

# Node-Level Topological Representation Learning on Point Clouds

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Topological Data Analysis (TDA) allows us to extract powerful topological, and  
 2 higher-order information on the global shape of a data set or point cloud. Tools  
 3 like Persistent Homology or the Euler Transform give a *single* complex description  
 4 of the *global structure* of the point cloud. However, common machine learning  
 5 applications like classification require *point-level* information and features to be  
 6 available. In this paper, we bridge this gap and propose a novel method to extract  
 7 node-level topological features from complex point clouds using discrete variants  
 8 of concepts from algebraic topology and differential geometry. We verify the  
 9 effectiveness of these topological point features (TOPF) on both synthetic and  
 10 real-world data and study their robustness under noise.

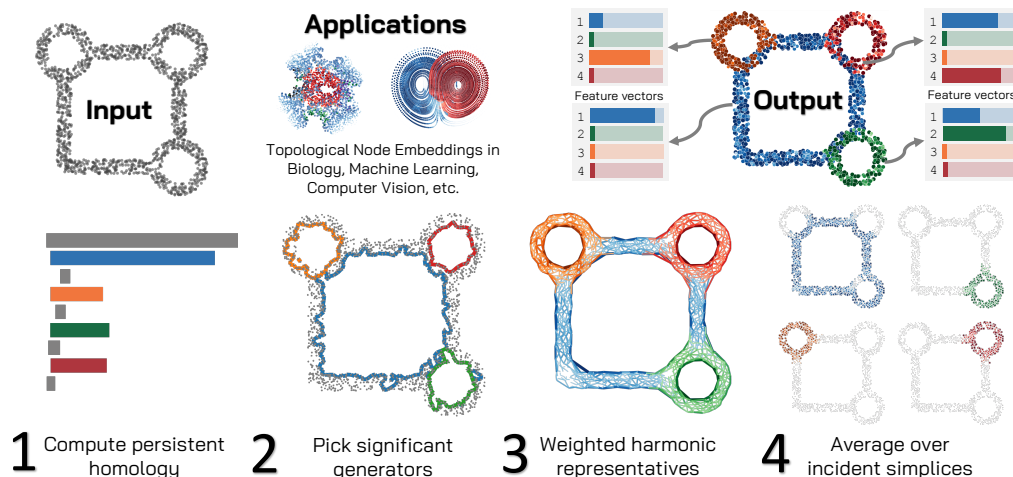


Figure 1: **Schematic of Computing Topological Point Features (TOPF).** **Input.** A point cloud  $X$  in  $n$ -dimensional space. **Step 1.** To extract global topological information, the persistent homology is computed on an  $\alpha/VR$ -filtration. The most significant topological features  $\mathcal{F}$  across all specified dimensions are selected. **Step 2.**  $k$ -homology generators associated to all features  $f_{i,k} \in \mathcal{F}$  are computed. For every feature, a simplicial complex is built at a step of the filtration where  $f_{i,k}$  is alive. **Step 3.** The homology generators are projected to the harmonic space of the simplices. **Step 4.** The vectors are normalised to obtain vectors  $\mathbf{e}_k^i$  indexed over the  $k$ -simplices. For every point  $x$  and feature  $f \in \mathcal{F}$ , we compute the mean of the entries of  $\mathbf{e}_k^i$  corresponding to simplices containing  $x$ . The output is a  $|X| \times |\mathcal{F}|$  matrix which can be used for downstream ML tasks. **Optional.** We weigh the simplicial complexes resulting in a topologically more faithful harmonic representative in **Step 3**.

# 1 Introduction

In modern machine learning [39], objects are described by feature vectors within a high-dimensional space. However, the coordinates of a single vector can often only be understood in relation to the entire data set: if the value  $x$  is small, average, large, or even an outlier depends on the remaining data. In a 1-dimensional (or low-dimensional) case this issue can be addressed simply by normalising the data points according to the global mean and standard deviation or similar procedures. We can interpret this as the most straight-forward way to construct *local* features informed by the *global* structure of the data set.

In the case where not all data dimensions are equally relevant, or contain correlated and redundant information, we can apply (sparse) PCA to project the data points to a lower dimensional space using information about the *global structure* of the point cloud [51]. For even more complex data, we may first have to learn the encoded structure itself: indeed, a typical assumption underpinning many unsupervised learning methods is the so-called “manifold hypothesis” which posits that real world data can be described well via submanifolds of  $n$ -dimensional space [36, 21]. Using eigenvectors of some Laplacian, we can then obtain a coordinate system intrinsic to the point cloud (see e.g. [47, 4, 15]). Common to all these above examples is the goal is to construct locally interpretable point-level features that encode *globally meaningful positional information* robust to local perturbations of the data. However, none of these approaches is able to represent higher-order topological information, making point clouds with these kind of structure inaccessible to point-level machine learning algorithms.

Instead of focussing on the interpretation of individual points, topological data analysis (TDA), [9], follows a different approach. TDA extracts a global description of the shape of data, which is typically considered in the form of a high-dimensional point cloud. This is done measuring topological features like persistence homology, which counts the number of generalised “holes” in the point cloud on multiple scales. Due to their flexibility and robustness these global topological features have been shown to contain relevant information in a broad range of application scenarios: In medicine, TDA has provided methods to analyse cancer progression [33]. In biology, persistent homology has been used to analyse knotted protein structures [5], and the spectrum of the Hodge Laplacian has been used for predicting protein behaviour [50].

This success of topological data analysis is a testament to the fact that relevant information is encoded in the global topological structure of point cloud data. Such higher-order topological information is however invisible to standard tools of data analysis like PCA or  $k$ -means clustering, and can also not be captured by graph models of the point cloud. We are now faced by a situation where **(i)** important parts of the global structure of a complex point cloud can only be described by the language of applied topology, however **(ii)** most standard methods to obtain positional point-level information are not sensitive to the higher-order topology of the point cloud.

**Contributions** We introduce TOPF (Figure 1), a novel method to compute node-level topological features relating individual points to global topological structures of point clouds. TOPF **(i)** *outperforms* other methods and embeddings for clustering downstream tasks on topologically structured data, returns **(ii)** *provably meaningful representations*, and is **(iii)** *robust to noise*. Finally, we introduce the topological clustering benchmark suite, the first benchmark for topological clustering.

**Related Work** The intersection of topological data analysis, topological signal processing and geometry processing has many interesting related developments in the past few years. On the side of homology and TDA, the authors in [16] and [41] use harmonic *cohomology* representatives to reparametrise point clouds based on circular coordinates. This implicitly assumes that the underlying structure of the point cloud is amenable to such a characterization. In [2, 26], the authors develop and use harmonic persistent homology for data analysis. However, among other differences their focus is not on providing robust topological point features. [24] uses the harmonic space of the Hodge Laplacians to cluster point clouds respecting topology, but is unstable against some form of noise, has no possibility for features selection across scales and is computationally far more expensive than TOPF. For a more in-depth review of related work, see Appendix A

**Organisation of the paper** In Section 2, we give an overview over the main ideas and concepts behind of TOPF. In Section 3, we describe how to compute TOPF. In Section 4, we give a theoretical

64 result guaranteeing the correctness of TOPF. Finally, we will apply TOPF on synthetic and real-world  
 65 data in Section 5. Furthermore, Appendix A contains a brief history of topology and a detailed  
 66 discussion of related work. Appendix B contains additional theoretical considerations, Appendix C  
 67 describes the novel topological clustering benchmark suite, Appendix D contains details on the  
 68 implementation and the choice of hyperparameters, Appendix E gives a detailed treatment of feature  
 69 selection, Appendix F discusses simplicial weights, and Appendix G discusses limitations in detail.

## 70 2 Main Ideas of TOPF

71 A main goal of algebraic topology is to capture the shape of spaces. Techniques from topology  
 72 describe globally meaningful structures that are indifferent to local perturbations and deformations.  
 73 This robustness of topological features to local perturbations is particularly useful for the analysis  
 74 of large-scale noisy datasets. To apply the ideas of algebraic topology in our TOPF pipeline, we  
 75 need to formalise and explain the notion of *topological features*. An important observation for  
 76 this is that high-dimensional point clouds and data may be seen as being sampled from topological  
 77 spaces — most of the time, even low-dimensional submanifolds of  $\mathbb{R}^n$  [21].

78 In this section we provide a broad overview over the most important concepts of topology and TDA  
 79 for our context, prioritising intuition over technical formalities. The interested reader is referred  
 80 to [7, 27, 49] for a complete technical account of topology and [38] for an overview over TDA.

81 **Simplicial Complexes** Spaces in topology are *continuous*, consist of *infinitely* many points, and  
 82 often live in *abstract space*. Our input data sets however consist of *finitely* many points embedded  
 83 in *real space*  $\mathbb{R}^n$ . In order to bridge this gap and open up topology to computational methods, we  
 84 need a notion of discretised topological spaces consisting of finitely many base points with finite  
 85 description length. A *Simplicial Complex* is the simplest discrete model that can still approximate  
 86 any topological space occurring in practice [43]:

87 **Definition 2.1** (Simplicial complexes). A *simplicial complex* (SC)  $\mathcal{S}$  consists of a set of vertices  $V$   
 88 and a set of finite non-empty subsets (simplices,  $S$ ) of  $V$  closed under taking non-empty subsets, such  
 89 that the union over all simplices  $\bigcup_{\sigma \in \mathcal{S}} \sigma$  is  $V$ . In the following, we will often identify  $\mathcal{S}$  with its set  
 90 of simplices  $S$  and denote by  $\mathcal{S}_k$  the set of simplices  $\sigma \in S$  with  $|\sigma| = k + 1$ , called *k-simplices*. We  
 91 say that  $\mathcal{S}$  is *n-dimensional*, where  $n$  is the largest  $k$  such that the set of  $k$ -simplices  $\mathcal{S}_k$  is non-empty.  
 92 The *k-skeleton* of SC contains the simplices of dimension at most  $k$ . If the vertices  $V$  lie in real space  
 93  $\mathbb{R}^n$ , we call the convex hull in  $\mathbb{R}^n$  of a simplex  $\sigma$  its *geometric realisation*  $|\sigma|$ . When doing this for  
 94 every simplex of  $\mathcal{S}$ , we call this the *geometric realisation of  $\mathcal{S}$* ,  $|\mathcal{S}| \subset \mathbb{R}^n$ .

95 Concretely, we can construct an  $n$ -dimensional SC  $\mathcal{S}$  in  $n + 1$  steps: First, we start with a set of  
 96 vertices  $V$  which we can identify with the 0-simplices  $\mathcal{S}_0$ . Second, we connect certain pairs of  
 97 vertices with edges, which constitute the set of 1-simplices. We can then choose to fill in some triples  
 98 of vertices which are fully connected by 1-simplices with triangles, i.e. 2-simplices. More generally,  
 99 in the  $k^{\text{th}}$  step, we can add a  $k$ -simplex for every set  $\sigma_k$  of  $k + 1$  vertices such that every  $k$ -element  
 100 subset  $\sigma_{k-1}$  of  $\sigma_k$  is already a  $(k - 1)$ -simplex.

101 **Vietoris–Rips and  $\alpha$ -complexes** We now need a way to construct a *simplicial complex* that  
 102 approximates the *topological structure* inherent in our data set  $X \subset \mathbb{R}^n$ . Such a construction will  
 103 always depend on the scale of the structures we are interested in. When looking from a very large  
 104 distance, the point cloud will appear as a singular connected blob in the otherwise empty and infinite  
 105 real space, on the other hand when we continue to zoom in, the point cloud will at some point appear  
 106 as a collection of individual points separated by empty continuous space; all interesting information  
 107 can be found in-between these two extreme scales where some vertices are joined by simplices and  
 108 others are not. Instead of having to pick a single scale, the *Vietoris–Rips (VR) filtration* and the  
 109  *$\alpha$ -filtration* take as input a point cloud and return a nested sequence of simplicial complexes indexed  
 110 by a scale parameter  $\varepsilon$  approximating the topology of the data across all possible scales.

111 **Definition 2.2** (VR complex). Given a finite point cloud  $X$  in a metric space  $(\mathcal{M}, d)$  and a non-  
 112 negative real number  $\varepsilon \in \mathbb{R}_{\geq 0}$ , the associated VR complex  $VR_\varepsilon(X)$  is given by the vertex set  $X$  and  
 113 the set of simplices  $S = \{\sigma \subset X \mid \sigma \neq \emptyset, \forall x, y \in \sigma : d(x, y) \leq \varepsilon\}$

114 Intuitively, a VR complex with parameter  $\varepsilon$  consists of all simplices  $\sigma$  where all vertices  $x \in \sigma$  have a  
 115 pair-wise distance of at most  $\varepsilon$ . For  $r \leq r'$ , we obtain the canonical inclusions  $i_{r,r'}(X) : VR_r(X) \hookrightarrow$

116  $VR_r(X)$ . The set of VR complexes on  $X$  for all possible  $r \in \mathbb{R}_{\geq 0}$  together with the inclusions then  
 117 form the *VR filtration* on  $X$ . For large point clouds, using the VR complex for computations becomes  
 118 expensive due to its large number of simplices. In contrast, the more sophisticated  $\alpha$ -complex  
 119 approximates the topology of a point cloud using far fewer simplices and thus we will make use of it.  
 120 For a complete account and definition of  $\alpha$ -complexes and our reason to use them, see Appendix B.

121 **Boundary matrices** So far, we have discussed a discretised version of topological spaces in the  
 122 form of SCs and a way to turn point clouds into a sequence of SCs indexed by a scale parameter.  
 123 However, we still need an *algebraic representation* of simplicial complexes that is capable of encoding  
 124 the structure of the SC and enables extraction of the *topological features*: The *boundary matrices*  
 125  $\mathcal{B}_k$  associated to an SC  $\mathcal{S}$  store all structural information of SC. The rows of  $\mathcal{B}_k$  are indexed by the  
 126  $k$ -simplices of  $\mathcal{S}$  and the columns are indexed by the  $(k+1)$ -simplices.

127 **Definition 2.3** (Boundary matrices). Let  $\mathcal{S}$  be a simplicial complex and  $\preceq$  a total order on its vertices  
 128  $V$ . Then, the  $i$ -th face map in dimension  $n$   $f_i^n: \mathcal{S}_n \rightarrow \mathcal{S}_{n-1}$  is given by

$$f_i^n: \{v_0, v_1, \dots, v_n\} \mapsto \{v_0, v_1, \dots, \widehat{v}_i, \dots, v_n\}$$

129 with  $v_0 \preceq v_1 \preceq \dots \preceq v_n$  and  $\widehat{v}_i$  denoting the omission of  $v_i$ . Now, the  $n$ -th *boundary operator*  
 130  $\mathcal{B}_n: \mathbb{R}[\mathcal{S}_{n+1}] \rightarrow \mathbb{R}[\mathcal{S}_n]$  with  $\mathbb{R}[\mathcal{S}_n]$  being the real vector space over the basis  $\mathcal{S}_n$  is given by

$$\mathcal{B}_n: \sigma \mapsto \sum_{i=0}^{n+1} (-1)^i f_i^{n+1}(\sigma).$$

131 When lexicographically ordering the simplex basis, we can view  $\mathcal{B}_n$  as a *matrix*. We call  $\mathbb{R}[\mathcal{S}_n]$  the  
 132 space of  $n$ -chains. Now,  $\mathcal{B}_0$  is the vertex-edge incidence matrix of the associated graph consisting of  
 133 the 0- and 1-simplices of  $\mathcal{S}$  and  $\mathcal{B}_1$  is the edge-triangle incidence matrix of  $\mathcal{S}$

134 **Betti Numbers and Persistent Homology** We now turn to the notion of  
 135 *topological features* and how to extract them. *Homology* is one of the main  
 136 algebraic invariants to capture the shape of topological spaces and SC. From  
 137 a technical point of view, the  $k$ -th homology module  $H_k(\mathcal{S})$  of an SC  $\mathcal{S}$   
 138 with boundary operators  $\mathcal{B}_k$  is defined as  $H_k(\mathcal{S}) := \ker \mathcal{B}_{k-1} / \text{Im } \mathcal{B}_k$ . The  
 139 *generator* or representative of a homology class is an element of the kernel  
 140  $\ker \mathcal{B}_{k-1}$ . In dimension 1, these are given by formal sums of 1-simplices  
 141 forming closed loops in the SC. Importantly, the rank  $\text{rk } H_k(\mathcal{S})$  is called  
 142 the  $k$ -th *Betti number*  $B_k$  of  $\mathcal{S}$ . In dimension 0,  $B_0$  counts the number of  
 143 connected components,  $B_1$  counts the number of loops around ‘holes’ of  
 144 the space,  $B_2$  counts the number of 3-dimensional voids with 2-dimensional  
 145 boundary, and so on.

146 If we are now given a filtration of simplicial complexes instead of a single  
 147 SC, we can track how the homology modules evolve as the simplicial  
 148 complex grows. The mathematical formalisation, *persistent homology*, thus  
 149 turns a point cloud via a simplicial filtration into an algebraic object summarising the topological  
 150 feature of the point cloud. For better computational performance, the computations are usually done  
 151 in one of the small finite fields  $\mathbb{Z}/p\mathbb{Z}$ . Because we will later be interested in the sign of numbers  
 152 to distinguish different simplex orientations, we will use  $\mathbb{Z}/3\mathbb{Z}$ -coefficients, with  $\mathbb{Z}/3\mathbb{Z}$  being the  
 153 smallest field being able to distinguish 1 and  $-1$ .

154 **The Hodge Laplacian and the Harmonic Space** In the previous part, we have introduced a  
 155 language to characterise the global shape of spaces and point clouds. However, we still need to find  
 156 a way to relate these *global characterisations* back to *local properties* of the point cloud. We will  
 157 do so by using ideas and concepts from differential geometry and topology: The simplicial Hodge  
 158 Laplacian is a discretisation of the Hodge–Laplace operator acting on differential forms of manifolds:

159 **Definition 2.4** (Hodge Laplacian). Given a simplicial complex  $\mathcal{S}$  with boundary operators  $\mathcal{B}_k$ , we  
 160 define the  $n$ -th Hodge Laplacian  $L_n: \mathbb{R}[\mathcal{S}_n] \rightarrow \mathbb{R}[\mathcal{S}_n]$  by setting

$$L_n := \mathcal{B}_{n-1}^\top \mathcal{B}_{n-1} + \mathcal{B}_n \mathcal{B}_n^\top.$$

161 The Hodge Laplacian gives rise to the Hodge decomposition theorem:

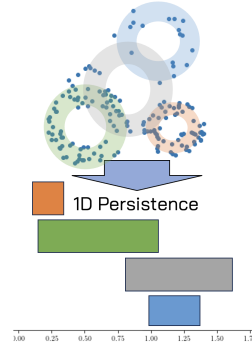


Figure 2: Sketch of Persistent Homology, [23]

149 turns a point cloud via a simplicial filtration into an algebraic object summarising the topological  
 150 feature of the point cloud. For better computational performance, the computations are usually done  
 151 in one of the small finite fields  $\mathbb{Z}/p\mathbb{Z}$ . Because we will later be interested in the sign of numbers  
 152 to distinguish different simplex orientations, we will use  $\mathbb{Z}/3\mathbb{Z}$ -coefficients, with  $\mathbb{Z}/3\mathbb{Z}$  being the  
 153 smallest field being able to distinguish 1 and  $-1$ .

154 **The Hodge Laplacian and the Harmonic Space** In the previous part, we have introduced a  
 155 language to characterise the global shape of spaces and point clouds. However, we still need to find  
 156 a way to relate these *global characterisations* back to *local properties* of the point cloud. We will  
 157 do so by using ideas and concepts from differential geometry and topology: The simplicial Hodge  
 158 Laplacian is a discretisation of the Hodge–Laplace operator acting on differential forms of manifolds:

159 **Definition 2.4** (Hodge Laplacian). Given a simplicial complex  $\mathcal{S}$  with boundary operators  $\mathcal{B}_k$ , we  
 160 define the  $n$ -th Hodge Laplacian  $L_n: \mathbb{R}[\mathcal{S}_n] \rightarrow \mathbb{R}[\mathcal{S}_n]$  by setting

$$L_n := \mathcal{B}_{n-1}^\top \mathcal{B}_{n-1} + \mathcal{B}_n \mathcal{B}_n^\top.$$

161 The Hodge Laplacian gives rise to the Hodge decomposition theorem:



---

**Algorithm 1** Topological Point Features (TOPF)
 

---

**Input:** Point cloud  $X \in \mathbb{R}^n$ , maximum homology dimension  $d \in \mathbb{N}$ , interpolation coeff.  $\lambda$ .

1. Compute persistent homology with generators in dimension  $k \leq d$ .
2. Select set of significant features  $(b_i, d_i, g_i)$  with birth, death, and generator in  $\mathbb{F}_3$  coordinates.
3. Embed  $g_i$  into real space and project into harmonic subspace of SC at step  $t = \lambda b_i + (1 - \lambda)d_i$ .
4. Normalise projections to  $e_i^k$  and compute  $F_k^i(x) := \text{avg}_{g \in \sigma} (e_i^k l(\sigma))$  for all points  $x \in X$ .

**Output:** Features of  $x \in X$

---

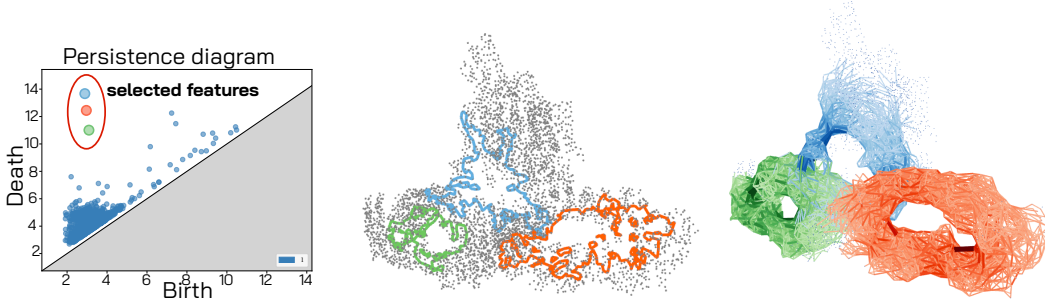


Figure 3: **TOPF pipeline applied to NALCN channelosome, a membran protein [32].** *Left:* Steps 1&2a, when computing persistent 1-homology, three classes are more prominent than the rest. *Centre:* Step 2b: The selected homology generators. *Right:* Step 3: The projections of the generators into (weighted) harmonic are now each supported on one of the three rings.

162 **Theorem 2.5** (Hodge Decomposition [34, 46, 44]). *For an SC  $\mathcal{S}$  with boundary matrices  $(\mathcal{B}_i)$  and*  
 163 *Hodge Laplacians  $(L_i)$ , we have in every dimension  $k$*

$$\mathbb{R}[\mathcal{S}_k] = \underbrace{\text{Im } \mathcal{B}_{k-1}^\top}_{\text{gradient space}} \oplus \underbrace{\ker L_k}_{\text{harmonic space}} \oplus \underbrace{\text{Im } \mathcal{B}_k}_{\text{curl space}}.$$

164 This, together with the fact that the  $k$ -th harmonic space is isomorphic to the  $k$ -th real-valued  
 165 homology group  $\ker L_k \cong H_k(\mathbb{R})$  means that we can associate a *unique harmonic representative*  
 166 to every homology class. The harmonic space encodes higher-order generalisations of smooth flow  
 167 around the holes of the simplicial complex. Intuitively, this means that for every abstract global  
 168 homology class of persistent homology from above we can now compute one unique harmonic  
 169 representative in  $\ker L_k$  that assigns every simplex a value based on how much it contributes to the  
 170 homology class. Thus, the Hodge Laplacian is a gateway between the *global topological features*  
 171 and the *local properties* of our SC. It is easy to show that the kernel of the Hodge Laplacian is the  
 172 intersection of the kernel of the boundary and the coboundary map  $\ker L_k = \ker \mathcal{B}_{n-1} \cap \ker \mathcal{B}_n^\top$ .  
 173 Because we have finite SCs we can identify the spaces of chains and cochains. This leads to another  
 174 characterisation of the harmonic space: The space of chains that are simultaneously homology and  
 175 cohomology representatives.

### 176 3 How to Compute Topological Point Features

177 In this section, we will combine the ideas and insights of the previous section to give a complete  
 178 account of how to compute Topological point features (TOPF). A pseudo-code version can be found  
 179 in Algorithm 1 and an overview in Figure 1. We start with a finite point cloud  $X \subset \mathbb{R}^n$ .

180 **Step 1: Computing the persistent homology** First, we need to determine the *most significant*  
 181 *persistent homology classes* which determine the shape of the point cloud. By doing this, we can  
 182 also extract the “interesting” scales of the data set. We will later use this to construct SCs to derive  
 183 local variants of the global homology features. Thus we first compute the persistent  $k$ -homology  
 184 modules  $P_k$  including a set of homology representatives  $R_k$  of  $X$  using an  $\alpha$ -filtration for  $n \leq 3$  and  
 185 a VR filtration for  $n > 3$ . We use  $\mathbb{Z}/3\mathbb{Z}$  coefficients to be sensitive to simplex orientations. In case we  
 186 have prior knowledge on the data set, we can choose a real number  $R \in \mathbb{R}_{>0}$  and only compute the  
 187 filtration and persistent homology connecting points up to a distance of at most  $R$ . In data sets like  
 188 protein atom coordinates, this might be useful as we have prior knowledge on what constitutes the

189 “interesting” scale, reducing computational complexity. See Figure 3 *left* for a persistent homology  
 190 diagram.

191 **Step 2: Selecting the relevant topological features** We now need to select the relevant *homology*  
 192 *classes* which carry the most important *global information*. The persistent homology  $P_k$  module in  
 193 dimension  $k$  is given to us as a list of pairs of birth and death times  $(b_i^k, d_i^k)$ . We can assume these  
 194 pairs are ordered in non-increasing order of the durations  $l_i^k = d_i^k - b_i^k$ . This list is typically very  
 195 long and consists to a large part of noisy homological features which vanish right after they appear.  
 196 In contrast, we are interested in connected components, loops, cavities, etc. that *persist* over a long  
 197 time, indicating that they are important for the shape of the point cloud. Distinguishing between the  
 198 relevant and the irrelevant features is in general difficult and may depend on additional insights on  
 199 the domain of application. In order to provide a heuristic which does not depend on any a-priori  
 200 assumptions on the number of relevant features we pick the smallest quotient  $q_i^k := l_{i+1}^k/l_i^k > 0$   
 201 as the point of cut-off  $N_k := \arg \min_i q_i^k$ . The only underlying assumption of this approach is that  
 202 the band of “relevant” features is separated from the “noisy” homological features by a drop in  
 203 persistence. If this assumption is violated, the only possible way to do meaningful feature selection  
 204 depends on application-specific domain knowledge. We found that our proposed heuristics work well  
 205 across a large scale of applications. See Figure 3 *left* and *centre* for an illustration and Appendix E  
 206 for more technical details and ways to improve and adapt the feature selection module of TOPF. We  
 207 call the chosen  $k$ -homology classes including  $k$ -homology generators in dimension  $f_k^i$ .

208 **Step 3: Projecting the features into harmonic space and normalising** In this step, we need to  
 209 relate the *global topology* extracted in the previous step to the simplices which we will use to compute  
 210 the *local* topological point feature. Every selected feature  $f_k^i$  of the previous step comes with a birth  
 211 time  $b_{i,k}$  and a death time  $d_{i,k}$ . This means that the homology class  $f_k^i$  is present in every SC of  
 212 the filtration between step  $\varepsilon = b_{i,k}$  and  $\varepsilon = d_{i,k}$  and we could choose any of the SCs for the next  
 213 step. Picking a *small*  $\varepsilon$  will lead to *fewer* simplices in the SC and thus to a very *localised* harmonic  
 214 representative. Picking a *large*  $\varepsilon$  will lead to *many* simplices in the SC and thus to a very *smooth*  
 215 and “blurry” harmonic representative with large support. Finding a middle ground between these  
 216 regimes returns optimal results. For the interpolation parameter  $\gamma \in (0, 1)$ , we will thus consider the  
 217 simplicial complex  $\mathcal{S}^{t_{i,k}}(X)$  at step  $t_{i,k} := b_{i,k}^{1-\gamma} d_{i,k}^\gamma$  for  $k > 0$  and at step  $t_{i,k} := \gamma d_{i,k}$  for  $k = 0$   
 218 of the simplicial filtration. At this point, the homology class  $f_k^i$  is still alive. We then consider the  
 219 real vector space  $\mathbb{R}[\mathcal{S}_k^{t_{i,k}}(X)]$  with formal basis consisting of the  $k$ -simplices of the SC  $\mathcal{S}^{t_{i,k}}$ . From  
 220 the persistent homology computation of the first step, we also obtain a generator of the feature  $f_k^i$ ,  
 221 consisting of a list  $\Sigma_k^i$  of simplices  $\hat{\sigma}_j \in \mathcal{S}_k^{b_{i,k}}$  and coefficients  $c_j \in \mathbb{Z}/3\mathbb{Z}$ . We need to turn this  
 222 formal sum of simplices with  $\mathbb{Z}/3\mathbb{Z}$ -coefficients into a vector in the real vector space  $\mathbb{R}[\mathcal{S}_k^{t_{i,k}}(X)]$ :  
 223 Let  $\iota: \mathbb{Z}/3\mathbb{Z} \rightarrow \mathbb{R}$  be the map induced by the canonical inclusion of  $\{-1, 0, 1\} \hookrightarrow \mathbb{R}$ . We can now define  
 224 an indicator vector  $e_k^i \in \mathbb{R}[\mathcal{S}_k^{t_{i,k}}(X)]$  associated to the feature  $f_k^i$ .

$$e_k^i(\sigma) := \begin{cases} \iota(c_j) & \exists \hat{\sigma}_j \in \Sigma_k^i : \sigma = \hat{\sigma}_j \\ 0 & \text{else} \end{cases}.$$

225 While this homology representative lives in a real vector space, it is not unique, has a small support,  
 226 and can differ largely between close simplices. All of these problems can be solved by projecting  
 227 the homology representative to the harmonic subspace  $\ker L_k$  of  $\mathbb{R}[\mathcal{S}_k^{t_{i,k}}(X)]$ . Rather than directly  
 228 projecting  $e_k^i$  to the harmonic subspace, we make use of the Hodge decomposition theorem (The-  
 229 orem 2.5) which allows us to compute the gradient and curl projections solving computationally  
 230 efficient least square problems:

$$e_{k,\text{grad}}^i := \mathcal{B}_{k-1}^\top \arg \min_{x \in \mathbb{R}[\mathcal{S}_{k-1}]} \|e_k^i - \mathcal{B}_{k-1}^\top x\|_2^2 \quad \text{and} \quad e_{k,\text{curl}}^i := \mathcal{B}_k \arg \min_{x \in \mathbb{R}[\mathcal{S}_{k+1}]} \|e_k^i - e_{k,\text{grad}}^i - \mathcal{B}_k x\|_2^2$$

231 and then setting  $\hat{e}_k^i := e_k^i - e_{k,\text{grad}}^i - e_{k,\text{curl}}^i$ . (Cf. Figure 3 *right* for a visualisation.) Because homology  
 232 representatives are gradient-free, we only need to consider the projection of  $e_k^i$  into the curl space.

233 **Step 4: Processing and aggregation at a point level** In the previous step, we have computed  
 234 a set of simplex-valued harmonic representatives of homology classes. However, these simplices  
 235 likely have no real-world meaning and the underlying simplicial complexes differ depending  
 236 on the birth and death times of the homology classes. Hence in this step, we will collect the

237 features on the point-level after performing some necessary preprocessing. Given a simplex-valued  
 238 vector  $\hat{e}_k^i$  and a hyperparameter  $\delta$ , we now construct  $e_k^i: \mathcal{S}_k^{t_i,k}(X) \rightarrow [0, 1]$  by setting  $e_k^i: \sigma \mapsto \in$   
 239  $\{|\hat{e}_k^i(\sigma)|/(\delta \max_{\sigma' \in \mathcal{S}_k^{t_i,k}(X)} |\hat{e}_k^i(\sigma')|), 1\}$  such that  $\hat{e}_k^i$  is normalised to  $[0, 1]$ , the values of  $[0, \delta]$  are  
 240 mapped linearly to  $[0, 1]$  and everything above is sent to 1. We found empirically that a thresholding  
 241 parameter of  $\delta = 0.07$  works best across at the range of applications considered below. However,  
 242 TOPF is not sensitive to small changes to  $\delta$  because entries of  $\hat{e}_k^i$  are concentrated around 0.

243 For every feature  $f_k^i$  in dimension  $k$  with processed simplicial feature vector  $e_k^i$  and simplicial  
 244 complex  $\mathcal{S}^{t_i,k}$ , we define the point-level feature map  $F_i^k: X \rightarrow \mathbb{R}$  mapping from the initial point  
 245 cloud  $X$  to  $\mathbb{R}$  by setting

$$F_i^k: v \mapsto \frac{\sum_{\sigma_k \in \mathcal{S}_k^{t_i,k}: v \in \sigma_k} e_k^i(\sigma_k)}{\max(1, |\{\sigma_k \in \mathcal{S}_k^{t_i,k}: v \in \sigma_k\}|)}.$$

246 For every point  $v$ , we can thus view the vector  $(F_i^k(v): f_k^i \in \mathcal{F})$  as a feature vector for  $v$ . We call  
 247 this collection of features *Topological Point Features* (TOPF). (Cf. Figure 4 for an example).

248 **Choosing Simplicial Weights** By default, the simplicial complexes of  $\alpha$ - and VR filtrations are  
 249 unweighted. However, the weights determine the entries of the harmonic representatives, increasing  
 250 and decreasing the influence of certain simplices and parts of the simplicial complex. We can use this  
 251 observation to increase the robustness of TOPF against the influence of heterogeneous point cloud  
 252 structure, which is present in virtually all real-world data sets. For a complete technical account of  
 253 how and why we do this, see Appendix F.

## 254 4 Theoretical guarantees

255 In this section, we prove the relationship between TOPF and actual topological structure in datasets:

256 **Theorem 4.1** (Topological Point Features of Spheres). *Let  $X$  consist of at least  $(n + 2)$  points*  
 257 *(denoted by  $S$ ) sampled uniformly at random from a unit  $n$ -sphere in  $\mathbb{R}^{n+1}$  and an arbitrary number*  
 258 *of points with distance of at least 2 to  $S$ . When we now consider the  $\alpha$ -filtration on this point*  
 259 *cloud, with probability 1 we have that (i) there exists an  $n$ -th persistent homology class generated*  
 260 *by the 2-simplices on the convex hull of  $S$ , (ii) the associated unweighted harmonic homology*  
 261 *representative takes values in  $\{0, \pm 1\}$  where the 2-simplices on the boundary of the convex hull are*  
 262 *assigned a value of  $\pm 1$ , and (iii) the support of the associated topological point feature (TOPF)  $\mathcal{F}_n^*$*   
 263 *is precisely  $S$ :  $\text{supp}(\mathcal{F}_n^*) = S$ . (iv) The same holds true for point clouds sampled from multiple*  
 264  *$n_i$ -spheres if the above conditions are met on each individual sphere.*

265 We will give a proof of this theorem in Appendix B.

266 *Remark 4.2.* In practice, datasets with topological structure consist in a majority of cases of points  
 267 sampled with noise from deformed  $n$ -spheres. The theorem thus guarantees that TOPF will recover  
 268 these structural information in an idealised setting. Experimental evidence suggests that this holds  
 269 under the addition of noise as well which is plausible as harmonic persistent homology is robust  
 270 against some noise [2].

## 271 5 Experiments

272 In this section, we conduct experiments on real world and synthetic data, compare the clustering  
 273 results with clustering by TPCC, other classical clustering algorithms, and other point features, and  
 274 demonstrate the robustness of TOPF against noise.

275 **Topological Point Cloud Clustering Benchmark** We introduce the topological clustering bench-  
 276 mark suite (Appendix C) and report running times and the accuracies of clustering based on TOPF  
 277 and other methods and point embeddings, see Table 1. We see that TOPF *outperforms* all classical  
 278 clustering algorithms on all but one dataset by a wide margin. We also see that TOPF closely matches  
 279 the performance of the only other higher-order topological clustering algorithm, TPCC on two datasets  
 280 with clear topological features, whereas TOPF *outperforms* TPCC on datasets with more complex  
 281 structure. In addition, TOPF has a consistently lower running time with better scaling for the more

Table 1: **Quantitative performance comparison of clustering with TOPF and other features/clustering algorithms.** Four 2D and three 3D data sets of the topological clustering benchmark suite (Appendix C, cf. Figure 6 for ground truth labels and Figure 7 for clustering results of TOPF). We ran each algorithm 20 times and list the mean adjusted rand index (ARI) with standard deviation  $\sigma$  and mean running time. We omit  $\sigma$  for algorithms with  $\sigma = 0$  on every dataset. TOPF consistently outperforms or almost matches the other algorithms while having significantly better run time than the second best performing algorithm TPCC. Spectral Clustering (SC), DBSCAN, and Agglomerative Clustering (AgC) are standard clustering algorithms, ToMATo is a topological clustering algorithm [11], Geo clusters using 12-dimensional point geometric features extracted by pgeof and the normal point coordinates, whereas node2vec [25] produces node embeddings on a  $k$ -nearest neighbour graph built upon an affinity matrix. We highlight all ARI scores within  $\pm 0.05$  of the best ARI score.

		TOPF (ours)	TPCC	SC	DBSCAN	AgC	ToMATo	Geo	node2vec
4spheres	ARI	<b>0.81</b>	0.52±0.17	0.37	0.00	0.45	0.32	0.20	0.00±0.00
	time (s)	14.5	23.3	0.2	0.0	0.0	0.0	0.2	48.4
Ellipses	ARI	<b>0.95</b>	0.47±0.04	0.25	0.19	0.52	0.29	0.81	0.02±0.00
	time (s)	12.7	14.4	0.1	0.0	0.0	0.0	0.1	11.2
Spheres+Grid	ARI	0.70	0.39±0.04	<b>0.90</b>	<b>0.92</b>	<b>0.89</b>	0.82	0.41	0.01±0.00
	time (s)	13.0	28.5	0.5	0.0	0.0	0.0	0.3	63.8
Halved Circle	ARI	<b>0.71</b>	0.18±0.12	0.24	0.00	0.20	0.16	0.08	0.00±0.01
	time (s)	12.2	14.3	0.1	0.0	0.0	0.0	0.1	18.2
2Spheres2Circles	ARI	<b>0.94</b>	<b>0.97±0.01</b>	0.70	0.00	0.51	0.87	0.12	0.00±0.00
	time (s)	38.9	1662.2	1.6	0.0	0.3	0.0	0.9	348.6
SphereinCircle	ARI	<b>0.97</b>	<b>0.98±0.0</b>	0.34	0.00	0.29	0.06	0.69	0.13±0.03
	time (s)	14.5	8.0	0.0	0.0	0.0	0.0	0.08	20.1
Spaceship	ARI	<b>0.92</b>	0.56±0.03	0.28	0.26	0.47	0.30	<b>0.87</b>	0.07±0.00
	time (s)	16.3	341.8	16.7	0.0	0.0	0.0	0.2	49.8
<b>mean</b>	ARI	<b>0.86</b>	0.58	0.44	0.16	0.48	0.40	0.45	0.03
	time (s)	17.5	298.9	0.4	0.0	0.0	0.0	0.3	80.0

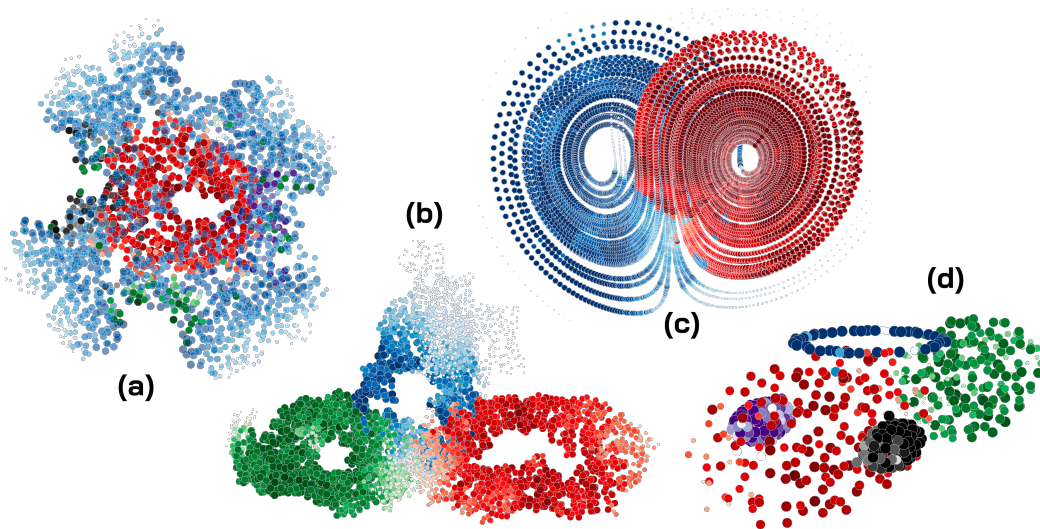


Figure 4: **TOPF on 3D real-world and synthetic point clouds.** For every point, we highlight the largest corresponding topological feature, where colour stands for the different features and saturation for the value of the feature. (a): Atoms of mutated Cys123 of E. coli [29]. We added auxiliary points on the convex hull and considered 2-homology, to detect the protein pockets which are crucial for protein-environment interactions (Cf. [40]). (b): Atoms of NALCN Channelosome [32] display three distinct loops. (c): Points sampled in the state space of a Lorentz attractor. The two features correspond to the two lobes of the attractor. (d): Point cloud spaceship of our newly introduced topological clustering benchmark suite (See Appendix C).

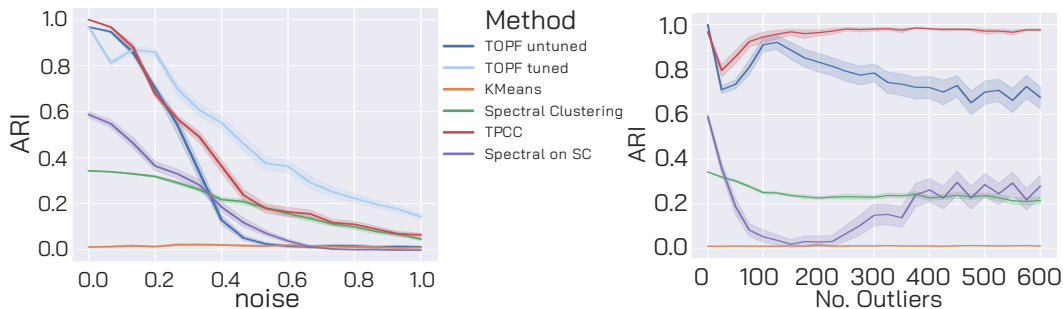


Figure 5: **Performance of Clustering based on TOPF features in increasing noise/outlier levels with 95% CI.** *Left:* We add i.i.d. Gaussian noise to every point with standard deviation indicated by the noise parameter. We see that even when compared with TPCC on a data set specifically crafted for TPCC, TOPF requires significantly less information and delivers almost equal performance. When tuned for datasets with a high noise level, the TOPF even outperform TPCC and drastically outperform all classical clustering algorithms. *Right:* We add outliers with the same standard deviation as the point cloud to the data set. We then measure the adjusted rand index obtained restricted on the original points. We see that even when compared with TPCC on a data set specifically crafted for TPCC, TOPF requires significantly less information and delivers matching to superior performance, significantly outperforming all other classical clustering algorithms.

282 complex datasets, while also not requiring prior knowledge on the best topological scale. As for the  
 283 other point embeddings, Node2Vec is not able to capture any meaningful topological information,  
 284 whereas the performance of clustering using geometric features depends on the data set.

285 **Feature Generation** In Figure 4, we show qualitatively that TOPF constructs meaningful topological  
 286 features on data sets from Biology and Physics, and synthetic data, corresponding to for example  
 287 rings and pockets in proteins or trajectories around different attractors in dynamical systems. (For  
 288 individual heatmaps see Figure 8)

289 **Robustness against noise** We have evaluated the robustness of TOPF against Gaussian noise on  
 290 the dataset introduced in [24] and compared the results against TPCC, Spectral Clustering, Graph  
 291 Spectral Clustering on the graph constructed by TPCC, and against  $k$ -means in Figure 5 *Left*. We have  
 292 also analysed the robustness of TOPF against the addition of outliers in Figure 5 *Right*. We see that  
 293 TOPF performs well in both cases, underlining our claim of robustness.

## 294 6 Discussion

295 **Limitations** TOPF can — by design — only produce meaningful output on point clouds with a  
 296 *topological structure* quantifiable by persistent homology. In practice it is thus desirable to combine  
 297 TOPF with some geometric or other point-level feature extractor. As TOPF relies on the computation of  
 298 persistent homology, its runtime increases on very large point clouds, especially in higher dimensions  
 299 where  $\alpha$ -filtrations are computationally infeasible. However, subsampling, either randomly or using  
 300 landmarks, usually preserves relevant topological features while improving run time [41]. Finally,  
 301 selection of the relevant features is a very hard problem. While our proposed heuristics work well  
 302 across a variety of domains and application scenarios, only domain- and problem-specific knowledge  
 303 makes correct feature selection feasible.

304 **Future Work** The integration of higher-order TOPF features into ML pipelines that require point-  
 305 level features potentially leads to many new interesting insights across the domains of biology, drug  
 306 design, graph learning and computer vision. Furthermore, efficient computation of simplicial weights  
 307 leading to the provably most faithful topological point features is an exciting open problem.

308 **Conclusion** We introduced point-level features TOPF founded on algebraic topology relating global  
 309 structural features to local information. We gave theoretical guarantees for the correctness of their  
 310 construction and evaluated them quantitatively and qualitatively on synthetic and real-world data sets.  
 311 Finally, we introduced the novel topological clustering benchmark suite and showed that clustering  
 312 using TOPF outperforms other available clustering methods and features extractors.

## References

- 313 [1] Michael Atiyah. *K-theory*. CRC press, 1989.
- 314 [2] Saugata Basu and Nathanael Cox. Harmonic persistent homology. In *2021 IEEE 62nd Annual*  
315 *Symposium on Foundations of Computer Science (FOCS)*, pages 1112–1123. IEEE, 2022.
- 316 [3] Ulrich Bauer. Ripser: efficient computation of vietoris–rips persistence barcodes. *Journal of*  
317 *Applied and Computational Topology*, 5(3):391–423, 2021.
- 318 [4] Mikhail Belkin and Partha Niyogi. Laplacian eigenmaps for dimensionality reduction and data  
319 representation. *Neural computation*, 15(6):1373–1396, 2003.
- 320 [5] Katherine Benjamin, Lamisah Mukta, Gabriel Moryoussef, Christopher Uren, Heather A  
321 Harrington, Ulrike Tillmann, and Agnese Barbensi. Homology of homologous knotted proteins.  
322 *Journal of the Royal Society Interface*, 20(201):20220727, 2023.
- 323 [6] A. K. Bousfield. The localization of spaces with respect to homology. *Topology*, 14(2):133–150,  
324 1975.
- 325 [7] G.E. Bredon, J.H. Ewing, F.W. Gehring, and P.R. Halmos. *Topology and Geometry*. Graduate  
326 Texts in Mathematics. Springer, New York, 1993.
- 327 [8] Peter Bubenik et al. Statistical topological data analysis using persistence landscapes. *J. Mach.*  
328 *Learn. Res.*, 16(1):77–102, 2015.
- 329 [9] Gunnar Carlsson and Mikael Vejdemo-Johansson. *Topological Data Analysis with Applications*.  
330 Cambridge University Press, 2021.
- 331 [10] Charu Chaudhry, Arthur L Horwich, Axel T Brunger, and Paul D Adams. Exploring the struc-  
332 tural dynamics of the e. coli chaperonin groel using translation-libration-screw crystallographic  
333 refinement of intermediate states. *Journal of molecular biology*, 342(1):229–245, 2004.
- 334 [11] Frédéric Chazal, Leonidas J. Guibas, Steve Y. Oudot, and Primoz Skraba. Persistence-based  
335 clustering in riemannian manifolds. *J. ACM*, 60(6), nov 2013.
- 336 [12] Frédéric Chazal and Bertrand Michel. An introduction to topological data analysis: fundamental  
337 and practical aspects for data scientists. *Frontiers in artificial intelligence*, 4:108, 2021.
- 338 [13] Yu-Chia Chen and Marina Meilă. The decomposition of the higher-order homology embedding  
339 constructed from the k-laplacian. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and  
340 J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34,  
341 pages 15695–15709. Curran Associates, Inc., 2021.
- 342 [14] Yu-Chia Chen, Marina Meilă, and Ioannis G Kevrekidis. Helmholtzian eigenmap: Topological  
343 feature discovery & edge flow learning from point cloud data. *arXiv preprint arXiv:2103.07626*,  
344 2021.
- 345 [15] Ronald R Coifman and Stéphane Lafon. Diffusion maps. *Applied and computational harmonic*  
346 *analysis*, 21(1):5–30, 2006.
- 347 [16] Vin De Silva and Mikael Vejdemo-Johansson. Persistent cohomology and circular coordinates.  
348 In *Proceedings of the twenty-fifth annual symposium on Computational geometry*, pages 227–  
349 236, 2009.
- 350 [17] Richard Dedekind. *Was sind und was sollen die Zahlen?* Verlag Friedrich Vieweg und Sohn,  
351 Braunschweig, 1888.
- 352 [18] Boris Delaunay et al. Sur la sphere vide. *Izv. Akad. Nauk SSSR, Otdelenie Matematicheskii i*  
353 *Estestvennyka Nauk*, 7(793-800):1–2, 1934.
- 354 [19] Stefania Ebli and Gard Spreemann. A notion of harmonic clustering in simplicial complexes.  
355 In *2019 18th IEEE International Conference On Machine Learning And Applications (ICMLA)*,  
356 pages 1083–1090, 2019.
- 357

- 358 [20] Samuel Eilenberg and Saunders MacLane. General theory of natural equivalences. *Transactions*  
359 *of the American Mathematical Society*, 58:231–294, 1945.
- 360 [21] Charles Fefferman, Sanjoy Mitter, and Hariharan Narayanan. Testing the manifold hypothesis.  
361 *Journal of the American Mathematical Society*, 29(4):983–1049, Oct 2016.
- 362 [22] David Chin-Lung Fong and Michael Saunders. Lsmr: An iterative algorithm for sparse least-  
363 squares problems. *SIAM Journal on Scientific Computing*, 33(5):2950–2971, 2011.
- 364 [23] Vincent P. Grande and Michael T Schaub. Non-isotropic persistent homology: Leveraging the  
365 metric dependency of ph. In *Learning on Graphs Conference*, pages 17–1. PMLR, 2023.
- 366 [24] Vincent P. Grande and Michael T. Schaub. Topological point cloud clustering. In *Proceedings*  
367 *of the 40th International Conference on Machine Learning, ICML’23*, 2023.
- 368 [25] Aditya Grover and Jure Leskovec. node2vec: Scalable feature learning for networks. In  
369 *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and*  
370 *data mining*, pages 855–864, 2016.
- 371 [26] Davide Gurnari, Aldo Guzmán-Sáenz, Filippo Utro, Aritra Bose, Saugata Basu, and Laxmi  
372 Parida. Probing omics data via harmonic persistent homology. *arXiv preprint arXiv:2311.06357*,  
373 2023.
- 374 [27] Allen Hatcher. *Algebraic Topology*. Cambridge University Press, Cambridge, 2002.
- 375 [28] Felix Hausdorff. Grundzüge einer theorie der geordneten mengen. *Mathematische Annalen*,  
376 65:435–505, 1908.
- 377 [29] Esther Hidber, Edward R Brownie, Koto Hayakawa, and Marie E Fraser. Participation of  
378 cys123 $\alpha$  of escherichia coli succinyl-coa synthetase in catalysis. *Acta Crystallographica*  
379 *Section D: Biological Crystallography*, 63(8):876–884, 2007.
- 380 [30] David Hilbert. *Grundlagen der Geometrie*. Wissenschaft und Hypothese. B. G. Teubner,  
381 Leipzig, 1899.
- 382 [31] Sze-tsen Hu. *Homotopy theory*. Academic press, 1959.
- 383 [32] Marc Kschonsak, Han Chow Chua, Claudia Weidling, Nouridine Chakouri, Cameron L. Noland,  
384 Katharina Schott, Timothy Chang, Christine Tam, Nidhi Patel, Christopher P. Arthur, Alexander  
385 Leitner, Manu Ben-Johny, Claudio Ciferri, Stephan Alexander Pless, and Jian Payandeh.  
386 Structural architecture of the human nalcn channelosome. *Nature*, 603(7899):180–186, Mar  
387 2022.
- 388 [33] Peter Lawson, Andrew B Sholl, J Quincy Brown, Brittany Terese Fasy, and Carola Wenk.  
389 Persistent homology for the quantitative evaluation of architectural features in prostate cancer  
390 histology. *Scientific reports*, 9(1):1139, 2019.
- 391 [34] Lek-Heng Lim. Hodge laplacians on graphs. *SIAM Review*, 62(3):685–715, 2020.
- 392 [35] Jacob Lurie. Stable infinity categories. *arXiv preprint math/0608228*, 2006.
- 393 [36] Yunqian Ma and Yun Fu. *Manifold learning theory and applications*, volume 434. CRC press  
394 Boca Raton, 2012.
- 395 [37] Facundo Mémoli, Zhengchao Wan, and Yusu Wang. Persistent laplacians: Properties, algorithms  
396 and implications. *SIAM Journal on Mathematics of Data Science*, 4(2):858–884, 2022.
- 397 [38] Elizabeth Munch. A user’s guide to topological data analysis. *Journal of Learning Analytics*,  
398 4(2):47–61, 2017.
- 399 [39] Kevin P. Murphy. *Probabilistic Machine Learning: An introduction*. MIT Press, 2022.
- 400 [40] Haruhisa Oda, Mayuko Kida, Yoichi Nakata, and Hiroki Kurihara. Novel definition and quantita-  
401 tive analysis of branch structure with topological data analysis. *arXiv preprint arXiv:2402.07436*,  
402 2024.



- 403 [41] Jose A. Perea. Sparse circular coordinates via principal  $\mathbb{Z}$ -bundles. In Nils A. Baas, Gunnar E.  
404 Carlsson, Gereon Quick, Markus Szymik, and Marius Thaule, editors, *Topological Data*  
405 *Analysis*, pages 435–458, Cham, 2020. Springer International Publishing.
- 406 [42] Henri Poincaré. Analysis situs. *J. de l'Ecole Poly.*, 1, 1895.
- 407 [43] Daniel G. Quillen. *Homotopical Algebra*, volume 43 of *Lecture Notes in Mathematics*. Springer,  
408 Berlin, 1967.
- 409 [44] T Mitchell Roddenberry, Nicholas Glaze, and Santiago Segarra. Principled simplicial neural  
410 networks for trajectory prediction. In *International Conference on Machine Learning*, pages  
411 9020–9029. PMLR, 2021.
- 412 [45] Michael T Schaub, Austin R Benson, Paul Horn, Gabor Lippner, and Ali Jadbabaie. Random  
413 walks on simplicial complexes and the normalized hodge 1-laplacian. *SIAM Review*, 62(2):353–  
414 391, 2020.
- 415 [46] Michael T. Schaub, Yu Zhu, Jean-Baptiste Seby, T. Mitchell Roddenberry, and Santiago Segarra.  
416 Signal processing on higher-order networks: Livin’ on the edge... and beyond. *Signal Processing*,  
417 187:108149, 2021.
- 418 [47] Jianbo Shi and Jitendra Malik. Normalized cuts and image segmentation. *IEEE Transactions*  
419 *on pattern analysis and machine intelligence*, 22(8):888–905, 2000.
- 420 [48] The GUDHI Project. *GUDHI User and Reference Manual*. GUDHI Editorial Board, 2015.
- 421 [49] Tammo tom Dieck. *Algebraic topology*, volume 8. European Mathematical Society, Zürich,  
422 2008.
- 423 [50] JunJie Wee, Jiahui Chen, Kelin Xia, and Guo-Wei Wei. Integration of persistent laplacian and  
424 pre-trained transformer for protein solubility changes upon mutation. *Computers in Biology*  
425 *and Medicine*, page 107918, 2024.
- 426 [51] Hui Zou, Trevor Hastie, and Robert Tibshirani. Sparse principal component analysis. *Journal*  
427 *of computational and graphical statistics*, 15(2):265–286, 2006.
- 428 [52] Matija Čufar. Ripserer.jl: flexible and efficient persistent homology computation in julia.  
429 *Journal of Open Source Software*, 5(54):2614, 2020.

## 430 A Extended Background

431 **A brief history of topology and machine learning** Algebraic topology is a discipline of Mathe-  
432 matics dating back roughly to the late 19<sup>th</sup> century [42]. Starting with Henri Poincaré and continuing  
433 in the early 20<sup>th</sup> century, the mathematical community became interested in developing a framework  
434 to capture the global shapes of manifolds and topological spaces in concise algebraic terms. This de-  
435 velopment was partly made possible by the push towards a formalisation of mathematics and analysis,  
436 in particular, which took place inside the mathematical community in the 1800’s and early 1900’s (e.g.  
437 [17, 30, 28]). The axiomatisation of analysis in the early 20<sup>th</sup> century is an important result of this  
438 process. These abstract ideas made it possible for Topologists to talk about the now common notions  
439 of Euler characteristics, Betti number, simplicial homology of manifolds, topological spaces, and  
440 simplicial and CW complexes. Over the course of the last 100 years, branching into many sub-areas  
441 like low-dimensional topology, differential topology, K-theory or homotopy theory [1, 31], algebraic  
442 topology has resolved many of the important questions and provides a comprehensive tool-box for  
443 the study of topological spaces. These achievements were tied to an abstraction and generalisation of  
444 concepts: topological spaces turned into spectra, diffeomorphism to homotopy equivalences and later  
445 weak equivalences, and Topologists turned to category theory [20], model categories [6] and recently  
446  $\infty$ -categories [35] as the language of choice.

447 The 21<sup>st</sup> century saw the advent and rise of topological data analysis (TDA, [8, 12]). In short,  
448 mathematicians realised that the same notions of shape and topology that their predecessors carefully  
449 defined a century earlier were now characterising the difference between healthy and unhealthy  
450 tissue, between normal and abnormal behaviour protein behaviour, or more general between different  
451 categories in their complex data sets.

452 **Related Work** The intersection of topological data analysis, topological signal processing and  
453 geometry processing has many interesting related developments in the past few years. On the side  
454 of homology and TDA, the authors in [16] and [41] use harmonic *co*homology representatives to  
455 reparametrise point clouds based on circular coordinates. This implicitly assumes that the underlying  
456 structure of the point cloud is amenable to such a characterization. Although circular coordinates are  
457 orthogonal to the core goal of TOPF, the approaches share many key ideas and insights. In [2, 26],  
458 the authors develop and use harmonic persistent homology and provide a way to pool features to  
459 the point-level. However, their focus is not on providing robust topological point features and their  
460 approach includes no tunable homology feature selection across dimensions, no support for weighted  
461 simplicial complexes, and they only construct the simplicial complex at birth. In their paper on  
462 topological mode analysis, [11] use persistent homology to cluster point clouds. However, they only  
463 consider 0-dimensional homology to base the clustering on densities and there is no clear way to  
464 generalise this to higher dimensions.

465 On the more geometric-centred side, [19] already provide a notion of harmonic clustering on simplices,  
466 [13, 14] analyse the notion of geometry and topology encoded in the Hodge Laplacian and its relation  
467 to homology decompositions, [45] study the normalised and weighted Hodge Laplacian in the context  
468 of random walks, and [24] use the harmonic space of the Hodge Laplacians to cluster point clouds  
469 respecting topology. Finally, a persistent variant of the Hodge Laplacian is used to study filtrations of  
470 simplicial complexes [37].

471 In [24], the authors have introduced TPCC, the first method to cluster a point cloud based on the  
472 higher-order topological features encoded in the data set. However, TPCC is **(i)** computationally  
473 expensive due to extensive eigenvector computations, **(ii)** depending on high-dimensional subspace  
474 clustering algorithms, which are prone to instabilities and errors, **(iii)** sensitive to the correct choice  
475 of hyperparameters, **(iv)** requiring the topological true features and noise to occur in different steps  
476 of the simplicial filtration, and it **(v)** solely focussed on clustering the points rather than extracting  
477 relevant node-level features. This paper solves all the above by completely revamping the TPCC  
478 pipeline, introducing several new ideas from applied algebraic topology and differential geometry.  
479 The core insight is: When you have the time to compute persistent homology with generators on a  
480 data set, you get the topological node features with similar computational effort.

## 481 B Theoretical Considerations

482 **More details on VR and  $\alpha$ -filtrations** Vietoris–Rips complexes are easy to define, approximate  
483 the topological properties of a point cloud across all scales and computationally easy to implement.  
484 However for moderately large  $r$ , the associated VR complex contains a large number of simplices —  
485 up to  $\binom{|X|}{n}$   $n$ -simplices for large enough  $r$  — leading to poor computational performance for any  
486 downstream task on some large point clouds. One way to see this is the following: After adding the  
487 first edge that connects two components or the final simplex that fills a hole in the simplicial complex  
488 the VR complex keeps adding more and more simplices in the same area that keep the topology  
489 unchanged. One way to mitigate this problem is to pre-compute a set of simplices that are able to  
490 express the entire topology of the point cloud. For a point cloud  $X \subset \mathbb{R}^n$ , the  $\alpha$ -filtration consists of  
491 the intersection of the simplicial complexes of the VR filtration on  $X$  with the (higher-dimensional)  
492 Delaunay triangulation of  $X$  in  $\mathbb{R}$ . Due to algorithmic reasons, the filtration value of a simplex is  
493 then the radius of the circumscribed sphere instead of the maximum pair-wise distance of vertices.  
494 This reduces the number of required simplices across all dimensions to  $O(|X|^{\lceil n/2 \rceil})$ . However, the  
495 Delaunay triangulation becomes computationally infeasible for larger  $n$ .

496 **Definition B.1** ( $n$ -dimensional Delaunay triangulation). Given a set of vertices  $V \in \mathbb{R}^n$ , a Delaunay  
497 triangulation  $DT(V)$  is a triangulation of  $V$  such that for any  $n$ -simplex  $\sigma_n \in DT(V)$  the interior  
498 of the circum-hypersphere of  $\sigma_n$  contains no point of  $DT(V)$ . A triangulation of  $V$  is a SC  $\mathcal{S}$  with  
499 vertex set  $V$  such that its geometric realisation covers the convex hull of  $V$   $\text{hull}(V) = |\mathcal{S}|$  and we  
500 have for any two simplices  $\sigma, \sigma'$  that the intersection of geometric realisations  $|\sigma| \cap |\sigma'|$  is either  
501 empty or the geometric realisation  $|\hat{\sigma}|$  of a common sub-simplex  $\hat{\sigma} \subset \sigma, \sigma'$ .

502 If  $V$  is in general position, the Delaunay triangulation is unique and guaranteed to exist [18].

503 **Definition B.2** ( $\alpha$ -complex of a point cloud). Given a finite point cloud  $X$  in real space  $\mathbb{R}^n$ , the  
504  $\alpha$ -complex  $\alpha_\varepsilon(X)$  is the subset of the  $n$ -dimensional Delaunay triangulation  $DT(X)$  consisting of  
505 all  $\sigma \in DT(X)$  with a radius  $r$  of its circumscribed sphere with  $r \leq \varepsilon$ .

506 **Proof of the main theorem** We will now give the proof of the theorem that guarantees that TOPF  
 507 works. First, let us recall Theorem 4.1:

508 **Theorem 4.1** (Topological Point Features of Spheres). *Let  $X$  consist of at least  $(n + 2)$  points  
 509 (denoted by  $S$ ) sampled uniformly at random from a unit  $n$ -sphere in  $\mathbb{R}^{n+1}$  and an arbitrary number  
 510 of points with distance of at least 2 to  $S$ . When we now consider the  $\alpha$ -filtration on this point  
 511 cloud, with probability 1 we have that (i) there exists an  $n$ -th persistent homology class generated  
 512 by the 2-simplices on the convex hull of  $S$ , (ii) the associated unweighted harmonic homology  
 513 representative takes values in  $\{0, \pm 1\}$  where the 2-simplices on the boundary of the convex hull are  
 514 assigned a value of  $\pm 1$ , and (iii) the support of the associated topological point feature (TOPF)  $\mathcal{F}_n^*$   
 515 is precisely  $S$ :  $\text{supp}(\mathcal{F}_n^*) = S$ . (iv) The same holds true for point clouds sampled from multiple  
 516  $n_i$ -spheres if the above conditions are met on each individual sphere.*

517 *Proof.* Assume that we are in the scenario of the theorem. Now because the  $n$ -volume of  $(n - 1)$ -  
 518 submanifolds is zero, we have that with probability 1 the points of  $S$  don't lie on a single  $(n - 1)$ -  
 519 sphere inside the  $n$ -sphere. Let us now look at the  $\alpha$ -filtration of the simplices in  $S$ : Recall that the  
 520 filtration values of a  $k$ -simplex is given by the radius of the  $(k - 1)$ -sphere determined by its vertices.  
 521 Because all of the  $(n + 1)$ -simplices  $\sigma_{n+1}$  with vertices  $V \subset S$  in  $S$  lie on the same unit  $n$ -sphere  $S_n$ ,  
 522 they all share the filtration value of  $\alpha(\sigma_{n+1}) = 1$ . By the same argument as above, with probability 1  
 523 there are no  $(n + 1)$  points in  $S$  that lie on an *unit*  $(n - 1)$ -sphere. Thus all of the  $n$ -simplices  $\sigma_n$  lie  
 524 on  $(n - 1)$ -spheres  $S_n$  with a radius  $r < 1$  smaller than 1 and hence have a filtration value  $\alpha(\sigma_n)$   
 525 smaller than 1. Let

$$b := \max(\{\alpha(\sigma_n) : \sigma_n \subset \partial \text{hull}(S)\})$$

526 be the maximum filtration value of an  $n$ -simplex on the boundary of the convex hull of  $S$ . Then, then  
 527 a linear combination  $g$  of the  $n$ -simplices of the boundary of the convex hull of  $S$  with coefficients in  
 528  $\pm 1$  is a generator of a persistent homology class with life time  $(b, 1)$  (this follows from the fact that  
 529  $n$ -spheres and their triangulations are orientable). This proves claim (i).

530 Because of the assumption that all points not contained in  $S$  have a distance of at least 2 to the points  
 531 in  $S$ , all  $(n + 1)$ -simplices  $\sigma_{n+1}$  with vertices both in  $S$  and its complement in  $X$  will have a filtration  
 532 value  $\alpha(\sigma_{n+1}) \geq 1$  of at least 1. Recall that all  $(n + 1)$ -simplices  $\sigma_{n+1} \subset S$  with vertices inside  $S$   
 533 have a filtration value of  $\alpha(\sigma_{n+1}) = 1$ . Thus the adjoint of the  $n$ -th boundary operator  $\mathcal{B}_n^\top$  is trivial  
 534 on the homology generator  $g$ . Thus, we have that for the  $n$ -th Hodge Laplacian

$$L_n g = \mathcal{B}_{n-1}^\top \mathcal{B}_{n-1} g + \mathcal{B}_n \mathcal{B}_n^\top g = 0 + 0 = 0$$

535 and hence  $g$  is a harmonic generator for the entire filtration range of  $(b, 1)$ , which proves claim (ii).  
 536 Claim (iii) and (iv) then follow from the construction of the TOPF values.  $\square$

## 537 C Topological Clustering Benchmark Suite

538 We introduce seven point clouds for topological point cloud clustering in the topological clustering  
 539 benchmark suite (TCBS). The ground truth and the point clouds are depicted in Figure 6. The point  
 540 clouds represent a mix between 0-, 1- and 2-dimensional topological structures in noiseless and noisy  
 541 settings in ambient 2-dimensional and 3-dimensional space. The results of clustering according to  
 542 TOPF can be found in Figure 7.

## 543 D Implementation

544 We will release an implementation of TOPF and the code and data required to reproduce  
 545 the experimental results of this paper under [https://anonymous.4open.science/r/topf\\_](https://anonymous.4open.science/r/topf_submission-5C40/)  
 546 [submission-5C40/](https://anonymous.4open.science/r/topf_submission-5C40/). In particular, we will release the topological clustering benchmark suite.

547 All experiments were run on a Apple M1 Pro chipset with 10 cores and 32 GB memory. TOPF  
 548 and the experiments are implemented in Python and Julia. For persistent homology computations,  
 549 we used GUDHI [48] (© The GUDHI developers, MIT license) and Ripserer [52] (© mtsch, MIT  
 550 license), which is a modified Julia implementation of [3]. For the least square problems, we used  
 551 the LSMR implementation of SciPy [22]. We used the Node2Vec python implementation <https://github.com/eliorc/node2vec>  
 552 <https://github.com/eliorc/node2vec> (© Elior Cohen, MIT License) based on the Node2Vec Paper  
 553 [25]. We used the pgeof Python package for computation of geometric features <https://github.com>

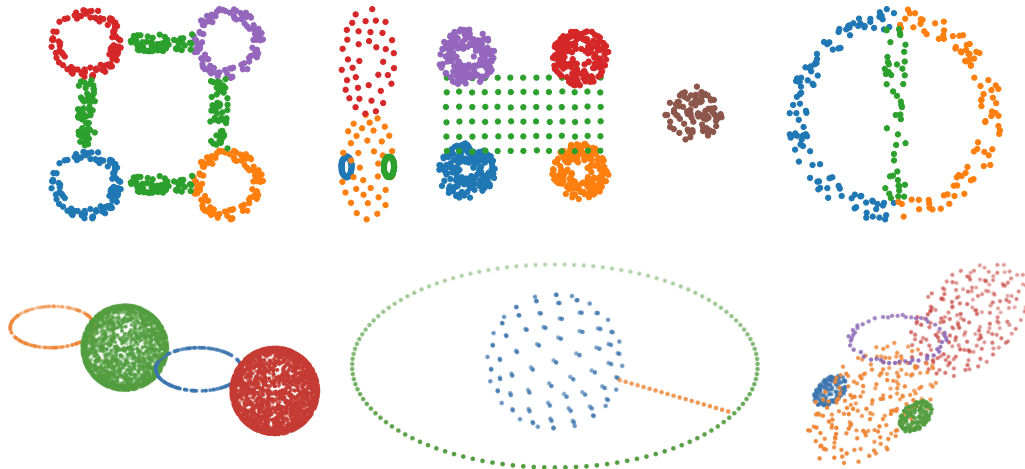


Figure 6: **Data sets of the Topological Clustering Benchmark Suite (TCBS) with true labels.** *Top: 2D data sets. From left to right: 4Spheres (656 points), Ellipses (158 points), Spheres+Grid (866 points), Halved Circle (249 points).* *Bottom: 3D data sets. From left to right: 2Spheres2Circles (4600 points), SphereinCircle (267 points), spaceship (650 points).*

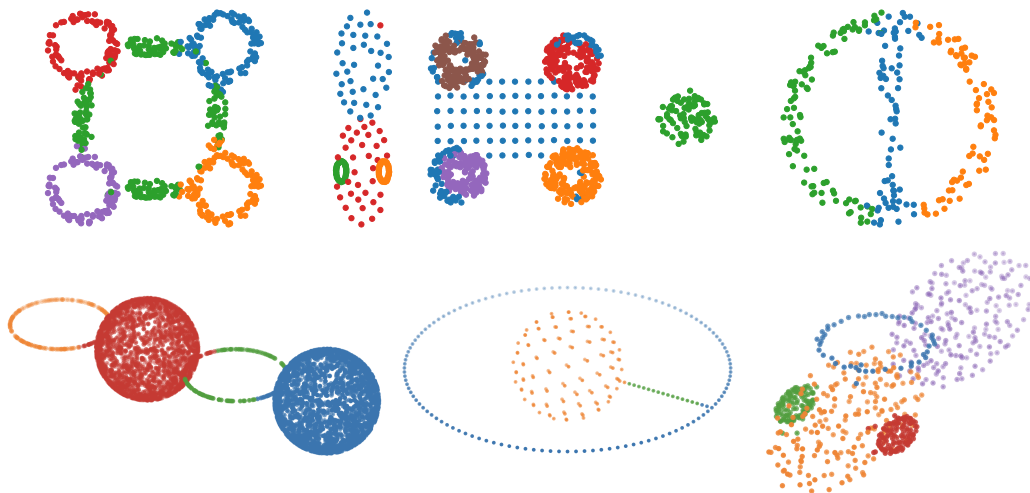


Figure 7: **Data sets of the Topological Clustering Benchmark Suite (TCBS) with labels generated by TOPF.** *Top: 2D data sets. From left to right: 4Spheres (0.81 ARI), Ellipses (0.95 ARI), Spheres+Grid (0.70 ARI), Halved Circle (0.71 ARI).* *Bottom: 3D data sets. From left to right: 2Spheres2Circles (0.94 ARI), SphereinCircle (0.97 ARI), spaceship (0.92 ARI).*

554 com/drprojects/point\_geometric\_features (© Damien Robert, Loic Landrieu, Romain Jan-  
 555 vier, MIT license). We use parts of the implementation of TPCC <https://git.rwth-aachen.de/netsci/publication-2023-topological-point-cloud-clustering> (© Computational  
 556 Network Science Group, RWTH Aachen University, MIT license).  
 557

## 558 D.1 Hyperparameters

559 All the relevant hyperparameters are already mentioned in their respective sections. However, for  
 560 convenience we gather and briefly discuss them in this section. We note that TOPF is robust and  
 561 applicable in most scenarios when using the default parameters without tuning hyperparameters. The  
 562 hyperparameters should more be thought of as an additional way where detailed domain-knowledge  
 563 can enter the TOPF pipeline.

564 **Maximum Homology Dimension  $d$**  The maximum homology dimension determines the dimen-  
565 sions of persistent homology the algorithm computes.

566 For the choice of the maximum homology degree  $d$  to be considered there are mainly three heuristics  
567 which we will list in decreasing importance (Cf. [24]):

568 I. In applications, we usually know which kind of topological features we are interested in, which  
569 will then determine  $d$ . This means that 1-dimensional homology and  $d = 1$  suffices when we  
570 are looking at loops of protein chains. On the other hand, if we are working with voids and  
571 cavities in 3d histological data, we need  $d = 2$  and thus compute 2-dimensional homology.

572 II. Algebraic topology tells us that there are no closed  $n$ -dimensional submanifolds of  $\mathbb{R}^n$ . Hence  
573 their top-homology will always vanish and all interesting homological activity will appear for  
574  $d < n$ .

575 III. In the vast majority of cases, the choice will be between  $d = 1$  or  $d = 2$  because empirically  
576 there are virtually no higher-dimensional topological features in practice.

577 In our quantitative experiments, we have always chosen  $d = n - 1$ .

578 **Thresholding parameter  $\delta$**  In step 4 of the algorithm, we normalise and threshold the harmonic  
579 representatives. After normalising, the entries of the vectors lie in the interval of  $[0, 1]$ . The  
580 thresholding parameter  $\delta$  now essentially determines an interval of  $[0, \delta]$  which we will linearly map  
581 to  $[0, 1]$ , while mapping all entries above  $\delta$  to 1 as well. This is necessary as most of the entries in  
582 the vector  $e_k^i$  are very close to 0 with a very small number of entries being close to 1. Without this  
583 thresholding, TOPF would now be almost entirely determined by these few large values. Thus this  
584 step limits the maximum possible influence of a single entry. However, because most of the entries of  
585  $e_k^i$  are concentrated around 0, small changes in  $\delta$  will not have a large effect and we chose  $\delta = 0.07$   
586 in all our experiments.

587 **Interpolation coefficient  $\lambda$**  The interpolation coefficient  $\lambda \in [0, 1)$  determines whether we build  
588 our simplicial complexes close to the birth or the death of the relevant homological features at time  
589  $t = b^{1-\lambda}d$ . This then in turns controls how localised or smooth the harmonic representative will  
590 be. In general, the noisier the ground data is the higher we should choose  $\lambda$ . However, TOPF is not  
591 sensitive to small changes in  $\lambda$ . We have picked  $\lambda = 0.3$  for all the quantitative experiments, which  
592 empirically represents a good choice for a broad range of applications.

593 **Feature selection factor  $\beta$**  Increasing  $\beta$  leads to TOPF preferring to pick a larger number of relevant  
594 topological features. Without specific domain-knowledge,  $\beta = 0$  represents a good choice.

595 **Feature selection quotients `max_total_quot`, `min_rel_quot`, and `min_0_ratio`** These are  
596 technical hyperparameters controlling the feature selection module of TOPF. For a technical account  
597 of them, see Appendix E. In most of the cases without domain knowledge, they do not have an effect  
598 on the performance of TOPF and should be kept at their default values.

599 **Simplicial Complex Weights** Although the simplicial weights are not technically a hyperparameter,  
600 there are many potential ways to weigh the considers SCs that can highlight or suppress different  
601 topological and geometric properties. In all our experiments, we use  $w_\Delta$  weights discussed in  
602 Appendix F.

## 603 E How to pick the most relevant topological features

604 **Simplified heuristic** The persistent homology  $P_k$  module in dimension  $k$  is given to us as a list of  
605 pairs of birth and death times  $(b_i^k, d_i^k)$ . We can assume these pairs are ordered in non-increasing order  
606 of the durations  $l_i^k = d_i^k - b_i^k$ . This list is typically very long and consists to a large part of noisy  
607 homological features which vanish right after they appear. In contrast, we are interested in connected  
608 components, loops, cavities, etc. that *persist* over a long time, indicating that they are important for  
609 the shape of the point cloud. Distinguishing between the relevant and the irrelevant features is in  
610 general difficult and may depend on additional insights on the domain of application. In order to  
611 provide a heuristic which does not depend on any a-priori assumptions on the number of relevant

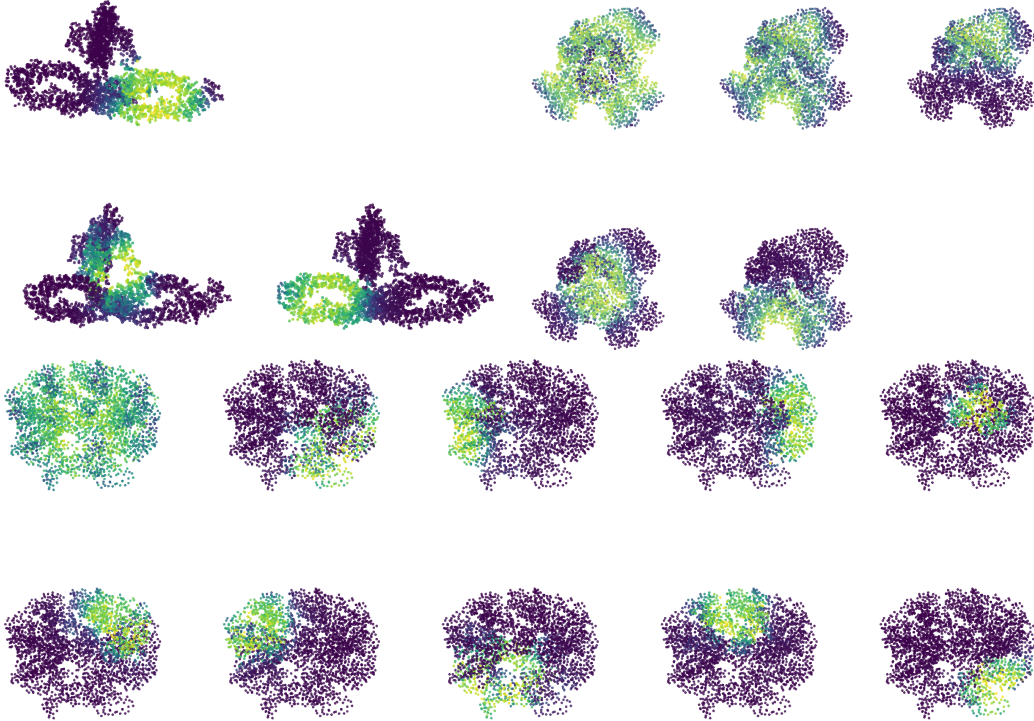


Figure 8: **TOPF heatmaps for three proteins.** *Top left* NALCN channelosome [32] *Top right:* Mutated Cys123 of E. coli [29], with convex hull added during computation, only 2-dimensional homology features *Bottom:* GroEL of E. coli [10] (Selected features).

612 features we pick the smallest quotient  $q_i^k := l_{i+1}^k / l_i^k > 0$  as the point of cut-off  $N_k := \arg \min_i q_i^k$ .  
 613 The only underlying assumption of this approach is that the band of “relevant” features is separated  
 614 from the “noisy” homological features by a drop in persistence.

615 **Advanced Heuristic** However, certain applications have a single very prominent feature, followed  
 616 by a range of still relevant features with significantly smaller life times, that are then followed by  
 617 the noisy features after another drop-off. This then could potentially lead the heuristic to find the  
 618 wrong drop-off. We propose to mitigate this issue by introducing a hyperparameter  $\beta \in \mathbb{R}_{>0}$ . We  
 619 then define the  $i$ -th importance-drop-off quotient  $q_i^k$  by

$$q_i^k := l_{i+1}^k / l_i^k (1 + \beta/i).$$

620 The basic idea is now to consider the most significant  $N_k$  homology classes in dimension  $k$  when  
 621 setting  $N_k$  to be

$$N_k := \arg \min_i q_i^k.$$

622 Increasing  $\beta$  leads the heuristic to prefer selections with more features than with fewer features.  
 623 Empirically, we still found  $\beta = 0$  to work well in a broad range of application scenarios and used  
 624 it throughout all experiments. There are only a few cases where domain-specific knowledge could  
 625 suggest picking a larger  $\beta$ .

626 To catch edge cases with multiple steep drops or a continuous transition between real features and  
 627 noise, we introduce two more checks: We allow a minimal  $q_i^k$  of `min_rel_quot` = 0.1 and a  
 628 maximal quotient  $q_1^k / q_i^k$  of `max_total_quot` = 10 between any homology dimensions. Because  
 629 features in 0-dimensional homology are often more noisy than features in higher dimensions, we add  
 630 a minimum zero-dimensional homology ratio of `min_0_ratio` = 5, i.e. every chosen 0-dimensional  
 631 feature needs to be at least `min_0_ratio` more persistent than the minimum persistence of the  
 632 higher-dimensional features. Because these hyperparameters only deal with the edge cases of  
 633 feature selection, TOPF is not very sensitive to them. For all our experiments, we used the above  
 634 hyperparameters. We advise to change them only in cases where one has in-depth domain knowledge  
 635 about the nature of relevant topological features.

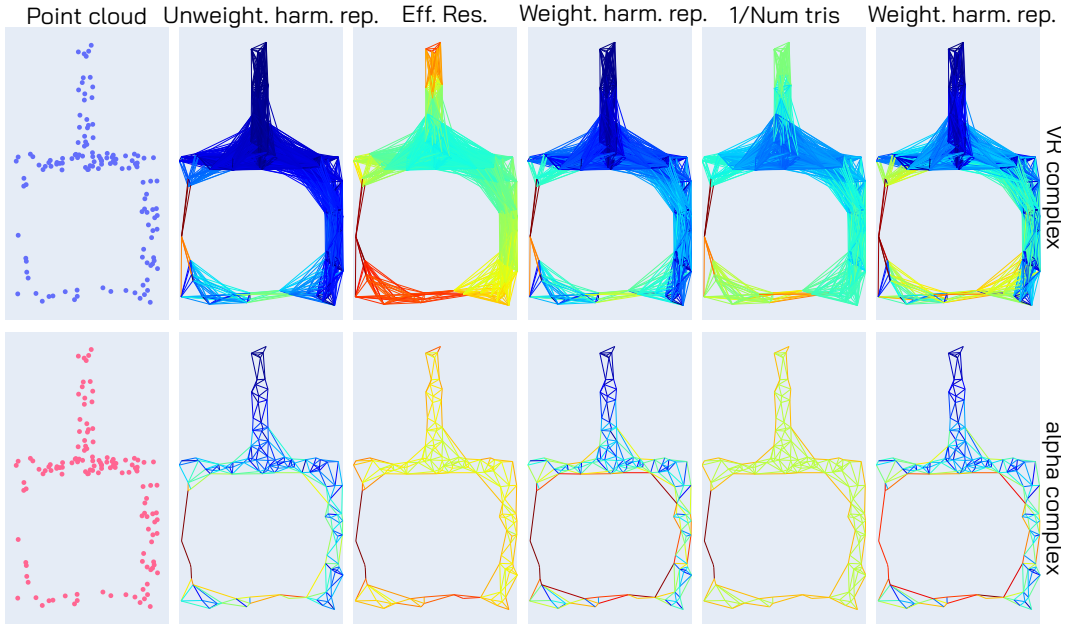


Figure 9: **Effect of weighing a simplicial complex on harmonic representatives.** *Top:* VR complex. *Bottom:*  $\alpha$ -complex *Left:* The base point cloud with different densities. *2<sup>nd</sup> Left:* Unweighted harmonic homology representative of the large loop. *3<sup>rd</sup> Right:* Effective resistance of the 1-simplices. *3<sup>rd</sup> Right:* Harmonic homology representative of the complex weighted by effective resistance. *2<sup>nd</sup> Right:* Inverse of number of incident triangles (Definition F.1). *Right:* Harmonic homology representative of the complex weighted by number of incident triangles. Up to a small threshold, the standard harmonic representative in the VR complex is almost exclusively supported in the low-density regions of the simplicial complex. This leads to poor and unpredictable classification performance in downstream tasks. In contrast, the harmonic homology representative of the weighted VR complex has a more homogenous support along the loop, while still being able to discriminate the edges not contributing to the loop. The  $\alpha$ -complex suffers less from this phenomenon (at least in dimension 2), and hence reweighing is not necessarily required.

## 636 F Simplicial Weights

637 In an ideal world, the harmonic eigenvectors in dimension  $k$  would be vectors assigning  $\pm 1$  to all  
638  $k$ -simplices contributing to  $k$ -dimensional homological feature, a 0 to all  $k$ -simplices not contributing  
639 or orthogonal to the feature, and a value in  $(-1, 1)$  for all simplices based on the alignment of the  
640 simplex with the boundary of the void. However, this is not the case: In dimension 1, we can for  
641 example imagine a total flow of 1 circling around the hole. This flow is then split up between all  
642 parallel edges which means *two* things: **I** Edges where the loop has a *larger diameter* have *smaller*  
643 *harmonic values* than edges in thin areas and **II** in VR complexes, which are the most frequently  
644 used simplicial complexes in TDA, edges in areas with a *high point density* have *smaller harmonic*  
645 *values* than edges in low-density areas. Point **II** is another advantage of  $\alpha$ -complexes: The expected  
646 number of simplices per point does not scale with the point density in the same way as it does in the  
647 VR complex, because only the simplices of the Delaunay triangulation can appear in the complex.

648 We address this problem by weighing the  $k$ -simplices of the simplicial complex. The idea behind this  
649 is to weigh the simplicial complex in such a way that it increases and decreases the harmonic values  
650 of some simplices in an effort to make the harmonic eigenvectors more homogeneous. For weights  
651  $w \in \mathbb{R}^{S_k}$ ,  $W = \text{diag}(w)$ , the symmetric weighted Hodge Laplacian [45] takes the form of

$$L_k^w = W^{1/2} \mathcal{B}_{k-1} \mathcal{B}_{k-1}^\top W^{1/2} + W^{-1/2} \mathcal{B}_k \mathcal{B}_k^\top W^{-1/2}.$$



652 Because we want the homology representative to lie in the weighted gradient space, we have to scale  
 653 its entries with the weight and set  $e_{k,w}^i := W^{-1/2}e_k^i$ . With this, we have that

$$\mathcal{B}_{k-1}^\top W^{1/2} e_{k,w}^i = \mathcal{B}_{k-1}^\top W^{1/2} W^{-1/2} e_k^i = \mathcal{B}_{k-1}^\top e_k^i = 0$$

654 We propose two options to weigh the simplicial complex. The first option is to weigh a  $k$ -simplex by  
 655 the square of the number of  $k + 1$ -simplices the simplex is contained in:

$$w_\Delta(\sigma_k) = 1/(|\{\sigma_{k+1} \in \mathcal{S}_{k+1}^t : \sigma_k \subset \sigma_{k+1}\}| + 1)^2$$

656 where the  $+1$  is to enforce good behaviour at simplices that are not contained in any higher-order  
 657 simplices. One of the advantages of the  $\alpha$ -complex is that we don't have large concentrations of  
 658 simplices in well-connected areas. The proposed weighting  $w_\Delta$  is computationally straightforward,  
 659 as it can be obtained as the column sums of the absolute value of the boundary matrix  $|\mathcal{B}_k|$ . The  
 660 weights also deal with the previously mentioned problem **II**: As the homology representative is scaled  
 661 inversely to the weight vector  $w$ , the simplices in high-density regions will be assigned a low weight  
 662 and thus their weighted homology representative will have a larger entry. By the projection to the  
 663 orthogonal complement of the curl space, this large entry is then diffused among the high-density  
 664 region of the SC with many simplices, whereas the lower entries of the simplices in low-density  
 665 regions are only diffused among fewer adjacent simplices.

666 However, the first weight is not able to incorporate the number of parallel simplices into the weighting.  
 667 This is why we propose a second simplicial weight function based on generalised effective resistance.

668 **Definition F.1** (Effective Hodge resistance weights). For a simplicial complex  $\mathcal{S}$  with boundary  
 669 matrices  $(\mathcal{B}_k)$ , we define the effective Hodge resistance weights  $w_R$  on  $k$ -simplices to be:

$$w_R := \text{diag}(\mathcal{B}_{k-1}^+ \mathcal{B}_{k-1})^2$$

670 where  $\text{diag}(-)$  denotes the vector of diagonal entries and  $(-)^+$  denotes taking the Moore–Penrose  
 671 inverse.

672 Intuitively for  $k = 1$ , we can assume that every edge has a resistance of 1 and then the effective  
 673 resistance coincides with the notion from Physics. Thus simplices with many parallel simplices are  
 674 assigned a small effective resistance, whereas simplices with few parallel simplices are assigned an  
 675 effective resistance close to 1. However, computing the Moore–Penrose inverse is computationally  
 676 expensive and only feasible for small simplicial complexes.

677 In Figure 9, we show that the weights  $w_\Delta$  are a good approximation of the effective resistance in  
 678 terms of the resulting harmonic representative. The standard form of TOPF used in all experiments  
 679 uses  $w_\Delta$ -weights.

## 680 G Limitations

681 **Topological features are not everywhere** The proposed topological point features take relevant  
 682 persistent homology generators and turn these into point-level features. As such, applying TOPF  
 683 only produces meaningful results on point clouds that have a topological structure. On these point  
 684 clouds, TOPF can extract structural information unobtainable by non-topological methods. Although  
 685 TDA has been successful in a wide range of applications, a large number of data sets does not  
 686 possess a meaningful topological structure. Applying TOPF in these cases will produce no additional  
 687 information. Other data sets require pre-processing before containing topological features. In Figure 4  
 688 *left*, the  $2d$  topological features characterising protein pockets of Cys123 only appear after artificially  
 689 adding points sampled on the convex hull of the point cloud (Cf [40]).

690 **Computing persistent homology can be computationally expensive** As TOPF relies on the  
 691 computation of persistent homology including homology generators, its runtime increases on very  
 692 large point clouds. This is especially true when using VR instead of  $\alpha$ -filtrations, which become  
 693 computationally infeasible for higher-dimensional point clouds. Persistent homology computations  
 694 for dimensions above 2 are only feasible for very small point clouds. Because virtually all discovered  
 695 relevant homological features in applications appear in dimension 0, 1, or 2, this does not present  
 696 a large problem. Despite these computational challenges, subsampling, either randomly or using  
 697 landmarks, usually preserves relevant topological features and thus extends the applicability of TDA  
 698 in general and TOPF even to very large point clouds.

699 **Automatic feature selection is difficult without domain knowledge** While the proposed heuristics  
700 works well across a variety of domains and application scenarios, only domain- and problem-specific  
701 knowledge makes truthful feature selection feasible.

702 **Experimental Evaluation** There are no benchmark sets for topological point features in the  
703 literature, which makes benchmarking TOPF not straightforward. On the level of clustering, we  
704 introduced the topological clustering benchmark suite to make quantitative comparisons of TOPF  
705 possible, and benchmarked TOPF on some of the point clouds of [24]. On both the level of point  
706 features and real-world data sets, it is however hard to establish what a *ground truth* of topological  
707 features would mean. Instead we chose to qualitatively report the results of TOPF on proteins and  
708 real-world data, see Figure 4.

## 709 **NeurIPS Paper Checklist**

### 710 **1. Claims**

711 Question: Do the main claims made in the abstract and introduction accurately reflect the  
712 paper's contributions and scope?

713 Answer: [\[Yes\]](#)

714 Justification: The claims about TOPF are supported by the theoretical background in Section 2  
715 and Section 3, quantitatively and qualitatively validated and benchmarked in Section 5.  
716 Furthermore, a theoretical guarantee can be found in Section 4.

717 Guidelines:

- 718 • The answer NA means that the abstract and introduction do not include the claims  
719 made in the paper.
- 720 • The abstract and/or introduction should clearly state the claims made, including the  
721 contributions made in the paper and important assumptions and limitations. A No or  
722 NA answer to this question will not be perceived well by the reviewers.
- 723 • The claims made should match theoretical and experimental results, and reflect how  
724 much the results can be expected to generalize to other settings.
- 725 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
726 are not attained by the paper.

### 727 **2. Limitations**

728 Question: Does the paper discuss the limitations of the work performed by the authors?

729 Answer: [\[Yes\]](#)

730 Justification: We believe that being open about limitations is crucial for the practice of  
731 doing good Science. We briefly discuss the main limitations in ??, and talk in detail about  
732 limitations in Appendix G. Finally, we are open about limitations when talking about the  
733 theoretical background and the algorithm in Section 2 and Section 3 and the remark in  
734 Section 4.

735 Guidelines:

- 736 • The answer NA means that the paper has no limitation while the answer No means that  
737 the paper has limitations, but those are not discussed in the paper.
- 738 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 739 • The paper should point out any strong assumptions and how robust the results are to  
740 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
741 model well-specification, asymptotic approximations only holding locally). The authors  
742 should reflect on how these assumptions might be violated in practice and what the  
743 implications would be.
- 744 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
745 only tested on a few datasets or with a few runs. In general, empirical results often  
746 depend on implicit assumptions, which should be articulated.
- 747 • The authors should reflect on the factors that influence the performance of the approach.  
748 For example, a facial recognition algorithm may perform poorly when image resolution  
749 is low or images are taken in low lighting. Or a speech-to-text system might not be  
750 used reliably to provide closed captions for online lectures because it fails to handle  
751 technical jargon.
- 752 • The authors should discuss the computational efficiency of the proposed algorithms  
753 and how they scale with dataset size.
- 754 • If applicable, the authors should discuss possible limitations of their approach to  
755 address problems of privacy and fairness.
- 756 • While the authors might fear that complete honesty about limitations might be used by  
757 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
758 limitations that aren't acknowledged in the paper. The authors should use their best  
759 judgment and recognize that individual actions in favor of transparency play an impor-  
760 tant role in developing norms that preserve the integrity of the community. Reviewers  
761 will be specifically instructed to not penalize honesty concerning limitations.

762 **3. Theory Assumptions and Proofs**

763 Question: For each theoretical result, does the paper provide the full set of assumptions and  
764 a complete (and correct) proof?

765 Answer: [Yes]

766 Justification: We provide a full set of assumptions for the theorem in Section 4 and a  
767 complete proof in Appendix B. We give references for all cited propositions and theorems  
768 exceeding basic common mathematical knowledge.

769 Guidelines:

- 770 • The answer NA means that the paper does not include theoretical results.
- 771 • All the theorems, formulas, and proofs in the paper should be numbered and cross-  
772 referenced.
- 773 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 774 • The proofs can either appear in the main paper or the supplemental material, but if  
775 they appear in the supplemental material, the authors are encouraged to provide a short  
776 proof sketch to provide intuition.
- 777 • Inversely, any informal proof provided in the core of the paper should be complemented  
778 by formal proofs provided in appendix or supplemental material.
- 779 • Theorems and Lemmas that the proof relies upon should be properly referenced.

780 **4. Experimental Result Reproducibility**

781 Question: Does the paper fully disclose all the information needed to reproduce the main ex-  
782 perimental results of the paper to the extent that it affects the main claims and/or conclusions  
783 of the paper (regardless of whether the code and data are provided or not)?

784 Answer: [Yes]

785 Justification: We list all the steps necessary to reproduce TOPF in Section 3, Appendix E, Ap-  
786 pendix F and talk in detail about the hyperparameter choices in Appendix D.1. Furthermore,  
787 we will both release the Topological Clustering Benchmark Suite and the code necessary to  
788 reproduce all experiments in this paper.

789 Guidelines:

- 790 • The answer NA means that the paper does not include experiments.
- 791 • If the paper includes experiments, a No answer to this question will not be perceived  
792 well by the reviewers: Making the paper reproducible is important, regardless of  
793 whether the code and data are provided or not.
- 794 • If the contribution is a dataset and/or model, the authors should describe the steps taken  
795 to make their results reproducible or verifiable.
- 796 • Depending on the contribution, reproducibility can be accomplished in various ways.  
797 For example, if the contribution is a novel architecture, describing the architecture fully  
798 might suffice, or if the contribution is a specific model and empirical evaluation, it may  
799 be necessary to either make it possible for others to replicate the model with the same  
800 dataset, or provide access to the model. In general, releasing code and data is often  
801 one good way to accomplish this, but reproducibility can also be provided via detailed  
802 instructions for how to replicate the results, access to a hosted model (e.g., in the case  
803 of a large language model), releasing of a model checkpoint, or other means that are  
804 appropriate to the research performed.
- 805 • While NeurIPS does not require releasing code, the conference does require all submis-  
806 sions to provide some reasonable avenue for reproducibility, which may depend on the  
807 nature of the contribution. For example
  - 808 (a) If the contribution is primarily a new algorithm, the paper should make it clear how  
809 to reproduce that algorithm.
  - 810 (b) If the contribution is primarily a new model architecture, the paper should describe  
811 the architecture clearly and fully.
  - 812 (c) If the contribution is a new model (e.g., a large language model), then there should  
813 either be a way to access this model for reproducing the results or a way to reproduce  
814 the model (e.g., with an open-source dataset or instructions for how to construct  
815 the dataset).

816 (d) We recognize that reproducibility may be tricky in some cases, in which case  
817 authors are welcome to describe the particular way they provide for reproducibility.  
818 In the case of closed-source models, it may be that access to the model is limited in  
819 some way (e.g., to registered users), but it should be possible for other researchers  
820 to have some path to reproducing or verifying the results.

## 821 5. Open access to data and code

822 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
823 tions to faithfully reproduce the main experimental results, as described in supplemental  
824 material?

825 Answer: [Yes]

826 Justification: We will release the full code necessary to reproduce all experimental results of  
827 this paper. Furthermore, we will release the topological clustering benchmark suite to the  
828 public.

829 Guidelines:

- 830 • The answer NA means that paper does not include experiments requiring code.
- 831 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/  
832 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 833 • While we encourage the release of code and data, we understand that this might not be  
834 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not  
835 including code, unless this is central to the contribution (e.g., for a new open-source  
836 benchmark).
- 837 • The instructions should contain the exact command and environment needed to run to  
838 reproduce the results. See the NeurIPS code and data submission guidelines ([https://  
839 nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 840 • The authors should provide instructions on data access and preparation, including how  
841 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 842 • The authors should provide scripts to reproduce all experimental results for the new  
843 proposed method and baselines. If only a subset of experiments are reproducible, they  
844 should state which ones are omitted from the script and why.
- 845 • At submission time, to preserve anonymity, the authors should release anonymized  
846 versions (if applicable).
- 847 • Providing as much information as possible in supplemental material (appended to the  
848 paper) is recommended, but including URLs to data and code is permitted.

## 849 6. Experimental Setting/Details

850 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-  
851 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the  
852 results?

853 Answer: [Yes]

854 Justification: We do not train neural networks in this paper. However, we will release  
855 the topological clustering benchmark suite. We talk in detail about how to reproduce the  
856 algorithm and the relevant choices of hyperparameters, and how we evaluate the experiments.

857 Guidelines:

- 858 • The answer NA means that the paper does not include experiments.
- 859 • The experimental setting should be presented in the core of the paper to a level of detail  
860 that is necessary to appreciate the results and make sense of them.
- 861 • The full details can be provided either with the code, in appendix, or as supplemental  
862 material.

## 863 7. Experiment Statistical Significance

864 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
865 information about the statistical significance of the experiments?

866 Answer: [Yes]

867 Justification: We provide standard deviations where applicable in Table 1, unless the  
868 standard deviation is 0, which we talk about in the caption of the table. In Figure 5 we give  
869 a confidence interval for all the experiments.

870 Guidelines:

- 871 • The answer NA means that the paper does not include experiments.
- 872 • The authors should answer "Yes" if the results are accompanied by error bars, confi-  
873 dence intervals, or statistical significance tests, at least for the experiments that support  
874 the main claims of the paper.
- 875 • The factors of variability that the error bars are capturing should be clearly stated (for  
876 example, train/test split, initialization, random drawing of some parameter, or overall  
877 run with given experimental conditions).
- 878 • The method for calculating the error bars should be explained (closed form formula,  
879 call to a library function, bootstrap, etc.)
- 880 • The assumptions made should be given (e.g., Normally distributed errors).
- 881 • It should be clear whether the error bar is the standard deviation or the standard error  
882 of the mean.
- 883 • It is OK to report 1-sigma error bars, but one should state it. The authors should  
884 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
885 of Normality of errors is not verified.
- 886 • For asymmetric distributions, the authors should be careful not to show in tables or  
887 figures symmetric error bars that would yield results that are out of range (e.g. negative  
888 error rates).
- 889 • If error bars are reported in tables or plots, The authors should explain in the text how  
890 they were calculated and reference the corresponding figures or tables in the text.

## 891 8. Experiments Compute Resources

892 Question: For each experiment, does the paper provide sufficient information on the com-  
893 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
894 the experiments?

895 Answer: [Yes]

896 Justification: We list the hardware used in Appendix D and list the required running times in  
897 our quantitative experiments, see Table 1. Because we did not train neural networks, the  
898 results are easily reproducible on any PC in reasonable time.

899 Guidelines:

- 900 • The answer NA means that the paper does not include experiments.
- 901 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,  
902 or cloud provider, including relevant memory and storage.
- 903 • The paper should provide the amount of compute required for each of the individual  
904 experimental runs as well as estimate the total compute.
- 905 • The paper should disclose whether the full research project required more compute  
906 than the experiments reported in the paper (e.g., preliminary or failed experiments that  
907 didn't make it into the paper).

## 908 9. Code Of Ethics

909 Question: Does the research conducted in the paper conform, in every respect, with the  
910 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

911 Answer: [Yes]

912 Justification: We have reviewed the NeurIPS Ethics guidelines to make sure our research  
913 complies with them. (It does comply.)

914 Guidelines:

- 915 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 916 • If the authors answer No, they should explain the special circumstances that require a  
917 deviation from the Code of Ethics.

- 918 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-  
919 eration due to laws or regulations in their jurisdiction).

## 920 10. Broader Impacts

921 Question: Does the paper discuss both potential positive societal impacts and negative  
922 societal impacts of the work performed?

923 Answer: [NA]

924 Justification: As the paper is of foundational nature, we do not foresee any direct societal  
925 impacts.

926 Guidelines:

- 927 • The answer NA means that there is no societal impact of the work performed.
- 928 • If the authors answer NA or No, they should explain why their work has no societal  
929 impact or why the paper does not address societal impact.
- 930 • Examples of negative societal impacts include potential malicious or unintended uses  
931 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations  
932 (e.g., deployment of technologies that could make decisions that unfairly impact specific  
933 groups), privacy considerations, and security considerations.
- 934 • The conference expects that many papers will be foundational research and not tied  
935 to particular applications, let alone deployments. However, if there is a direct path to  
936 any negative applications, the authors should point it out. For example, it is legitimate  
937 to point out that an improvement in the quality of generative models could be used to  
938 generate deepfakes for disinformation. On the other hand, it is not needed to point out  
939 that a generic algorithm for optimizing neural networks could enable people to train  
940 models that generate Deepfakes faster.
- 941 • The authors should consider possible harms that could arise when the technology is  
942 being used as intended and functioning correctly, harms that could arise when the  
943 technology is being used as intended but gives incorrect results, and harms following  
944 from (intentional or unintentional) misuse of the technology.
- 945 • If there are negative societal impacts, the authors could also discuss possible mitigation  
946 strategies (e.g., gated release of models, providing defenses in addition to attacks,  
947 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from  
948 feedback over time, improving the efficiency and accessibility of ML).

## 949 11. Safeguards

950 Question: Does the paper describe safeguards that have been put in place for responsible  
951 release of data or models that have a high risk for misuse (e.g., pretrained language models,  
952 image generators, or scraped datasets)?

953 Answer: [NA]

954 Justification: We do not foresee any such risks.

955 Guidelines:

- 956 • The answer NA means that the paper poses no such risks.
- 957 • Released models that have a high risk for misuse or dual-use should be released with  
958 necessary safeguards to allow for controlled use of the model, for example by requiring  
959 that users adhere to usage guidelines or restrictions to access the model or implementing  
960 safety filters.
- 961 • Datasets that have been scraped from the Internet could pose safety risks. The authors  
962 should describe how they avoided releasing unsafe images.
- 963 • We recognize that providing effective safeguards is challenging, and many papers do  
964 not require this, but we encourage authors to take this into account and make a best  
965 faith effort.

## 966 12. Licenses for existing assets

967 Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
968 the paper, properly credited and are the license and terms of use explicitly mentioned and  
969 properly respected?

970 Answer: [Yes]



971 Justification: We credit the creators and owners of code used in the model, and state the  
972 licenses.

973 Guidelines:

- 974 • The answer NA means that the paper does not use existing assets.
- 975 • The authors should cite the original paper that produced the code package or dataset.
- 976 • The authors should state which version of the asset is used and, if possible, include a  
977 URL.
- 978 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 979 • For scraped data from a particular source (e.g., website), the copyright and terms of  
980 service of that source should be provided.
- 981 • If assets are released, the license, copyright information, and terms of use in the  
982 package should be provided. For popular datasets, `paperswithcode.com/datasets`  
983 has curated licenses for some datasets. Their licensing guide can help determine the  
984 license of a dataset.
- 985 • For existing datasets that are re-packaged, both the original license and the license of  
986 the derived asset (if it has changed) should be provided.
- 987 • If this information is not available online, the authors are encouraged to reach out to  
988 the asset's creators.

### 989 13. New Assets

990 Question: Are new assets introduced in the paper well documented and is the documentation  
991 provided alongside the assets?

992 Answer: [\[Yes\]](#)

993 Justification: The code and benchmark suite which we will release with the paper are  
994 described and documented in the paper.

995 Guidelines:

- 996 • The answer NA means that the paper does not release new assets.
- 997 • Researchers should communicate the details of the dataset/code/model as part of their  
998 submissions via structured templates. This includes details about training, license,  
999 limitations, etc.
- 1000 • The paper should discuss whether and how consent was obtained from people whose  
1001 asset is used.
- 1002 • At submission time, remember to anonymize your assets (if applicable). You can either  
1003 create an anonymized URL or include an anonymized zip file.

### 1004 14. Crowdsourcing and Research with Human Subjects

1005 Question: For crowdsourcing experiments and research with human subjects, does the paper  
1006 include the full text of instructions given to participants and screenshots, if applicable, as  
1007 well as details about compensation (if any)?

1008 Answer: [\[NA\]](#)

1009 Justification: This paper does not involve crowdsourcing nor research with human subjects.

1010 Guidelines:

- 1011 • The answer NA means that the paper does not involve crowdsourcing nor research with  
1012 human subjects.
- 1013 • Including this information in the supplemental material is fine, but if the main contribu-  
1014 tion of the paper involves human subjects, then as much detail as possible should be  
1015 included in the main paper.
- 1016 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,  
1017 or other labor should be paid at least the minimum wage in the country of the data  
1018 collector.

### 1019 15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 1020 Subjects

1021 Question: Does the paper describe potential risks incurred by study participants, whether  
1022 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
1023 approvals (or an equivalent approval/review based on the requirements of your country or  
1024 institution) were obtained?

1025 Answer: [NA]

1026 Justification: This paper does not involve crowdsourcing nor research with human subjects.

1027 Guidelines:

- 1028 • The answer NA means that the paper does not involve crowdsourcing nor research with  
1029 human subjects.
- 1030 • Depending on the country in which research is conducted, IRB approval (or equivalent)  
1031 may be required for any human subjects research. If you obtained IRB approval, you  
1032 should clearly state this in the paper.
- 1033 • We recognize that the procedures for this may vary significantly between institutions  
1034 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
1035 guidelines for their institution.
- 1036 • For initial submissions, do not include any information that would break anonymity (if  
1037 applicable), such as the institution conducting the review.