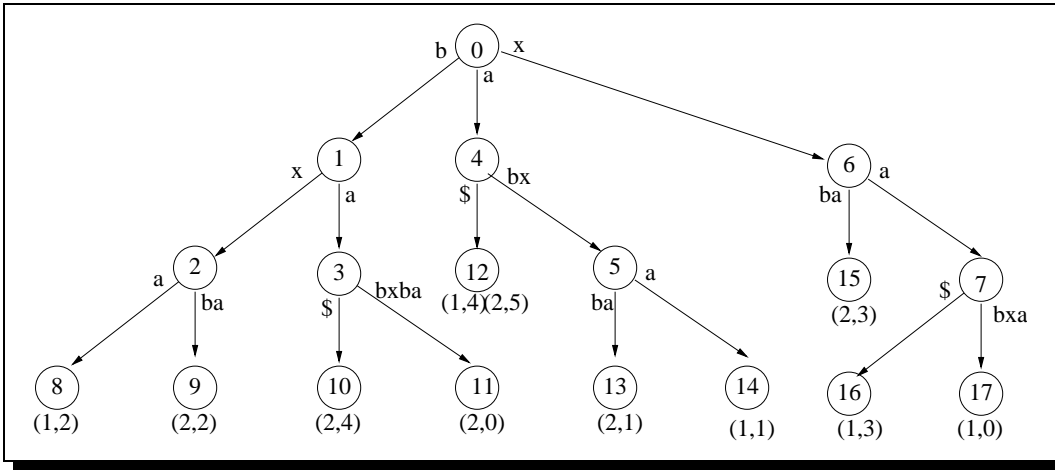


Fig. 4. GST for sequences  $S_1 = xabxa$  and  $S_2 = babxba$

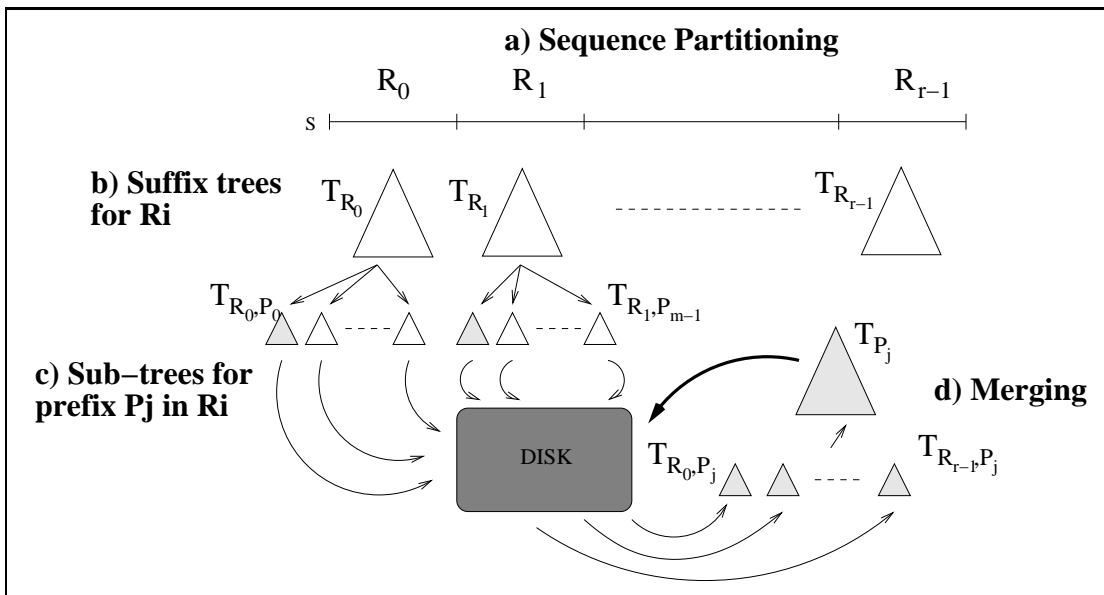


Fig. 5. Overview of External-Memory Suffix Tree

<b>Input</b>	: query Node $N_q$ , database Node $N_d$ , distance $\epsilon$ , length threshold $l$
<b>Output</b>	: maximal matches set ( $MMSet$ )
<b>Initialization:</b>	$MMSet = \emptyset$
<b>Procedure: MMS(<math>N_q, N_d, \epsilon, l</math>)</b>	
<b>foreach</b> edge $E_q$ out of $N_q$ <b>do</b>	
	<b>foreach</b> edge $E_d$ out of $N_d$ <b>do</b>
	NodeSearch( $E_q, 0, E_d, 0, \epsilon, l$ ).

Fig. 6. Maximal Match Search (MMS) Algorithm

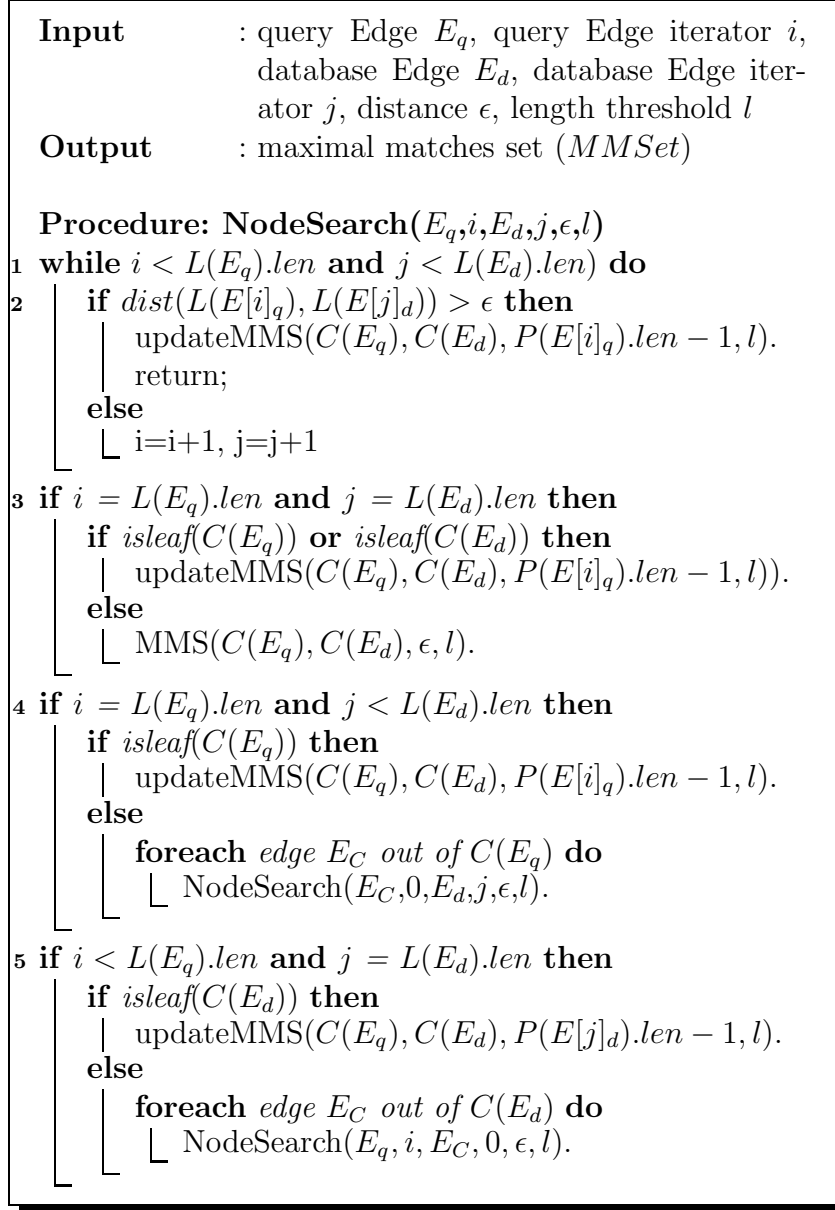


Fig. 7. Node Search Algorithm

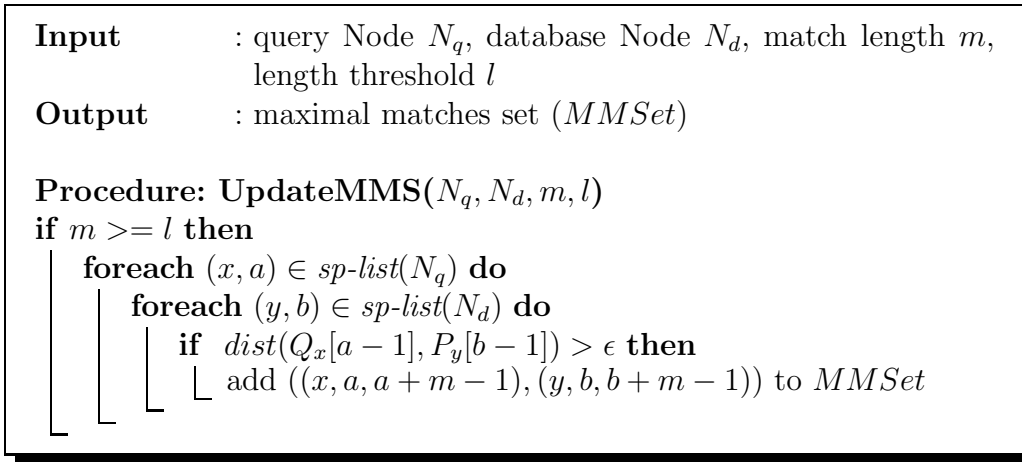


Fig. 8. Update Maximal Match Set Algorithm

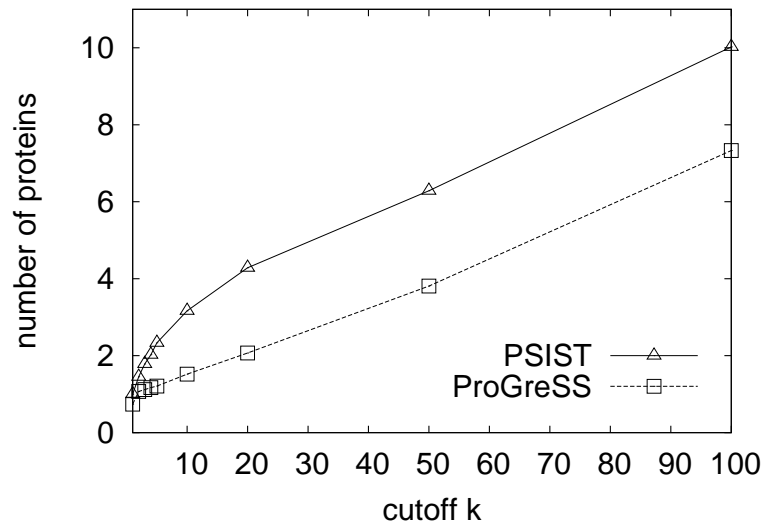


Fig. 9. Number of proteins found from the same topology for different top- $k$  value ( $w = 3$ ,  $b = 2$ ,  $\epsilon = 0$  and  $l = 10$ ).

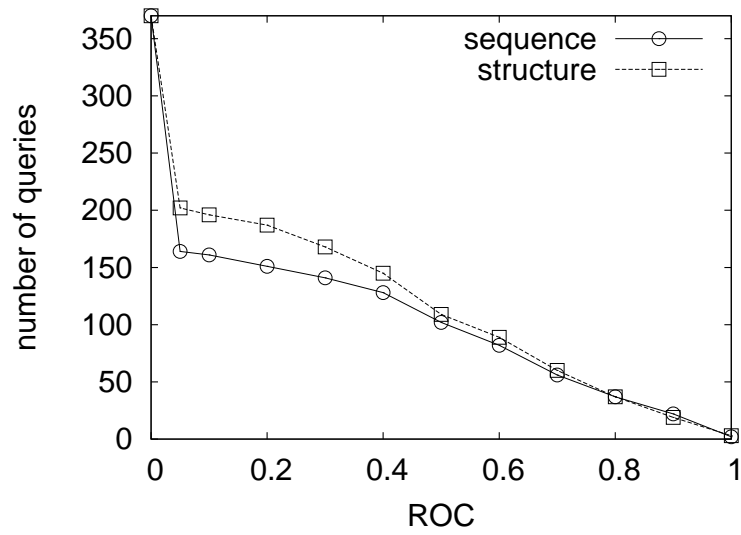


Fig. 10. Relative performance