

# 1 Protein Classification with Improved Topological 2 Data Analysis

3 **Tamal Dey**

4 Department of Computer Science and Engineering, The Ohio State University, Columbus, USA  
5 dey.8@osu.edu

6  <http://web.cse.ohio-state.edu/~dey.8/>

7 **Sayan Mandal**

8 Department of Computer Science and Engineering, The Ohio State University, Columbus, USA  
9 mandal.25@osu.edu

10  <http://web.cse.ohio-state.edu/~mandal.25/>

## 11 — Abstract —

---

12 Automated annotation and analysis of protein molecules have long been a topic of interest due to  
13 immediate applications in medicine and drug design. In this work, we propose a topology based,  
14 fast, scalable, and parameter-free technique to generate protein signatures.

15 We build an initial simplicial complex using information about the protein's constituent atoms,  
16 including its radius and existing chemical bonds, to model the hierarchical structure of the mo-  
17 lecule. Simplicial collapse is used to construct a filtration which we use to compute persistent  
18 homology. This information constitutes our signature for the protein. In addition, we demon-  
19 strate that this technique scales well to large proteins. Our method shows sizable time and  
20 memory improvements compared to other topology based approaches. We use the signature to  
21 train a protein domain classifier. Finally, we compare this classifier against models built from  
22 state-of-the-art structure-based protein signatures on standard datasets to achieve a substantial  
23 improvement in accuracy.

24 **2012 ACM Subject Classification** Applied Computing → Life and medical sciences

25 **Keywords and phrases** topological data analysis, persistent homology, simplicial collapse, super-  
26 vised learning, topology based protein feature vector, protein classification

27 **Digital Object Identifier** 10.4230/LIPIcs.WABI.2018.6

28 **Supplement Material** <http://web.cse.ohio-state.edu/~dey.8/proteinTDA>

29 **Acknowledgements** This work has been supported by NSF grants CCF-1318595, CCF-1526513,  
30 and CCF-1733798.

## 31 **1** Introduction

32 Proteins are by far the most anatomically intricate and functionally sophisticated molecules  
33 known. The benchmarking and classification of unannotated proteins have been done by  
34 researchers for quite a long time. This effort has direct influence in understanding behavior of  
35 unknown proteins or in more advanced tasks as genome sequencing. Since the sheer volume  
36 of protein structures is huge, up till the last decade, it had been a cumbersome task for  
37 scientists to manually evaluate and classify them. For the last decade, several works aiming  
38 at automating the classification of proteins have been developed. The majority of annotation  
39 and classification techniques are based on sequence comparisons (for example in BLAST [19],  
40 HHblits [2] and [18]) that try to align protein sequences to find homologs or a common  
41 ancestor. However, since those methods focus on finding sequence similarity, they are more



© Tamal K. Dey and Sayan Mandal;

licensed under Creative Commons License CC-BY

18th International Workshop on Algorithms in Bioinformatics (WABI 2018).

Editors: Laxmi Parida and Esko Ukkonen; Article No. 6; pp. 6:1–6:13

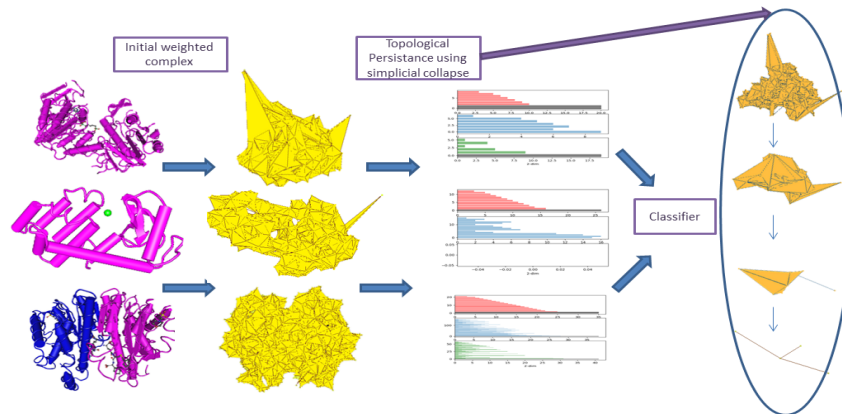
Leibniz International Proceedings in Informatics



LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

## 6:2 Protein Classification with Improved Topological Data Analysis

42 efficient in finding close homologs. Some domains such as remote homologs are known to  
43 have less than 25% sequential similarity and yet have common ancestors and are functionally  
44 similar. So, we miss out important information on structural variability while classifying  
45 proteins solely based on sequences. Even though, sometimes, homology is established by  
46 comparing structural alignment [14], accurate and fast structural classification techniques for  
the rapidly expanding Protein Data Bank remains a challenge.



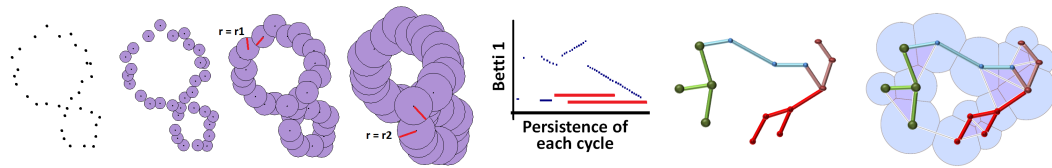
■ **Figure 1** Workflow of our technique

47

48 Several works on the classification of protein structures exist in the literature. The main  
49 intuition behind these works draws upon a heuristic that generates a signature for each  
50 protein strand so that structurally close proteins have similar signatures. Essentially, the  
51 signature alignment quantifies the similarity between two protein structures. The problem,  
52 however, remains with the speed of computing these signatures and the degree of their  
53 representative power. We want a fingerprint for the protein that can be computed fast and  
54 can tell whether two proteins are dissimilar or even marginally similar.

55 Some works use vector of frequencies to describe structural features while others take vari-  
56 ous physical properties into account such as energy, surface area, volume, flexibility/rigidity  
57 or use other features from geometric modeling. The "Bag-Of-Word" (BOW) representation  
58 to describe an object has been used in computer vision, natural language processing and  
59 various other fields. The work by Budowski-Tal [3] have described protein structure using a  
60 fragment library in a similar context. Since we use this work for comparison, we shall discuss  
61 its details later.

62 Topological data analysis [10], a newly developed data analysis technique has been  
63 shown to give some encouraging results in protein structure analysis. Topological signatures,  
64 particularly based on Persistent Homology, enjoy some nice theoretical properties including  
65 their robustness and scale invariance. These features are global and more resilient to local  
66 perturbations. Moreover, they are invariant to scaling and any isometric transformation of  
67 the input. The authors in [23] extract topological fingerprints based on the alignment of  
68 atoms and molecules in three dimensional space. Their work shows the impact of persistent  
69 homology in the modeling of protein flexibility which is ultimately used in protein B-factor  
70 analysis. This work also characterizes the evolution of topology during protein folding and  
71 thereby predicts its stability. For this task, the authors have introduced a coarse grain (CG)  
72 representation of proteins by considering an amino acid molecule as an atom  $C_\alpha$ . This helps  
73 them describe the higher level protein structures using the topological fingerprint perfectly.  
74 However, since the CG homology may be inconsistent due to ambiguity in choosing the CG



■ **Figure 2** Persistence of a point cloud in  $\mathbb{R}^2$  and its corresponding barcode

■ **Figure 3** Weighted Alpha complex for protein structure

75 particle, we present a similar study on secondary structures using our signature and show that  
 76 our method does not require such a representation as it is inherently scaling independent.

77 The authors in [4] have used persistent homology to generate feature vector in the context  
 78 of machine learning algorithms applied to protein structure explorations. We explore further  
 79 to improve upon the technique to eliminate its deficiencies. First, the approach in [4] does not  
 80 differentiate between atoms belonging to different elements. Also, it does not account for the  
 81 existing chemical bonds between the atoms while building the signature. Most importantly  
 82 it uses Vietoris Rips(VR) complex to generate the topological features for protein complex  
 83 which suffers from the well-known problem of scalability. As we will describe later, the VR  
 84 complex developed in the early 20th century grows rapidly in size even for moderate size  
 85 protein structures. Current state-of-the-art techniques, which have addressed the problem to  
 86 some extent, are still very cumbersome and slow especially for structures having about 30,000  
 87 atoms on an average. Among the several methods that generate persistence signature from a  
 88 point cloud, the PHAT toolbox [1] based on several efficient matrix reduction strategies and  
 89 GUDHI [22] library based on some compression techniques have been popular because of  
 90 their space and time efficiencies. A recent software called SimBa [8] published last year, has  
 91 been shown to work faster for large datasets. Yet, for our application, SimBa falls short as  
 92 we shall see later.

93 The algorithm that we present here is a fast technique to generate a topological signature  
 94 for protein structures. We build our signature based on the coordinates of the atoms in  $\mathbb{R}^3$   
 95 using their radius as weights. Since we also consider existing chemical bonds between the  
 96 atoms while building the signature, we believe that the hierarchical convoluted structure of  
 97 protein is captured in our features. Finally, we have developed a new technique to generate  
 98 persistence that is much quicker and uses less space than even the current state-of-the-art  
 99 such as SimBa. It helps us generate the signature even for reasonably large protein structures.  
 100 In sum, in this paper, we focus on three problems: (1) effectively map a protein structure into  
 101 a suitable complex; (2) develop a technique to generate fast persistent signature from this  
 102 complex; (3) use this signature to train a machine learning model for classification and compare  
 103 against other techniques. Our entire method is summarized in figure 1. We also illustrate  
 104 this method using a supplementary video available at <https://youtu.be/yf9UWgdTo>.

## 105 2 Methods

106 We use the theory of topological persistence to generate features for protein structures. These  
 107 topological features serve as a distinct signature for each protein strand. In this section, we  
 108 give some background on persistent homology followed by how we construct our signature.

109 **2.1 Persistence signature of point cloud data**

110 We start with a point cloud data in any  $n$ -dimensional Euclidean space. These will essentially  
 111 be the centers of protein atoms in the three dimensional space. However, to illustrate the  
 112 theory of persistent homology, we consider a toy example of taking a set of points in two  
 113 dimensions sampled uniformly from a two-hole structure (Fig. 2). We start growing balls  
 114 around each point, increasing their radius  $r$  continually and tracking the behavior of the  
 115 union of these growing balls. If we start increasing  $r$  from zero, we notice that at  $r = r_1$   
 116 (third from left in Fig 2) both holes are prominent in the union of ball structure. Further  
 117 increasing  $r$  to  $r_2$ , leads to filling of the smaller hole (fourth figure from left). This continues  
 118 till the value of  $r$  is large enough for the union of balls to fill the entire structure. During the  
 119 change in the structure of the union of balls due to increase in radius, the larger of the two  
 120 holes ‘*persists*’ for a larger range of  $r$  compared to the smaller one. Hence features that are  
 121 more prominent are expected to persist for longer periods of increasing  $r$ . This is the basic  
 122 intuition for topological persistence. The holes in this example are captured by calculating a  
 123 set of *birth-death* pairs of homology cycle classes that indicate at which value of  $r$  the class is  
 124 born and where it dies. The persistence is visualized in  $\mathbb{R}^2$  using horizontal line segments  
 125 that connect two points whose  $x$ -coordinates coincide with the birth and death values of  
 126 the homology classes. These collection of line segments, as shown in Figure 2, are called  
 127 barcodes [5]. The length of each line segment corresponds to the persistence of a cycle in the  
 128 structure. Hence, the short blue line segments correspond to the tiny holes that are formed  
 129 intermittently as the radius increases. The two long red line segments correspond to the two  
 130 holes in the structure, the largest being the bigger hole. For computational purposes, the  
 131 growing sequence of the union of balls is converted to a growing sequence of triangulations,  
 132 simplicial complexes in general, called a *filtration*. In some cases, some cycles called the  
 133 ‘*essential cycles*’ persists till the end of the filtration.

134 The rank of the persistent homology group called the persistent Betti numbers capture  
 135 the number of persistent features. For  $n$ -dimensional homology group, we denote this number  
 136 as  $\beta_n$ . This means  $\beta_0$  counts the number of connected components that arise in the filtration.  
 137 Similarly,  $\beta_1$  counts the number of *circular* holes being born as we proceed through the  
 138 filtration. It is due to this fact that all the folds in the tertiary structure, as well as the helix  
 139 and strands in the secondary structure of proteins, are recorded in our signature.

140 With the above technique, difficulties are faced as  $r$  increases. An average protein in  
 141 a database such as CATH [20] has 20,000~30,000 atoms, thus creating a point cloud of  
 142 the same size in  $\mathbb{R}^3$ . Furthermore, the initial complex including 3-simplices (or tetrahedra)  
 143 becomes quite large. On an average, this complex size grows to  $(50\sim 100)\times 10^4$  simplices of  
 144 dimension upto 4 and becomes quite difficult to process. Building a filtration using this  
 145 growing sequence of balls is thus not scalable. We attack the problem with two strategies: (1)  
 146 we only consider simplices on the boundary of the entire simplicial complex in our algorithm  
 147 and (2) compute a new filtration technique that is based on collapsing simplices rather than  
 148 growing their numbers by addition.

149 **Topological persistence**

150 Traditionally, given a point cloud, its persistence signature is calculated by building  
 151 a filtration over a simplicial complex called *Vietoris-Rips*(VR). This technique is also  
 152 used in [4] which takes the 3D position of the centers of the atoms as points in the  
 153 point cloud. Given a parameter  $\alpha$ , we can define VR complex over a point cloud  $\mathbf{P}$  as:  
 154  $\mathcal{VR}^\alpha(\mathbf{P}) = \{\sigma \mid \mathbf{d}(\mathbf{p}, \mathbf{q}) < \alpha \forall \mathbf{p}, \mathbf{q} \in \sigma\}$ .

155 As the value of  $\alpha$  increases, more edges and higher order simplices are introduced, and  
 156 a *filtration* is obtained. Finally, the persistence of this *filtration* is computed. For a better  
 157 representation of protein molecules, we take into account the radius of different atoms as  
 158 weight of the points. So, we replace each point  $p \in P$  with a tuple  $\hat{p} = (p, r_p)$  where  $r_p$  is the  
 159 radius of the atom represented by  $p$ . For the resulting weighted point cloud  $\hat{P} = \{(p, r_p)\}$ ,  
 160 we consider the weighted VR complex:  $\mathcal{VR}^\alpha(\hat{P}) = \{\sigma \mid \mathbf{d}(\mathbf{p}, \mathbf{q}) < \alpha(\mathbf{r}_p + \mathbf{r}_q) \forall \mathbf{p}, \mathbf{q} \in \sigma\}$ .

161 The VR complex is easy to implement, but its size can become a hindrance for an even a  
 162 moderate size protein molecule. Thus, instead of a VR complex, we use the (weighted) alpha  
 163 complex that is sparser and has been used to model molecules in earlier works [11].

164  
 165 **Alpha complex  $AC(\alpha)$ :** For a given value of  $\alpha$ , a simplex  $\sigma \in AC(\alpha)$  if:

- 166 ■ The circumball of  $\sigma$  is empty and has radius  $< \alpha$ , or
- 167 ■  $\sigma$  is a face of some other higher dimensional simplex in  $AC(\alpha)$ .

168  
 169 **Weighted Alpha Complex  $WAC_{\hat{P}}(\alpha)$ :** Let  $B_k(\hat{p})$  be a  $k$ -dimensional closed ball with  
 170 center  $p$ , and weight  $r_p$ . It is orthogonal or sub-orthogonal to a weighted point  $(p', r_{p'})$  iff  
 171  $\|\mathbf{p} - \mathbf{p}'\|^2 = \mathbf{r}_p^2 + \mathbf{r}_{p'}^2$  or  $\|\mathbf{p} - \mathbf{p}'\|^2 < \mathbf{r}_p^2 + \mathbf{r}_{p'}^2$ , respectively.

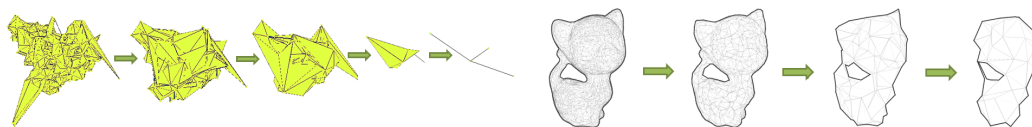
172 An *orthoball* of a  $k$ -simplex  $\sigma = \{\hat{p}_0, \dots, \hat{p}_k\}$  is a  $k$ -dimensional ball that is orthogonal to  
 173 every vertex  $p_i$ . A simplex is in the weighted alpha complex  $WAC_{\hat{P}}(\alpha)$  iff its orthoball has  
 174 radius less than  $\alpha$  and is suborthogonal to all other weighted points in  $\hat{P}$ .

## 175 2.2 Collapse-induced persistent homology from point clouds

176 The following procedure computes a topological signature for a weighted point cloud  
 177  $\hat{P} = \{p, r_p\}$  using subsamples and subsequent collapses:

- 178 1. Compute a weighted alpha complex  $\mathbb{C}^0$  on the point set  $\hat{P} = \{p, r_p\}$  using the algorithm  
 179 described in [22]. Let  $V^0$  be the vertex set of  $\mathbb{C}^0$ .
- 180 2. Compute a sequence of subsamples  $V^0 \supset V^1 \supset \dots \supset V^k$  of the initial vertex set  $V^0$  based  
 181 on the Morton Ordering as discussed later. (For every  $V^i$ , we remove every  $n^{th}$  point  
 182 in the Morton Ordering from  $V^i$  to form  $V^{i+1}$ . We choose ‘n’ based on the number of  
 183 initial points).
- 184 3. This sequence of subsets of  $V^i$  allows us to define a simplicial map between any two  
 185 adjacent subsets  $V^i$  and  $V^{i+1}$ :  $f^i(p) = \begin{cases} p & \text{if } p \in V^{i+1} \\ \operatorname{argmin}_{v \in V^{i+1}} d(p, v) & \text{otherwise} \end{cases}$
- 186 4. This vertex map  $f^i : V^i \rightarrow V^{i+1}$  in turn generates a sequence of collapsed complexes:  
 187  $\mathbb{C}^0, \mathbb{C}^1, \dots, \mathbb{C}^n$ . Each vertex map induces a simplicial map  $f^i : \mathbb{C}^{i-1} \rightarrow \mathbb{C}^i$  that associates  
 188 simplices in  $\mathbb{C}^{i-1}$  to simplices in  $\mathbb{C}^i$  (see Figure 4)
- 189 5. Compute the persistence for the simplicial maps in the sequence  $\mathbb{C}^0 \xrightarrow{f_1} \mathbb{C}^1 \xrightarrow{f_2} \dots \xrightarrow{f_k} \mathbb{C}^k$   
 190 to generate the topological signature of the point set  $\hat{P}$ .

191 In step 1 of the procedure, weighted points alone lead to disconnected weighted atoms in  
 192  $\mathbb{C}^0$  rather than capturing the actual protein structure. To sidestep this difficulty, we increase  
 193 the weights of these points based on the existence of covalent or ionic bonds in the structure.  
 194 That is, if there exists a chemical bond between two atoms (which we get from the input  
 195 .pdb file), we scale-up the weight of each point so that they are connected in the weighted  
 196 alpha complex  $WAC_{\hat{P}}(\alpha)$  (see Fig. 3). We determine a global multiplying factor  $\rho \geq 1$  for  
 197 this purpose. As mentioned earlier, we take the boundary of this weighted complex which  
 198 forms our initial simplicial complex  $\mathbb{C}^0$ .



■ **Figure 4** (a) Collapse of weighted alpha complex generated from protein structure via simplicial map. (b) Same algorithm applied to a kitten model in  $\mathbb{R}^3$

In step 2, in order to generate the sequence of subsamples, we pick vertices uniformly from the simplicial complex to be collapsed to their respective nearest neighbors. To choose a subsample that respects local density, we use a space curve generation technique called Morton Ordering [15]. The Morton curve generates a total ordering on the point set  $V^0$ . This ordering is explicitly defined by the Morton Ordering map  $M : \mathbb{Z}^N \mapsto \mathbb{Z}$  given by:

$$\mathbf{M}(\mathbf{p}) = \bigvee_{\mathbf{b}=0}^{\mathbf{B}} \bigvee_{\mathbf{i}=0}^{\mathbf{N}} \mathbf{x}_2^{\mathbf{i},\mathbf{b}} \ll \mathbf{N}(\mathbf{b} + \mathbf{1}) - (\mathbf{i} + \mathbf{1}),$$

where  $x_2^{i,b}$ :  $b^{\text{th}}$  bit value of the binary representation of the  $i^{\text{th}}$  component of  $x$ .

This map merely interleaves bits of the different components of  $p$ . Application of  $M$  to  $V^0$  yields a total ordering on our initial point set. To generate a new subset  $V^1 \subset V^0$ , we simply choose a value  $n$  such that  $1 < n \leq \|V^0\|$ . Then,  $V^{i+1}$  is taken as:

$$\mathbf{V}^{i+1} = \{\mathbf{x}_j \mid \mathbf{x}_j \in \mathbf{V}^i, \mathbf{j} \not\equiv \mathbf{0} \pmod{\mathbf{n}}\},$$

199 where  $x_j$  is the  $j^{\text{th}}$  vertex in the Morton Ordering of  $V^i$ . We choose  $n = 12$  as it has  
 200 procured good results for the datasets we experimented on (having 20,000~30,000 atoms on  
 201 an average). Following this approach, the process can be repeated to create a sequence of  
 202 subsets  $V^0 \supset V^1 \supset \dots \supset V^n, \|V^n\| \leq k$  as done in step 2 of our procedure above.

203 Finally, as described in step 3, instead of constructing the filtration by increasing the  
 204 value of  $\alpha$ , we perform a series of successive collapses starting with the initial simplicial  
 205 complex. This leads to a sequence of complexes that decreases in size instead of growing as  
 206 we proceed forward. Effectively, it generates a sequence called *tower* of simplicial complexes  
 207 where successive complexes are connected by *simplicial maps*. These maps which are the  
 208 counterpart of continuous maps for the combinatorial setting extend maps on vertices (vertex  
 209 maps) to simplices (see [16] for details). In our case, collapses of vertices generate these  
 210 simplicial maps between a simplicial complex in the tower to the one it is collapsed to.  
 211 Persistence for towers under simplicial maps can be computed by the algorithm introduced  
 212 in [7]. We use the package called Simpers that the authors have reported for the same.

213 To summarize, the algorithm generates an initial weighted alpha complex. It then  
 214 proceeds by recursively choosing vertices based on Morton Ordering to be collapsed to their  
 215 nearest neighbors resulting in vertex maps. These vertex maps are then extended to higher  
 216 order simplices (such as triangles and tetrahedra) using the simplicial map. Finally given  
 217 the simplicial map, we generate the persistence and get the barcodes for the zero and one  
 218 dimensional homology groups.

## 219 2.3 Feature vector generation

220 We discuss how we generate a feature vector given a protein structure. We take protein data  
 221 bank (\*.pdb) files as input to extract protein structures. It contains the coordinates of every  
 222 atom, their name, chemical bond with neighboring atoms and other meta-data such as helix,  
 223 sheet and peptide residue information. We introduce a weighted point for each atom in the  
 224 protein where the point is the center of the atom and its weight is the specified radius. For

225 instance, for a Nitrogen atom in the amino acid, we assign a weight equal to its covalent  
 226 radius of 71(pm). On this weighted point cloud  $\hat{p} = (p, r_p)$ , if two atoms  $\hat{p}$  and  $\hat{q}$  are involved  
 227 in a chemical bond, we increase their weights so that  $p$  and  $q$  get connected in the alpha  
 228 complex. We compute the persistence by generating the initial alpha complex and undergoing  
 229 a series of collapses as described in the previous section. For computational efficiency, we  
 230 only consider the barcodes in zero and one dimensional homology groups. Note that some of  
 231 the barcodes can have death time equal to infinity indicating an essential feature. For finite  
 232 barcodes, shorter lengths (*death - birth*) indicate noise. Elimination of these intermittent  
 233 features serves some interesting purpose as we will see in section 3. To find relatively long  
 234 barcodes, we sort them in descending order of their lengths. Let  $\{l_1, l_2, \dots, l_k\}$  be this sorted  
 235 sequence. Consider the sequence  $\{l'_1, l'_2, \dots, l'_{k-1}\}$  where  $l'_i = l_{i+1} - l_i$  and let  $l'_m$  be a  
 236 maximal element for  $1 \leq m \leq k - 1$ . All barcodes with the lengths  $[l_1..l_m]$  form part of the  
 237 feature vector. Essentially we remove all barcodes whose lengths are shorter than the largest  
 238 gap between two consecutive barcodes when sorted according to their lengths. A similar  
 239 technique used in [13] has shown improved results in image segmentation over other heuristics  
 240 and parameterizations. Since the feature vector needs to be of a fixed length for feeding into  
 241 a classifier, we compute the index  $m$  of  $l'_m$  over all protein structures and take an average.  
 242 The feature vector also includes the number of *essential* zero and one dimensional cycles.  
 243 Therefore, we have a feature vector of length  $2 \times m + 2$  :  $\{l_1^0, l_2^0, \dots, l_m^0, l_1^1, l_2^1, \dots, l_m^1, c_{\beta_0}, c_{\beta_1}\}$ .  
 244 Here  $l_i^0$  and  $l_i^1$  are the lengths of zero and one dimensional homology cycles respectively  
 245 whereas  $c_{\beta_i}$  are the total number of essential cycles in  $i$ -dimensional homology.

### 246 3 Experiments and results

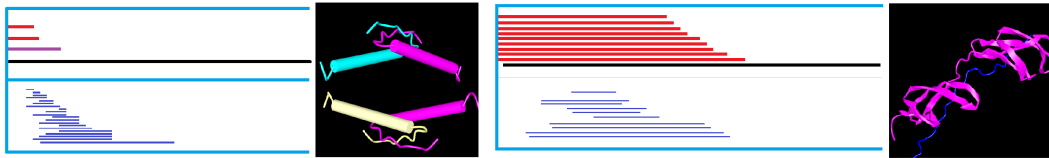
247 We perform several experiments to establish the utility of the generated topological signature.  
 248 First, we show how our feature vector captures various connections in the single strands  
 249 of secondary structures and compare them against the signatures obtained in [23]. Then  
 250 we investigate if there is a correlation between the count of such secondary structures and  
 251 our feature vector. Next, we describe the topological feature vector obtained from two  
 252 macromolecular proteins structures. We also compare the size and time needed by our  
 253 algorithm (software) over the other commonly used persistence software (as in [4]). Lastly,  
 254 we show the effectiveness of our approach in classifying protein structures using machine  
 255 learning models.

#### 256 Topological description of alpha helix and beta sheet

257 It is known that barcodes can explain the structure of an alpha helix and a beta sheet [23].  
 258 The authors in [23] use a coarse-grain(CG) representation of the protein by replacing each  
 259 amino acid molecule with a single atom. This representation removes the short barcodes



■ **Figure 5** (a) Left: Alpha helix from PCB 1C26 , Middle: Barcode of [23], Right: Our Barcode,  
 (b) Left: Beta sheet from PCB 2JOX, Middle: Barcode of [23], Right: Our Barcode. Each segment  
 of the barcodes shows  $\beta_0$ (top) and  $\beta_1$ (bottom)



■ **Figure 6** Barcode and Ribbon diagram of (Left): PDB: 1c26. (Right): PDB: 1o9a. Diagram courtesy NCBI [17]

260 corresponding to the edges and cycles of the chemical bonds inside the amino acid molecule.  
 261 We do not need this CG representation as our procedure can implicitly determine a threshold  
 262  $l_m$  and therefore delete all barcodes of length shorter than the largest gap between two  
 263 consecutive barcodes (as described in section 2.3). So, we get a barcode that describes the  
 264 essential features of the secondary structures without including noise or short lived cycles  
 265 from the amino acids. For a fair comparison, we compute our barcodes on the same alpha  
 266 helix residue as in [23] with 19 residues extracted from the protein strand having PDB ID  
 267 1C26 (see figure 5). Analogous to the barcode of [23] (as shown in the middle diagram of  
 268 figure 5a), we have 19 bars in the zero-dimensional homology for the alpha helix representing  
 269 the nineteen initial residues. These components die as edges are introduced in the weighted  
 270 alpha complex which gets them connected. For one-dimensional homology, an initial ring  
 271 with 4 residues is formed followed by additional rings resulting from the growing connections  
 272 in each amino acid. These cycles eventually die by the collapse operations in our algorithm.

273 The same process is followed for beta sheets after we extract two parallel beta sheet  
 274 strands from the protein structure with PDB ID 2JOX. The zero-dimensional homology  
 275 cycles are killed when individual disconnected amino acid residues belonging to the same  
 276 beta sheet strand are connected by edges, as represented in the top 17 barcodes (leftmost  
 277 figure of 5b). However, other than these barcodes and the longest bar corresponding to the  
 278 *essential cycle*, there is one bar in the zero-dimensional homology which is longer than the  
 279 top 17 bars. This bar represents the component which is killed by joining the closest adjacent  
 280 amino acid molecules from the two parallel beta strands. The one dimensional homology  
 281 bars are formed as more adjacent amino acid molecules are connected and killed once the  
 282 collapse operation starts. Note that the two barcodes shown in figure 5 comparing our work  
 283 with [23] are not to scale. This is because, in contrast to [23], the barcodes in our figure are  
 284 not plotted against Euclidean distance rather the step at which each insertion and collapse  
 285 operation occurs.

## 286 A caveat

287 Our aim is to compute signatures that capture discriminating structural information  
 288 useful for classifying proteins. Even though we can use our signature to describe secondary  
 289 structures, we do not want our signature to be directly correlated to the *number of* alpha  
 290 helix or beta sheet as it would mean they are redundant. We generate a  $2 \times 12$  matrix where  
 291 each cell contains the correlation value between beta-sheet(top row) and alpha-helix(bottom  
 292 row) with each individual component in the feature vector:  $\{l_1^0, l_2^0, \dots, l_m^0, l_1^1, l_2^1, \dots, l_m^1, c_{\beta 0}, c_{\beta 1}\}$ .  
 293 We use proteins in the PCB00020 record of the PCBC database to compute this matrix and  
 294 depict it by a heatmap (Fig 7). Essentially, we first generate two vectors  $v_\alpha$  and  $v_\beta$  of the  
 295 number of alpha helices and beta sheets respectively in each protein over all entries in the  
 296 database. Similarly, we produce a vector for each value in the feature vector:  $\{v_{l_1^0}, \dots, v_{c_{\beta 1}}\}$ .  
 297 Now we populate the matrix by calculating the correlation between each of these individual



Data	Dim	Size			Time (in sec)		
		VR	SimBa	Our	VR	SimBa	Our
CATH	3	-	1422	443	-	1.75	0.35
Soneko	3	324802	10188	576	32	6.77	2.05
Surv-1	150	-	$3.1 \times 10^6$	$1.09 \times 10^6$	-	$5.08 \times 10^3$	884
PrCl-I	25	-	$10.2 \times 10^6$	$0.22 \times 10^6$	-	585	141.3

■ **Table 1** Time comparison of our algorithm against SimBa [8] and VR complex.

Class	SVM			KNN		
	FB	Cang	Our	FB	Cang	Our
Architecture	91.08	89.07	92.36	86.01	86.40	86.39
Topology	92.19	94.87	96.71	91.54	94.02	96.20
Homology	93.33	94.06	94.17	90.28	91.11	93.30

■ **Table 2** Accuracy comparison with Frag-Bag and Cang

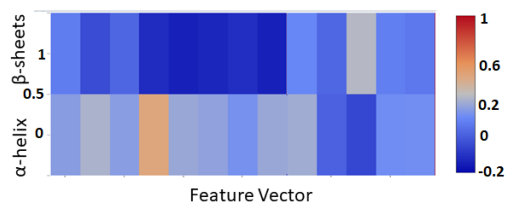
298 vectors with  $v_\alpha$  and  $v_\beta$ . For example, row 1 and column 1 of the matrix contain the cor-  
 299 relation value between the vectors  $v_{l_1^0}$  and  $v_\beta$ . The heat map color ranges from blue for  
 300 zero correlation to dark-red for complete correlation. As we can see from the figure, almost  
 301 all matrix entries have a blue tinge indicating low correlation. This shows that our feature  
 302 vector is non-redundant over the frequency of secondary structures.

303

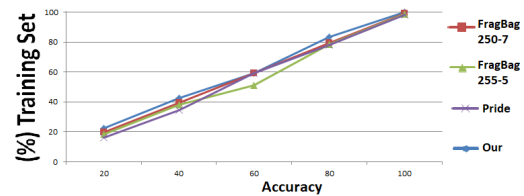
### 304 Topological Description of macromolecular structures

305 In the previous section, we use our signature to describe the secondary structures and  
 306 compare it with the work in [23]. In this section, we further show how our signature works  
 307 by describing two macromolecular protein structures that are built on multiple secondary  
 308 structures. We start by describing the tetrameric protein: 1C26. The ribbon diagram  
 309 and associated barcode after noise removal is given in figure 6 . It essentially contains four  
 310 monomers, associated pairwise to form two dimers. These two dimers, in turn, join across a  
 311 distinct parallel helix-helix interface to form the tetramer. When we build the filtration on  
 312 this protein structure, two monomers on opposite sides are killed first by connecting to their  
 313 adjacent monomers to form two distinct dimers. This is evident as there are two short bars  
 314 in the zero dimensional barcode (Fig. 6 right: shown in red). We now have two dimers, one  
 315 of which is killed when it joins with the other to form a third slightly longer non-essential  
 316 barcode (shown in purple). The second dimer lives on as the tetramer and forms an essential  
 317 barcode (shown in black). Next, if we look into the one dimensional homology (shown as blue  
 318 lines), we notice that the most notable feature for the protein is the tetramer structure which  
 319 contains a large loop when the two dimers are connected. This is evident in our 1D-barcode  
 320 as there is a distinct long bar representing the large one dimensional cycle. Note that the  
 321 birth time of this cycle in 1D corresponds with the death time of the non-essential dimer in  
 322 0D.

323 Next, we consider the protein structure 1O9A. The structure contains several antiparallel  
 324 beta-strands and is an example of a tandem beta-zipper. As we can see from the ribbon  
 325 diagram in Fig. 6, there are six beta sheets on one side and five on another, connected  
 326 together to form a fibronectin. This is evident as there are ten non essential and one essential  
 327 bar in the zero dimensional homology owing to the six beta sheets on one side and five on  
 328 the other. Each component is killed as the beta sheets join with another as the filtration  
 329 proceeds. Note that the last connected component after joining all beta sheets forms an  
 330 essential bar. Moreover, since there is no distinct cycle in the structure, we do not get any  
 331 distinct long bar in the one dimensional homology. The presence of multiple one dimensional  
 332 bars of similar size are probably due to the antiparallel beta-strands on either side which  
 333 form a ring once joined. Thus, we can see that using the same signature generation method,  
 334 we can describe secondary structures (as in the previous section) as well as macromolecular  
 335 proteins without any change in the parameter. It is therefore evident that our signature is



■ **Figure 7** Heatmap correlating secondary structure against our feature vector. Each column in the heatmap is the feature vector.



■ **Figure 8** Plot showing accuracy against varying training data size. 100(%) indicates the entire training and test data.

336 intrinsic and scale independent.

337

338

339

### Time and space comparison with VR-complex and SimBa

340 The method in [4] uses persistent homology as feature vectors for machine learning.  
 341 However, as mentioned earlier, the use of Veitoris-Rips (VR) complex leads to a size blow up  
 342 that not only increases runtime but also in most cases, causes the operating system to abort  
 343 the process due to space requirements. Results in [4] procure good results as the datasets  
 344 are of moderate size, but the same could not be reported for larger and real life protein  
 345 structures. In table 1, we show a size and time comparison of our approach with the original  
 346 feature generation technique used in [4]. We also tabulate the size and time to generate the  
 347 same feature vector in [4] using a state-of-the-art persistence algorithm called SimBa [8].  
 348 Table 2 contains a mix of protein databases and other higher dimensional datasets. As we  
 349 see in the table, our algorithm is faster even when the features in [4] are generated with SimBa.

350

351

### 3.1 Supervised learning based classification models

352 **Classification model.** For the purpose of protein classification, we train two  
 353 classifiers: an SVM model and a k-nn model on some protein databases. Once the model is  
 354 trained, we test it to find accuracy, precision, and recall. The reason behind choosing Sup-  
 355 port Vector Machine and k-nn based supervised learning technique over other sophisticated  
 356 and state-of-the-art classifiers is their basic nature. Results obtained from basic learning  
 357 techniques prove the effectiveness of the feature vectors rather than that of the classifier. We  
 358 can further improve the classification accuracy for proteins using some advanced supervised  
 359 learning or Neural Network based classifiers using our proposed features.

360

361

362

363

364

365

366

367

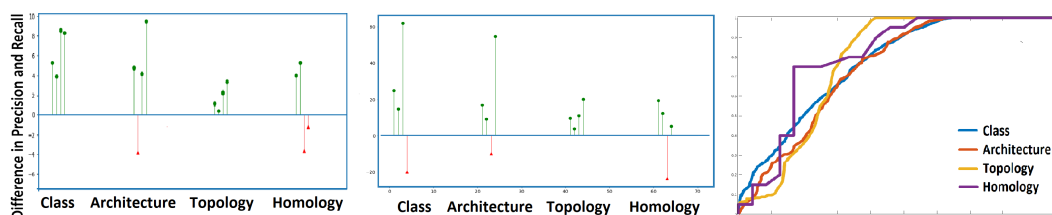
368

369

370

371

**Benchmark techniques.** In order to test the effectiveness of our protein signature,  
 we need to compare it against some of the state-of-the-art protein structure classification  
 techniques. We generate feature vectors through these techniques and train and test the  
 same classification models as before. The first technique, known as PRIDE [9], classifies  
 proteins at the Homology level in the CATH database. It represents protein structures by  
 distance distributions and compares two sets of distributions via contingency table analysis.  
 The more recent work by Budowski-Tal et al. [3], which has achieved significant improvement  
 in protein classification is used as our second benchmark technique. Their work, known  
 as FragBag mimics the bag of word representation used in natural language processing  
 and computer vision. They maintain protein fragments as a benchmark. Given a large  
 protein strand, they generate a vector which is essentially a count of the number of times



■ **Figure 9** Left:a) Difference in precision and recall from FragBag. Middle: b) Difference in precision and recall from [4]. Right: c) ROC curve for SVM classification of our algorithm

372 each fragment approximates a segment in this strand. This vector now acts as a signature  
 373 for the protein structure and that is what forms the basis for their feature vector which  
 374 we use to train and test our classifier. The protein fragment benchmark is available from  
 375 the library [12]. We choose 250 protein fragments of length 7. The third work that we  
 376 test against is the topological approach to generate a protein fingerprint [4]. However, as  
 377 we saw earlier, it is not possible to generate all the protein signatures using the original  
 378 algorithm used by the authors. Therefore, we replace the Vietoris-Rips filtration by the  
 379 state-of-the-art SimBa and generate feature vectors the same way as mentioned in their paper.

380  
 381 **Database.** The database that we use is called Protein Classification Benchmark  
 382 Collection (PCBC) [21]. It has 20 classification tasks, each derived from either the SCOP  
 383 or CATH databases, each containing about 12000 protein structures. The classification  
 384 experiment involves the group of the domain as positive set, and the positive test set is the  
 385 sub-domain group. The negative set includes the rest of the database outside the superfamily  
 386 divided into a negative test or negative training set. The result for some of the classification  
 387 tasks for the database is given in Table 3. As evident from the table, the accuracy obtained  
 388 by using our signature has a considerable improvement over the state of the art techniques.  
 389 The only classification task in which our algorithm under-performs is with the protein domain  
 390 CATH95\_Class\_5fold\_Filtered (fourth row of table 3). The class domain is randomly  
 391 sub-divided into 5 subclasses in this task. Since the class is divided randomly into subclasses,  
 392 we believe some proteins belonging to different sub-classes have generated a similar initial  
 393 complex resulting in a similar filtration and ultimately a decrease in performance.

394 The PCBC dataset, even though suitable for learning algorithms, suffers from being  
 395 skewed as the number of negative instances in any classification is much larger than the  
 396 number of positive instances, leading to probable incorrect classifications. Therefore, we test  
 397 on one of the most popular protein databases known as CATH [6]. The CATH database  
 398 contains proteins divided into different domains (C: class; A: architecture; T: topology; H:  
 399 homologous superfamily). For each domain, we get protein structures and their labels in  
 400 accordance with the sub-domain they belong to. For any classification task, we randomly  
 401 choose positive instances from one sub-domain and the same number of negative instances  
 402 sampled equally from the other sub-domains. Each such task, on average has 400 protein  
 403 structures containing approximately 30,000 atoms each. We then divide this into 80%-training  
 404 and 20%-test set. The result of classification on the CATH database averaged over several  
 405 such randomly chosen sub-domains as positive classes, are illustrated in table 2. We see  
 406 yet again that for each case, there is an improvement of about 3-4% over the benchmark  
 407 techniques.

408

	SVM				k-NN			
	Pride	Fragbag	Cang	Our	Pride	Fragbag	Cang	Our
SC Sf Fm F	90.09	93.01	93.39	95.24	89.58	87.31	89.83	91.66
CA T 5f	94.23	92.97	94.87	99.53	90.96	91.16	94.57	97.87
CA T H F	90.15	89.89	95.06	98.80	84.98	81.11	86.65	95.51
CA C 5f F	85.09	84.76	80.98	82.36	80.18	84.74	83.83	78.81
CA H Si F	98.60	95.89	98.24	99.05	95.45	91.11	79.469	97.56
CA A T F	87.56	91.58	74.58	90.95	67.47	89.00	68.90	87.00

■ **Table 3** Classification accuracy for different techniques on Protein dataset. SC: SCOP95, CA: CATH95, Sf: Superfamily, Fm: Family, F: Filtered, T: Topology, H: Homology, C: Class, 5f: 5fold, A: Architecture, Si: Similarity

### 409 3.1.1 Classification result

410 We have listed our main results in tables 2 and 3 showing the improvement in accuracy  
 411 using our method over the state-of-the-art techniques of FragBag, PRIDE and the preceding  
 412 work on topology by Cang et al. [4]. We provide further evidence of the efficiency of our  
 413 algorithm by comparing the precision and recall in figures 9a and b. In these plots, we show  
 414 the difference between the precision and recall obtained using our algorithm against that of  
 415 FragBag(9a) and Cang(9b). A green bar indicates that our algorithm performed better and  
 416 the difference is positive while a red bar suggests the opposite. This experiment is done on  
 417 the CATH database and the figure shows the precision and recall for each domain: class(C),  
 418 architecture(A), topology(T) and homology(H). Notice that, since the classification is binary,  
 419 we get two precision and two recall for every class in each domain. Thus, there are four  
 420 bars for each of C,A,T,H. Yet again, other than a few marginal cases, our algorithm largely  
 421 performs better. Finally, we calculate the ROC curve using SVM on a subset of the CATH  
 422 dataset, the result of which is shown in figure 9c. The ROC curve is a plot of the true  
 423 positive rate against false positive rate obtained by changing the input size and parameter.  
 424 This means that the further the lines are away from the diagonal, the better is the classifier.

425 For the positive test cases, we investigate further the trend of the output. We try to see  
 426 the correlation of accuracy with the change in training set size. We therefore change the  
 427 training and test set sizes by taking a fraction of the entire dataset and trace the accuracy in  
 428 each case. This is done over all the test cases shown in Table 3 and the average is shown in  
 429 Fig 8. We have plotted the output of our algorithm in blue with two instances of FragBag  
 430 with (fragment, library) sizes (5,225) and (7,250) in red and green respectively. In addition,  
 431 we have plotted the output of PRIDE as well. Ideally, the accuracy should decrease uniformly  
 432 with a decrease in training set size and we should get a straight line across the diagonal.  
 433 In this case, all the trendlines are almost close to the diagonal and hence we can say that  
 434 they are correlated. Moreover, we observe that even as the training data size decreases,  
 435 the accuracy of our algorithm remains better or comparable to the other algorithms. This  
 436 indicates that topological features work better with a lower number of samples as well.

## 437 4 Conclusion

438 We present a practical topological technique to generate signatures for protein molecules  
 439 that can be used as feature vectors for its classification. Since we investigated the descriptive  
 440 power of our signature, we believe it can be used for other purposes such as protein energy  
 441 computation, or finding protein B-factor. We believe that this signature can be extended to  
 442 other biomolecular data such as DNA or enzymes.

443 ——— **References** ———

- 444 **1** Ulrich Bauer, Michael Kerber, Jan Reininghaus, and Hubert Wagner. Phat - persistent  
445 homology algorithms toolbox. *J. Symb. Comput.*, 78(C):76–90, January 2017.
- 446 **2** Juliana Bernardes, Gerson Zaverucha, Catherine Vaquero, Alessandra Carbone, and Levitt  
447 Michael. Improvement in protein domain identification is reached by breaking consensus,  
448 with the agreement of many profiles and domain co-occurrence. 12, 07 2016.
- 449 **3** Inbal Budowski-Tal, Yuval Nov, and Rachel Kolodny. Fragbag, an accurate representa-  
450 tion of protein structure, retrieves structural neighbors from the entire pdb quickly and  
451 accurately. *PNAS*, 107(8):3481–3486, February 2010.
- 452 **4** Zixuan Cang, Lin Mu, Kedi Wu, Kristopher Opron, Kelin Xia, and Guo-Wei Wei. A  
453 topological approach for protein classification. *MBMB*, Nov 2015.
- 454 **5** Gunnar Carlsson, Afra Zomorodian, Anne Collins, and Leonidas Guibas. Persistence bar-  
455 codes for shapes. In *Proceedings of the 2004 Eurographics/ACM SIGGRAPH Symposium*  
456 *on Geometry Processing, SGP '04*, pages 124–135. ACM, 2004.
- 457 **6** Natalie Dawson, Tony E Lewis, Sayoni Das, Jonathan Lees, David Lee, Paul Ashford,  
458 Christine Orengo, and Ian Sillitoe. Cath: An expanded resource to predict protein function  
459 through structure and sequence. 45, 11 2016.
- 460 **7** Tamal K. Dey, Fengtao Fan, and Yusu Wang. Computing topological persistence for sim-  
461 plicial maps. *Symposium on Computational Geometry*, pages 345–354, june 2014.
- 462 **8** Tamal K. Dey, Dayu Shi, and Yusu Wang. Simba: An efficient tool for approximating  
463 rips-filtration persistence via simplicial batch-collapse. In *ESA*, volume 57, 2016.
- 464 **9** Zoltán Gáspári, Kristian Vlahovick, and Sándor Pongor. Efficient recognition of folds in  
465 protein 3d structures by the improved pride algorithm. *Bioinformatics*, 21(15), 2005.
- 466 **10** Edelsbrunner Herbert and John Harer. *Computational topology: an introduction*. 2010.
- 467 **11** Liang J, Edelsbrunner H, Fu P, Sudhakar PV, and Subramaniam S. Analytical shape  
468 computation of macromolecules: Ii. molecular area and volume through alpha shape. In  
469 *Proteins*, volume 33, pages 18–29, 1998.
- 470 **12** Rachel Kolodny, Patrice Koehl, Leonidas Guibas, and Michael Levitt. Small libraries of  
471 protein fragments model native protein structures accurately. *JMB*, 323, 2002.
- 472 **13** Vitaliy Kurlin. A fast persistence-based segmentation of noisy 2D clouds with provable  
473 guarantees. 83:3–12, 2015.
- 474 **14** Holm Liisa and Rosenström Päivi. Dali server: conservation mapping in 3d. *Nucleic Acids*  
475 *Research*, 38:W545–W549, 2010. doi:10.1137/070711669.
- 476 **15** G. M. Morton. A computer oriented geodetic data base; and a new technique in file  
477 sequencing. *International Business Machines Co.*, 1966.
- 478 **16** J. R. Munkres. *Elements of Algebraic Topology*, chapter 1. 1 edition, 1984.
- 479 **17** USA National Institutes of Health, 1988. URL: <https://www.ncbi.nlm.nih.gov/>.
- 480 **18** M Remmert, A Biegert, and Söding J. Hauser A. Hhblits: lightning-fast iterative protein  
481 sequence searching by hmm-hmm alignment. *Nature Methods*, 9, Dec 2011.
- 482 **19** Altschul S.F., Gish W., Miller W., Myers E.W., and Lipman D.J.n. Basic local alignment  
483 search tool. *Journal of Molecular Biology*, 215:403–410, 1990.
- 484 **20** Ian Sillitoe, Tony E Lewis, and et al. Cath: Comprehensive structural and functional  
485 annotations for genome sequences. 43, 01 2015.
- 486 **21** Paolo Sonogo, Mircea Pacurar, Somdutta Dhir, Attila Kertesz-Farkas, András Kocsor,  
487 Zoltán Gáspári, Jack A M Leunissen, and Sándor Pongor. A protein classification bench-  
488 mark collection for machine learning. 35:D232–6, 02 2007.
- 489 **22** The GUDHI Project. *GUDHI User and Reference Manual*. GUDHI Editorial Board, 2015.  
490 URL: <http://gudhi.gforge.inria.fr/doc/latest/>.
- 491 **23** Kelin Xia and Guo-Wei Wei. Persistent homology analysis of protein structure, flexibility  
492 and folding. *IJNMBE*, 30(8), 2014. URL: doi:10.1002/cnm.2655.