ECLAYR: FAST AND ROBUST TOPOLOGICAL LAYER BASED ON DIFFERENTIABLE EULER CHARACTERISTIC CURVE

Anonymous authors

Paper under double-blind review

ABSTRACT

In the realm of Topological Data Analysis, persistent homology has traditionally served as a primary tool for extracting topological features. However, approaches relying on persistent homology often encounter practical challenges due to their high computational costs. To address this issue, we propose a computationally efficient novel topological layer tailored for general deep learning architectures, leveraging the Euler Characteristic Curve (ECC). Unlike methods based on persistent homology, ECC offers computational advantages by circumventing the need for persistent homology calculation, while still allowing access to crucial information about the underlying topological structure. The proposed layer can readily adapt to diverse data modalities by allowing appropriate filtration according to the user's preference, enabling its application across various learning problems without data preprocessing. We present a novel technique for stable backpropagation that effectively mitigates the vanishing gradient problems commonly encountered in existing methods, allowing for seamless integration of our layer into deep learning models. We go on to present stability analysis, showing that the proposed layer is robust against noise and outliers. We apply our method to topological autoencoders, showing that the standard loss function can effectively regularize topological structures of the latent space. Through classification experiments across various datasets, we illustrate the benefits of our approach in mitigating information loss under conditions of data scarcity or data contamination.

031 032

006

008 009 010

011 012 013

014

015

016

017

018

019

021

024

025

026

027

028

029

033

1

INTRODUCTION

034

035

In recent years, machine learning communities have witnessed increasing efforts to incorporate Topological Data Analysis (TDA) into deep learning workflows, an emerging paradigm known as 037 topological deep learning (Carlsson & Gabrielsson, 2020; Papamarkou et al., 2024). Topological deep learning integrates tools from TDA to exploit essential topological features within the data that elude conventional methods, or to enhance understanding and control of computational models. 040 Persistent homology (PH), a primary tool in TDA, captures multi-scale topological features of the 041 underlying data structure by tracking the birth and death of homology features, thereby producing 042 topological summaries such as persistence diagrams or barcodes (Chazal & Michel, 2021). Given 043 that PH is a multiset by nature, a number of strategies have been proposed to transform these PH-044 based topological summaries into alternative representations that are more suitable for subsequent machine learning tasks (e.g., Bubenik et al., 2015; Adams et al., 2017; Umeda, 2017) (see, for example, Hensel et al. (2021) for a review). 046

Recent efforts have shed light on the possibility of incorporating PH-based topological summaries as input features for neural networks, enhancing their ability to learn from the intrinsic geometric structure of the data. Hofer et al. (2017); Carrière et al. (2020) introduced topological layers aimed at learning vector embeddings of persistence diagrams using a particular parametrization, yet lacking differentiability required for enabling gradient backpropagation. Hofer et al. (2019); Gabrielsson et al. (2020); Carriere et al. (2021); Leygonie et al. (2022) explored the differentiability aspects of the PH-based functions or losses, highlighting the potential of incorporating topological insights into deep learning frameworks. Kim et al. (2020) were the first to propose a generic differentiable

topological layer allowing backpropagation, offering flexibility in its integration within arbitrary network architectures.

Despite its popularity, computations involving PH can demand significant computational resources, 057 rendering them impractical for large-scale deep learning applications. Its time complexity generally scales poorly with the size and dimensionality of the data; when the number of simplices is given by N, the time complexity of the standard PH computation algorithm is $O(N^3)$ (Otter et al., 060 2017). As a result, there has recently been a growing need for alternative features capable of captur-061 ing topological information in a computationally efficient manner. The Euler Characteristic Curve 062 (ECC) is one such feature that can be computed without the need for PH calculation. Due to their 063 ability to drastically boost computational efficiency, ECC-based descriptors have recently received 064 increased attention (e.g., Beltramo et al., 2021; Chen et al., 2022; Dłotko & Gurnari, 2023; Hacquard & Lebovici, 2023; Malott & Wilsey, 2023; Jiang et al., 2023; Richardson & Werman, 2014; Laky 065 & Zavala, 2024). However, these descriptors have been primarily utilized in a static manner within 066 the framework of feature engineering. A more recent contribution in this field is the Differentiable 067 Euler Characteristic Transformation (DECT) (Röell & Rieck, 2024), though detailed analytical ex-068 plorations were not conducted. As indicated by its name, DECT utilizes ECT; a collection of ECCs 069 computed from various directions (Turner et al., 2014).

In this paper, we aim to develop a novel computationally efficient ECC-based topological layer that 071 facilitates integration with general deep learning models via stable backpropagation. The preceding 072 work most closely related to ours is the recent development of Kim et al. (2020) and Röell & Rieck 073 (2024). Kim et al. (2020) adopted persistence landscapes to construct a differentiable topological 074 layer. Notwithstanding its merits, persistence landscapes inherit the high computational complexity 075 of PH, and their gradients can often be highly sparse, lacking substantial information (see Figure 076 5 in Appendix). The proposal for a topological layer utilizing DECT (Röell & Rieck, 2024) relies 077 specifically on the height filtration, which works best with graphs and meshes, but not necessarily with other data structures; it may be applicable to point clouds, for example, yet with a compromise 079 in connectivity information. Moreover, Röell & Rieck (2024) achieve differentiability of ECT by employing a sigmoid approximation, which may result in inconsistent gradients or even vanishing 081 gradient problems when evaluated at discretized points; this will be further discussed in Section 4.

082 This work proposes a novel fast and robust ECC-based topological layer, ECLayr, designed to ad-083 dress all the aforementioned drawbacks comprehensively. Our proposed method obviates the need 084 for PH calculations, thereby significantly enhancing computational efficiency while preserving the 085 capability to extract key topological information from underlying data structures. ECLayr is capable of utilizing generic filtrations, exhibiting versatility across various data modalities without 087 necessitating data preprocessing or resorting to a particular filtration. Importantly, we introduce a 880 novel approach for stable backpropagation with respect to the layer input, addressing the inconsistent/vanishing gradient issues associated with the sigmoid approximation in DECT. We also provide 089 a stability analysis, showing that the proposed layer is robust against noise and outliers. Our exper-090 imental analysis show that ECLayr delivers performance comparable to state-of-the-art PH-based 091 topological layers, while being significantly faster by several orders of magnitude. Using our pro-092 posed backpropagation algorithm, ECLayr exhibits improved performance and greater efficiency 093 than DECT. We further demonstrate its versatility through the application on topological autoen-094 coders.

095 096

098

2 MATHEMATICAL BACKGROUND

This section provides a brief overview of the essential tools in TDA used throughout the development, as well as some notations. For further information, see, for example, Hatcher (2002);
Edelsbrunner & Harer (2010); Chazal & Michel (2021); Kaczynski et al. (2004).

Simplex and Simplicial Complex. Let u_0, \ldots, u_k be affinely independent points in \mathbb{R}^d . A *k*simplex is the convex hull of the k + 1 points, $\sigma_k = \operatorname{conv}\{u_0, \ldots, u_k\}$ (e.g., 0-simplex is a vertex, 1-simplex is an edge, 2-simplex is a triangle, etc.). The *dimension* of σ_k is k. τ is a *face* of σ_k if it is a convex hull constructed from any non-empty subset of the k+1 points of σ_k . A *simplicial complex* K is a finite collection of simplices such that (i) the face of any simplex in K is also in K, and (ii) the intersection of two simplices in K is either empty or a face of both simplices. Commonly used simplicial complexes include the Vietoris-Rips complex and the Alpha complex (see Appendix A).



Figure 1: Computation of ECC using sublevel set filtration in filtered cubical complex.

Filtration. A filtration $\mathcal{F} = \{K(a) \subset K | a \in \mathbb{R}\}$ is a collection of nested simplicial complexes that satisfy $K(a) \subset K(b)$ whenever $a \leq b$. A typical way of constructing a filtration is to use a monotonic filtration function $f: K \to \mathbb{R}$. f is monotonic in the sense that $f(\tau) \leq f(\sigma)$ whenever τ is a face of σ . By defining $K(a) \coloneqq f^{-1}(-\infty, a]$, we have $K(a) \subset K(b)$ whenever $a \leq b$.

Cubical Complex. Cubical complex is an analogy of simplicial complex that consists of k-cubes 129 (e.g., vertices, edges, squares, cubes, etc.). It provides a suitable framework for analyzing data that 130 is naturally aligned with a grid structure (e.g., digital images). An elementary interval is an interval 131 of form I = [l, l+1] or I = [l, l] for some $l \in \mathbb{Z}$, where the former interval is called *nondegenerate* 132 and the latter degenerate. An elementary cube is the finite product of elementary intervals, i.e., 133 $Q = I_1 \times I_2 \times \cdots \times I_n$. The *dimension* of Q is the number of nondegenerate elementary intervals in the product. P is a face of Q if $P \subset Q$ where P and Q are both elementary cubes. A *cubical* 134 *complex* K is a finite collection of elementary cubes such that the face of any cube in K is also in 135 K. A *filtered* cubical complex can be constructed by assigning a filtration value to each of the cubes 136 (see Appendix B for details). 137

Euler Characteristic. The *Euler characteristic* is a topological invariant that provides a single number summarizing the essential topological¹ features of data. Given a simplicial or cubical complex K, it is defined as the alternating sum of the number of k-simplices or k-cubes in K. It can equivalently be defined as an alternating sum of Betti numbers.

$$\zeta(K) = \sum_{k=0}^{\infty} (-1)^k |K^k| = \sum_{k=0}^{\infty} (-1)^k \beta_k,$$
(1)

where K^k is the set of k-dimensional simplices or cubes in K, $|K^k|$ is its cardinality, and β_k is the k-th Betti number of K. We can obtain an *Euler Characteristic Curve* (ECC) $C : \mathbb{R} \to \mathbb{R}$ by computing the Euler characteristic along a filtration, where the x-axis corresponds to the filtration values and y-axis corresponds to the Euler characteristic of the subcomplex at a given filtration value, i.e., for $t \in \mathbb{R}$, $C(t) = \chi(K(t))$ (see Figure 1).

150 151 152

153

142 143

144 145

146

147

148

149

122 123

3 LAYER CONSTRUCTION

)

The construction of ECLayr involves two steps: (i) computing the ECC from input data, and (ii) passing the ECC through a differentiable map. To compute ECC, we consider an alternative representation of the Euler characteristic. Let us denote K(t) as the subcomplex of K at a given filtration value t. Then, the Euler characteristic of K(t) in equation 1 can be equivalently defined as

$$\chi(K(t)) = \sum_{k=0}^{\infty} (-1)^k \sum_{\sigma \in K^k} \mathbb{1}\left[f_{\sigma} \le t\right],\tag{2}$$

¹Depending on the filtration, both PH and ECC can capture geometric information as well.

162 **Algorithm 1:** Computation of Euler Characteristic Curve: $X \to \mathcal{E}$ 163 1 Hyperparameters: T_{min}, T_{max}, v 164 ² Input: X 165 ³ Choose a simplicial complex suitable for the input data and build a filtration 166 4 Set $tseq = \{t_1, \ldots, t_v\}$, a sequence of v evenly-spaced discretized points from T_{min} to T_{max} 167 s Initialize $\mathcal{E} = (0, \dots, 0) \in \mathbb{R}^v$, which are values corresponding to locations in *tseq* 168 6 for $\sigma \in K$ do 169 if $f_{\sigma} > T_{max}$ then 7 170 continue 8 171 $t^* \leftarrow \min\{t_i \in tseq | t_i > f_\sigma\}$ 9 172 $\mathcal{E}(t^*) \leftarrow \mathcal{E}(t^*) + (-1)^{\dim(\sigma)}$ 10 173 11 $\mathcal{E} \leftarrow \text{cumsum}(\mathcal{E})$ 174 12 **Output**: $\mathcal{E} \in \mathbb{R}^v$ 175

176

177 178 179

where f_{σ} is the filtration value of σ . The equivalence between equation 1 and equation 2 is straightforward, as the sum of indicator functions is identical to the number of k-simplices in the subcomplex K(t). To simplify notation, let X, \mathcal{E} , and \mathcal{O}_{θ} represent the input, vectorized ECC, and output of our layer, respectively.

182 183

185 186

181

3.1 COMPUTATION OF ECC: $X \rightarrow \mathcal{E}$

Before calculating ECC from input data, a filtration must be defined by choosing an appropriate sim-187 plicial complex K and a function $f: K \to \mathbb{R}$. This is often data-dependent; Vietoris-Rips or Alpha 188 complexes are commonly used for point clouds while sub/superlevel set filtrations on filtered cubical 189 complexes are a natural choice for images. Upon constructing a filtration, we proceed to obtain the 190 vectorized approximation of ECC based on equation 2. First, we set a closed interval $[T_{min}, T_{max}]$ 191 and sample v evenly-spaced grid points ranging from T_{min} to T_{max} . We denote these discretized 192 points as $tseq = \{t_1, \ldots, t_v\}$, where $t_1 = T_{min}$ and $t_v = T_{max}$. Our objective is to derive a vector 193 \mathcal{E} containing the Euler Characteristics of each subcomplex $K(t_i)$; $\mathcal{E} = (\chi(K(t_1)), \dots, \chi(K(t_v)))$. 194 This vector serves as a finite sample approximation of the ECC function C. To compute \mathcal{E} , we begin 195 by initializing \mathcal{E} as a zero vector of size v. Next, we iterate over all simplices $\sigma \in K$ and perform 196 the following steps: (i) find $t^* = \min\{t_i \in tseq | t_i > f_{\sigma}\}$, which denotes the smallest grid point that is larger than the filtration value of σ , and (ii) add $(-1)^{\dim(\sigma)}$ to $\mathcal{E}(t^*)$. If f_{σ} exceeds the upper 197 bound T_{max} and t^* cannot defined, we proceed to the subsequent simplex in the iteration. Once the 199 iteration is terminated, we return the cumulative sum of \mathcal{E} for each point t_i . The resulting output \mathcal{E} is a vector in \mathbb{R}^{v} . The procedure is summarized in Algorithm 1. 200

Time Complexity. Given a filtration, computation of ECC (Steps 6 to 12 in Algorithm 1) requires O(N+v) time, where N is the number of simplices and v is the number of grid points. This shows a substantial improvement over the $O(N^3)$ of the standard PH computation. In our experimental analysis in Section 6.1, we empirically demonstrate that our proposed layer reduces the computational complexity up to several orders of magnitude compared to PH.

206 207 208

209

3.2 Computation of Layer Output: $\mathcal{E} ightarrow \mathcal{O}_{ heta}$

By definition, ECC depends on the number of generators (Betti numbers), even if they are small (noise) generators. This implies that ECC may potentially contain some noise information. Thus, we do not use ECC directly, but employ a differentiable parametrized map g_{θ} to project ECC to a learnable task-optimal representation. Given $\mathcal{E} \in \mathbb{R}^v$ and an output dimension m, the map g_{θ} : $\mathbb{R}^v \to \mathbb{R}^m$ takes \mathcal{E} as input and outputs $\mathcal{O}_{\theta} \in \mathbb{R}^m$. There are no restrictions regarding the structure of g_{θ} as long differentiability with respect to θ is guaranteed. In this paper, we use a sequence of fully connected layers and Relu nonlinearity functions to construct g_{θ} .

4 **STABLE BACKPROPAGATION**

We first present the differentiability result for ECC. By applying chain rule, we decompose the derivative of Euler characteristic into two elements: (i) derivative of the filtration value with respect to input X, and (ii) derivative of the indicator function with respect to the filtration value, as shown below.

216

217 218

219

220

221

225

226 227

228

229

230 231

232

The term (i) depends on the specific choice of filtration. While our framework allows the use of arbitrary differentiable filtration, here we focus on three widely-used filtrations; we provide results for Vietoris-Rips, Alpha, and sub/superlevel set filtration on filtered cubical complex in Appendix C.

 $\frac{\partial \chi(K(t))}{\partial X} = \sum_{k=0}^{\infty} (-1)^k \sum_{\sigma \in K^k} \underbrace{\frac{\partial f_{\sigma}}{\partial X}}_{(i)} \underbrace{\frac{\partial \mathbb{1}(f_{\sigma} \le t)}{\partial f_{\sigma}}}_{(ii)}.$

4.1**GRADIENT INCONSISTENCY IN SIGMOID APPROXIMATION**

233 The key barrier in achieving differentiability arises from the discontinuous nature of the indicator 234 function $\mathbb{1}(f_{\sigma} \leq t)$ in the term (ii). To bypass the need for direct differentiation, Röell & Rieck 235 (2024) adopted a smooth approximation by substituting the indicator function with a sigmoid func-236 tion $S(\lambda(t-f_{\sigma}))$, where the extra hyperparameter λ controls the precision of approximation. De-237 spite being theoretically sound, such smoothing-based approximation poses a practical problem in 238 backpropagation procedures as values are evaluated over a finite set of discretized grid points rather 239 than a continuous domain. When a filtration value f_{σ} lies between grid points, the gradient with respect to f_{σ} is not evaluated at the precise location, but at its neighboring grid points. Thus, the 240 magnitude of gradient varies depending on the proximity of f_{σ} to its adjacent grid points. This will 241 be referred to as *gradient inconsistency* (see Figure 2). 242

243 Especially, we show that the sigmoid approximation is prone to gradient vanishing problems when 244 insufficient v leads to excessive spacing between grid points, or when λ is too large (see Figure 245 2-(c)). To formally state this issue, we suppose the gradients of the indicator function $\mathbb{1}(f_{\sigma} \leq t)$ with respect to f_{σ} are approximated on a fixed grid $tseq = \{t_1, \ldots, t_v\}$ with $\Delta t := \frac{t_{i+1} - t_i}{2}$ being equal. Let $S'^{tseq}_{\lambda, f_{\sigma}} \in \mathbb{R}^v$ be the gradient vector of the sigmoid function $S(\lambda(t - f_{\sigma}))$ computed at 246 247 248 t_1, \ldots, t_v , i.e., 249 $S'^{tseq}_{\lambda,f_{\sigma}} = \frac{\partial S(\lambda(t-f_{\sigma}))}{\partial f_{\sigma}}|_{t=t_1,\ldots,t_v}.$

250

251 252

253

The next proposition shows that the local gradient of the sigmoid approximation can approach arbitrarily close to zero, regardless of its true value.

Proposition 4.1. When $\frac{\partial S(\lambda(t-f_{\sigma}))}{\partial f_{\sigma}}$ is viewed as a function of t, then its L_{∞} norm is computed as

$$\left\|\frac{\partial S(\lambda(t-f_{\sigma}))}{\partial f_{\sigma}}\right\|_{\infty} = \frac{\lambda}{4}$$

while its discretization over tseq is L_{∞} bounded as

$$\left\|S'^{tseq}_{\lambda,f_{\sigma}}\right\|_{\infty} \leq \lambda S(\lambda d(f_{\sigma}, tseq)) \left[1 - S(\lambda d(f_{\sigma}, tseq))\right].$$

So in particular when $\lambda \exp(-\lambda \Delta t) \rightarrow 0$,

$$\inf_{f_{\sigma} \in [t_1 - \Delta t, t_v + \Delta t)} \left\| S'^{tseq}_{\lambda, f_{\sigma}} \right\|_{\infty} \to 0$$

The issue of diminishing gradients during the training is particularly vexing in deep learning. As 268 the downstream gradient is computed via a series of multiplications involving local gradients, the 269 sigmoid approximation in ECC will eventually impede effective backpropagation.



Figure 2: An illustration of gradient inconsistency issues when using the sigmoid approximation with $\lambda = 200$. The values are assessed at 16 evenly-spaced points over [-1, 1]. In (a), the sigmoid approximation performs well when the filtration value precisely aligns with one of the grid points. In (b), however, a slight shift in the filtration value results in a significant change in the gradient. A further slight adjustment eventually leads to gradient vanishing, as shown in (c). This indicates that the gradients vary significantly when the sigmoid approximation is used for the ECC layer, depending on the position of the filtration value relative to the grid points.

4.2 STABLE BACKPROPAGATION VIA DISTRIBUTIONAL DERIVATIVES

287 To resolve the gradient inconsistency issue, here we propose an alternative approach in which we 288 approximate the gradient rather than the indicator function itself. In order to compute the gradient, we resort to distributional derivatives; $\frac{\partial \mathbb{I}[f_{\sigma} \leq t]}{\partial f_{\sigma}} = -\delta(t - f_{\sigma})$ where $\delta(x) = \lim_{\beta \to 0} \frac{1}{|\beta|\sqrt{\pi}} e^{-(x/\beta)^2}$ is the dirac delta, a function that has a single impulse at x = 0 and zero elsewhere. Since the height of this single impulse is infinite, we proceed with approximation $\max_{x} \frac{1}{|\beta|\sqrt{\pi}} e^{-(x/\beta)^2} = \frac{1}{|\beta|\sqrt{\pi}}$ 289 290 291 292 where β is a hyperparameter that determines the height of the spike. Namely, the estimated gradient 293 will always be $\frac{1}{|\beta|\sqrt{\pi}}$ at the impulse point and zero elsewhere. It is also critical to ensure that the gradient does not leak; while the approximated gradient has a single impulse at $t = f_{\sigma}$, it may not 295 necessarily correspond to the predefined positions *tseq*. Thus, we shift the location of impulse so 296 that it aligns with one of the points in *tseq*. Recall from Section 3.1 that for a given simplex σ , 297 $t^* = \min\{t_i \in tseq | t_i > f_{\sigma}\}$ is the grid point where the jump is reflected during the forward 298 pass. Consequently, we backpropagate the gradient to the identical location t^* . In this formulation, 299 the gradient invariably traverses one of the discretized locations, unless it was ignored during the 300 forward pass. As a result, we can assure that a consistent gradient value is properly backpropagated 301 to the preceding layer. The following proposition shows that our proposed method prevents the 302 gradient vanishing issues associated with the sigmoid approximation.

Proposition 4.2. Let $\hat{\delta}_{\beta,f_{\sigma}}^{tseq} \in \mathbb{R}^{v}$ be our gradient approximation of $\frac{\partial \mathbb{I}(f_{\sigma} \leq t)}{\partial f_{\sigma}}$ computed at t_{1}, \ldots, t_{v} , so $\hat{\delta}_{\beta,f_{\sigma}}^{tseq}$ is jth element is $-\frac{1}{\beta\sqrt{2\pi}}$ if $f_{\sigma} \in [t_{j-1}, t_{j})$, and other elements are 0. Then the L_{∞} norm of $\hat{\delta}_{\beta,f_{\sigma}}^{tseq}$ is given as

$$\left\|\hat{\delta}^{tseq}_{\beta,f_{\sigma}}\right\|_{\infty} = \frac{1}{\beta\sqrt{2\pi}}.$$

Aside from the issue of diminishing gradients, we have shown that our proposed approaches can achieve much lower errors in approximating true gradient values. Specifically, with $\Delta t := \frac{t_{i+1}-t_i}{2}$, $\forall i$, we show that by letting $\beta = \frac{\sqrt{\pi}}{2\Delta t}$, our proposed methods may attain consistency. We refer to Appendix D for detailed theoretical results.

Time Complexity of Sigmoid Approximation. Computation of ECC via sigmoid approximation requires O(vN) time, as the sigmoid function must be applied to every $t_i \in tseq$ during each iteration across all simplices $\sigma \in K$. Our backpropagation method enables the use of Algorithm 1 during forward pass, achieving enhanced efficiency of O(N + v).

319

303

304

305 306

307 308

310

311

312

313 314

5 STABILITY THEOREM

320 321

An essential benefit of using a topological layer is its robustness against noise. Extending the results of Dłotko & Gurnari (2023), we can establish a stability property for the layer output with respect to changes in the input. For notation, let X, X' be two distinct inputs, and $f_X, f_{X'}$ be corresponding 324 filtration functions on fixed simplicial complexes K, K', respectively. Let $\mathcal{D}_k(X), \mathcal{D}_k(X')$ be cor-325 responding k-dimensional persistence diagrams, and let $\mathcal{C}_X, \mathcal{C}_{X'} : \mathbb{R} \to \mathbb{R}$ be corresponding ECC 326 functions. See Appendix A for the definition of persistence diagrams and Wasserstein distance. 327

We first see the relation between the final layer output and ECC functions.

Proposition 5.1. Let $t_1^* < t_2^* < \cdots < t_w^*$ be unique values of all births and deaths in 329 $\{\mathcal{D}_k(X), \mathcal{D}_k(X') : k \ge 0\}$, and let $tseq = (t_1, \ldots, t_v)$. Suppose there exists $\Delta t > 0$ satisfy-330 ing that $\Delta t < t_{j+1} - t_j$ and $\Delta t < t_{j+1}^* - t_j^*$. Let g_{θ} be L-Lipschitz with respect to $\|\cdot\|_1$ -norm, i.e., 331 $||g_{\theta}(x) - g_{\theta}(y)||_{1} \leq L ||x - y||_{1}$. Then 332

$$\left\|\mathcal{O}_{\theta}(X) - \mathcal{O}_{\theta}(X')\right\|_{1} \leq \frac{2L}{\Delta t} \left\|\mathcal{C}_{X} - \mathcal{C}_{X'}\right\|_{1}.$$

335 Hence what we really need to establish is the stability of ECC functions. We first address the most 336 general stability result with respect to the 1-Wasserstein distance of the persistence diagrams of 337 input, which is directly from Dłotko & Gurnari (2023). 338

Proposition 5.2 (Dłotko & Gurnari (2023), Proposition 3.2). 339

$$\left\|\mathcal{C}_X - \mathcal{C}_{X'}\right\|_1 \le 2\sum_{k=0}^{\infty} W_1(\mathcal{D}_k(X), \mathcal{D}_k(X')).$$

343 The behavior of the 1-Wasserstein distance $W_1(\mathcal{D}_k(X), \mathcal{D}_k(X'))$ is in general complicated and 344 difficult to analyze. It is possible to further upper bound this by the difference of the filtration 345 functions f_X and $f_{X'}$. The difference is represented as L_{∞} distance below, but there is a more 346 general version of Theorem 5.3 as well. 347

Theorem 5.3. Suppose K = K' and is a finite simplicial complex or cubical complex. Then there exists a constant C_K only depending on K such that

$$\left|\mathcal{C}_{X} - \mathcal{C}_{X'}\right\|_{1} \le C_{K} \left\|f_{X} - f_{X'}\right\|_{\infty}$$

Theorem 5.3 provides a stability result whose relation to the difference of the input is clear, and also applicable to general filtration functions. Since we use DTM functions in Section 6, we present a specific result for DTM.

Corollary 5.4. Suppose K is a finite cubical complex, and f_X , $f_{X'}$ are restrictions of DTM functions $d_{P_X,m_0}, d_{P_{X'},m_0}$ to K, where $P_X, P_{X'}$ are empirical distributions on X and X', respectively. (for detailed meaning, see Appendix G.) Then

$$\|\mathcal{C}_X - \mathcal{C}_{X'}\|_1 \le \frac{C_K}{\sqrt{m_0}} W_2(P_X, P_{X'})$$

Due to the inherent reliance of Euler characteristics on even small generators, we note that the above stability results in terms of the Wasserstein distance are less strict than those bounded by the 362 Bottleneck distance in Kim et al. (2020). Thus, Euler characteristic-based descriptors compromise 363 stability in order to attain computational efficiency over PH-based descriptors. 364

365

361

328

333 334

340 341 342

348

349 350 351

352

353

354

355

356

6 **EXPERIMENTS**

366 367

To showcase the versatility and effectiveness of our layer, we conduct a series of experiments. First, 368 we demonstrate the computational efficiency of our approach by measuring runtime metrics across 369 different datasets. Next, we proceed to implement a topological autoencoder using point clouds to 370 illustrate an application of our layer in imposing topological constraints on the latent space. Fi-371 nally, we perform classification tasks on two image datasets: MNIST and Br35H. The first image 372 classification task shows that our layer can effectively mitigate information loss under conditions of 373 data scarcity or data contamination. The subsequent experiment highlights the distinct advantages 374 of our layer by performing operations on moderately high-dimensional data, which would otherwise 375 necessitate intensive computation for PH, rendering it impractical for real-world applications. All 376 experiments are implemented using GUDHI (The GUDHI Project, 2021) and Pytorch. Here, we present only a partial summary of the experimental findings; for comprehensive results and detailed 377 descriptions of the architecture and hyperparameter selection, please refer to Appendix I.

 Model
 Data (Number of samples)

 MNIST (60000)
 Br35H (209)
 Synth. (1000)

 ECC
 3.129 sec
 0.458 sec
 2.17 sec

 PH
 33.700 sec
 11.033 sec
 59.288 sec



Table 1: Average runtime performance per iteration (in seconds).

Figure 3: Latent representations of Spheres data.

6.1 COMPUTATIONAL EFFICIENCY

In this section, we analyze the empirical time complexity of our method in comparison to PH. The time complexity of each topological descriptor is assessed by measuring the runtime for a complete iteration through the training dataset, averaged over 10 repetitions. PH computes the persistence diagram using the GUDHI package, while ECC computes the vectorized approximation of ECC using Algorithm 1. In order to investigate how each descriptor scales with increasing data dimensions, we additionally generate a synthetic dataset containing 1000 samples of size 224×224 , where each pixel is randomly sampled from a uniform distribution. The experiment results for different datasets are provided in Table 1. We can observe that PH scales poorly as the dimension of the data increases. Considering the additional computation often required to transform persistence diagrams into alternative representations better suited for machine learning, our approach offers a significant benefit over all PH-based metrics in terms of computation, both in theory and in practice.

6.2 TOPOLOGICAL AUTOENCODER

The idea of imposing topological constraints on the latent space was first explored by Hofer et al. (2019); Moor et al. (2020). Whereas existing works rely on a topology-based loss term to regularize the latent space, our formulation allows for the utilization of standard loss functions, such as Mean Squared Error (MSE) or Mean Absolute Error (MAE) to achieve a similar goal. Inspired by the stability results regarding L_1 distance in Section 5, we employ the MAE loss between ECC of input and ECC of latent representation as our topological constraint. The respective ECCs are computed using Vietoris-Rips filtration, with maximum dimension set to 1. For the experiment, we use the synthetic Spheres dataset from Moor et al. (2020). The dataset consists of ten 100-spheres with radius r = 5 enclosed by one larger 100-sphere with radius = 25, all embedded in 101-dimension. The ten smaller spheres are shifted in random directions according to Gaussian noise.

Result. We discover that our approach effectively preserves the underlying shape of the encompass-ing sphere (yellow points in Figure 3), in contrast to the vanilla autoencoder, which loses this shape. Moreover, it constrains the smaller spheres to remain on the boundary of the encompassing sphere, whereas in the vanilla autoencoder, numerous smaller circles lie far beyond the boundaries of the encompassing circle. However, with this simplistic architecture, its capacity to comprehensively articulate the nested relationship inherent in the data was somewhat restricted. While our method demonstrates capability of regularizing the latent space, we do not claim superiority over alternative approaches. Rather, we present it as a motivating example of how topological characterization in the latent space can be promoted via simple standard loss functions.

6.3 CLASSIFICATION AGAINST DATA SCARCITY AND DATA CONTAMINATION

Our primary interest in this section is to demonstrate that our layer can effectively mitigate information loss under conditions of data scarcity and data contamination. For such purpose, we
consider two scenarios on the MNIST dataset. In the first scenario, we restrict the training data
to 100, 300, 500, 700 and 1000 samples to observe how model performance changes with data
size. In the second scenario, we consider a corruption and noise process where the pixels are randomly omitted and subsequently contaminated by random noise between 0 and 1 with probability 0.05, 0.1, 0.15, and 0.2.



Figure 4: MNIST test accuracy: (a) Performance across different sample sizes; (b) Performance in the presence of noise.

Table 2: Average runtime per epoch over 1000 MNIST data without noise (in seconds).

449 **Experimental Setup.** To impartially illustrate the advantages of our layer, we purposefully retain 450 a simple experimental setting. The base model consists of two CNN layers followed by two fully 451 connected layers. We compare the performance of our proposed layer with a base model, and two other topological layers applicable to image datasets: PersLay (Carrière et al., 2020) and PLLay 452 (Kim et al., 2020). For the data scarcity scheme, we additionally implement an ECLayr using the 453 sigmoid approximation (denoted as CNN + DECT) previously applied by (Röell & Rieck, 2024). 454 For all topological layers, we place a parallel layer at the beginning of the network (referred to as 455 CNN + EC(i), Pers, PL(i)). For topological layers that allow backpropagation, we add an additional 456 layer after the last convolutional layer (referred to as CNN + EC, DECT, PL). We implement su-457 perlevel cubical filtration for the experiment with varying data size. In the experiment involving 458 different noise levels, we employ the DTM filtration, a tool used in TDA to robustly extract topo-459 logical features in the presence noise (see Appendix A for further details). As DTM can control the 460 level of locality when extracting topological information, we place two parallel topological layers 461 with different scales at the beginning of the network when using DTM filtration. Utilizing a very 462 simple model on limited training samples, we observed random failures across all models with outliers significantly affecting the outcome. To remove the influence of outliers and solely evaluate 463 model performance, we repeat each experiment 15 times and select the top 10 test accuracies for as-464 sessment. 30% of the training data is used as a validation set, while model performance is evaluated 465 on the full test set. 466

467 **Result.** In Figure 4, we observe that by utilizing topological information, the performance of ECLayr consistently surpasses the baseline in all scenarios. Surprisingly, we notice that ECLayr 468 outperforms PH-based models despite the fact that PH is more informative than Euler Characteris-469 tics. We speculate that this phenomenon stems from an optimization process, coupled with the con-470 sideration that solely macroscopic topological features are adequate for this uncomplicated dataset. 471 PH provides multiple summaries for each homology dimension, which complicates optimization 472 in scenarios with limited data, whereas ECC yields a single summary for all dimensions. Conse-473 quently, the simplicity of ECC renders our layer more appropriate for scenarios with insufficient 474 data. We also observe that our model outperforms ECLayr with sigmoid approximation, supporting 475 the use of our proposed stable backpropagation method. Furthermore, ECLayr exhibits resistance 476 to approximately $5 \sim 10\%$ of data contamination compared to the baseline model. Nevertheless, the 477 inherent dependence of ECC on even small generators result results in our layer exhibiting reduced 478 noise resistance compared to, for instance, PersLay. Runtime metrics are provided in Table 2. Our 479 method scales approximately 20 to 30 times faster compared to PH-based methods, highlighting the significant improvement in computational efficiency. The high runtime of CNN + DECT results 480 from the increased time complexity of O(vN) when using sigmoid approximation. 481

482 483

484

445

446

447 448

6.4 CLASSIFICATION ON MODERATELY HIGH-DIMENSIONAL DATA

485 A significant drawback of PH is that its time complexity generally scales poorly with the dimension of data, rendering it impractical for high-dimensional data applications. Conversely, the computa-

487	Model	Test Accuracy	Runtime
488		81 705	
489	ResNet	(± 2.899)	0.110 sec
490		82.936	
491	ResNet + $EC(i)$	(± 2.520)	0.415 sec
492		84.351	
493	ResNet + EC	(± 3.308)	0.470 sec
101			

196

496 497

Table 3: Br35H test accuracy and runtime per epoch.

tional efficiency of ECC enables the use of moderately high-dimensional data without compromising
 significant computational costs. We demonstrate that our layer can enhance model performance by
 effectively exploiting topological information using a real-world dataset with dimension that would
 normally require intense computation for PH applications.

502 **Experimental Setup.** We conduct a binary classification task of detecting brain tumors on the Br35H dataset. The Br35H dataset consists of 3000 brain MRI images that have different size for 504 each dimension and a varying number of channels. We preprocess the data by cropping along the shorter dimension, resizing it to 112×112 , and converting it to grayscale. ResNet18 (He et al., 2016) 505 is employed as a baseline model, with an additional fully connected layer of size 64 appended at the 506 end of the network. For ResNet + EC(i), we add a parallel ECLayr at the beginning of the network 507 and concatenate the output with the residual layer output before feeding to the fully connected layer. 508 For ResNet + EC, we place an additional ECLayr before the first residual layer. As the task is a 509 simple binary classification problem, we only use 10% of the data as training samples and 30% of 510 training data is used for validation. Our training scheme utilizing limited data mirrors real-world 511 challenges, as access to medical data is often limited and costly. Each simulation is repeated 10 512 times, with the average test accuracy and average runtime per epoch reported in Table 3.

Result. The results in Table 3 show that our layer can enhance model performance while maintaining manageable computational costs on moderately high dimensional data. Furthermore, it suggests that our layer can be effectively integrated with large models such as ResNet for practical usage. Another interesting observation is that using an additional ECLayr before the first residual layer yields further improvement in performance. Fully exploiting the computational efficiency of ECC, our layer facilitates operations on moderately high-dimensional data that would be impractical for PH, highlighting the significance of ECLayr for real-world applications.

7 DISCUSSION

521 522 523

520

ECLayr is a novel topological layer that offers computationally efficiency and stable backprop-524 agation, allowing for seamless integration into a wide range of deep learning architectures while 525 enhancing both robustness and convergence behavior. Our proposed layer can be used generically 526 for an extensive variety of data structures as long as the filtration is differentiable with respect to 527 the input data. Nonetheless, there are some important caveats and limitations which should be 528 addressed. First, while ECCs offer computational efficiency, PH-based summaries provide more detailed, multi-scale topological information. Understanding this tradeoff is essential. Therefore, 529 our proposed ECLayr is particularly well-suited for applications where computational efficiency 530 is prioritized over detailed topological insights. Moreover, as discussed in greater detail in Section 531 5, ECCs are topologically weaker invariants compared to PHs. Consequently, the ECC-based layer 532 generally exhibits less robustness than the PH-based layers. Next, as with other topological layers, 533 further research is necessary to achieve successful systematic hyperparameter exploration. Finally, 534 extending our analysis to other filtrations, such as the clique complex of a multigraph, and applying 535 ECLayr to time-series embeddings (Kim et al., 2018; Umeda, 2017) would be a valuable direction 536 for future research, which could further demonstrate the versatility of our proposed methods.

- 537 538
- 530

540 REFERENCES 541

549

554

558

559

560

561

567

568

569

573

- Henry Adams, Tegan Emerson, Michael Kirby, Rachel Neville, Chris Peterson, Patrick Shipman, 542 Sofya Chepushtanova, Eric Hanson, Francis Motta, and Lori Ziegelmeier. Persistence images: A 543 stable vector representation of persistent homology. Journal of Machine Learning Research, 18 544 (8):1-35, 2017.
- 546 Hirokazu Anai, Frédéric Chazal, Marc Glisse, Yuichi Ike, Hiroya Inakoshi, Raphaël Tinarrage, and 547 Yuhei Umeda. Dtm-based filtrations. In Topological Data Analysis: The Abel Symposium 2018, 548 pp. 33-66. Springer, 2020.
- Gabriele Beltramo, Rayna Andreeva, Ylenia Giarratano, Miguel O Bernabeu, Rik Sarkar, and Pri-550 moz Skraba. Euler characteristic surfaces. arXiv preprint arXiv:2102.08260, 2021. 551
- 552 Peter Bubenik et al. Statistical topological data analysis using persistence landscapes. J. Mach. 553 Learn. Res., 16(1):77-102, 2015.
- Dmitri Burago, Yuri Burago, and Sergei Ivanov. A course in metric geometry, volume 33 of Grad-555 uate Studies in Mathematics. American Mathematical Society, Providence, RI, 2001. ISBN 556 0-8218-2129-6.
 - Gunnar Carlsson and Rickard Brüel Gabrielsson. Topological approaches to deep learning. In Topological Data Analysis: The Abel Symposium 2018, pp. 119–146. Springer, 2020.
- Mathieu Carrière, Frédéric Chazal, Yuichi Ike, Théo Lacombe, Martin Royer, and Yuhei Umeda. Perslay: A neural network layer for persistence diagrams and new graph topological signatures. 562 In International Conference on Artificial Intelligence and Statistics, pp. 2786–2796. PMLR, 2020. 563
- 564 Mathieu Carriere, Frédéric Chazal, Marc Glisse, Yuichi Ike, Hariprasad Kannan, and Yuhei Umeda. 565 Optimizing persistent homology based functions. In International conference on machine learn-566 ing, pp. 1294–1303. PMLR, 2021.
 - Frédéric Chazal and Bertrand Michel. An introduction to topological data analysis: fundamental and practical aspects for data scientists. Frontiers in artificial intelligence, 4:108, 2021.
- 570 Frédéric Chazal, David Cohen-Steiner, Marc Glisse, Leonidas J Guibas, and Steve Y Oudot. Prox-571 imity of persistence modules and their diagrams. In Proceedings of the twenty-fifth annual sym-572 posium on Computational geometry, pp. 237-246. ACM, 2009.
- Frédéric Chazal, David Cohen-Steiner, and Quentin Mérigot. Geometric inference for probability 574 measures. Foundations of Computational Mathematics, 11:733–751, 2011. 575
- 576 Frédéric Chazal, Vin de Silva, and Steve Oudot. Persistence stability for geometric complexes. 577 Geom. Dedicata, 173:193-214, 2014. ISSN 0046-5755. doi: 10.1007/s10711-013-9937-z. URL 578 https://doi.org/10.1007/s10711-013-9937-z.
- Frédéric Chazal, Vin de Silva, Marc Glisse, and Steve Oudot. The structure and stability of per-580 sistence modules. SpringerBriefs in Mathematics. Springer, [Cham], 2016a. ISBN 978-3-319-581 42543-6; 978-3-319-42545-0. doi: 10.1007/978-3-319-42545-0. URL https://doi.org/ 582 10.1007/978-3-319-42545-0. 583
- 584 Frédéric Chazal, Pascal Massart, and Bertrand Michel. Rates of convergence for robust geometric inference. 2016b. 585
- 586 Yuzhou Chen, Ignacio Segovia-Dominguez, Baris Coskunuzer, and Yulia Gel. Tamp-s2gcnets: coupling time-aware multipersistence knowledge representation with spatio-supra graph convolu-588 tional networks for time-series forecasting. In International Conference on Learning Representa-589 tions, 2022. 590
- David Cohen-Steiner, Herbert Edelsbrunner, John Harer, and Yuriy Mileyko. Lipschitz functions have L_p -stable persistence. Found. Comput. Math., 10(2):127–139, 2010. ISSN 592 doi: 10.1007/s10208-010-9060-6. URL https://doi.org/10.1007/ 1615-3375. s10208-010-9060-6.

- 594 Paweł Dłotko and Davide Gurnari. Euler characteristic curves and profiles: a stable shape invariant for big data problems. *GigaScience*, 12:giad094, 2023. 596 H. Edelsbrunner and J. Harer. Computational Topology: An Introduction. Applied Mathematics. 597 American Mathematical Society, 2010. ISBN 9780821849255. 598 Herbert Edelsbrunner and Ernst P Mücke. Three-dimensional alpha shapes. ACM Transactions On 600 Graphics (TOG), 13(1):43-72, 1994. 601 Rickard Brüel Gabrielsson, Bradley J. Nelson, Anjan Dwaraknath, and Primoz Skraba. A topology 602 layer for machine learning. In Silvia Chiappa and Roberto Calandra (eds.), Proceedings of the 603 Twenty Third International Conference on Artificial Intelligence and Statistics, volume 108 of 604 Proceedings of Machine Learning Research, pp. 1553–1563. PMLR, 26–28 Aug 2020. URL 605 https://proceedings.mlr.press/v108/gabrielsson20a.html. 606 Marcio Gameiro, Yasuaki Hiraoka, and Ippei Obayashi. Continuation of point clouds via persistence 607 diagrams. *Physica D: Nonlinear Phenomena*, 334:118–132, 2016. 608 609 Olympio Hacquard and Vadim Lebovici. Euler characteristic tools for topological data analysis. 610 arXiv preprint arXiv:2303.14040, 2023. 611 A. Hatcher. Algebraic Topology. Algebraic Topology. Cambridge University Press, 2002. ISBN 612 9780521795401. 613 614 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-615 nition. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 616 770–778, 2016. 617 Felix Hensel, Michael Moor, and Bastian Rieck. A survey of topological machine learning methods. 618 Frontiers in Artificial Intelligence, 4:681108, 2021. 619 620 Christoph Hofer, Roland Kwitt, Marc Niethammer, and Andreas Uhl. Deep learning with topologi-621 cal signatures. Advances in neural information processing systems, 30, 2017. 622 Christoph Hofer, Roland Kwitt, Marc Niethammer, and Mandar Dixit. Connectivity-optimized rep-623 resentation learning via persistent homology. In International conference on machine learning, 624 pp. 2751–2760. PMLR, 2019. 625 Shengli Jiang, Nanqi Bao, Alexander D Smith, Shraddha Byndoor, Reid C Van Lehn, Manos 626 Mavrikakis, Nicholas L Abbott, and Victor M Zavala. Scalable extraction of information from 627 spatiotemporal patterns of chemoresponsive liquid crystals using topological descriptors. The 628 Journal of Physical Chemistry C, 127(32):16081–16098, 2023. 629 630 Tomasz Kaczynski, Konstantin Michael Mischaikow, and Marian Mrozek. Computational homol-631 ogy, volume 157. Springer, 2004. 632 Kwangho Kim, Jisu Kim, and Alessandro Rinaldo. Time series featurization via topological data 633 analysis. arXiv preprint arXiv:1812.02987, 2018. 634 635 Kwangho Kim, Jisu Kim, Manzil Zaheer, Joon Kim, Frédéric Chazal, and Larry Wasserman. Pllay: 636 Efficient topological layer based on persistent landscapes. Advances in Neural Information Pro-637 cessing Systems, 33:15965–15977, 2020. 638 Daniel J Laky and Victor M Zavala. A fast and scalable computational topology framework for the 639 euler characteristic. Digital Discovery, 2024. 640 641 Jacob Leygonie, Steve Oudot, and Ulrike Tillmann. A framework for differential calculus on persistence barcodes. Foundations of Computational Mathematics, pp. 1–63, 2022. 642 643 Nicholas O Malott and Philip A Wilsey. Scalable homology classification through decomposed 644 euler characteristic curves. In 2023 IEEE International Conference on Big Data (BigData), pp. 645 768-777. IEEE, 2023. 646
- 647 Michael Moor, Max Horn, Bastian Rieck, and Karsten Borgwardt. Topological autoencoders. In *International conference on machine learning*, pp. 7045–7054. PMLR, 2020.

648 649	Nina Otter, Mason A Porter, Ulrike Tillmann, Peter Grindrod, and Heather A Harrington. A roadmap for the computation of persistent homology. <i>EPJ Data Science</i> , 6:1–38, 2017.
650 651 652 653 654 655 656 657 658	Theodore Papamarkou, Tolga Birdal, Michael M. Bronstein, Gunnar E. Carlsson, Justin Curry, Yue Gao, Mustafa Hajij, Roland Kwitt, Pietro Lio, Paolo Di Lorenzo, Vasileios Maroulas, Nina Mi- olane, Farzana Nasrin, Karthikeyan Natesan Ramamurthy, Bastian Rieck, Simone Scardapane, Michael T Schaub, Petar Veličković, Bei Wang, Yusu Wang, Guowei Wei, and Ghada Zamzmi. Position: Topological deep learning is the new frontier for relational learning. In Ruslan Salakhut- dinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), <i>Proceedings of the 41st International Conference on Machine Learning</i> , vol- ume 235 of <i>Proceedings of Machine Learning Research</i> , pp. 39529–39555. PMLR, 21–27 Jul 2024. URL https://proceedings.mlr.press/v235/papamarkou24a.html.
660 661	Eitan Richardson and Michael Werman. Efficient classification using the euler characteristic. <i>Pattern Recognition Letters</i> , 49:99–106, 2014.
662 663 664 665	Ernst Röell and Bastian Rieck. Differentiable euler characteristic transforms for shape classi- fication. In International Conference on Learning Representations, 2024. URL https: //openreview.net/forum?id=MO632iPq3I.
666 667	The GUDHI Project. <i>GUDHI User and Reference Manual</i> . GUDHI Editorial Board, 3.4.1 edition, 2021. URL https://gudhi.inria.fr/doc/3.4.1/.
668 669	Katharine Turner, Sayan Mukherjee, and Doug M Boyer. Persistent homology transform for model- ing shapes and surfaces. <i>Information and Inference: A Journal of the IMA</i> , 3(4):310–344, 2014.
671 672	Yuhei Umeda. Time series classification via topological data analysis. <i>Information and Media Technologies</i> , 12:228–239, 2017.
673	
675	
676	
677	
678	
679	
680	
681	
682	
683	
684	
685	
686	
687	
688	
689	
690	
691	
692	
693	
694	
695	
696	
697	
698	
699	
700	
701	

703 704 705

706

708

709

714

723

742 743 744

751 752 753

APPENDIX

MORE BACKGROUNDS IN TOPOLOGICAL DATA ANALYSIS А

We briefly review basic concepts in Topological Data Analysis that are needed to develop stability results in Section 5 of this paper, mainly coming from Kim et al. (2020). We refer interested readers to Chazal & Michel (2021); Hatcher (2002); Edelsbrunner & Harer (2010); Chazal et al. (2009; 710 2016a) for details and formal definitions.

711 **Vietoris-Rips Complex.** Let X be a finite set of points in \mathbb{R}^d . For r > 0, the Vietoris-Rips complex 712 is a collection of simplices where the distance between any two vertices is smaller than 2r: 713

$$\operatorname{Rips}(r) = \{ \sigma \subset X | d(u_i, u_j) < 2r, \forall u_i, u_j \in \sigma \}.$$

715 Notice that $\operatorname{Rips}(r_1) \subset \operatorname{Rips}(r_2)$ when $r_1 \leq r_2$. Thus, we can build a filtration on the Vietoris-Rips 716 complex by monotonically increasing r.

717 Alpha Complex. Let X be a finite set of points in \mathbb{R}^d . For each $u_i \in X$, the Voronoi cell of u_i is 718 the set of points that are closest to u_i ; $V_{u_i} = \{x \in \mathbb{R}^d | d(u_i, x) \le d(u_j, x), \forall u_j \in X, u_j \ne u_i\}$. For 719 r > 0 and each $u_i \in X$, let us denote the closed *r*-ball with center u_i and radius *r* as $B_{u_i}(r)$. Then, 720 we define $R_{u_i}(r) = B_{u_i}(r) \cap V_{u_i}$, which is the intersection of each r-ball with its corresponding 721 Voronoi cell. The Alpha complex is a collection of simplices such that all $R_{u_i}(r)$ of the vertices in 722 the simplex have an intersection:

 $Alpha(r) = \{ \sigma \subset X | \cap_{u_i \in \sigma} R_{u_i}(r) \neq \emptyset \}.$

724 Similar to the Vietoris-Rips complex, we can build a filtration on the Alpha complex by monotoni-725 cally increasing r. 726

Persistent Homology and Persistence Diagram. Persistent homology is a multiscale approach 727 to represent the topological features of the complex K, and can be represented in the persistence 728 diagram. For a filtration \mathcal{F} and for each nonnegative k, we keep track of when k-dimensional 729 homological features (e.g., 0-dimension: connected component, 1-dimension: loop, 2-dimension: 730 cavity,...) appear and disappear in the filtration. If a homological feature α_i appears at b_i and 731 disappears at d_i , then we say α_i is born at b_i and dies at d_i . By considering these pairs (b_i, d_i) as 732 points in the plane, one obtains the persistence diagram defined as follows. 733

Definition A.1. Let $\mathbb{R}^2_* := \{(b, d) \in (\mathbb{R} \cup \infty)^2 : d > b\}$. A persistence diagram \mathcal{D} is a finite multiset 734 of $\{(b_i, d_i) : (b_i, d_i) \in \mathbb{R}^2_*\}$. 735

736 Wasserstein Distance. We suggest two versions of Wasserstein distances, one is for persistence 737 diagrams and the other is for probability measures.

738 We first start with Wasserstein distance for persistence diagrams. A *matching* between two persis-739 tence diagrams \mathcal{D}_1 and \mathcal{D}_2 , is a subset $m \subset \mathcal{D}_1 \times \mathcal{D}_2$ such that every off-diagonal point in \mathcal{D}_1 and 740 \mathcal{D}_2 only appears once in m. The p-Wasserstein distance between persistence diagrams is defined by 741

$$W_p(\mathcal{D}_1, \mathcal{D}_2) = \inf_{\text{matching } m} \left(\sum_{(x, y) \in m} \|x - y\|_{\infty}^p \right)^{1/p}$$

745 Now we see Wasserstein distance for persistence diagrams. Let P and Q be probability measures 746 on \mathcal{X} , and let $\mathcal{J}(P,Q)$ denote all joint distributions J for $\mathcal{X} \times \mathcal{X}$ that have marginals P and Q. In 747 other words, $(\Pi_1)_{\#}J = P$ and $(\Pi_2)_{\#}J = Q$ where $\Pi_1(x, y) = x$ and $\Pi_2(x, y) = y$, and $T_{\#}P$ is a 748 push-forward measure of P, i.e., $T_{\#}P(A) = P(\{x : T(x) \in A\}) = P(T^{-1}(A))$. For $p \ge 1$, the 749 Kantorovich, or Wasserstein, distance is 750

$$W_p(P,Q) = \left(\inf_{J \in \mathcal{J}(P,Q)} \int_{\mathcal{X} \times \mathcal{X}} ||x - y||^p dJ(x,y)\right)^{1/p}$$

Gromov-Hausdorff distance. The Hausdorff distance is on sets embedded in the same metric 754 spaces. This distance measures how two sets are close to each other in the embedded metric space. 755 When $S \subset \mathbb{X}$, we denote by S^r the r-neighborhood of a set S in \mathbb{R}^d , i.e. $S^r = \bigcup_{x \in S} B_x(r)$.

756 **Definition A.2** (Hausdorff distance (Burago et al., 2001, Definition 7.3.1)). Let $X, Y \subset X$ be subsets of \mathbb{R}^d . The Hausdorff distance between X and Y, denoted by $d_H(X,Y)$, is defined as 758

$$d_H(X,Y) \coloneqq \inf \left\{ r > 0 : X \subset Y^r \text{ and } Y \subset X^r \right\}.$$

760 The notion of the Hausdorff distance can be generalized to the comparison of any pair of metric 761 spaces. The Gromov-Hausdorff distance measures how two sets are far from being isometric to each 762 other. 763

Definition A.3 ((Burago et al., 2001, Definition 7.3.10)). Let X and Y be two metric spaces. The 764 Gromov-Hausdorff distance between X and Y, denoted by $d_{GH}(X, Y)$, is defined as 765

$$d_{GH}(X,Y) \coloneqq \inf\{d_H(X',Y') : \text{there exists a metric space } Z \text{ and } X',Y' \subset Z$$

with X,Y isometric to X',Y', respectively.}

769 Distance to measure. Distance to measure (DTM) (Chazal et al., 2011; 2016b; Anai et al., 2020) 770 is a distance-like function² that is robust to outliers. For a probability measure μ and parameters $m_0 \in [0,1)$ and $r \ge 1$ (default is r = 2), the DTM function $\hat{d}_{\mu,m_0} : \mathbb{R}^d \to \mathbb{R}$ is defined as 771

$$d_{\mu,m_0}(x) = \left(\frac{1}{m_0} \int_0^{m_0} \delta_{\mu,m}^r(x) dm\right)^{1/r}$$

775 where $\delta_{\mu,m}(x) = \inf\{t > 0 | \mu(B_x(t)) > m\}$ and $B_x(t)$ is a closed t-ball centered at x. In practice, 776 an empirical DTM is used. If input data X is considered as weights corresponding to fixed points Y, 777

$$\hat{d}_{m_0}(x) = \left(\frac{\sum_{Y_i \in N_k(x)} X_i' \|Y_i - x\|^r}{m_0 \sum_{i=1}^n X_i}\right)^{1/r},\tag{3}$$

where $N_k(x)$ is a subset of Y containing the k nearest neighbors of x. k is such that satisfies 781 $\sum_{Y_i \in N_{k-1}(x)} X_i < m_0 \sum_{i=1}^n X_i \le \sum_{Y_i \in N_k(x)} X_i, \text{ and } X'_i = \sum_{Y_j \in N_k(x)} X_j - m_0 \sum_{j=1}^n X_j \text{ if at least one of } Y_i \text{'s is in } N_k(x) \text{ and } X'_i = X_i \text{ otherwise (see Figure 5 (b)).}$ 782 783

When input data is considered as empirical data points, the empirical DTM becomes 785

$$\hat{d}_{m_0}(x) = \left(\frac{\sum_{X_i \in N_k(x)} w_i' \|X_i - x\|^r}{m_0 \sum_{i=1}^n w_i}\right)^{1/r}$$

where $N_k(x)$ is a subset of X containing the k nearest neighbors of x. k is such that satisfies 789 $\sum_{X_i \in N_{k-1}(x)} w_i < m_0 \sum_{i=1}^n w_i \le \sum_{X_i \in N_k(x)} w_i$, and $w'_i = \sum_{X_j \in N_k(x)} w_j - m_0 \sum_{j=1}^n w_j$ if at 790 least one of X_i 's is in $N_k(x)$ and $w'_i = w_i$ otherwise. 791

792 The parameter m_0 determines how much local/global structures should be extracted, with smaller 793 m_0 corresponding to more local structures. The DTM function is differentiable (Kim et al., 2020), 794 and adopting a sublevel or superlevel set filtration on the DTM transformed data yields a DTM filtration that is robust to outliers.

796 797 798

759

766 767 768

772 773 774

778

779

784

786 787 788

CONSTRUCTING FILTERED CUBICAL COMPLEXES FROM IMAGE DATA В

799 Let $X \in \mathbb{R}^{H \times W}$ be a 2D image. There are two methods of constructing a filtered cubical complex: 800 T-construction and V-construction. 801

T-construction In *T-construction*, each pixel in the image is mapped to a top-dimensional cell in 802 the cubical complex, which is a square in case of 2D images. The filtration value of each square 803 is assigned as the intensity of its corresponding pixel, and these filtration values are recursively 804 extended to lower dimensional cubes. The filtration value of each edge is assigned as the minimum 805 of the filtration values of its neighboring squares. Similarly, the filtration value of each vertex is 806 assigned as the minimum of the filtration values its neighboring edges. 807

⁸⁰⁸ ²This distance function is not the distance function giving a metric between two input points such as l_p 809 distance, but rather measures a distance between a single input point and the support set of a probability distribution.

V-construction In V-construction, each pixel in the image is mapped to a vertex in the cubical complex. The filtration value of each vertex is assigned as the intensity of its corresponding pixel, and these filtration values are recursively extended to higher dimensional cubes. The filtration value of each edge is assigned as the maximum of the filtration values of its neighboring vertices. Similarly, the filtration value of each square is assigned as the maximum of the filtration values of its neighboring edges.

For both constructions, a sublevel set at a given filtration value t defines a subcomplex $K(t) := \{\sigma \in K | f(\sigma) \le t\}$; the collection of cubes with filtration value less than or equal to t. Consequently, a sublevel set filtration can be built by monotonically increasing t. A superlevel set filtration can also be obtained by applying the sublevel set filtration to a cubical complex constructed from -X rather than X. In case of 2D images with multiple channels, such as color images represented by RGB channels where $X \in \mathbb{R}^{3 \times H \times W}$, cubical complexes are constructed independently for each channel.

822 823

824 825

826

832 833

834

835 836

838

839

840

841

842

843

844

845 846

847

851 852 853 C DERIVATIVE OF FILTRATION VALUE WITH RESPECT TO INPUT X: $\frac{\partial f(\sigma)}{\partial X}$

C.1 VIETORIS-RIPS FILTRATION

827 Assume Vietoris-Rips general position for a point cloud X: (i) all points in X are unique, and (ii) 828 the length of all attaching edges are unique. The filtration value of a simplex σ in the Vietoris-Rips 829 filtration is half the length of the longest edge in σ . This edge is the *attaching edge* of σ , denoted as 830 τ_{σ} . Letting x_i and x_j be the vertices of τ_{σ} , the derivatives of filtration value $f(\sigma) = \frac{\|x_i - x_j\|}{2}$ with 831 respect to the points x_i and x_j are given by (Gameiro et al., 2016):

$$\frac{\partial f(\sigma)}{\partial x_i} = \frac{1}{2} \frac{x_i - x_j}{\|x_i - x_j\|}, \quad \frac{\partial f(\sigma)}{\partial x_j} = \frac{1}{2} \frac{x_j - x_i}{\|x_i - x_j\|}.$$
(4)

The derivatives with respect to points other than x_i and x_j are all zero.

837 C.2 Alpha Filtration

Assume Alpha general position of a point cloud X: (i) general position in the sense of Edelsbrunner & Mücke (1994), and (ii) filtration values of all attaching simplices are unique. In Alpha filtration, all simplices are either an attaching simplex, or a simplex attached by another simplex of higher dimension. In the latter case, filtration value of the attached simplex is given by the filtration value of its attaching simplex. The filtration value of an attaching simplex σ is the radius of the smallest circumcircle of σ (Edelsbrunner & Mücke, 1994; Gameiro et al., 2016) and it can be differentiated with respect to the coordinates of each of the vertices.

C.3 SUB/SUPERLEVEL SET FILTRATION ON FILTERED CUBICAL COMPLEXES

Let us treat a 2D image $X \in \mathbb{R}^{H \times W}$ as a vector $x = (x_1, \dots, x_{HW}) \in \mathbb{R}^{HW}$, where the elements of the vector are arranged in row-major order. Then, the derivative of the filtration value with respect to the input data can be written as

$$\frac{\partial f(\sigma)}{\partial x} = \left(\frac{\partial f(\sigma)}{\partial x_1}, \dots, \frac{\partial f(\sigma)}{\partial x_{HW}}\right)$$

Given that the filtration value varies depending on the construction used, we provide differentiability results for both T-construction and V-constructions. For simplicity of notation, we denote $\mathcal{I} = \{1, 2, \dots, HW\}$ as the index set.

T-construction. In T-construction, each pixel is mapped to a square, with the pixel intensity serving as the filtration value of the corresponding square. Thus, we first explore the scenario where σ is a square, and then extend our analysis to lower dimensional cubes.

(i) Assume σ is a square, i.e., dim $(\sigma) = 2$. Let $j \in \mathcal{I}$ denote the index of the pixel in x that corresponds to σ . Then, $f(\sigma) = x_j$ and thus,

$$\frac{\partial f(\sigma)}{\partial x_i} = \begin{cases} 1, & \text{if } i = j \\ 0, & \text{otherwise} \end{cases}$$

for all $i \in \mathcal{I}$.

870

871

872

873 874

878 879

880

891

892

893

894

897

899 900

901

902

903

904

905 906

907

908

909

914 915

916

(ii) Assume σ is an edge, i.e., dim $(\sigma) = 1$. Recall that $f(\sigma)$ is assigned as the minimum filtration value of its neighboring squares, which in turn is equivalent to the minimum pixel intensity of the pixels corresponding to those neighboring squares. Thus, we can identify the pixel associated with σ by

- 1. find neighboring squares of σ
- 2. determine the neighboring square with minimum filtration value
- 3. identify the pixel that corresponds to the square found in (2)

In step 2, multiple neighboring squares may have the same minimum filtration value. In this case, we identify the set of pixels that corresponds to all such squares. Letting $J \subset \mathcal{I}$ denote an index set labeling the members of such set of pixels,

$$\frac{\partial f(\sigma)}{\partial x_i} = \begin{cases} 1/|J|, & \text{if } i \in J\\ 0, & \text{otherwise} \end{cases}$$

for all $i \in \mathcal{I}$. Observe that when multiple pixels contribute to σ , we distribute the gradient evenly between those pixels.

(iii) Assume σ is a vertex, i.e., dim $(\sigma) = 0$. Recall that $f(\sigma)$ is assigned as the minimum filtration value of its neighboring edges. Therefore, once we find the neighboring edge(s) with minimum filtration value, we can repeat the process in (ii) to identify the set of pixels associated with σ . Letting $J \subset \mathcal{I}$ denote an index set labeling the members of such set of pixels,

$$\frac{\partial f(\sigma)}{\partial x_i} = \begin{cases} 1/|J|, & \text{if } i \in J\\ 0, & \text{otherwise} \end{cases}$$

for all $i \in \mathcal{I}$.

V-construction. In V-construction, each pixel is mapped to a vertex, with the pixel intensity serving as the filtration value of the corresponding vertex. Thus, we first explore the scenario where σ is a vertex, and then extend our analysis to higher dimensional cubes.

(i) Assume σ is a vertex, i.e., dim $(\sigma) = 0$. Let $j \in \mathcal{I}$ be the index of the pixel in x that corresponds to σ . Then, $f(\sigma) = x_j$ and thus,

$$\frac{\partial f(\sigma)}{\partial x_i} = \begin{cases} 1, & \text{if } i = j\\ 0, & \text{otherwise} \end{cases}$$

for all $i \in \mathcal{I}$.

(ii) Assume σ is an edge, i.e., dim $(\sigma) = 1$. Recall that $f(\sigma)$ is assigned as the maximum filtration value of its neighboring vertices, which in turn is equivalent to the maximum pixel intensity of the pixels corresponding to those neighboring vertices. Thus, we can identify the pixel associated with σ by

- 1. find neighboring vertices of σ
- 2. determine the neighboring vertex with maximum filtration value
- 3. identify the pixel that corresponds to the vertex found in (2)

In step 2, multiple neighboring vertices may have the same maximum filtration value. In this case, we identify the set of pixels that corresponds to all such vertices. Letting $J \subset \mathcal{I}$ denote an index set labeling the members of such set of pixels,

$$\frac{\partial f(\sigma)}{\partial x_i} = \begin{cases} 1/|J|, & \text{if } i \in J\\ 0, & \text{otherwise} \end{cases}$$

for all $i \in \mathcal{I}$. Observe that when multiple pixels contribute to σ , we distribute the gradient evenly between those pixels.



Figure 5: (b) is the DTM transformation of (a) using $m_0 = 0.05$ on a 28×28 unit grid. (c) & (d) visualize the respective gradients of ECC and persistence landscape with respect to (b). The gradient of ECC provides more detailed and interpretable information compared to the gradient of persistence landscapes, which is very sparse.

(iii) Assume σ is a square, i.e., dim $(\sigma) = 2$. Recall that $f(\sigma)$ is assigned as the maximum filtration value of its neighboring edges. Therefore, once we find the neighboring edge(s) with maximum filtration value, we can repeat the process in (ii) to identify the set of pixels associated with σ . Letting $J \subset \mathcal{I}$ denote an index set labeling the members of such set of pixels,

$$\frac{\partial f(\sigma)}{\partial x_i} = \begin{cases} 1/|J|, & \text{if } i \in J\\ 0, & \text{otherwise} \end{cases}$$

for all $i \in \mathcal{I}$.

D APPROXIMATION OF GRADIENTS

For more detailed theoretical analysis of the approximations of the gradients, we suppose that the algorithm is to approximate the gradients of the indicator function $\mathbb{I}(f_{\sigma} \leq t)$ with respect to f_{σ} on a fixed grid $t \in tseq = \{t_1, \ldots, t_v\}$, with $\Delta t := \frac{t_{i+1}-t_i}{2}$ being equal. Suppose the algorithm outputs approximations of gradients $\frac{\partial \mathbb{I}(f_{\sigma} \leq t)}{\partial f_{\sigma}}|_{t=t_1,\ldots,t_v}$ as g_1, \ldots, g_v , we treat that the gradient $\frac{\partial \mathbb{I}(f_{\sigma} \leq t)}{\partial f_{\sigma}}$ at t is approximated as g_1 on $t \in [t_1 - \Delta t, t_1 + \Delta t)$, g_2 on $t \in [t_2 - \Delta t, t_2 + \Delta t)$, and g_v on $t \in [t_v - \Delta t, t_v + \Delta t)$, where g_j can depend on f_{σ} . Hence the corresponding approximation $g : [t_1 - \Delta t, t_v + \Delta t)$ is

$$q(t) = g_j, \quad \text{for } t \in [t_j - \Delta t, t_j + \Delta t)$$

If g is a good approximation of $\frac{\partial \mathbb{I}(f_{\sigma} \leq t)}{\partial f_{\sigma}}$, then

$$\int g(t) df_{\sigma} \approx \int \frac{\partial \mathbb{I}(f_{\sigma} \leq t)}{\partial f_{\sigma}} df_{\sigma} \qquad \text{for each } t \in (t_1 - \Delta t, t_v + \Delta t),$$

and

$$\int g(t)dt \approx \int \frac{\partial \mathbb{I}(f_{\sigma} \leq t)}{\partial f_{\sigma}} dt \qquad \text{for each } f_{\sigma} \in (t_1 - \Delta t, t_v + \Delta t)$$

We will analyze the approximations of the gradients based on these criteria.

For given $f_{\sigma} \in (t_1 - \Delta t, t_v + \Delta t)$, let the sigmoid approximation be $S'_{\lambda, f_{\sigma}} : [t_1 - \Delta t, t_v + \Delta t)$ as $S'_{\lambda, f_{\sigma}}(t) = (S'^{tseq}_{\lambda, f_{\sigma}})_j = -\lambda \cdot S(\lambda(t_j - f_{\sigma})) [1 - S(\lambda(t_j - f_{\sigma}))], \quad \text{for } t \in [t_j - \Delta t, t_j + \Delta t).$ Similarly, let our gradient approximation be $\hat{\delta}_{\beta, f_{\sigma}} : [t_1 - \Delta t, t_v + \Delta t)$ as

$$\hat{\delta}_{\beta,f_{\sigma}}(t) = \left(\hat{\delta}_{\beta,f_{\sigma}}^{tseq}\right)_{j} = \begin{cases} -\frac{1}{\beta\sqrt{\pi}}, & \text{if } f_{\sigma} \in [t_{j-1}, t_{j}), \\ 0, & \text{otherwise}, \end{cases} \quad \text{for } t \in [t_{j} - \Delta t, t_{j} + \Delta t). \end{cases}$$

Proposition D.1. For all $t \in (t_1 - \Delta t, t_v + \Delta t)$,

$$\left| \int_{t_1 - \Delta t}^{t_v + \Delta t} \hat{S'}_{\lambda, f_\sigma}(t) df_\sigma - \int_{t_1 - \Delta t}^{t_v + \Delta t} \frac{\partial \mathbb{I}(f_\sigma \le t)}{\partial f_\sigma} df_\sigma \right| \ge 2S(-\lambda(2v - 1)\Delta t).$$
(5)

972 And if $f_{\sigma} = t_j - \Delta t$ for some j, then

$$\left| \int_{t_1 - \Delta t}^{t_v + \Delta t} \hat{S'}_{\lambda, f_\sigma}(t) dt - \int_{t_1 - \Delta t}^{t_v + \Delta t} \frac{\partial \mathbb{I}(f_\sigma \le t)}{\partial f_\sigma} dt \right| \ge \left| 1 - 2v\lambda \Delta t \exp(-\lambda \Delta t) \right|.$$

Proposition D.2. For all $t \in (t_1 - \Delta t, t_v + \Delta t)$,

$$\int_{t_1-\Delta t}^{t_v+\Delta t} \hat{\delta}_{\beta,f_\sigma}(t) df_\sigma - \int_{t_1-\Delta t}^{t_v+\Delta t} \frac{\partial \mathbb{I}(f_\sigma \le t)}{\partial f_\sigma} df_\sigma = -\frac{2\Delta t}{\beta\sqrt{\pi}} + 1$$

and for all $f_{\sigma} \in (t_1 - \Delta t, t_v + \Delta t)$,

$$\int_{t_1-\Delta t}^{t_v+\Delta t} \hat{\delta}_{\beta,f_\sigma}(t)dt - \int_{t_1-\Delta t}^{t_v+\Delta t} \frac{\partial \mathbb{I}(f_\sigma \le t)}{\partial f_\sigma}dt = -\frac{2\Delta t}{\beta\sqrt{\pi}} + 1.$$

Suppose the grid is fixed, so Δt and v is fixed. Then for equation 5 to go to $0, \lambda \to \infty$ should hold. However, as $\lambda \to \infty$, the lower bound of equation 6 converges to 1, which means that the integral of the sigmoid approximation $\int_{t_1-\Delta t}^{t_v+\Delta t} \hat{S}'_{\lambda,f_{\sigma}}(t) dt$ becomes inconsistent. This is already expected from the vanishing gradient behavior. However, Proposition equation D.2 suggests that when β is appropriately chosen as $\beta = \frac{\sqrt{\pi}}{2\Delta t}$, the gradient approximation becomes consistent for the integral with respect to both f_{σ} and t.

E PROOFS FOR SECTION 4

Proof for Proposition 4.1. First, note that the sigmoid function $S(x) = \frac{1}{1 + \exp(-x)}$ satisfies

$$\frac{dS}{dx}(x) = \frac{\exp(-x)}{(1 + \exp(-x))^2} = S(x)(1 - S(x)),$$

and hence

$$\frac{\partial S(\lambda(t-f_{\sigma}))}{\partial f_{\sigma}} = -\lambda S(\lambda(t-f_{\sigma}))(1-S(\lambda(t-f_{\sigma}))$$

Since $x \mapsto |x(1-x)|$ on [0,1] is maximized when $x = \frac{1}{2}$, so

$$\left\|\frac{\partial S(\lambda(t-f_{\sigma}))}{\partial f_{\sigma}}\right\|_{\infty} = \frac{\lambda}{4}$$

1007 Meanwhile,

$$\left| \frac{\partial S(\lambda(t-f_{\sigma}))}{\partial f_{\sigma}} \right|_{t=t_{j}} \right| = \lambda S(\lambda(t_{j}-f_{\sigma})) \left[1 - S(\lambda(t_{j}-f_{\sigma})) \right]$$
$$\leq \lambda S(\lambda d(f_{\sigma}, tseq)) \left[1 - S(\lambda d(f_{\sigma}, tseq)) \right].$$

1012 Hence

$$\begin{split} \left\| S'^{tseq}_{\lambda,f_{\sigma}} \right\|_{\infty} &= \left\| \frac{\partial S(\lambda(t-f_{\sigma}))}{\partial f_{\sigma}} |_{t=t_{1},\dots,t_{v}} \right\|_{\infty} \\ &\leq \lambda S(\lambda d(f_{\sigma},tseq)) \left[1 - S(\lambda d(f_{\sigma},tseq)) \right]. \end{split}$$

1017 Hence, if $f_{\sigma} = t_j - \Delta t$ for some j, then

1018
1019
1020
1021

$$\left\|S'^{tseq}_{\lambda,f_{\sigma}}\right\|_{\infty} \leq \lambda S(\lambda \Delta t) \left[1 - S(\lambda \Delta t)\right] \leq \lambda \exp(-\lambda \Delta t),$$

1022 and therefore when $\lambda \exp(-\lambda \Delta t) \to 0$,

$$\inf_{\substack{f_{\sigma} \in [t_1 - \Delta t, t_v + \Delta t)}} \left\| S^{\prime tseq}_{\lambda, f_{\sigma}} \right\|_{\infty} \to 0.$$

(6)

Proof for Proposition 4.2. $\hat{\delta}^{tseq}_{\beta,f_{\sigma}}$ always has the form

$$\left(0,\ldots,0,-\frac{1}{\beta\sqrt{2\pi}},0,\ldots,0
ight)$$

1031 Therefore,

 $\left\| \hat{\delta}^{tseq}_{\beta,f_{\sigma}} \right\|_{\infty} = \frac{1}{\beta \sqrt{2\pi}}.$

	L
	L
	L

1038 F PROOFS FOR APPENDIX D

1040 Before beginning the proofs of Appendix D, we would first like to emphasize that for all $t \in (t_1 - \Delta t, t_v + \Delta t)$,

$$\int_{t_1 - \Delta t}^{t_v + \Delta t} \frac{\partial \mathbb{I}(f_\sigma \le t)}{\partial f_\sigma} df_\sigma = \mathbb{I}(t_v + \Delta t \le t) - \mathbb{I}(t_1 - \Delta t \le t) = -1.$$

and for all $f_{\sigma} \in (t_1 - \Delta t, t_v + \Delta t)$,

$$\int \frac{\partial \mathbb{I}(f_{\sigma} \leq t)}{\partial f_{\sigma}} dt = -\int \frac{\partial \mathbb{I}(f_{\sigma} \leq t)}{\partial t} dt = -\mathbb{I}(f_{\sigma} \leq t_v + \Delta t) + \mathbb{I}(f_{\sigma} \leq t_1 - \Delta t) = -1.$$

Proof for Proposition D.1. For given $t \in (t_1 - \Delta t, t_v + \Delta t)$, let $t_j \in tseq$ be such that $t \in [t_j - \Delta t, t_j + \Delta t)$. Then

$$\int \hat{S}'_{\lambda,f_{\sigma}}(t)df_{\sigma} = \int_{t_1-\Delta t}^{t_v+\Delta t} \frac{\partial S(\lambda(t_j-f_{\sigma}))}{\partial f_{\sigma}}df_{\sigma}$$
$$= S(\lambda(t_j-t_v-\Delta t)) - S(\lambda(t_j-t_1+\Delta t)).$$

This is minimized when t_j is close to $\frac{t_1+t_v}{2}$. Hence,

$$\int \hat{S}'_{\lambda, f_{\sigma}}(t) df_{\sigma} \ge S(-\lambda((2v-1)\Delta t) - S(\lambda((2v-1)\Delta t),$$

and hence

$$\left| \int \hat{S'}_{\lambda, f_{\sigma}}(t) df_{\sigma} - \int \frac{\partial \mathbb{I}(f_{\sigma} \le t)}{\partial f_{\sigma}} df_{\sigma} \right| \ge 1 - S(\lambda(2v - 1)\Delta t) + S(-\lambda(2v - 1)\Delta t)$$
$$= 2S(-\lambda(2v - 1)\Delta t).$$

Also, note that from the calculation in the proof of Proposition 4.1, if $f_{\sigma} = t_j - \Delta t$ for some j, then

$$\begin{split} \left\| S^{\prime tseq}_{\lambda,f_{\sigma}} \right\|_{\infty} &\leq \lambda S(\lambda \Delta t) \left[1 - S(\lambda \Delta t) \right]. \\ &\leq \lambda \exp(-\lambda \Delta t). \end{split}$$

And

$$\left|\int \hat{S'}_{\lambda,f_{\sigma}}(t)dt\right| \leq \left\|S'^{tseq}_{\lambda,f_{\sigma}}\right\|_{\infty} 2\Delta tv \leq 2\lambda\Delta tv \exp(-\lambda\Delta t).$$

1076 Therefore,

$$\left|\int \hat{S}'_{\lambda, f_{\sigma}}(t)dt - \int \frac{\partial \mathbb{I}(f_{\sigma} \leq t)}{\partial f_{\sigma}}dt\right| \geq |1 - 2\lambda\Delta tv \exp(-\lambda\Delta t)|$$

 $\int \hat{\delta}_{\beta, f_{\sigma}}(t) df_{\sigma} = \int_{t_{j-1}}^{t_j} -\frac{1}{\beta \sqrt{\pi}} df_{\sigma} = -\frac{2\Delta t}{\beta \sqrt{\pi}}.$

1086 And for given $f_{\sigma} \in (t_1 - \Delta t, t_v + \Delta t)$, let $t_j \in tseq$ be such that $f_{\sigma} \in [t_{j-1}, t_j)$. Then $\hat{\delta}_{\beta, f_{\sigma}}(t)$ is 1087 nonzero if and only if $t \in [t_j - \Delta t, t_j + \Delta t)$, and hence

$$\int \hat{\delta}_{\beta, f_{\sigma}}(t) dt = \int_{t_j - \Delta t}^{t_j + \Delta t} - \frac{1}{\beta \sqrt{\pi}} dt = -\frac{2\Delta t}{\beta \sqrt{\pi}}.$$

		L
		L
		L
		L

G PROOFS FOR SECTION 5

1096 Proof for Proposition 5.1. Since $\mathcal{O}_{\theta}(X) = g_{\theta}(\mathcal{C}_X(tseq))$ and $\mathcal{O}_{\theta}(X') = g_{\theta}(\mathcal{C}_{X'}(tseq))$,

$$\begin{aligned} \|\mathcal{O}_{\theta}(X) - \mathcal{O}_{\theta}(X')\|_{1} &= \|g_{\theta}(\mathcal{C}_{X}(tseq)) - g_{\theta}(\mathcal{C}_{X'}(tseq))\|_{1} \\ &\leq L \|\mathcal{C}_{X}(tseq) - \mathcal{C}_{X'}(tseq)\|_{1}. \end{aligned}$$

1101 Now, note that ECC C_X can be expanded using persistence diagrams $\{\mathcal{D}_k(X) : k \ge 0\}$ as follows: 1102 if $\mathcal{D}_k(X) = \{(b_{ki}, d_{ki}) : 1 \le i \le n_k\}$, then

$$\mathcal{C}_X(t) = \sum_{k=0}^{\infty} (-1)^k \mathbb{I}(b_{ki} \le t < d_{ki})$$

Since $b_{ki}, d_{ki} \in \{t_i^*\}$, there exists $a_1, \ldots, a_m \in \mathbb{Z}$ and $n_1 < \cdots < n_{2m}$ such that $\mathcal{C}_X(t) - \mathcal{C}_{X'}(t)$ can be expressed as

$$\mathcal{C}_X(t) - \mathcal{C}_{X'}(t) = \sum_{i=1}^m a_i \mathbb{I}(t_{n_{2i}}^* \le t < t_{n_{2i+1}}^*)$$

1110 Then $\|\mathcal{C}_X(tseq) - \mathcal{C}_{X'}(tseq)\|_1$ and $\|\mathcal{C}_X - \mathcal{C}_{X'}\|_1$ is expanded as

$$\|\mathcal{C}_X(tseq) - \mathcal{C}_{X'}(tseq)\|_1 = \sum_{j=1}^v \sum_{i=1}^m |a_i| \,\mathbb{I}(t_{n_{2i}}^* \le t_j < t_{n_{2i+1}}^*)$$

1115 and

$$\left\|\mathcal{C}_{X} - \mathcal{C}_{X'}\right\|_{1} = \sum_{i=1}^{m} |a_{i}| \left(t_{n_{2i+1}}^{*} - t_{n_{2i}}^{*}\right)$$

1118 Now for each i = 1, ..., m, $\sum_{j=1}^{v} \mathbb{I}(t_{n_{2i}}^* \leq t_j < t_{n_{2i+1}}^*)$ is the number of t_j 's that falls within the 1119 interval $[t_{n_{2i}}^*, t_{n_{2i+1}}^*)$. But since $t_{j+1} - t_j \geq \Delta t$, such number is at most $\left[\frac{(t_{n_{2i+1}}^* - t_{n_{2i}}^*)}{\Delta t}\right]$, and also 1121 from $t_{n_{2i+1}}^* - t_{n_{2i}}^* \geq \Delta t$,

$$\sum_{j=1}^{v} \mathbb{I}(t_{n_{2i}}^* \le t_j < t_{n_{2i+1}}^*) \le \left\lceil \frac{(t_{n_{2i+1}}^* - t_{n_{2i}}^*)}{\Delta t} \right\rceil \le \frac{2(t_{n_{2i+1}}^* - t_{n_{2i}}^*)}{\Delta t}.$$

1125 Hence $\|\mathcal{C}_X(tseq) - \mathcal{C}_{X'}(tseq)\|_1$ can be correspondingly upper bounded as

1127
1128
$$\|\mathcal{C}_X(tseq) - \mathcal{C}_{X'}(tseq)\|_1 = \sum_{i=1}^m |a_i| \sum_{j=1}^v \mathbb{I}(t_{n_{2i}}^* \le t_j \le t_{n_{2i+1}}^*)$$
1129

$$<\frac{2}{m}\sum_{i=1}^{m}|a_{i}|(t_{m}^{*},\ldots,-t_{m}^{*})$$

$$-\Delta t \sum_{i=1}^{+1} (u_{2i+1} - u_{2i})$$

1133
$$= \frac{2}{\Delta t} \left\| \mathcal{C}_X - \mathcal{C}_{X'} \right\|_1.$$

And correspondingly,

1136

1137

1138 1139 When K = K', $||f_X - f_{X'}||_{\infty} = \sup_{\sigma \in K} |f_X(\sigma) - f_{X'}(\sigma)|$, and if $K \neq K'$, $||f_X - f_{X'}||_{\infty}$ is not 1140 well defined. However, there is a general distance between two filtration functions f_X and $f_{X'}$ even 1141 when base simplicial complexes (or cubical complexes) are different; it is the interleaving distance 1142 $d_I(f_X, f_{X'})$. For the definition, see Section 5.1 from Chazal et al. (2016a). When K = K', there is 1143 a bound

1144

1149

1156 1157

$$d_I(f_X, f_{X'}) \le \|f_X - f_{X'}\|_{\infty}$$

 $\left\|\mathcal{O}_{\theta}(X) - \mathcal{O}_{\theta}(X')\right\|_{1} \leq \frac{2L}{\Delta t} \left\|\mathcal{C}_{X} - \mathcal{C}_{X'}\right\|_{1}.$

1145 Hence we have a general version of Theorem 5.3 as follows.

Theorem G.1. Suppose K is a finite simplicial complex or cubical complex. Then there exists a constant C_K only depending on K such that

$$\left\|\mathcal{C}_{X} - \mathcal{C}_{X'}\right\|_{1} \le C_{K} d_{I}(f_{X}, f_{X'}).$$

Proof for Theorem G.1 is in a similar manner from the proof of Wasserstein Stability Theorem of Cohen-Steiner et al. (2010).

Proof for Theorem G.1. From Proposition 5.2, it is sufficient to show that there exists C'_K depending only on K such that

$$\sum_{k=0}^{\infty} W_1(\mathcal{D}_k(X), \mathcal{D}_k(X')) \le C'_K d_I(f_X, f_{X'}).$$

1158 Fix $k \ge 0$, and let $\epsilon_k := W_{\infty}(\mathcal{D}_k(X), \mathcal{D}_k(X'))$ be the bottleneck distance between two diagrams 1159 $\mathcal{D}_k(X)$ and $\mathcal{D}_k(X')$. Let $\gamma_k : \mathcal{D}_k(X) \to \mathcal{D}_k(X')$ be the bijection that realizes the bottleneck 1160 distance, i.e., for any $p \in \mathcal{D}_k(X)$,

1161
$$\|p - \gamma_k(p)\|_{\infty} \le \epsilon_k.$$

1162 Then 1-Wasserstein distance $W_1(\mathcal{D}_k(X), \mathcal{D}_k(X'))$ satisfies

1163
1164
1165
1166
1167
1168

$$W_1(\mathcal{D}_k(X), \mathcal{D}_k(X')) = \inf_{\gamma} \sum_{x \in \mathcal{D}_k(X)} \|x - \gamma(x)\|_{\infty}$$

$$\leq \sum_{x \in \mathcal{D}_k(X)} \|x - \gamma_k(x)\|_{\infty}$$

$$\leq \epsilon_1 \|\mathcal{D}_k(X)\|$$

1169
$$\leq \epsilon_k |\nu_k(X)|.$$

1170 And hence if we let $\epsilon \coloneqq \sup_{k \ge 0} {\epsilon_k}$, then summing over $k \ge 0$ gives

1171
1172
1173
$$\sum_{k=0}^{\infty} W_1(\mathcal{D}_k(X), \mathcal{D}_k(X')) \le \sum_{k=0}^{\infty} \epsilon_k |\mathcal{D}_k(X)|$$
1173

1174
1175
1176
$$\leq \epsilon \sum_{k=0}^{\infty} |\mathcal{D}_k(X)|$$

1177 Now $\sum_{k=0}^{\infty} |\mathcal{D}_k(X)|$ is the number of points in persistence diagrams of all homological dimensions 1178 on K. This can be bounded by some constant C'_K that depends only on K: one rough bound can 1179 be as $|\{\sigma : \sigma \in K\}|^2$, since each point in persistence diagrams has a unique pair (σ_b, σ_d) of a birth 1180 simplex σ_b and a death simplex σ_d . And therefore, 1181

1182
1183
$$\sum_{k=0}^{\infty} W_1(\mathcal{D}_k(X), \mathcal{D}_k(X')) \le C'_K \epsilon$$

Now from f_X and $f_{X'}$ being on a finite simplicial complex K, they are q-tame (see Section 3.8 from Chazal et al. (2016a)). So from the bottleneck stability theorem (see e.g., Section 5.1 and Theorem 5.23 from Chazal et al. (2016a)), for all $k \ge 0$,

$$\epsilon_k \le d_I(f_X, f_{X'})$$

And hence, $\sum_{k=0}^{\infty} W_1(\mathcal{D}_k(X), \mathcal{D}_k(X')) \le C'_K d_I(f_X, f_{X'}).$ Before proving Corollary 5.4, we explain what an empirical distribution and a 'restriction of a DTM function' mean. We say P_X is an empirical distribution on $X = \{X_1, \ldots, X_n\}$, when $P_X = \{X_1, \ldots, X_n\}$ $\frac{1}{n}\sum_{i=1}^{n}\delta_{X_i}$, where δ_{X_i} is a Dirac measure on X_i , i.e., $\delta_{X_i}(A) = \mathbb{I}(X_i \in A)$. And suppose $\{\sigma_i\} \subset K$ be vertices of K for V-construction, or top dimensional cells of K for T-construction. Then we say f_X is a 'restrictions of a DTM function' d_{P_X,m_0} , if there exists a grid $\mathcal{G} = \{x_i\} \subset \mathbb{R}^d$, with $f_X(\sigma_i) = d_{P_X,m_0}(x_i)$.

Proof for Corollary 5.4. From Theorem 5.3,

$$\left\|\mathcal{C}_X - \mathcal{C}_{X'}\right\|_1 \le C_K \left\|f_X - f_{X'}\right\|_{\infty}$$

1206 Now we further bound $||f_X - f_{X'}||_{\infty}$. Note that

1207
1208
1209
1210
1210
1210
1211

$$\|f_X - f_{X'}\|_{\infty} = \max_{\sigma \in K} |f_X(\sigma) - f_{X'}(\sigma)|$$

$$= \max_{x \in \mathcal{G}} |d_{P_X,m_0}(x) - d_{P_{X'},m_0}(x)|$$

$$\leq \|d_{P_X,m_0} - d_{P_{X'},m_0}\|_{\infty}.$$

1213 And from Chazal et al. (2011)[Theorem 3.5],

$$d_{P_X,m_0} - d_{P_{X'},m_0} \Big\|_{\infty} \le \frac{1}{\sqrt{m_0}} W_2(P_X, P_{X'})$$

Hence putting these things together gives

$$\|\mathcal{C}_X - \mathcal{C}_{X'}\|_1 \le \frac{(d+1)DC_K}{\sqrt{m_0}}W_2(P_X, P_{X'}).$$

			1	
			L	
			L	
-	-	-	2	

1224 Let d_{GH} be the Gromov-Hausdorff distance.

1226 Corollary G.2. Suppose f_X , $f_{X'}$ are Vietoris-Rips filtrations of X and X', respectively. Then

$$\left\|\mathcal{C}_X - \mathcal{C}_{X'}\right\|_1 \le C_K d_{GH}(X, X').$$

1230 Proof for Corollary 5.4. From Theorem G.1,

$$\left\|\mathcal{C}_{X} - \mathcal{C}_{X'}\right\|_{1} \le C_{K} d_{I}(f_{X}, f_{X'})$$

¹²³³ 1234 Then from Lemma 4.3 of Chazal et al. (2014),

 $d_I(f_X, f_{X'}) \le d_{GH}(X, X').$

1237 Hence putting these things together gives

- 1239 $\|\mathcal{C}_X - \mathcal{C}_{X'}\|_1 \le C_K d_{GH}(X, X').$





(a) 30 sampled points from unit circle with small noise.

(b) ECC of (a) calculated using Alpha filtration with 32 and 2000 grid points.

Figure 6: Within the interval [0.06, 0.7], both ECCs capture the Euler characteristic of the underlying
loop structure, which is 0. The ECC with 2000 grid points exhibits more noise than the ECC with
32 grid points, as the dense discretization captures even the small (noise) generators that do not
represent the global structure of data.

H CHOICE OF TDA HYPERPARAMETERS

1267 1268

1265 1266

1259

1260

1269 In this section, we discuss the choice of several TDA hyperparameters in ECLayr: filtration, T_{min} , 1270 T_{max} , v, and β .

1271 Choice of filtration. Although numerous filtration options exist, certain filtrations are commonly favored for specific data modalities and training contexts. For example, Vietoris-Rips and Alpha filtrations are extensively utilized for point clouds, while a sub/superlevel filtration on a filtered cubical complex is a natural choice for data with grid structure. DTM filtration provides robustness against outliers, and thereby preferable in scenarios of data contamination. Despite not being discussed here, other choices of filtrations are also available. Nevertheless, the fundamental idea remains the same; choose a filtration that can mostly effectively extract topological features from the given data.

1278 **Choice of** $[T_{min}, T_{max}]$. A naive and convenient approach is to assign T_{min} and T_{max} as the 1279 minimum and maximum of possible filtration values, respectively. For instance, one can set $T_{min} =$ 1280 0 and $T_{max} = 1$ for image data with min-max normalized pixel values. An alternative method is to select $[T_{min}, T_{max}]$ as a tighter interval within the range of possible filtration values, focusing on 1281 regions of the filtration that contain meaningful topological and geometrical information. Such an 1282 interval can be identified via hyperparameter search, or chosen intuitively by examining the ECC 1283 computed for some data. For example, the ECCs in Figure 6-(b) reveal that [0.06, 0.7] is a suitable 1284 interval for capturing the underlying loop structure depicted in Figure 6-(a). 1285

1286 **Choice of** v. Spacing between grid points is important as our vectorized ECC does not account 1287 for cycles that are born and dead between t_i and t_{i+1} , where $t_i, t_{i+1} \in tseq$. Provided that the discretization is not overly sparse, the uncaptured cycles are often small (noise) generators with life 1288 span shorter than $t_{i+1} - t_i$. This implies that with appropriate discretization, noise can be partially 1289 filtered by design. Therefore, using a highly dense discretization is not necessarily beneficial, as 1290 it captures even the small (noise) generators (see Figure 6-(b)). Conversely, using an excessively 1291 sparse discretization may jeopardize the capturing of essential global features. The optimal choice 1292 of v is not always evident; we recommend hyperparameter search using cross validation to determine 1293 the adequate v that balances the two circumstances. 1294

Choice of β . The hyperparameter β regulates the gradient's magnitude, with smaller values of β yielding larger gradients. However, the optimal choice of β is somewhat ambiguous. Unfortunately,

we do not have a clear rationale for choosing β ; it is contingent upon numerous factors, including model architecture and the specific task at hand. Therefore, we recommend conducting hyperparameter search via cross validation to select an appropriate β .

¹³⁰⁰ 1301 I Experiment Details

All experiments were conducted with NVIDIA RTX A6000 GPU, with the exception of empirical time complexity analysis, which was run on Apple M2.

1306 I.1 COMPUTATIONAL EFFICIENCY

The runtime metrics are computed on the training set of MNIST, Br35H, and synthetic data. MNIST contains 60000 training samples of size 28×28 . For Br35H, we use the same training set as the experiment conducted in Section 6.4, which is 209 samples of size 112×112 . We generate 1000 samples of size 224×224 for the synthetic data, where each pixel is randomly sampled from a uniform distribution. Both ECC and PH use superlevel set filtration on V-constructed cubical complex. PH is computed using the GUDHI package, while Algorithm 1 is used to compute ECC with v = 32 on interval [0, 1].

1314

1316

1305

1315 I.2 TOPOLOGICAL AUTOENCODER

The encoder and decoder network each consists of three fully connected layers, with input dimension 1317 size 101, hidden dimension size 32, and latent dimension size 2. BatchNorm and ReLu nonlinearity 1318 is used after each layer, except for the latent dimension. We place one ECLayr at the beginning of 1319 the encoder and one ECLayr at the latent dimension in order to compute the MAE loss between 1320 ECC of input point cloud and ECC of latent representation. This MAE loss acts as a topology 1321 regularizing term, with lambda=0.001 controlling its magnitude. Vietoris-Rips filtration is employed 1322 to compute ECC and max dimension is restricted to 1. v is set to 1000 over interval [0, 2], where 1323 filtration values indicate edge length. β , which controls the magnitude of the gradient, is assigned as 1324 0.01 and we do not use q_{θ} for this experiment. Adam optimizer is used for training with batch size 1325 32 and learning rate 0.0001. We run for 100 epochs and adopt early stopping after patience 10.

1326

1327 I.3 MNIST DATASET 1328

1329 The MNIST dataset contains 60000 training data and 10000 test data of handwritten digits from 0 to 9. We implement a 4 layer baseline model, consisting of two CNN layers followed by two 1330 fully connected layers, with ReLU nonlinearity between every layer. Both CNN layers use 3×3 1331 kernels with stride 1 and padding 1. Each CNN layer has channel size 32 and 1, respectively. The 1332 output of the CNN network is flattened and passed to the subsequent fully connected layers with 1333 hidden dimension of 64. For training, we employ the Adam optimizer with learning rate 0.001 and 1334 batch size 32. The learning rate is decayed by a factor of 0.1 when the validation loss plateaus for 1335 10 epochs. While we use a maximum of 1000 epochs for training, early stopping is implemented 1336 to stop training after the validation loss plateaus for 25 epochs. Cross-entropy loss is used for 1337 classification. We compare the performance of our proposed layer with the base model, and two 1338 other topological layers applicable to image datasets: PersLay (Carrière et al., 2020) and PLLay 1339 (Kim et al., 2020). Utilizing a very simple model on limited training samples, we observed random 1340 failures across all models with outliers significantly affecting the outcome. To remove the influence 1341 of outliers and solely evaluate model performance, we repeat each experiment 15 times and select the top 10 test accuracies for assessment. 30% of the training data is used as a validation set and 1342 model performance is evaluated on the complete test data. 1343

Data Scarcity. To observe how model performance changes with data size, we sample training data of size 100, 300, 500, 700 and 1000 with equal proportion for each label. For all topological layers, we place a parallel layer at the beginning of the network (referred to as CNN + EC(i), CNN + Pers, and CNN + PL(i)). For ECLayr and PLLay, which allow backpropagation, we add an additional layer after the last convolutional layer (referred to as CNN + EC, CNN + DECT and CNN + PL). $\beta = 0.01$ is used to control the gradient intensity for ECLayr. For all topological layers, we implement superlevel set filtrations on T-constructed cubical complex and use v = 32

1351

1352

1353

1354

1355

1356

1357

1358

1359

1363

1365

1367

1380

1382

1384 1385

1386

1387 1388

Data Size Models 100 300 500 700 1000 69.075 83.282 85.346 86.868 89.066 CNN (± 0.305) (± 1.049) (± 0.566) (± 0.397) (± 0.518) 86.549 87.649 90.659 73.417 88.642 CNN + EC(i) (± 0.718) (± 1.050) (± 0.251) (± 0.335) (± 0.397) 73.652 86.551 88.335 88.755 90.361 CNN + EC (± 0.983) (± 0.400) (± 1.035) (± 0.458) (± 0.319) 70.872 84.004 85.702 86.940 89.421 CNN + DECT (± 1.484) (± 0.364) (± 0.502) (± 0.343) (± 0.234) 64.982 81.256 84.239 86.200 88.444 CNN + Pers (± 2.225) (± 1.376) (± 0.617) (± 0.786) (± 0.474) 66.425 83.176 86.110 87.996 90.072 CNN + PL(i) (± 2.078) (± 1.156) (± 0.598) (± 0.476) (± 0.364) 69.656 83.382 86.251 88.123 90.164 CNN + PL (± 2.673) (± 1.764) (± 1.200) (± 0.231) (± 0.713)

Table 4: Test accuracy of models trained on different sizes of MNIST dataset. For each data size, the best accuracy is highlighted in bold.

Models	Corruption & Noise Probability					
Widdels	0.00	0.05	0.10	0.15	0.20	
CNINI	89.066	86.556	85.711	83.288	81.357	
CININ	(± 0.518)	(± 0.495)	(± 0.546)	(± 0.735)	(± 0.464)	
CNN + EC(i)	90.659	88.315	86.844	84.341	82.122	
CININ + EC(I)	(± 0.397)	(± 0.542)	(± 0.564)	(± 0.753)	(± 0.791)	
CNN + EC	90.361	88.164	86.762	84.310	81.976	
CNN + EC	(± 0.400)	(± 0.700)	(± 0.700)	(± 0.800)	(± 0.647)	
CNN + Dara	88.444	87.686	86.680	84.385	81.845	
CININ + Fels	(± 0.474)	(± 1.229)	(± 0.304)	(± 0.559)	(± 0.483)	
CNN + DI (i)	90.072	85.723	84.742	82.303	80.593	
CNN + PL(1)	(± 0.364)	(± 0.740)	(± 0.641)	(± 0.971)	(± 0.493)	
CNN + DI	90.164	85.665	84.9	82.668	80.852	
CININ + FL	(± 0.713)	(± 0.665)	(± 0.637)	(± 0.595)	(± 0.926)	

Table 5: Test accuracy of models trained on 1000 MNIST data with different corruption & noise probability. For each corruption and noise probability, the best accuracy is highlighted in bold.

over interval [0, 1]. PersLay uses line point transform and a 10×10 unit grid for learnable weights over interval $[0, 1] \times [0, 1]$. *top*2 function is used as the permutation invariant operation for PLLay and PersLay. After the topological descriptor is computed from the respective layers, it is fed into a fully connected layer g_{θ} with output dimension of 32 to compute the final output of each topological layer. We concatenate the output of the topological layer with the output of the CNN network to feed it to the subsequent fully connected layer. See Table 4 for the full experiment results

1395 Data Contamination. To observe how robust models are against data contamination, we consider a 1396 corruption and noise process where the pixels are randomly omitted and subsequently contaminated by random noise between 0 and 1 with probability 0.05, 0.1, 0.15, and 0.2. We apply different levels 1398 of noise to 1000 training samples with equal proportion for each label. The overall architecture 1399 and training scheme remain the same as before, except that we now place two parallel layers at the 1400 beginning of the network using DTM filtration. For each DTM filtration, we align the data on a unit grid and use $m_0 = 0.05$ and $m_0 = 0.2$ to examine the topological structure at different local/global 1401 scales. We use the interval [0.02, 0.28] and [0.06, 0.29] for $m_0 = 0.05$ and $m_0 = 0.2$ respectively. 1402 For CNN + EC and CNN + PL, we place an additional layer using DTM filtration with $m_0 = 0.05$ 1403 after the last convolutional layer. β is set at 0.01. See Table 5 for the full experiment results





Filtration			Interval $[T_{min}, T_{max}]$		
T, sub	V, sup	V, sub	[0.1, 1]	[0, 0.8]	[0.1, 0.8]
88.927	89.205	88.817	89.249	89.659	89.467
(± 0.262)	(± 0.291)	(± 0.219)	(± 0.230)	(± 0.430)	(± 0.134)
Discretization v			Gra	dient Contr	ol β
16	64	128	0.1	0.001	0.0001
89.299	89.271	88.979	89.228	89.206	88.659
(± 0.453)	(± 0.274)	(± 0.228)	(± 0.176)	(± 0.332)	(± 0.657)

Table 6: Test accuracy on 1000 MNIST data for different choice of hyperparameters.

1471 I.5 HYPERPARAMETER INFLUENCE

1473 ECLayr introduces four hyperparameters: filtration, interval $[T_{min}, T_{max}]$, discretization v, and 1474 gradient control β . To evaluate the influence of each hyperparameter on performance, we provide 1475 an ablation study on 1000 noiseless samples from the MNIST dataset. We vary a single hyperpa-1476 rameter while keeping all others consistent with the experiment in Section 6.3. The test accuracies 1477 of EC + CNN model are presented in Table 6. "T", "V", "sub", and "sup" for different choices of 1478 filtration refer to T-construction, V-construction, sublevel set filtration, and superlevel set filtration,