

# 000 001 002 003 004 005 006 007 008 009 010 011 012 013 014 015 016 017 018 019 020 021 022 023 024 025 026 027 028 029 030 031 032 033 034 035 036 037 038 039 040 041 042 043 044 045 046 047 048 049 050 051 052 053 TOPOFORMER: TOPOLOGY MEETS ATTENTION FOR GRAPH LEARNING

Anonymous authors

Paper under double-blind review

## ABSTRACT

We introduce TOPOFORMER, a lightweight and scalable framework for graph representation learning that encodes topological structure into attention-friendly sequences. At the core of our method is *Topo-Scan*, a novel module that decomposes a graph into a short, ordered sequence of topological tokens by slicing over node or edge filtrations. These sequences capture multi-scale structural patterns, from local motifs to global organization, and are processed by a Transformer to produce expressive graph-level embeddings. Unlike traditional persistent homology pipelines, *Topo-Scan* is parallelizable, avoids costly diagram computations, and integrates seamlessly with standard deep learning architectures. We provide theoretical guarantees on the stability of our topological encodings and demonstrate state-of-the-art performance across graph classification and molecular property prediction benchmarks. Our results show that TOPOFORMER matches or exceeds strong GNN and topology-based baselines while offering predictable and efficient compute. This work opens a new path for parallelizable and unifying approaches to graph representation learning that integrate topological inductive biases into attention frameworks.

## 1 INTRODUCTION

Graphs are powerful data structures for modeling relational data in biology, chemistry, and social networks. While recent advances in graph learning have produced strong task-specific models, most architectures lack the generalization of foundation models in vision and language (Radford et al., 2021; Bubeck et al., 2023). Achieving such general-purpose capability in graphs is difficult due to their irregular, non-Euclidean structure (Wu et al., 2020), which complicates the design of transferable inductive biases.

Topological Data Analysis (TDA) provides a principled approach by encoding global and local structure in a way that is stable to perturbations and insensitive to node identity (Hensel et al., 2021; Pham et al., 2025). In principle, Persistent Homology (PH) offers a canonical summary of how connectivity and cycles evolve across scales, and has proven useful across domains (Skaf et al., 2022; Obayashi et al., 2022; Shultz, 2023). In practice, however, PH pipelines depend on persistence diagrams, which require expensive global reductions and a subsequent vectorization step (e.g., images, landscapes, curves) whose design can materially affect downstream performance. On graphs, common sublevel/superlevel filtrations also tend to *early-saturate*, high-valued vertices activate early, quickly filling the complex and suppressing late-emerging features. These computational and modeling frictions have slowed the adoption of PH in graph representation learning, despite the clear promise of topological signals for multi-resolution structure.

To overcome these barriers, we develop a lightweight yet expressive alternative that bypasses full persistence diagrams while retaining multi-resolution topological information in a form consumable by transformers. We introduce TOPOFORMER, a scalable framework that integrates topological descriptors with attention architectures. At its core is Topo-Scan, a module that converts a graph into a short, ordered sequence of topological tokens across multiple resolutions. These sequences are directly consumable by attention mechanisms (Vaswani et al., 2017), enabling efficient graph-level representations within the same token-based interface used by large-scale transformer models. We therefore view TOPOFORMER as a step toward topology-aware graph foundation models, rather than a full foundation model itself, and leave large-scale pretraining on heterogeneous graph corpora to future work. TOPOFORMER achieves strong performance on graph classification and molecular

054 property prediction under unified evaluation protocols, with theoretical guarantees on the stability of  
 055 its topological encodings. Our Contributions are as follows:  
 056

- 057 • We introduce a scalable method for turning topological structure into attention-ready sequences, en-  
 058 abling transformers to process graphs without relying on node embeddings or heavy preprocessing.
- 059 • We propose a new framework that bridges topological data analysis and deep learning, capturing  
 060 both local and global graph structure through a unified attention mechanism.
- 061 • We conduct a comprehensive evaluation across diverse graph learning tasks, demonstrating strong  
 062 performance on both graph classification and molecular property prediction benchmarks.
- 063 • We provide theoretical guarantees on the robustness of our representations and show that our  
 064 approach offers predictable and efficient compute, making it practical for large-scale applications.

## 066 2 MOTIVATION AND BACKGROUND

068 This section reviews recent work and highlights the need to integrate advanced topological methods  
 069 with modern ML to overcome limitations in graph representation learning.

070 **Persistent Homology for Graphs.** Persistent Homology (PH) was first defined for filtered  
 071 simplicial complexes in the early 2000s (Edelsbrunner et al., 2002; Zomorodian & Carlsson, 2005).  
 072 Early applications centered on point clouds  $\mathcal{X} \subset \mathbb{R}^N$ , where Vietoris–Rips filtrations generate  
 073 nested complexes  $\Delta_1(\mathcal{X}) \subset \Delta_2(\mathcal{X}) \subset \dots$ , allowing topological features to be tracked across  
 074 scales (Carlsson, 2009). The persistence diagram  $\text{PD}_k(\mathcal{X}) = \{[b_i, d_i]\}$  records births and deaths of  
 075  $k$ -dimensional features, with longer intervals  $(d_i - b_i)$  interpreted as more persistent and thus more  
 076 structurally significant (Dey & Wang, 2022).

077 PH has since been applied to graphs and images. For graphs, two principal approaches are used. *Power*  
 078 (*distance*) *filtrations* treat nodes as a point cloud with graph distances as pairwise distances, then  
 079 build a Rips filtration (Aktas et al., 2019), which is typically computationally heavy. A more practical  
 080 alternative in graph learning is *sublevel filtrations*, where a scalar node/edge function  $f$  induces  
 081 nested subgraphs that are lifted to simplicial complexes via cliques (upper-star extension is standard).  
 082 A key interpretability difference follows: in power filtrations, bar lengths reflect geometric scale; in  
 083 sublevel filtrations they reflect *differences in f* rather than physical size, so “long bars  $\Rightarrow$  important  
 084 features” need not hold universally. Poorly chosen  $f$  may yield many short-lived features or early  
 085 saturation, while task-aligned or *learnable* filtrations can mitigate these effects (Hofer et al., 2020).  
 086 Rather than viewing this as an intrinsic weakness, we take it as motivation to design fixed-budget,  
 087 stable summaries that integrate smoothly with modern ML (See App. C.8 for discussion).

088 The standard PH pipeline for graphs has three main steps (Coskunuzer & Akçora, 2024): *filtration*,  
 089 *persistence computation*, and *vectorization*. Given a graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , a function  $f : \mathcal{V} \rightarrow \mathbb{R}$   
 090 with thresholds  $\{\alpha_i\}_{i=1}^N$  induces subgraphs  $\mathcal{G}_1 \subseteq \dots \subseteq \mathcal{G}_N$ , where  $\mathcal{G}_i$  contains vertices  $\{v \in$   
 091  $\mathcal{V} \mid f(v) \leq \alpha_i\}$ . Lifting each  $\mathcal{G}_i$  to its clique complex  $\widehat{\mathcal{G}}_i$  yields a filtration  $\{\widehat{\mathcal{G}}_i\}$ . Persistence  
 092 diagrams  $\text{PD}_k(\mathcal{G}, f) = \{(b_j, d_j)\}$  record births and deaths of  $H_k(\widehat{\mathcal{G}}_i)$  and are typically vectorized  
 093 via persistence images, landscapes, or Betti curves (Ali et al., 2022).

094 In recent years, the ML community has increasingly recognized the value of topological encodings  
 095 for graph-level tasks, with PH-based methods showing strong results across domains (Immonen  
 096 et al., 2024; Demir et al., 2022; Verma et al., 2024; Loiseaux et al., 2024; Horn et al., 2021; Chen  
 097 et al., 2024b). Despite this promise, two bottlenecks hinder broader adoption: (1) the computational  
 098 overhead of persistence computations in large pipelines (Otter et al., 2017), and (2) the difficulty of  
 099 choosing vectorizations that align with downstream objectives (Ali et al., 2022). Our TOPOFORMER  
 100 framework addresses both by producing a compact *sequence* of stable, low-cost topological tokens  
 101 that feed directly into attention layers, thereby bypassing full persistence diagrams and bespoke  
 102 vectorizations while remaining compatible with efficient graph-specific computations.

103 **Recent methods learn neural approximations of persistence based topological features to reduce the**  
 104 **cost of exact PH.** RipsNet (de Surrel et al., 2022) estimates Rips persistence diagrams for point clouds  
 105 directly from raw data, while Yan et al. (2022) approximate graph topological features with a GNN.  
 106 Our approach is complementary: instead of approximating persistence diagrams, Topo Scan bypasses  
 107 global PH and directly builds short interlevel topological sequences tailored to Transformer encoders.

108 **Transformers.** Transformers (Vaswani et al., 2017) underpin transferable models in language and  
 109 vision (Devlin, 2018; Dosovitskiy et al., 2020) by learning from *ordered* token sequences with long-  
 110 range dependencies. On graphs, adapting attention is challenging due to variable size, permutation  
 111 invariance, and irregular connectivity. Our design sidesteps these issues: Topo-Scan yields a short,  
 112 1D *ordered* sequence of topological tokens with a fixed channel width, so positional encodings  
 113 and attention operate in their native regime, without graph-specific heavy machinery. This makes  
 114 transformers a natural, efficient backend for multi-resolution structural signals. **By contrast, recent**  
 115 **graph transformer models such as Graphomer** (Ying et al., 2021), **GPS** (Rampášek et al., 2022),  
 116 and related architectures operate directly on node tokens and inject structure via shortest-path or  
 117 Laplacian-based positional encodings and attention biases. In TOPOFORMER, each graph is first  
 118 compressed into a short sequence of topological tokens, so attention runs on a fixed-length, purely  
 119 topological sequence rather than on all nodes of the original graph.

120 **Molecular Property Prediction.** Molecular property prediction (MPP) is central to drug discovery  
 121 (ADMET). Classical pipelines use engineered fingerprints with RF/SVMs (Cereto-Massagué et al.,  
 122 2015); deep learning extends to MLPs on fingerprints, sequence models on SMILES (Rong et al.,  
 123 2020), and GNNs on molecular graphs (Wieder et al., 2020), with recent 3D methods trading accuracy  
 124 for higher compute and sensitivity to rotations (Gasteiger et al., 2021; Li et al., 2022b). Despite  
 125 progress, DL does not always surpass strong classical baselines on realistic benchmarks (Janela  
 126 & Bajorath, 2022; Valsecchi et al., 2022), motivating transformer variants (Sultan et al., 2024),  
 127 geometric models (Liu et al., 2022b), and topological approaches (Demir et al., 2022; Loiseaux  
 128 et al., 2023). Evaluation protocols also vary (e.g., scaffold vs. random splits), affecting reported  
 129 generalization. Our approach unifies robust topological structure with a scalable attention backend,  
 130 providing an effective, split-agnostic representation for MPP.

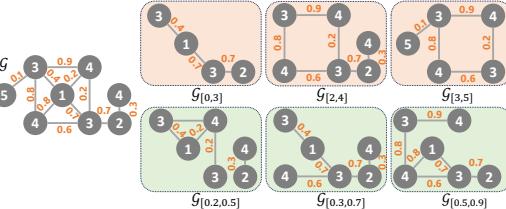
### 131 3 TOPO-TRANSFORMERS

132 TDA captures multi-scale structural patterns while offering robustness to noise, making it attractive for  
 133 representation learning. However, the standard Persistent Homology pipeline, consisting of filtration  
 134 construction, persistence diagram computation, and vectorization, introduces inefficiencies and lacks  
 135 adaptability, particularly in graph settings. While persistence diagram computation is standardized,  
 136 vectorization remains ad hoc and significantly impacts model performance Ali et al. (2022). Our goal  
 137 is to develop an efficient and scalable alternative to this workflow for graph representation learning.

138 Our first insight is that the *strict nested-*  
 139 *ness condition required in PH is not always*  
 140 *necessary for graphs.* Unlike point clouds,  
 141 graphs inherently encode structural relation-  
 142 ships that permit more flexible and direct  
 143 extraction of topological features. Building  
 144 on this observation, we bypass persistence  
 145 diagrams and vectorization by directly ex-  
 146 tracting topological sequences from struc-  
 147 tured graph slices. This shift enables effi-  
 148 cient and adaptable pattern extraction and  
 149 forms the foundation of a scalable learning  
 150 framework.

151 **Topo-Scan.** Traditional sublevel filtrations on graphs often saturate rapidly, and once most nodes  
 152 join the subgraph at low thresholds, little new structure emerges and important patterns at larger scales  
 153 are lost. Topo-Scan overcomes this by first imposing a directional hierarchy via a scalar function  
 154  $f: \mathcal{V} \rightarrow \mathbb{R}$ , then slicing the graph into a sequence of overlapping subgraphs along increasing values  
 155 of  $f$ . Rather than waiting for a single threshold to engulf the entire graph, each slice captures fresh  
 156 topological information, such as connectivity changes and emerging loops, without early collapse.  
 157 We compute basic invariants (e.g. Betti numbers) on each slice to form a compact, ordered signature  
 158 sequence. Feeding these ordered descriptors into a transformer lets the model attend to structure at  
 159 every scale, ensuring no signal is lost to premature saturation (See Fig. 2 for a toy example).

160 Let  $f: \mathcal{V} \rightarrow \mathbb{R}$  be a filtration function defined on the vertices, with thresholds  $\alpha_0 = \min_{v \in \mathcal{V}} f(v) <$   
 161  $\alpha_1 < \dots < \alpha_N = \max_{v \in \mathcal{V}} f(v)$ . In most cases, the thresholds are selected either as evenly spaced



162 **Figure 1: Topo-Scan.** Topo-Scan decomposes a graph into  
 163 sequential topological slices via node and edge filtrations.  
 164 The top row shows node-based filtrations; the bottom row,  
 165 edge-based ones.

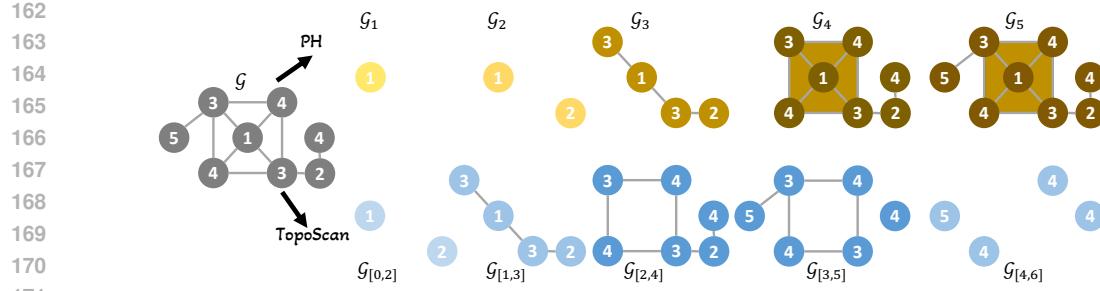


Figure 2: **Topo-Scan vs. PH.** This toy example highlights the key differences between the Topo-Scan filtration and standard PH filtration, where node values indicate filtration function values. In PH, early-activated nodes quickly saturate the graph, suppressing the emergence of later topological features. As shown, PH yields relatively uninformative barcodes with  $\beta_0 = \langle 1, 2, 1, 1, 1 \rangle$  and  $\beta_1 = \langle 0, 0, 0, 0, 0 \rangle$ , while Topo-Scan captures richer topological dynamics, producing  $\beta_0 = \langle 2, 1, 1, 2, 4 \rangle$  and  $\beta_1 = \langle 0, 0, 1, 1, 0 \rangle$ .

values or based on quintiles. Next, for each  $\alpha_i$ , we define  $\mathcal{V}_i = \{v_r \in \mathcal{V} \mid \alpha_i \leq f(v_r) \leq \alpha_{i+m}\}$  and  $\mathcal{G}_i$  as the induced subgraph  $\mathcal{G}_i = (\mathcal{V}_i, \mathcal{E}_i)$ , where  $\mathcal{E}_i = \{e_{rs} \in \mathcal{E} \mid v_r, v_s \in \mathcal{V}_i\}$ . The clique complex of  $\mathcal{G}_i$ , denoted  $\widehat{\mathcal{G}}_i$ , forms a sequence  $\{\widehat{\mathcal{G}}_i\}$  called *slicing*. We call this process *Topo-Scan*, which decomposes graphs into topological slices, similar to medical scans revealing structural layers. Leveraging a hierarchical structure, it adapts to node and edge filtrations, weighted graphs, and diverse relations, capturing local and global topological patterns for robust, scalable representation learning. It remains robust by tracking short-lived features effectively and is scalable through parallelized slice extraction.

The resolution ( $N$ ) determines the number of slices, while the thickness ( $m$ ) specifies the range of nodes included in each slice. After constructing  $\{\widehat{\mathcal{G}}_i\}$ , we compute four outputs for each slice:  $\beta_0(\widehat{\mathcal{G}}_i)$  (Betti-0, connected components),  $\beta_1(\widehat{\mathcal{G}}_i)$  (Betti-1, cycles/holes),  $|\mathcal{V}_i|$  (node count), and  $|\mathcal{E}_i|$  (edge count). These outputs form ordered sequences of size  $N$ , such as  $\widehat{\beta}_k(\mathcal{G}) = \{\beta_k(\widehat{\mathcal{G}}_i)\}_{i=1}^N$  for  $k = 0, 1$ . While  $\widehat{\beta}_k(\mathcal{G})$  are the primary topological outputs,  $\{|\mathcal{V}_i|\}$  and  $\{|\mathcal{E}_i|\}$  serve as normalization factors (see Figure 1). These sequences are concatenated into a sequence (vector)  $\Gamma(\mathcal{G})$  of length  $4N$  where  $N$  is the number of slices.

A key distinction from PH lies in activation: PH includes all nodes up to a threshold, causing early saturation in dense graphs, while Topo-Scan uses a sliding window to preserve late-emerging features and capture fine structure (see Fig. 3 and App. C.8). Slice thickness  $m$  controls locality, allowing flexibility across datasets. Its localized design ensures robustness to noise and enables parallelization, making it ideal for scalable ML workflows.

**Vectorization Choice.** Among many possible vectorizations, we deliberately use a very low-dimensional token per slice,  $(\beta_0, \beta_1, |\mathcal{V}_i|, |\mathcal{E}_i|)$ . Global vectorizations such as persistence images or landscapes aggregate information over the entire filtration into a single feature vector, which largely destroys the *sequential* structure that Topo-Scan is designed to expose. In contrast, our Betti-based tokens preserve how components and cycles evolve across slices; richer per-slice invariants could be plugged in, but we focus on this minimal choice to isolate the benefit of the sequential representation.

**TOPOFORMER.** We use the *ordered sequences* of topological features in Transformers, which excel at capturing sequential structures and complex dependencies through self-attention mechanisms, making them well-suited for tasks requiring order and contextual understanding. While traditional PH processes a sequence of simplicial complexes  $\{\Delta_i\}$ , this sequential structure is often lost during the persistence diagram and vectorization stages, where outputs are transformed into unordered vectors.

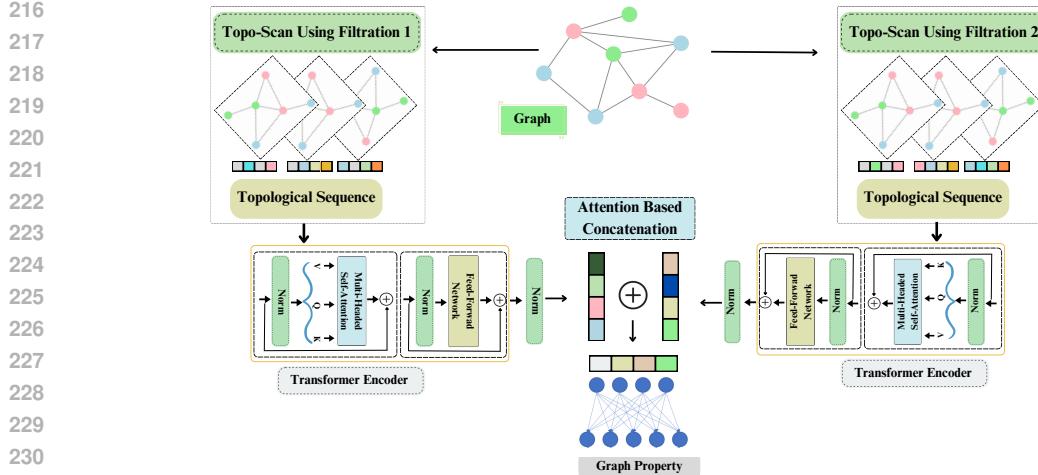


Figure 4: **TopoFormer Flowchart**: Given an input graph  $\mathcal{G}$ , sequential substructures are extracted via Topo-Scan. Each substructure is encoded into a four-dimensional topological signature. These sequences are processed by a transformer model, and outputs from multiple filtration functions are fused using attention-based concatenation. A final prediction layer maps the representation to the target graph property.

Topo-Scan preserves the sequential nature of topological features and aligns them with transformers' ability to model positional relationships.

**ML Model.** Our transformer architecture (Fig. 4) consists of an embedding layer that processes input sequences, a transformer encoder that captures hierarchical dependencies through self-attention mechanisms, and a fully connected classification head that maps learned representations to output predictions. To enhance generalization and mitigate overfitting, we integrate regularization techniques, such as dropout and weight decay, ensuring robustness across diverse graph learning tasks. Formally, given an input sequence  $\mathbf{x} \in \mathbb{R}^{m \times T \times D}$ , where  $m$  is the number of graphs,  $T$  the sequence length, and  $D$  the token dimensionality, the sequence is embedded via  $\mathbf{E}$ , with positional encoding  $\mathbf{P}$  added. This processed sequence is then passed through a multi-layer transformer encoder, producing an output representation  $\mathbf{z}$ , which is flattened and normalized before being classified via a fully connected layer.

Expanding this model, we introduce a dual-transformer framework with an integrated multi-layer perceptron (MLP) classifier to handle diverse input modalities. The model processes three distinct inputs:  $\mathbf{X}_1$  and  $\mathbf{X}_2$  through independent transformers  $\mathcal{T}_1$  and  $\mathcal{T}_2$ , and  $\mathbf{X}_3$  through an MLP  $\mathcal{M}$ . Their respective outputs  $\mathbf{z}_1$ ,  $\mathbf{z}_2$ , and  $\mathbf{z}_3$  are combined using a learnable weighted sum, allowing the model to adaptively balance feature contributions:

$$\mathbf{z}_{\text{combined}} = \alpha \cdot \mathbf{z}_1 + \beta \cdot \mathbf{z}_2 + (1 - \alpha - \beta) \cdot \mathbf{z}_3$$

where  $\alpha$  and  $\beta$  are learned during training. The aggregated feature vector is then batch-normalized and passed through a fully connected layer to produce the final classification output:  $\hat{y} = \mathbf{FC}(\mathbf{z}_{\text{combined}})$  where  $\hat{y} \in \mathbb{R}^C$  represents the class probabilities, with  $C$  being the number of output classes. Model details are given in Appendix C.

### 3.1 STABILITY OF TOPO-SCAN SEQUENCES

A useful graph vectorization should be robust: small changes in the filtration signal should not cause large changes in the output sequence. We formalize this for Topo-Scan on the *fixed clique complex*  $\widehat{G}$  of  $G = (V, E)$  using upper-star extensions of node functions.

**Setup.** Let  $f, g : V \rightarrow \mathbb{R}$  be filtration functions, extended to  $\widehat{G}$  by the upper-star rule  $\widehat{h}(\sigma) = \max_{v \in \sigma} h(v)$ . Fix a shared threshold grid  $\alpha_0 < \dots < \alpha_N$ , window width  $m$ , stride  $s$ , and windows  $I_t = [\alpha_{ts}, \alpha_{ts+m}]$  for  $t = 0, \dots, T-1$ , where  $T = \lfloor (N-m)/s \rfloor + 1$ . For  $k \in \{0, 1\}$ , the  $t$ -th Topo-Scan token is the interlevel Betti number

$$\beta_k^h(t) := \dim H_k((\widehat{G})_{I_t}^h), \quad h \in \{f, g\}.$$

270 **Theorem 3.1** (Discrete  $\ell_1$  stability of Topo-Scan). *There exists  $C = C(\widehat{G}, \{\alpha_i\}, m, s)$  such that for*  
 271  *$k \in \{0, 1\}$ ,*

$$272 \quad \|\widehat{\beta}_k(G, f) - \widehat{\beta}_k(G, g)\|_1 \leq C d_B(M_k^f, M_k^g),$$

273 *where  $M_k^h$  denotes the  $k$ -dimensional interlevel (level-set) persistence module of the upper-star  
 274 filtration induced by  $h$  on  $\widehat{G}$ , and  $d_B$  is the bottleneck distance between such modules.*

275 **Corollary 3.2.** *For upper-star filtrations on a fixed complex, interlevel modules satisfy*  
 276  *$d_B(M_k^f, M_k^g) \leq \|f - g\|_\infty$ . Hence  $\|\widehat{\beta}_k(G, f) - \widehat{\beta}_k(G, g)\|_1 \leq C \|f - g\|_\infty$ .*

277 **Outline.** Each token counts classes surviving exactly over  $I_t$ ; under a  $\delta$  bottleneck matching, only  
 278 classes within  $\delta$  of the interval boundary can change their contribution, so per-window changes are  
 279  $O(\delta)$  and summing over windows yields the discrete  $\ell_1$  bound with a constant depending on the  
 280 window schedule and the (finite) bar complexity of  $\widehat{G}$ . Full details and references are in Appendix B.

281 **Beyond stability,** the Topo-Scan sequences  $\widehat{\beta}_k(G, h)$  are closely related to classical PH invariants: they  
 282 can be viewed as a discrete sampling of the rank invariant / Betti curve of the interlevel module  $M_k^h$   
 283 along our window schedule. Thus, Topo-Scan provides a coarse but structured, Transformer-ready  
 284 discretization of the same homological information underlying barcodes and stable-rank summaries;  
 285 see Remark B.3 for further discussion.

## 286 4 EXPERIMENTS

### 287 4.1 EXPERIMENTAL SETUP

288 **Datasets.** We report the TOPOFORMER\* performance in two graph learning tasks: **graph classification**  
 289 and **molecular property prediction (MPP)**.

290 **Graph Classification Datasets.** We use nine graph classification benchmark datasets: (i) molecular  
 291 graphs from BZR, MUTAG and COX2 (Kriegel et al., 2012); (ii) biological graphs PROTEINS (Borg-  
 292 wardt et al., 2005); and (iii) social graphs, including IMDB-Binary, IMDB-Multi, REDDIT-Binary,  
 293 and REDDIT-Multi-5K (Yanardag et al., 2015). We also include a large-scale dataset, OGBG-  
 294 MOLHIV, from Open Graph Benchmark (Hu et al., 2020b). Dataset statistics are provided in Table 1.

295 **MPP Datasets.** For molecular prop-  
 296 erty prediction (MPP), we employ seven  
 297 datasets from MoleculeNet (Wu et al.,  
 298 2018): BBBP (blood-brain barrier pene-  
 299 tration), Tox21, ToxCast, ClinTox (toxic-  
 300 ity prediction), SIDER (adverse drug re-  
 301 actions), HIV (replication inhibition), and  
 302 BACE ( $\beta$ -secretase 1 inhibitors). Dataset  
 303 statistics are provided in Table 3 (top rows).

304 **Model Setup.** We use *Topo-Scan* to  
 305 generate topological signature sequences,  
 306 which are fed to Transformer classifiers. Each filtration (20 thresholds, width 2) yields four sequences  
 307 of length 19 (Betti-0, Betti-1, node count, edge count), giving 76 features per filtration. For graph  
 308 classification, we use Ollivier–Ricci curvature and Heat Kernel Signature; and for molecular property  
 309 prediction including the OGBG-MOLHIV dataset, we use atomic weight and Ollivier–Ricci curvature.  
 310 Independent Transformers process each filtration, and their outputs are fused by attention before a  
 311 final linear layer.

312 For MPP, we use TOPOFORMER with the standard molecular fingerprints, processed by a two-layer  
 313 MLP and combined with topological features via attention, yielding TOPOFORMER\*. We report  
 314 10-fold CV accuracy on graph classification, scaffold-split AUCs over three runs for MPP (Fang et al.,  
 315 2023a), and use the standard split for OGBG-MOLHIV.

316 **Hyperparameters.** For model optimization, we employed the Adam optimizer with a learning  
 317 rate of 0.001. We also use standard regularization techniques—dropout (0.5), weight decay (1e-4),  
 318 and batch normalization—commonly employed in transformer training. The transformer model  
 319 architecture was designed with a hidden dimension of 32. Hyperparameters are given in App. C.5.

Table 1: Graph classification datasets.

| Datasets    | #Graphs | $ \mathcal{V} $ | $ \mathcal{E} $ | Classes |
|-------------|---------|-----------------|-----------------|---------|
| BZR         | 405     | 35.75           | 38.36           | 2       |
| COX2        | 467     | 41.22           | 43.45           | 2       |
| MUTAG       | 188     | 17.93           | 19.79           | 2       |
| PROTEINS    | 1113    | 39.06           | 72.82           | 2       |
| IMDB-B      | 1000    | 19.77           | 96.53           | 2       |
| IMDB-M      | 1500    | 13.00           | 65.94           | 3       |
| REDDIT-B    | 2000    | 429.63          | 497.75          | 2       |
| REDDIT-5K   | 4999    | 508.52          | 594.87          | 5       |
| OGBG-MOLHIV | 41127   | 25.5            | 27.5            | 2       |

324 Table 2: **Graph Classification.** Accuracy on eight benchmark datasets using 10-fold CV. Baseline results are  
325 taken from the respective papers using the same setting. We mark the 1st (blue), 2nd (purple), and 3rd (green)  
326 per column. The last two columns report the average deviation (AvD) from the best-performing model and the  
327 average rank (AvR) across all datasets.

| Model                           | BZR                              | COX2                             | MUTAG                            | PROTEINS                         | IMDB-B                           | IMDB-M                           | REDDIT-B                         | REDDIT-5K                        | AvD↓       | AvR↓       |
|---------------------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|------------|------------|
| 6 GNNs (Errica et al., 2020)    | —                                | —                                | 80.42 $\pm$ 2.07                 | 75.80 $\pm$ 3.70                 | 71.20 $\pm$ 3.90                 | 49.10 $\pm$ 3.50                 | 89.90 $\pm$ 1.90                 | 56.10 $\pm$ 1.60                 | 6.0        | 8.5        |
| PersLay (Carrière et al., 2020) | —                                | 80.90 $\pm$ NA                   | 89.80 $\pm$ NA                   | 74.80 $\pm$ NA                   | 71.20 $\pm$ NA                   | 48.80 $\pm$ NA                   | —                                | 55.60 $\pm$ NA                   | 5.2        | 8.7        |
| DMP (Bodnar et al., 2021)       | —                                | —                                | 84.00 $\pm$ 8.60                 | 75.30 $\pm$ 3.30                 | 73.80 $\pm$ 4.50                 | 50.90 $\pm$ 2.50                 | 86.20 $\pm$ 6.80                 | 51.90 $\pm$ 2.10                 | 6.1        | 8.3        |
| FC-V (O’Bray et al., 2021)      | 85.61 $\pm$ 0.59                 | 81.01 $\pm$ 0.88                 | 87.31 $\pm$ 0.66                 | 74.54 $\pm$ 0.48                 | 73.84 $\pm$ 0.36                 | 46.80 $\pm$ 0.37                 | 89.41 $\pm$ 0.24                 | 52.36 $\pm$ 0.37                 | 5.7        | 9.2        |
| SubMix (Yoo et al., 2022)       | 86.34 $\pm$ 2.00                 | <b>84.68<math>\pm</math>3.70</b> | 80.99 $\pm$ 0.60                 | 67.80 $\pm$ 2.00                 | 70.30 $\pm$ 1.40                 | 46.47 $\pm$ 2.50                 | —                                | —                                | 8.4        | 11.5       |
| G-Mix (Han et al., 2022)        | 84.15 $\pm$ 2.30                 | 83.83 $\pm$ 2.10                 | 81.96 $\pm$ 0.60                 | 66.28 $\pm$ 1.10                 | 69.40 $\pm$ 1.10                 | 46.40 $\pm$ 2.70                 | —                                | —                                | 9.1        | 12.8       |
| RGCL (Li et al., 2022a)         | 84.54 $\pm$ 1.67                 | 79.31 $\pm$ 0.68                 | 87.66 $\pm$ 1.01                 | 75.03 $\pm$ 0.43                 | 71.85 $\pm$ 0.84                 | 49.31 $\pm$ 0.42                 | 90.34 $\pm$ 0.58                 | 56.38 $\pm$ 0.40                 | 5.2        | 8.0        |
| AutoGCL (Yin et al., 2022)      | 86.27 $\pm$ 0.71                 | 79.31 $\pm$ 0.70                 | 88.64 $\pm$ 1.08                 | 75.80 $\pm$ 0.36                 | 72.32 $\pm$ 0.93                 | 50.60 $\pm$ 0.80                 | 88.58 $\pm$ 1.49                 | <b>56.75<math>\pm</math>0.18</b> | 4.7        | 7.0        |
| WWLS (Fang et al., 2023b)       | 88.02 $\pm$ 0.61                 | 81.58 $\pm$ 0.91                 | 88.30 $\pm$ 1.23                 | 75.35 $\pm$ 0.74                 | <b>75.08<math>\pm</math>0.31</b> | 51.61 $\pm$ 0.62                 | —                                | —                                | 4.5        | 5.2        |
| PGOT (Qian et al., 2024)        | 87.32 $\pm$ 3.90                 | 82.98 $\pm$ 5.21                 | <b>92.63<math>\pm</math>2.58</b> | 73.21 $\pm$ 2.59                 | 62.90 $\pm$ 3.05                 | 51.33 $\pm$ 1.76                 | —                                | —                                | 6.1        | 7.5        |
| EMP (Chen et al., 2024a)        | —                                | —                                | 88.79 $\pm$ 0.63                 | 72.78 $\pm$ 0.54                 | 74.44 $\pm$ 0.45                 | 48.01 $\pm$ 0.42                 | <b>91.03<math>\pm</math>0.22</b> | 54.41 $\pm$ 0.32                 | 4.8        | 7.5        |
| EPIC (Heo et al., 2024)         | <b>88.78<math>\pm</math>2.30</b> | <b>85.53<math>\pm</math>1.60</b> | 82.44 $\pm$ 0.70                 | 69.06 $\pm$ 1.00                 | 71.70 $\pm$ 1.00                 | 47.93 $\pm$ 1.30                 | —                                | —                                | 6.9        | 9.0        |
| MP-HSM (Loiseaux et al., 2024)  | —                                | 77.10 $\pm$ 3.00                 | 85.60 $\pm$ 3.30                 | 74.60 $\pm$ 2.10                 | 74.80 $\pm$ 2.50                 | 47.90 $\pm$ 3.20                 | —                                | —                                | 6.9        | 10.1       |
| TopoGCL (Chen et al., 2024b)    | 87.17 $\pm$ 0.83                 | 81.45 $\pm$ 0.55                 | 90.09 $\pm$ 0.93                 | <b>77.30<math>\pm</math>0.89</b> | 74.67 $\pm$ 0.32                 | <b>52.81<math>\pm</math>0.31</b> | <b>90.40<math>\pm</math>0.53</b> | —                                | <b>3.5</b> | <b>4.3</b> |
| DASP (Ye et al., 2025)          | <b>89.40<math>\pm</math>3.10</b> | <b>84.80<math>\pm</math>4.60</b> | <b>91.90<math>\pm</math>8.60</b> | <b>77.20<math>\pm</math>3.10</b> | <b>81.40<math>\pm</math>3.60</b> | 51.20 $\pm$ 2.20                 | —                                | <b>57.60<math>\pm</math>1.60</b> | <b>1.6</b> | <b>2.8</b> |
| TOPOFORMER                      | <b>92.36<math>\pm</math>4.11</b> | 83.93 $\pm$ 4.03                 | <b>94.68<math>\pm</math>4.30</b> | <b>77.64<math>\pm</math>3.64</b> | <b>78.90<math>\pm</math>3.31</b> | <b>55.40<math>\pm</math>4.78</b> | <b>91.50<math>\pm</math>1.89</b> | <b>57.99<math>\pm</math>1.94</b> | <b>0.5</b> | <b>1.5</b> |

340 **Computational Complexity.** Classical PH requires global boundary–matrix reductions with  
341 cubic worst-case cost and poor parallelism (Otter et al., 2017). TOPOFORMER *skips persistence*  
342 *diagrams entirely*: instead of global reductions, it computes  $\beta_0$  and  $\beta_1$  *per slice* on the clique-complex  
343 2-skeleton.  $\beta_0$  uses union–find on the 1-skeleton, while  $\beta_1$  is derived from sparse edge–triangle  
344 operators after triangle enumeration (no cycle-rank identity due to clique complexes). This yields  
345  $\mathcal{O}(|V_t| + |E_t| + T_t)$  per slice, aggregated as  $\mathcal{O}(L \sum_t (|V_t| + |E_t| + T_t))$  across  $k$  slices and  $L$   
346 filtrations. Since slices are independent, Betti computations are fully parallelizable. By bypassing  
347 PD computation and vectorization, TOPOFORMER achieves multi-fold runtime and memory gains  
348 while retaining task-relevant topological features (Appendix C.3). In Appendix C.6, we further show  
349 that TOPOFORMER consistently outperforms classical PH across multiple filtration functions and  
350 vectorization schemes in the graph classification task.

351 **Implementation and Runtime.** We implemented our approach in Python and conducted ex-  
352 periments on a 12th Gen Intel Core i7-1270P vPro processor (E-cores up to 3.50 GHz, P-cores  
353 up to 4.80 GHz) with 32GB LPDDR5-6400MHz RAM. Topo-Scan feature extraction took 269.38  
354 seconds for OGBG-MOLHIV/HIV and 29.51 seconds for REDDIT-5K; other datasets were faster.  
355 The remaining model runtime was negligible. More timeruns and a comparison with PH can be found  
356 at Appendix C.3. Our code is available at the link <sup>1</sup>.

## 357 4.2 RESULTS

359 **Graph Classification Baselines.** We evaluate our method against 20 state-of-the-art baselines  
360 spanning several categories. These include: *GNN-based models* such as GCN, DGCNN, DiffPool,  
361 ECC, GIN, and GraphSAGE (with the best results reported by Errica et al. (2020)); *topological*  
362 *methods* including PersLay, DMP, FC-V, WWLS, MP-HSM, and EMP; *GNNs with data augmentation*  
363 such as SubMix, G-Mix, and EPIC; *contrastive learning methods* including RGCL, AutoGCL, and  
364 TopoGCL; and *prototype-based methods* such as PGOT. We further include the recent graph kernel  
365 method DASP (Ye et al., 2025). A complete list of baselines is provided in Table 2.

366 **Graph Classification Results.** In graph classification, TOPOFORMER attains the *best or second-best*  
367 *accuracy on 7 out of 8 benchmarks* (Table 2). Aggregating across datasets, it achieves an  
368 *average deviation (AvD)* of 0.5 from the best model and an *average rank (AvR)* of 1.5, demonstrating  
369 consistent top-tier performance. Notably, TOPOFORMER establishes new state-of-the-art on BZR,  
370 MUTAG, PROTEINS, IMDB-M, REDDIT-B, and REDDIT-5K, while remaining highly competitive  
371 elsewhere. It also surpasses common pooling-based methods on these datasets (see Table 12). On  
372 the large-scale OGBG-MOLHIV benchmark (Table 4), TOPOFORMER\* reaches an AUC within  
373  $\sim$ 2 points of the strong Graphomer baseline, underscoring both its scalability and the strength  
374 of topological signals as an inductive bias in graph learning. For this table, we restrict baselines  
375 to peer-reviewed published methods reported in the literature, rather than including unpublished  
376 leaderboard entries in (Hu et al., 2020b).

377 <sup>1</sup><https://anonymous.4open.science/r/TOPOFORMER-B0E3>

378 **Table 3: SOTA MPP Models.** ROC AUC comparison on molecular property prediction with scaffold splitting.  
379 We mark the 1st (blue), 2nd (purple), and 3rd (green) per column. The last two columns report the average  
380 deviation (AvD) from the best-performing model and the average rank (AvR) across all datasets.

| Model                        | BBBP                           | Tox21                          | ToxCast                        | SIDER                          | ClinTox                        | BACE                           | HIV                            | AvD $\downarrow$ | AvR $\downarrow$ |
|------------------------------|--------------------------------|--------------------------------|--------------------------------|--------------------------------|--------------------------------|--------------------------------|--------------------------------|------------------|------------------|
| # Molecules                  | 2,039                          | 7,831                          | 8,577                          | 1,427                          | 1,480                          | 1,513                          | 41,913                         |                  |                  |
| # Task                       | 1                              | 12                             | 617                            | 27                             | 2                              | 1                              | 1                              |                  |                  |
| N-GRAM (Liu et al., 2019)    | <b>91.2<math>\pm</math>3.0</b> | 76.1 $\pm$ 2.7                 | —                              | 63.2 $\pm$ 0.5                 | 87.5 $\pm$ 2.7                 | 79.1 $\pm$ 1.3                 | 78.7 $\pm$ 0.4                 | 8.5              | 8.4              |
| PT-GNN (Hu et al., 2020a)    | 70.8 $\pm$ 1.5                 | 78.7 $\pm$ 0.4                 | 65.7 $\pm$ 0.6                 | 62.7 $\pm$ 0.8                 | 72.6 $\pm$ 1.5                 | 84.5 $\pm$ 0.7                 | 79.9 $\pm$ 0.7                 | 12.4             | 8.7              |
| CMPNN (Song et al., 2020)    | <b>92.7<math>\pm</math>1.7</b> | 80.3 $\pm$ 1.3                 | <b>70.8<math>\pm</math>1.3</b> | 61.0 $\pm$ 3.6                 | 89.8 $\pm$ 0.8                 | 86.7 $\pm$ 0.2                 | 78.2 $\pm$ 2.2                 | <b>6.1</b>       | 6.2              |
| MGSSL (Zhang et al., 2021)   | 70.5 $\pm$ 1.1                 | 74.0 $\pm$ 1.4                 | 64.1 $\pm$ 0.7                 | 59.2 $\pm$ 0.6                 | 80.7 $\pm$ 2.1                 | 79.7 $\pm$ 0.8                 | 79.5 $\pm$ 1.1                 | 13.5             | 11.7             |
| GEM (Fang et al., 2022)      | 70.5 $\pm$ 2.0                 | 78.1 $\pm$ 0.6                 | 68.6 $\pm$ 0.2                 | 63.2 $\pm$ 1.5                 | 90.3 $\pm$ 0.7                 | 87.9 $\pm$ 1.0                 | <b>81.3<math>\pm</math>0.3</b> | 8.9              | 6.6              |
| GROVER (Rong et al., 2020)   | 86.8 $\pm$ 2.2                 | 82.0 $\pm$ 1.6                 | 56.8 $\pm$ 3.4                 | 61.2 $\pm$ 2.5                 | 70.3 $\pm$ 13.7                | 82.8 $\pm$ 3.6                 | 68.2 $\pm$ 1.1                 | 13.5             | 9.9              |
| GraphMVP (Liu et al., 2022a) | 72.4 $\pm$ 1.6                 | 76.5 $\pm$ 0.4                 | 63.1 $\pm$ 0.4                 | 63.9 $\pm$ 1.2                 | 79.1 $\pm$ 2.8                 | 81.2 $\pm$ 0.9                 | 77.0 $\pm$ 1.2                 | 12.7             | 10.4             |
| MolCLR (Wang et al., 2022)   | 72.6 $\pm$ 1.3                 | 77.2 $\pm$ 0.6                 | 65.9 $\pm$ 2.1                 | 61.3 $\pm$ 6.6                 | 89.8 $\pm$ 2.7                 | <b>88.5<math>\pm</math>2.2</b> | 77.4 $\pm$ 0.6                 | 9.9              | 7.8              |
| MolCLR-2 (Wang et al., 2022) | 72.4 $\pm$ 0.7                 | 78.4 $\pm$ 0.6                 | 69.1 $\pm$ 1.2                 | 59.7 $\pm$ 3.4                 | 88.0 $\pm$ 4.0                 | 85.0 $\pm$ 2.4                 | 77.8 $\pm$ 5.5                 | 10.2             | 8.6              |
| KANO (Fang et al., 2023a)    | <b>96.0<math>\pm</math>1.6</b> | <b>83.7<math>\pm</math>1.3</b> | <b>73.2<math>\pm</math>1.6</b> | <b>65.2<math>\pm</math>0.8</b> | 94.4 $\pm$ 0.3                 | <b>93.1<math>\pm</math>2.1</b> | <b>85.1<math>\pm</math>2.2</b> | <b>1.6</b>       | <b>2.0</b>       |
| MV-Mol (Luo et al., 2024)    | 73.6 $\pm$ 0.2                 | 80.3 $\pm$ 0.6                 | 70.0 $\pm$ 0.4                 | <b>67.3<math>\pm</math>0.0</b> | <b>95.6<math>\pm</math>1.6</b> | 88.2 $\pm$ 0.4                 | <b>81.4<math>\pm</math>0.3</b> | 6.5              | <b>3.6</b>       |
| MolFuse (Zheng et al., 2024) | 74.3 $\pm$ 1.3                 | 77.6 $\pm$ 0.4                 | 64.1 $\pm$ 0.3                 | <b>69.5<math>\pm</math>1.0</b> | <b>95.5<math>\pm</math>3.3</b> | 87.2 $\pm$ 1.3                 | 78.6 $\pm$ 0.9                 | 7.9              | 6.2              |
| MORE (Son et al., 2025)      | 71.9 $\pm$ 0.9                 | 75.6 $\pm$ 0.5                 | 64.6 $\pm$ 0.6                 | 60.9 $\pm$ 0.6                 | 81.0 $\pm$ 0.7                 | 82.8 $\pm$ 1.3                 | 77.0 $\pm$ 0.7                 | 12.6             | 11.1             |
| TOPOFORMER*                  | 89.5 $\pm$ 1.3                 | <b>82.7<math>\pm</math>0.5</b> | <b>75.3<math>\pm</math>0.5</b> | 63.1 $\pm$ 0.7                 | <b>96.5<math>\pm</math>0.6</b> | <b>95.9<math>\pm</math>0.3</b> | 81.2 $\pm$ 0.8                 | <b>2.5</b>       | <b>2.8</b>       |

395 **MPP Baselines.** We compare against strong  
396 supervised, self-supervised, and contrastive meth-  
397 ods for molecular property prediction (MPP). *Su-  
398 pervised:* CMPNN (message passing on molec-  
399 ular graphs). *Predictive self-supervision:* N-  
400 GRAM, PT-GNN, GROVER, MGSSL, GEM. *Con-  
401 trastive/augmentation and 3D:* GraphMVP (with  
402 3D), MolCLR, MolCLR-2. *Knowledge-aware /  
403 prompts:* KANO. *Recent multi-view/fusion mod-  
404 els:* MV-Mol (multi-view molecular representations),  
405 MORE (modality-aware molecular representation  
406 learning), and MolFuse (fusion of heterogeneous  
407 molecular signals). See Table 3 for full references.

408 **MPP Results.** On molecular property prediction,  
409 TOPOFORMER shows strong adaptability when  
410 paired with Extended Connectivity Fingerprints.  
411 Against state-of-the-art supervised, contrastive, and  
412 fusion baselines, TOPOFORMER\* achieves the *best ROC AUC* on *ToxCast*, *ClinTox*, and *BACE*, and  
413 is the *runner-up* on *Tox21* (Table 3). It remains competitive on *HIV*, trailing the leader by only a small  
414 margin. Aggregating across all seven benchmarks, TOPOFORMER attains the *second-lowest average*  
415 *deviation* from the column best ( $AvD = 2.5$ ) and the *second-lowest average rank* ( $AvR = 2.8$ ),  
416 confirming consistent top-tier performance alongside recent SOTA models such as KANO, MV-Mol,  
417 and MolFuse. We also benchmarked against hybrid classical (HC) models (Appendix A.3), where  
418 TOPOFORMER achieves the best result on four out of seven datasets and highly competitive results on  
419 others (Table 8). These findings highlight that transforming topology into compact, attention-ready  
420 tokens yields a robust and adaptable molecular predictor.

421 We further report *Hybrid Classical* baselines combining fingerprints, SMILES, and graph features  
422 with standard learners in Table 8. See Appendix A.3 for details of these models.

### 4.3 ABLATION STUDIES

423 We conduct **four ablation studies, as follows.**

424 **TOPOFORMER vs. PH** (Table 5): We compare TOPOFORMER with two persistent homology models  
425 using the same filtration functions and thresholds. *PH-MLP* uses sublevel filtrations with Betti  
426 vectorization followed by an MLP, while *PH-TR* replaces the MLP with a Transformer, treating  
427 Betti vectors as sequences. TOPOFORMER instead uses our proposed *Topo-Scan* to directly extract  
428 topological sequences. PH-TR outperforms PH-MLP, showing that sequential encodings preserve  
429 richer information than static features. TOPOFORMER further improves on PH-TR, indicating that  
430 Topo-Scan captures more expressive structure than standard PH filtrations.

Table 4: ROC AUC results for OGBG-MOLHIV dataset.

| Model                              | ROC AUC                          |
|------------------------------------|----------------------------------|
| GIN-VN (Xu et al., 2018)           | 77.80 $\pm$ 1.82                 |
| HGK-WL (Togninalli et al., 2019)   | 79.05 $\pm$ 1.30                 |
| WWL (Borgwardt et al., 2020)       | 75.58 $\pm$ 1.40                 |
| PNA (Corso et al., 2020)           | 79.05 $\pm$ 1.32                 |
| DGN (Beaini et al., 2021)          | 79.70 $\pm$ 0.97                 |
| GraphSNN (Wijesinghe et al., 2021) | 79.72 $\pm$ 1.83                 |
| GCN-GNorm (Cai et al., 2021)       | 78.83 $\pm$ 1.00                 |
| Graphomer (Ying et al., 2021)      | <b>80.51<math>\pm</math>0.53</b> |
| Cy2C-GCN (Choi et al., 2022)       | 78.02 $\pm$ 0.60                 |
| GAWL (Nikolentzos et al., 2023)    | 78.34 $\pm$ 0.39                 |
| LLM-GIN (Zhong et al., 2024)       | 79.22 $\pm$ NA                   |
| GMoE-GIN (Wang et al., 2024)       | 76.90 $\pm$ 0.90                 |
| TopER (Tola et al., 2025)          | <b>80.21<math>\pm</math>0.15</b> |
| TOPOFORMER*                        | 78.19 $\pm$ 0.19                 |

432 **Table 5: TOPOFORMER vs. PH:** Accuracy results for three topological models using degree centrality,  
 433 Ollivier-Ricci and HKS filtrations. The PH-MLP model utilizes Betti vectors derived from regular sublevel  
 434 filtrations combined with an MLP, while PH-TR applies transformers to the same vectors. The TOPOFORMER  
 435 uses Betti sequences generated via the Topo-Scan on the same filtration function and applies transformers.

| Filtration | Model      | BZR                              | COX2                             | MUTAG                            | PROTEINS                         | IMDB-B                           | IMDB-M                           | REDDIT-B                         |
|------------|------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|
| Degree     | PH-MLP     | 82.71 $\pm$ 6.51                 | 76.44 $\pm$ 5.39                 | 84.06 $\pm$ 4.65                 | 68.37 $\pm$ 3.97                 | 65.70 $\pm$ 4.03                 | 45.07 $\pm$ 2.59                 | 89.50 $\pm$ 2.87                 |
|            | PH-TR      | 86.43 $\pm$ 4.33                 | 78.15 $\pm$ 5.19                 | 86.11 $\pm$ 5.23                 | <b>77.54<math>\pm</math>2.64</b> | <b>75.00<math>\pm</math>2.11</b> | <b>50.67<math>\pm</math>3.57</b> | <b>92.30<math>\pm</math>1.77</b> |
|            | TOPOFORMER | <b>91.10<math>\pm</math>5.14</b> | <b>80.27<math>\pm</math>5.24</b> | <b>92.54<math>\pm</math>5.12</b> | 77.45 $\pm$ 4.02                 | 74.20 $\pm$ 5.01                 | 50.33 $\pm$ 1.52                 | 89.75 $\pm$ 2.18                 |
| O.Ricci    | PH-MLP     | 85.45 $\pm$ 3.36                 | 78.16 $\pm$ 5.09                 | 84.06 $\pm$ 5.21                 | 65.50 $\pm$ 4.26                 | 68.00 $\pm$ 3.55                 | 44.87 $\pm$ 3.65                 | 85.65 $\pm$ 2.62                 |
|            | PH-TR      | 88.62 $\pm$ 5.40                 | 78.16 $\pm$ 5.73                 | 87.61 $\pm$ 5.70                 | 77.27 $\pm$ 5.08                 | 72.20 $\pm$ 6.24                 | 48.00 $\pm$ 4.33                 | 90.65 $\pm$ 1.08                 |
|            | TOPOFORMER | <b>90.38<math>\pm</math>5.50</b> | <b>80.72<math>\pm</math>6.44</b> | <b>92.54<math>\pm</math>4.47</b> | <b>77.90<math>\pm</math>3.17</b> | <b>74.70<math>\pm</math>4.95</b> | <b>51.53<math>\pm</math>3.49</b> | <b>91.90<math>\pm</math>2.73</b> |
| HKS        | PH-MLP     | 84.96 $\pm$ 4.42                 | 78.19 $\pm$ 4.34                 | 84.09 $\pm$ 3.72                 | 70.80 $\pm$ 4.70                 | 71.10 $\pm$ 5.28                 | 47.93 $\pm$ 3.20                 | 88.10 $\pm$ 1.67                 |
|            | PH-TR      | 89.60 $\pm$ 5.84                 | 79.89 $\pm$ 4.66                 | 94.12 $\pm$ 5.42                 | 77.18 $\pm$ 3.15                 | 76.80 $\pm$ 3.97                 | 53.60 $\pm$ 3.31                 | 87.25 $\pm$ 1.95                 |
|            | TOPOFORMER | <b>90.62<math>\pm</math>4.91</b> | <b>83.95<math>\pm</math>2.99</b> | <b>95.32<math>\pm</math>5.58</b> | <b>77.35<math>\pm</math>2.86</b> | <b>77.90<math>\pm</math>5.72</b> | <b>54.07<math>\pm</math>2.54</b> | <b>90.05<math>\pm</math>2.41</b> |

445 **Effect of molecular fingerprints** (Table 7): We evaluate TOPOFORMER and Extended-Connectivity  
 446 Fingerprints (ECFPs) both independently and in combination, including integration with PubChem  
 447 descriptors. While topological and fingerprint models perform moderately on their own, their  
 448 combination consistently outperforms individual baselines, suggesting that topological features  
 449 complement domain-specific descriptors.

450 **Sensitivity to width parameter** (Table 6): We analyze how the sliding window size influences the  
 451 performance of Topo-Scan. See Appendix C.5 for further details.

452 **Table 6: Width Parameter.** Performance comparison for different window width parameters across  
 453 datasets.

|                   |         | BZR                              | COX2                             | MUTAG                            | PROTEINS                         | IMDB-B                           | IMDB-M                           | REDDIT-B                         |
|-------------------|---------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|
| Degree Centrality | $m = 2$ | <b>89.89<math>\pm</math>3.74</b> | 78.36 $\pm$ 4.93                 | 92.02 $\pm$ 7.24                 | <b>77.28<math>\pm</math>5.93</b> | <b>74.20<math>\pm</math>3.36</b> | <b>51.53<math>\pm</math>3.34</b> | 86.60 $\pm$ 2.97                 |
|                   | $m = 3$ | 88.64 $\pm$ 5.30                 | <b>78.38<math>\pm</math>5.04</b> | 90.41 $\pm$ 5.53                 | 76.92 $\pm$ 3.62                 | 73.20 $\pm$ 3.39                 | 51.13 $\pm$ 3.08                 | <b>86.90<math>\pm</math>2.31</b> |
|                   | $m = 4$ | 88.86 $\pm$ 4.33                 | 78.16 $\pm$ 6.07                 | <b>92.57<math>\pm</math>5.63</b> | 76.91 $\pm$ 3.28                 | 74.10 $\pm$ 3.93                 | 49.67 $\pm$ 5.36                 | 85.85 $\pm$ 2.85                 |
| O. Ricci          | $m = 2$ | <b>90.60<math>\pm</math>3.69</b> | <b>78.60<math>\pm</math>4.79</b> | 89.91 $\pm$ 3.86                 | 77.26 $\pm$ 4.29                 | <b>79.10<math>\pm</math>3.78</b> | <b>54.53<math>\pm</math>3.52</b> | <b>91.40<math>\pm</math>1.24</b> |
|                   | $m = 3$ | 89.14 $\pm$ 4.70                 | 78.15 $\pm$ 5.73                 | <b>91.02<math>\pm</math>6.45</b> | <b>77.72<math>\pm</math>3.36</b> | 78.80 $\pm$ 3.79                 | 53.73 $\pm$ 4.06                 | 89.95 $\pm$ 2.24                 |
|                   | $m = 4$ | 88.39 $\pm$ 6.44                 | 78.17 $\pm$ 5.05                 | 89.85 $\pm$ 6.40                 | 77.35 $\pm$ 4.05                 | 78.10 $\pm$ 3.14                 | 53.87 $\pm$ 5.27                 | 89.65 $\pm$ 2.37                 |
| HKS               | $m = 2$ | 90.62 $\pm$ 4.91                 | 83.95 $\pm$ 2.99                 | <b>95.32<math>\pm</math>5.58</b> | 77.35 $\pm$ 2.86                 | <b>77.90<math>\pm</math>5.72</b> | <b>54.07<math>\pm</math>2.54</b> | <b>90.05<math>\pm</math>2.41</b> |
|                   | $m = 3$ | <b>91.09<math>\pm</math>4.28</b> | <b>85.01<math>\pm</math>4.84</b> | 95.23 $\pm$ 3.89                 | <b>78.17<math>\pm</math>4.54</b> | 76.90 $\pm$ 3.48                 | 53.60 $\pm$ 3.30                 | 88.80 $\pm$ 1.86                 |
|                   | $m = 4$ | 90.63 $\pm$ 4.09                 | 83.75 $\pm$ 5.09                 | 95.12 $\pm$ 6.05                 | 78.07 $\pm$ 2.84                 | 77.00 $\pm$ 4.62                 | 53.40 $\pm$ 3.51                 | 89.00 $\pm$ 1.80                 |

467 **Single vs. multiple filtration functions** (Table 13): We test several node based and edge based  
 468 functions to study how filtration choice affects performance. We observe that single-filtration  
 469 TopoFormer (for example, using only HKS or only Ollivier–Ricci) already achieves strong results,  
 470 while combining filtrations yields modest but consistent improvements on some datasets. This  
 471 indicates that multiple filtrations are a flexible way to incorporate complementary structural signals  
 472 rather than a requirement for good performance.

473 **Discussion.** TOPOFORMER delivers consistently strong performance across a broad range of  
 474 graph classification benchmarks, outperforming state-of-the-art baselines and achieving the best  
 475 overall accuracy on most datasets. These results demonstrate the model’s ability to extract essential  
 476 structural information through topological patterns while producing *fixed-size sequential representations*.  
 477 Such representations are particularly well-suited for Graph Foundation Models, which require  
 478 consistent and transferable embeddings across graphs of varying sizes and domains. Table 2 further  
 479 reveals that among the six topological baselines (PersLay, DMP, FC-V, EMP, MP-HSM, TopoGCL),  
 480 TOPOFORMER achieves the best performance, despite being architecturally simpler and more computa-  
 481 tionally lightweight. This supports our design philosophy that robust topological summaries, when  
 482 properly structured, can outperform more complex pipelines.

483 Crucially, TOPOFORMER departs from the standard GNN paradigm of first learning node embeddings  
 484 followed by global pooling. While effective, this node-centric strategy treats graphs as unstructured  
 485 point clouds in latent space, requiring repeated updates as embeddings evolve, often at the cost of  
 coherence and efficiency (Mesquita et al., 2020; Liu et al., 2023). In contrast, topological models

486 treat graphs as structured wholes and directly encode global patterns. By bypassing intermediate  
 487 node embeddings, TOPOFORMER provides a streamlined and principled approach for learning stable,  
 488 transferable graph-level representations.

489 **Limitations and future work.** Our focus in this work is on a streamlined, graph-level instantiation  
 490 of TOPOFORMER, which also suggests several natural extensions. We restrict attention to low-  
 491 dimensional homology ( $H_0, H_1$ ) on a fixed clique complex with a small set of standard filtrations  
 492 (degree, curvature, HKS); incorporating richer per-slice invariants or learnable filtrations could  
 493 further boost expressivity while keeping the same Topo-Scan + Transformer backbone. Likewise, we  
 494 concentrate on widely used graph-classification and molecular benchmarks, leaving node-/edge-level  
 495 tasks and more heterogeneous settings (e.g., temporal or citation graphs) to future work. Finally,  
 496 Topo-Scan is designed as a lightweight, scan-style summary that complements rather than replaces full  
 497 persistent homology, and we see developing additional theory and applications for such summaries as  
 498 an interesting direction for the TDA community.

## 500 5 CONCLUSION

501 Fixed-size, transferable representations remain a central challenge in graph learning. We introduce  
 502 TOPOFORMER, a scalable framework that encodes multi-scale topological structure into attention-  
 503 ready sequences. By replacing full persistence diagrams with lightweight, slice-wise invariants via  
 504 Topo-Scan, our method integrates seamlessly with transformer architectures while offering theoretical  
 505 stability guarantees. TOPOFORMER achieve state-of-the-art results across graph classification and  
 506 molecular property prediction tasks, with predictable compute and strong generalization. Looking  
 507 ahead, we aim to extend this framework toward graph foundation models by combining topological  
 508 and spectral signals through large-scale self-supervised pretraining, and by adapting to dynamic and  
 509 heterogeneous graphs via learnable filtrations.

## 512 REFERENCES

514 Henry Adams, Tegan Emerson, Michael Kirby, Rachel Neville, Chris Peterson, Patrick Shipman,  
 515 Sofya Chepushtanova, Eric Hanson, Francis Motta, and Lori Ziegelmeier. Persistence images: A  
 516 stable vector representation of persistent homology. *Journal of Machine Learning Research*, 18,  
 517 2017.

518 Mehmet E Aktas, Esra Akbas, and Ahmed El Fatmaoui. Persistence homology of networks: methods  
 519 and applications. *Applied Network Science*, 4(1):1–28, 2019.

521 Dashti Ali, Aras Asaad, Maria-Jose Jimenez, Vudit Nanda, Eduardo Paluzo-Hidalgo, and Manuel  
 522 Soriano-Trigueros. A survey of vectorization methods in topological data analysis. *arXiv preprint*  
 523 *arXiv:2212.09703*, 2022.

524 Håvard Bakke Bjerkevik. On the stability of interval decomposable persistence modules. *Discrete &*  
 525 *Computational Geometry*, 66(1):92–121, 2021.

527 Ulrich Bauer and Michael Lesnick. Induced matchings of barcodes and the algebraic stability of  
 528 persistence. In *Proceedings of the thirtieth annual symposium on Computational geometry*, pp.  
 529 355–364, 2014.

531 Dominique Beaini, Saro Passaro, Vincent Létourneau, Will Hamilton, Gabriele Corso, and Pietro Liò.  
 532 Directional graph networks. In *ICML*, pp. 748–758. PMLR, 2021.

533 Filippo Maria Bianchi, Daniele Grattarola, and Cesare Alippi. Spectral clustering with graph neural  
 534 networks for graph pooling. In Hal III Daumé and Aarti Singh (eds.), *Proceedings of the 37th*  
 535 *International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning*  
 536 *Research*, pp. 874–883. PMLR, July 2020. URL <https://proceedings.mlr.press/v119/bianchi20a.html>.

538 Cristian Bodnar et al. Deep graph mapper: Seeing graphs through the neural lens. *Frontiers in big*  
 539 *Data*, 4:680535, 2021.

540 Karsten Borgwardt et al. Graph kernels: State-of-the-art and future challenges. *Foundations and*  
 541 *Trends® in Machine Learning*, 13(5-6):531–712, 2020.  
 542

543 Karsten M Borgwardt, Cheng Soon Ong, Stefan Schönauer, SVN Vishwanathan, Alex J Smola, and  
 544 Hans-Peter Kriegel. Protein function prediction via graph kernels. *Bioinformatics*, 21(suppl\_1):  
 545 i47–i56, 2005.

546 Magnus Botnan and Michael Lesnick. Algebraic stability of zigzag persistence modules. *Algebraic*  
 547 & *geometric topology*, 18(6):3133–3204, 2018.  
 548

549 Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar,  
 550 Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general intelligence:  
 551 Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*, 2023.

552 Chen Cai and Yusu Wang. Understanding the power of persistence pairing via permutation test. *arXiv*  
 553 *preprint arXiv:2001.06058*, 2020.  
 554

555 Tianle Cai, Shengjie Luo, Keyulu Xu, Di He, Tie-yan Liu, and Liwei Wang. Graphnorm: A principled  
 556 approach to accelerating graph neural network training. In *ICML*, pp. 1204–1215. PMLR, 2021.  
 557

558 Gunnar Carlsson. Topology and data. *Bulletin of the American Mathematical Society*, 46(2):255–308,  
 559 2009.

560 Mathieu Carrière, Frédéric Chazal, Yuichi Ike, Théo Lacombe, Martin Royer, and Yuhei Umeda.  
 561 Perslay: A neural network layer for persistence diagrams and new graph topological signatures. In  
 562 *AISTATS*, pp. 2786–2796, 2020.  
 563

564 Adrià Cereto-Massagué, María José Ojeda, Cristina Valls, Miquel Mulero, Santiago Garcia-Vallvé,  
 565 and Gerard Pujadas. Molecular fingerprint similarity search in virtual screening. *Methods*, 71:  
 566 58–63, 2015.

567 Frédéric Chazal, Brittany Terese Fasy, Fabrizio Lecci, Alessandro Rinaldo, and Larry Wasserman.  
 568 Stochastic convergence of persistence landscapes and silhouettes. In *SoCG*, pp. 474–483, 2014.  
 569

570 Yuzhou Chen, Ignacio Segovia-Dominguez, Cuneyt Gurcan Akcora, Zhiwei Chen, Murat Kantar-  
 571 cioglu, Yulia Gel, and Baris Coskunuzer. Emp: Effective multidimensional persistence for graph  
 572 representation learning. In *Learning on Graphs Conference*, pp. 24–1. PMLR, 2024a.

573 Yuzhou Chen et al. Topogcl: Topological graph contrastive learning. In *AAAI*, volume 38, pp.  
 574 11453–11461, 2024b.  
 575

576 Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger  
 577 Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder–decoder for  
 578 statistical machine translation. In *Proceedings of the 2014 Conference on Empirical Methods in*  
 579 *Natural Language Processing (EMNLP)*, pp. 1724–1734, 2014.

580 Yun Young Choi, Sun Woo Park, Youngho Woo, and U Jin Choi. Cycle to clique (cy2c) graph neural  
 581 network: A sight to see beyond neighborhood aggregation. In *The Eleventh ICLR*, 2022.  
 582

583 David Cohen-Steiner, Herbert Edelsbrunner, and John Harer. Stability of persistence diagrams.  
 584 *Discrete & Computational Geometry*, 37(1):103–120, 2007. doi: 10.1007/s00454-006-1276-5.  
 585

586 Gabriele Corso, Luca Cavalleri, Dominique Beaini, Pietro Liò, and Petar Veličković. Principal  
 587 neighbourhood aggregation for graph nets. *NeurIPS*, 33:13260–13271, 2020.

588 Baris Coskunuzer and Cüneyt Gürcan Akçora. Topological methods in machine learning: A tutorial  
 589 for practitioners. *arXiv preprint arXiv:2409.02901*, 2024.  
 590

591 Thibault de Surrel, Felix Hensel, Mathieu Carrière, Théo Lacombe, Yuichi Ike, Hiroaki Kurihara,  
 592 Marc Glisse, and Frédéric Chazal. Ripsnet: a general architecture for fast and robust estimation  
 593 of the persistent homology of point clouds. In *Topological, algebraic and geometric learning*  
 workshops 2022, pp. 96–106. PMLR, 2022.

594 Andac Demir, Baris Coskunuzer, Yulia Gel, Ignacio Segovia-Dominguez, Yuzhou Chen, and Bulent  
595 Kiziltan. Todd: Topological compound fingerprinting in computer-aided drug discovery. *NeurIPS*,  
596 35:27978–27993, 2022.

597 Jacob Devlin. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv*  
598 *preprint arXiv:1810.04805*, 2018.

600 Tamal Krishna Dey and Yusu Wang. *Computational Topology for Data Analysis*. Cambridge  
601 University Press, 2022.

602 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas  
603 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, and Sylvain Gelly. An image is  
604 worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*,  
605 2020.

606 Edelsbrunner, Letscher, and Zomorodian. Topological persistence and simplification. *Discrete &*  
607 *computational geometry*, 28(4):511–533, 2002.

609 Federico Errica, Marco Podda, Davide Bacciu, and Alessio Micheli. A fair comparison of graph  
610 neural networks for graph classification. In *ICLR*, 2020.

612 Xiaomin Fang, Lihang Liu, Jieqiong Lei, Donglong He, Shanzhuo Zhang, Jingbo Zhou, Fan Wang,  
613 Hua Wu, and Haifeng Wang. Geometry-enhanced molecular representation learning for property  
614 prediction. *Nature Machine Intelligence*, 4(2):127–134, 2022.

615 Yin Fang, Qiang Zhang, Ningyu Zhang, Zhuo Chen, Xiang Zhuang, Xin Shao, Xiaohui Fan, and  
616 Huajun Chen. Knowledge graph-enhanced molecular contrastive learning with functional prompt.  
617 *Nature Machine Intelligence*, 5(5):542–553, 2023a.

618 Zhongxi Fang, Jianming Huang, Xun Su, and Hiroyuki Kasai. Wasserstein graph distance based on  
619 11-approximated tree edit distance between weisfeiler-lehman subtrees. In *AAAI*, volume 37, pp.  
620 7539–7549, 2023b.

622 Hongyang Gao and Shuiwang Ji. Graph u-nets. In Kamalika Chaudhuri and Ruslan Salakhutdinov  
623 (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of  
624 *Proceedings of Machine Learning Research*, pp. 2083–2092. PMLR, June 2019. URL <https://proceedings.mlr.press/v97/gao19a.html>.

626 Johannes Gasteiger, Florian Becker, and Stephan Günnemann. Gemnet: Universal directional graph  
627 neural networks for molecules. *NeurIPS*, 34:6790–6802, 2021.

629 Xiaotian Han, Zhimeng Jiang, Ninghao Liu, and Xia Hu. G-mixup: Graph data augmentation for  
630 graph classification. In *ICML*, pp. 8230–8248. PMLR, 2022.

631 Felix Hensel, Michael Moor, and Bastian Rieck. A survey of topological machine learning methods.  
632 *Frontiers in Artificial Intelligence*, 4:52, 2021.

634 Jaeseung Heo, Seungbeom Lee, Sungsoo Ahn, and Dongwoo Kim. Epic: Graph augmentation with  
635 edit path interpolation via learnable cost. In *IJCAI*, 2024.

636 Christoph Hofer, Florian Graf, Bastian Rieck, Marc Niethammer, and Roland Kwitt. Graph filtration  
637 learning. In *ICML*, pp. 4314–4323. PMLR, 2020.

638 Shion Honda, Shoi Shi, and Hiroki R Ueda. Smiles transformer: Pre-trained molecular fingerprint  
639 for low data drug discovery. *arXiv preprint arXiv:1911.04738*, 2019.

641 Max Horn, Edward De Brouwer, Michael Moor, Yves Moreau, Bastian Rieck, and Karsten Borgwardt.  
642 Topological graph neural networks. In *ICLR*, 2021.

643 W Hu, B Liu, J Gomes, M Zitnik, P Liang, V Pande, and J Leskovec. Strategies for pre-training  
644 graph neural networks. In *ICLR*, 2020a.

646 Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele Catasta,  
647 and Jure Leskovec. Open graph benchmark: Datasets for machine learning on graphs. *NeurIPS*,  
33:22118–22133, 2020b. URL <https://ogb.stanford.edu>.

648 Johanna Immonen, Amauri Souza, and Vikas Garg. Going beyond persistent homology using  
 649 persistent homology. *NeurIPS*, 36, 2024.  
 650

651 Tiago Janela and Jürgen Bajorath. Simple nearest-neighbour analysis meets the accuracy of compound  
 652 potency predictions using complex machine learning models. *Nature Machine Intelligence*, 4(12):  
 653 1246–1255, 2022.

654 Dejun Jiang, Zhenxing Wu, Chang-Yu Hsieh, Guangyong Chen, Ben Liao, Zhe Wang, Chao Shen,  
 655 Dongsheng Cao, Jian Wu, and Tingjun Hou. Could graph neural networks learn better molecular  
 656 representation for drug discovery? a comparison study of descriptor-based and graph-based models.  
 657 *Journal of cheminformatics*, 13:1–23, 2021.

658

659 Talia B Kimber, Maxime Gagnebin, and Andrea Volkamer. Maxsmi: maximizing molecular property  
 660 prediction performance with confidence estimation using smiles augmentation and deep learning.  
 661 *Artificial Intelligence in the Life Sciences*, 1:100014, 2021.

662 Nils Kriege et al. Subgraph matching kernels for attributed graphs. In *ICML*, pp. 291–298, 2012.  
 663

664 Junhyun Lee, Inyeop Lee, and Jaewoo Kang. Self-attention graph pooling. In *ICML*, pp. 3734–3743.  
 665 PMLR, 2019.

666 Sihang Li, Xiang Wang, An Zhang, Yingxin Wu, Xiangnan He, and Tat-Seng Chua. Let invariant  
 667 rationale discovery inspire graph contrastive learning. In *ICML*, pp. 13052–13065. PMLR, 2022a.

668

669 Zhen Li, Mingjian Jiang, Shuang Wang, and Shugang Zhang. Deep learning methods for molecular  
 670 representation and property prediction. *Drug Discovery Today*, 27(12):103373, 2022b.

671 Chuang Liu, Yibing Zhan, Jia Wu, Chang Li, Bo Du, Wenbin Hu, Tongliang Liu, and Dacheng Tao.  
 672 Graph pooling for graph neural networks: progress, challenges, and opportunities. In *Proceedings  
 673 of the Thirty-Second International Joint Conference on Artificial Intelligence*, pp. 6712–6722,  
 674 2023.

675

676 Shengchao Liu, Hanchen Wang, Weiyang Liu, Joan Lasenby, Hongyu Guo, and Jian Tang. Pre-  
 677 training molecular graph representation with 3d geometry. In *ICLR Workshop on Geometrical and  
 678 Topological Representation Learning*, 2022a.

679

680 Shengchao Liu et al. N-gram graph: Simple unsupervised representation for graphs, with applications  
 681 to molecules. *NeurIPS*, 32, 2019.

682

683 Yi Liu, Limei Wang, Meng Liu, Yuchao Lin, Xuan Zhang, Bora Oztekin, and Shuiwang Ji. Spherical  
 684 message passing for 3d molecular graphs. In *ICLR*, 2022b.

685

686 David Loiseaux, Luis Scoccola, Mathieu Carrière, Magnus Bakke Botnan, and Steve Oudot. Stable  
 687 vectorization of multiparameter persistent homology using signed barcodes as measures. *NeurIPS*,  
 2023.

688

689 David Loiseaux, Luis Scoccola, Mathieu Carrière, Magnus Bakke Botnan, and Steve Oudot. Stable  
 690 vectorization of multiparameter persistent homology using signed barcodes as measures. *NeurIPS*,  
 36, 2024.

691

692 Yizhen Luo, Kai Yang, Massimo Hong, Xing Yi Liu, Zikun Nie, Hao Zhou, and Zaiqing Nie. Learning  
 693 multi-view molecular representations with structured and unstructured knowledge. In *Proceedings  
 694 of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 2082–2093,  
 2024.

695

696 Yao Ma, Suhang Wang, Charu C. Aggarwal, and Jiliang Tang. Graph convolutional networks with  
 697 eigenpooling. arXiv:1904.13107, 2019.

698

699 Diego Mesquita, Amauri Souza, and Samuel Kaski. Rethinking pooling in graph neural networks.  
 700 *NeurIPS*, 33:2220–2231, 2020.

701 Giannis Nikolentzos et al. Graph alignment kernels using weisfeiler and leman hierarchies. In  
 702 *International Conference on Artificial Intelligence and Statistics*, pp. 2019–2034. PMLR, 2023.

702 Ippei Obayashi, Takenobu Nakamura, and Yasuaki Hiraoka. Persistent homology analysis for  
 703 materials research and persistent homology software: Homcloud. *Journal of the physical society of*  
 704 *japan*, 91(9):091013, 2022.

705

706 Leslie O’Bray et al. Filtration curves for graph representation. In *Proceedings of the 27th ACM*  
 707 *SIGKDD Conference on Knowledge Discovery & Data Mining*, pp. 1267–1275, 2021.

708 Nina Otter, Mason A Porter, Ulrike Tillmann, Peter Grindrod, and Heather A Harrington. A roadmap  
 709 for the computation of persistent homology. *EPJ Data Science*, 6:1–38, 2017.

710

711 Phu Pham, Quang-Thinh Bui, Ngoc Thanh Nguyen, Robert Kozma, Philip S Yu, and Bay Vo.  
 712 Topological data analysis in graph neural networks: Surveys and perspectives. *IEEE Transactions*  
 713 *on Neural Networks and Learning Systems*, 2025.

714 Chen Qian, Huayi Tang, Hong Liang, and Yong Liu. Reimagining graph classification from a  
 715 prototype view with optimal transport: Algorithm and theorem. In *Proceedings of the 30th ACM*  
 716 *SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 2444–2454, 2024.

717

718 Alec Radford et al. Learning transferable visual models from natural language supervision. In *ICML*,  
 719 pp. 8748–8763. PMLR, 2021.

720 Ladislav Rampášek, Michael Galkin, Vijay Prakash Dwivedi, Anh Tuan Luu, Guy Wolf, and Do-  
 721 minique Beaini. Recipe for a general, powerful, scalable graph transformer. *Advances in Neural*  
 722 *Information Processing Systems*, 35:14501–14515, 2022.

723

724 David Rogers and Mathew Hahn. Extended-connectivity fingerprints. *Journal of Chemical Informa-*  
 725 *tion and Modeling*, 50(5):742–754, 2010.

726

727 Yu Rong, Yatao Bian, Tingyang Xu, Weiyang Xie, Ying Wei, Wenbing Huang, and Junzhou Huang.  
 728 Self-supervised graph transformer on large-scale molecular data. *NeurIPS*, 33:12559–12571, 2020.

729

730 Christopher Shultz. Applications of topological data analysis in economics. Available at *SSRN*  
 4378151, 2023.

731

732 Yara Skaf et al. Topological data analysis in biomedicine: A review. *Journal of Biomedical*  
 733 *Informatics*, 130:104082, 2022.

734

735 Yeongyeong Son, Dasom Noh, Gyoungyoung Heo, Gyoung Jin Park, and Sunyoung Kwon. More:  
 736 Molecule pretraining with multi-level pretext task. In *Proceedings of the AAAI Conference on*  
 737 *Artificial Intelligence*, volume 39, pp. 20531–20539, 2025.

738

739 Ying Song, Shuangjia Zheng, Zhangming Niu, Zhang-Hua Fu, Yutong Lu, and Yuedong Yang.  
 740 Communicative representation learning on attributed molecular graphs. In *IJCAI*, volume 2020,  
 pp. 2831–2838, 2020.

741

742 Afnan Sultan, Jochen Sieg, Miriam Mathea, and Andrea Volkamer. Transformers for molecular  
 743 property prediction: Lessons learned from the past five years. *Journal of Chemical Information*  
 744 *and Modeling*, 64(16):6259–6280, 2024.

745

746 Matteo Togninalli, Elisabetta Ghisu, Felipe Llinares-López, Bastian Rieck, and Karsten Borgwardt.  
 Wasserstein weisfeiler-lehman graph kernels. In *NeurIPS*, pp. 6439–6449, 2019.

747

748 Astrit Tola, Funmilola Mary Taiwo, Cuneyt Gurcan Akcora, and Baris Coskunuzer. Toper: Topologi-  
 749 cal embeddings in graph representation learning. *NeurIPS*, 2025.

750

751 Cecile Valsecchi, Magda Collarile, Francesca Grisoni, Roberto Todeschini, Davide Ballabio, and  
 752 Viviana Consonni. Predicting molecular activity on nuclear receptors by multitask neural networks.  
 753 *Journal of Chemometrics*, 36(2):e3325, 2022.

754

755 A Vaswani et al. Attention is all you need. *NeurIPS*, 2017.

Yogesh Verma, Amauri H Souza, and Vikas Garg. Topological neural networks go persistent,  
 equivariant, and continuous. In *ICML*, 2024.

756 Haotao Wang et al. Graph mixture of experts: Learning on large-scale graphs with explicit diversity  
 757 modeling. *NeurIPS*, 36, 2024.

758

759 Yuyang Wang, Jianren Wang, Zhonglin Cao, and Amir Barati Farimani. Molecular contrastive  
 760 learning of representations via graph neural networks. *Nature Machine Intelligence*, 4(3):279–287,  
 761 2022.

762 Yu Guang Wang, Ming Li, Zheng Ma, Guido Montúfar, Xiaosheng Zhuang, and Yanan Fan. Haar  
 763 graph pooling. In Hal III Daumé and Aarti Singh (eds.), *Proceedings of the 37th International  
 764 Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*,  
 765 pp. 9952–9962. PMLR, July 2020. URL <https://proceedings.mlr.press/v119/wang20m.html>.

766

767 Oliver Wieder, Stefan Kohlbacher, Mélaine Kuenemann, Arthur Garon, Pierre Ducrot, Thomas  
 768 Seidel, and Thierry Langer. A compact review of molecular property prediction with graph neural  
 769 networks. *Drug Discovery Today: Technologies*, 37:1–12, 2020.

770

771 Daniel S Wigh, Jonathan M Goodman, and Alexei A Lapkin. A review of molecular representation in  
 772 the age of machine learning. *Wiley Interdisciplinary Reviews: Computational Molecular Science*,  
 773 12(5):e1603, 2022.

774

775 Asiri Wijesinghe et al. A new perspective on " how graph neural networks go beyond weisfeiler-  
 776 lehman? ". In *ICLR*, 2021.

777

778 Zhengin Wu, Bharath Ramsundar, Evan N Feinberg, Joseph Gomes, Caleb Geniesse, Aneesh S  
 779 Pappu, Karl Leswing, and Vijay Pande. Moleculenet: a benchmark for molecular machine learning.  
*Chemical science*, 9(2):513–530, 2018.

780

781 Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and Philip S Yu. A  
 782 comprehensive survey on graph neural networks. *IEEE transactions on neural networks and  
 783 learning systems*, 32(1):4–24, 2020.

784

785 Jun Xia, Lecheng Zhang, Xiao Zhu, Yue Liu, Zhangyang Gao, Bozhen Hu, Cheng Tan, Jiangbin  
 786 Zheng, Siyuan Li, and Stan Z Li. Understanding the limitations of deep models for molecular  
 787 property prediction: Insights and solutions. *NeurIPS*, 36:64774–64792, 2023.

788

789 Zhaoping Xiong, Dingyan Wang, Xiaohong Liu, Feisheng Zhong, Xiaozhe Wan, Xutong Li, Zhaojun  
 790 Li, Xiaomin Luo, Kaixian Chen, Hualiang Jiang, et al. Pushing the boundaries of molecular  
 791 representation for drug discovery with the graph attention mechanism. *Journal of medicinal  
 792 chemistry*, 63(16):8749–8760, 2019.

793

794 Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural  
 795 networks? *ICLR*, 2018.

796

797 Zuoyu Yan, Tengfei Ma, Liangcai Gao, Zhi Tang, Yusu Wang, and Chao Chen. Neural approximation  
 798 of graph topological features. *Advances in neural information processing systems*, 35:33357–33370,  
 799 2022.

800

801 Pinar Yanardag et al. Deep graph kernels. In *Proceedings of the 21th ACM SIGKDD international  
 802 conference on knowledge discovery and data mining*, pp. 1365–1374, 2015.

803

804 Jingbo Yang, Yiyang Cai, Kairui Zhao, Hongbo Xie, and Xiujie Chen. Concepts and applications of  
 805 chemical fingerprint for hit and lead screening. *Drug discovery today*, 27(11):103356, 2022.

806

807 Chun Wei Yap. Padel-descriptor: An open source software to calculate molecular descriptors and  
 808 fingerprints. *Journal of Computational Chemistry*, 32(7):1466–1474, 2011.

809

810 Wei Ye, Shuhao Tang, Hao Tian, and Qijun Chen. Beyond histogram comparison: Distribution-aware  
 811 simple-path graph kernels. *IEEE Transactions on Artificial Intelligence*, 2025.

812

813 Yihang Yin, Qingzhong Wang, Siyu Huang, Haoyi Xiong, and Xiang Zhang. Autogcl: Automated  
 814 graph contrastive learning via learnable view generators. In *AAAI*, volume 36, pp. 8892–8900,  
 815 2022.

810 Chengxuan Ying et al. Do transformers really perform badly for graph representation? *NeurIPS*, 34:  
811 28877–28888, 2021.

812

813 Zhitao Ying, Jiaxuan You, Christopher Morris, Xiang Ren, William Hamilton, and Jure Leskovec.  
814 Hierarchical graph representation learning with differentiable pooling. In *Advances in Neural*  
815 *Information Processing Systems*, volume 31, pp. 4800–4810. Curran Associates, Inc., 2018.

816

817 Jaemin Yoo et al. Model-agnostic augmentation for accurate graph classification. In *Proceedings of*  
818 *the ACM Web Conference 2022*, pp. 1281–1291, 2022.

819

820 Zaixi Zhang, Qi Liu, Hao Wang, Chengqiang Lu, and Chee-Kong Lee. Motif-based graph self-  
821 supervised learning for molecular property prediction. *NeurIPS*, 34:15870–15882, 2021.

822

823 Yan Zheng, Song Wu, Junyu Lin, Yazhou Ren, Jing He, Xiaorong Pu, and Lifang He. Cross-view  
824 contrastive fusion for enhanced molecular property prediction. In *Proceedings of the Thirty-Third*  
825 *International Joint Conference on Artificial Intelligence*, volume 2, 2024.

826

827 Zhiqiang Zhong, Kuangyu Zhou, and Davide Mottin. Benchmarking large language models for  
828 molecule prediction tasks. *arXiv:2403.05075*, 2024.

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

## A TOPOFORMER\*: TOPOFORMER FOR MPP

868

869

### A.1 MOLECULAR FINGERPRINTS

870

871

872

873

874

875

Molecular fingerprints are widely used in computational chemistry and machine learning to represent molecular structures as fixed-length numerical vectors (Cereto-Massagué et al., 2015). They encode features such as atomic connectivity and substructural patterns, enabling efficient similarity search and predictive modeling. Popular methods include ECFP (Extended Connectivity Fingerprints) and PubChemFP, both extensively applied in drug discovery, virtual screening, and bioinformatics (Yang et al., 2022).

876

877

878

879

880

881

**ECFP Fingerprints.** Extended Connectivity Fingerprints (ECFP) capture structural features by iteratively hashing local atomic environments up to a specified radius (Rogers & Hahn, 2010). Unlike traditional hashed fingerprints, ECFP preserves substructural detail, making it effective for similarity search, QSAR modeling, and property prediction. It is invariant to atom ordering while retaining connectivity, enabling fine-grained molecular feature analysis. For a recent overview of ECFP fingerprints and their role in modern biochemical ML pipelines, see (Wigh et al., 2022).

882

883

### A.2 TOPOFORMER\* MODEL

884

885

886

887

888

889

890

891

892

893

For Molecular Property Prediction Task, we employ a hybrid model, TOPOFORMER\*, combining ECFP Fingerprints and our TOPOFORMER Model. This hybrid model shows the versatility of our TOPOFORMER model on its effective integration with complementary information (Table 7). We give the flowchart of our hybrid model in Figure 5. In our hybrid model, we used the same experimental setup for the TOPOFORMER component. For the MLP component, we employed a two-layer MLP with a hidden dimension of 200, ensuring that its output dimension matches the output dimension of the TOPOFORMER model. The model was optimized using the Adam optimizer with a learning rate of 0.01 and a weight decay of 1e-4. Both the MLP and TOPOFORMER components were trained in an end-to-end manner, allowing the model to leverage both topological signatures and complementary graph information, ultimately leading to improved performance.

894

895

896

Table 7: Performance comparison of standalone models (TOPOFORMER and FP-MLP) and the hybrid model (TOPOFORMER\*) in random splitting (8:1:1).

897

|         | PH-TR            | TOPOFORMER                       | FP-MLP           | PH+ECFP+TR                       | TOPOFORMER*                      |
|---------|------------------|----------------------------------|------------------|----------------------------------|----------------------------------|
| BACE    | 72.41 $\pm$ 3.15 | 83.29 $\pm$ 2.14                 | 90.29 $\pm$ 2.67 | 90.60 $\pm$ 2.99                 | <b>91.60<math>\pm</math>1.73</b> |
| HIV     | 69.29 $\pm$ 1.65 | 75.81 $\pm$ 0.23                 | 83.26 $\pm$ 1.01 | 83.97 $\pm$ 1.51                 | <b>85.10<math>\pm</math>0.49</b> |
| BBBP    | 83.37 $\pm$ 3.90 | <b>94.54<math>\pm</math>1.01</b> | 89.68 $\pm$ 3.46 | 93.47 $\pm$ 2.53                 | <b>95.90<math>\pm</math>0.28</b> |
| ClinTox | 75.89 $\pm$ 6.60 | <b>83.42<math>\pm</math>2.33</b> | 76.34 $\pm$ 6.54 | 82.04 $\pm$ 7.12                 | <b>86.20<math>\pm</math>3.83</b> |
| SIDER   | 62.91 $\pm$ 3.49 | 62.10 $\pm$ 1.44                 | 65.30 $\pm$ 0.99 | <b>66.99<math>\pm</math>1.85</b> | 66.80 $\pm$ 0.29                 |
| Tox21   | 68.24 $\pm$ 1.60 | 80.87 $\pm$ 0.19                 | 77.89 $\pm$ 1.54 | 79.03 $\pm$ 1.21                 | <b>81.50<math>\pm</math>1.85</b> |
| ToxCast | 64.74 $\pm$ 2.29 | 73.37 $\pm$ 1.42                 | 74.69 $\pm$ 1.33 | 75.73 $\pm$ 1.59                 | <b>78.40<math>\pm</math>1.57</b> |

904

905

906

### A.3 HYBRID CLASSICAL MPP BASELINES

907

908

909

910

911

912

913

914

915

916

917

We refer to the classical models combined with modern ML models as *Hybrid Classical (HC) Models*. The first family of HC baseline models consists of *Fingerprinting models* (Jiang et al., 2021), which use vectorized molecular fingerprints as input to traditional machine learning models, including SVM, XGB, RF, and MLP. The input fingerprints are a concatenation of 881-dimensional PubChem fingerprints (PubChemFP), 307-dimensional substructure fingerprints (SubFP), and 206-dimensional MOE 1-D and 2-D descriptors (Yap, 2011). The second family of baseline models comprises *SMILES models*, which treat SMILES strings as sequential input to 1D CNN (Kimber et al., 2021), a 3-layer bidirectional GRU (Cho et al., 2014), and a pre-trained SMILES transformer (TRSF) (Honda et al., 2019). The third family is *GNN models* which use 2D graph-based representations of compounds, where atom and bond features are encoded using one-hot schemes and fed into GCN, MPNN, GAT, and AFP models (Xiong et al., 2019). Another baseline is the SPN model, using SphereNet (Liu et al., 2022b), which employs 3D graphs of compounds as input.

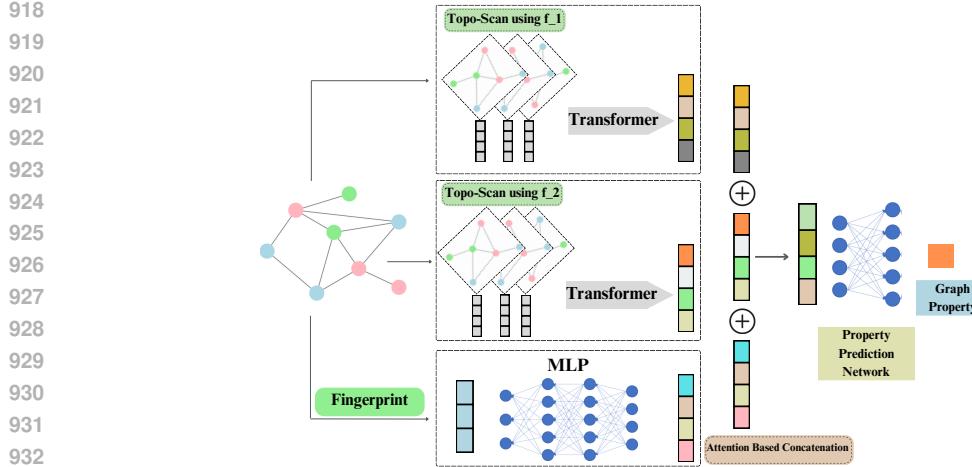


Figure 5: TOPOFORMER\*: To successfully integrate complementary graph information, such as ECFP, with our TOPOFORMER model, we employ a MLP. The MLP output is combined with the TOPOFORMER model using an attention mechanism, and the combined representation is then passed through a graph prediction network to perform the final prediction task.

Table 8: **Hybrid Classical MPP Models.** The ROC AUC results of ML models for molecular property prediction tasks with random splitting (8:1:1). The baseline results are reported from (Xia et al., 2023). The best and the second best performances are given in **bold**, and underlined, respectively.

| Dataset | Fingerprinting Models |             |             |      | SMILES Models |             |      | GNN Models  |             |             |      |      | Ours        |  |
|---------|-----------------------|-------------|-------------|------|---------------|-------------|------|-------------|-------------|-------------|------|------|-------------|--|
|         | SVM                   | XGB         | RF          | CNN  | RNN           | TRSF        | MLP  | GCN         | MPNN        | GAT         | AFP  | SPN  | TOPOFORMER* |  |
| BBBP    | 91.3                  | <u>92.6</u> | 92.3        | 89.7 | 76.0          | 69.3        | 89.7 | 91.8        | 91.5        | 87.2        | 90.2 | 90.5 | <b>96.6</b> |  |
| Tox21   | 82.0                  | 83.7        | 83.1        | 81.2 | 73.7          | 76.8        | 79.9 | <u>84.6</u> | 82.1        | <u>84.5</u> | 82.7 | 82.5 | 81.5        |  |
| ToxCast | 72.5                  | <u>78.5</u> | 77.8        | 73.5 | 67.8          | 78.0        | 78.1 | 76.7        | <u>78.8</u> | 77.2        | 76.8 | 77.2 | 78.4        |  |
| SIDER   | 62.6                  | 63.8        | 64.4        | 59.1 | 51.5          | <u>64.1</u> | 61.7 | 62.3        | 60.3        | 62.0        | 61.3 | 61.3 | <b>66.8</b> |  |
| ClinTox | 87.9                  | 91.9        | <u>93.0</u> | 88.8 | 68.5          | <u>96.3</u> | 93.0 | 88.9        | 86.8        | 89.8        | 87.9 | 91.2 | 86.2        |  |
| BACE    | 88.6                  | <u>89.6</u> | 89.0        | 81.5 | 55.9          | 83.5        | 88.7 | 88.0        | 84.6        | 88.6        | 87.9 | 88.2 | <b>91.6</b> |  |
| HIV     | 81.7                  | <u>83.9</u> | 82.0        | 82.6 | 73.3          | 74.8        | 79.1 | 83.4        | 81.4        | 81.2        | 81.8 | 81.8 | <b>85.1</b> |  |

## B PROOFS OF STABILITY THEOREMS

We work on the fixed clique complex  $\widehat{G}$  of  $G = (V, E)$ . For a node function  $h : V \rightarrow \mathbb{R}$ , we use the upper-star extension  $\widehat{h}(\sigma) = \max_{v \in \sigma} h(v)$  and the associated sublevel filtration on  $\widehat{G}$ . Throughout,  $k \in \{0, 1\}$  is the homological dimension used in our tokens.

**Preliminaries.** For  $a \leq b$ , define the interlevel (level-set) subcomplex  $\widehat{G}_{[a,b]}^h := \{\sigma \in \widehat{G} : a \leq \min_{v \in \sigma} h(v) \text{ and } \max_{v \in \sigma} h(v) \leq b\}$ . The associated pointwise finite-dimensional *interlevel persistence module* is the functor  $M_k^h : (a, b) \mapsto H_k(\widehat{G}_{[a,b]}^h)$ . Given a shared grid  $\alpha_0 < \dots < \alpha_N$ , window width  $m$ , and stride  $s$ , the Topo-Scan token at window  $t$  is

$$\widehat{\beta}_k^h(t) = \dim M_k^h(\alpha_{ts}, \alpha_{ts+m}), \quad t = 0, \dots, T-1, \quad T = \left\lfloor \frac{N-m}{s} \right\rfloor + 1.$$

We write  $d_B(M_k^f, M_k^g)$  for the bottleneck distance between the interval decompositions (barcodes) of the interlevel modules  $M_k^f$  and  $M_k^g$ .

### Two stability lemmas.

**Lemma B.1** (Interlevel stability). (Botnan & Lesnick, 2018, Thm 1.1 & 1.2) For  $k \geq 0$ , the interlevel modules of the upper-star filtrations induced by  $f, g : V \rightarrow \mathbb{R}$  on the fixed clique complex  $\widehat{G}$  satisfy

$$d_B(M_k^f, M_k^g) \leq \|f - g\|_\infty.$$

972 **Lemma B.2** (Lipschitzness of interval rank). *(Bauer & Lesnick, 2014; Bakke Bjerkevik, 2021)* Let  
 973  $M, N$  be interval-decomposable, p.f.d. modules with  $d_B(M, N) \leq \delta$ . For any interval  $I = [a, b]$ ,  
 974

$$975 \quad |\dim M(a, b) - \dim N(a, b)| \leq \mathcal{B}_M(I, \delta) + \mathcal{B}_N(I, \delta),$$

976 where  $\mathcal{B}_M(I, \delta)$  counts bars in  $\text{Bar}(M)$  whose endpoints lie within  $\delta$  of the boundary  $\{a, b\}$  (and  
 977 similarly for  $N$ ).

978 **Theorem 3.1.** *With the setup above, there exists  $C = C(\widehat{G}, \{\alpha_i\}, m, s)$  such that*

$$980 \quad \sum_{t=0}^{T-1} |\widehat{\beta}_k^f(t) - \widehat{\beta}_k^g(t)| \leq C d_B(M_k^f, M_k^g).$$

983 *Proof of Theorem 3.1.* Fix  $t$  and write  $I_t = [\alpha_{ts}, \alpha_{ts+m}]$ . By Lemma B.2, there exists a finite  
 984 constant  $C_0(\widehat{G}, I_t)$  such that  $|\dim M_k^f(I_t) - \dim M_k^g(I_t)| \leq C_0(\widehat{G}, I_t) d_B(M_k^f, M_k^g)$ . Summing  
 985 over  $t$  gives

$$987 \quad \sum_{t=0}^{T-1} |\widehat{\beta}_k^f(t) - \widehat{\beta}_k^g(t)| \leq \left( \sum_{t=0}^{T-1} C_0(\widehat{G}, I_t) \right) d_B(M_k^f, M_k^g) := C d_B(M_k^f, M_k^g).$$

990 On a fixed finite complex and fixed grid, the  $C_0(\widehat{G}, I_t)$  are finite and can be uniformly bounded,  
 991 yielding  $C = T C_0$ .  $\square$

992 **Corollary 3.2.** *For upper-star filtrations on a fixed complex,  $d_B(M_k^f, M_k^g) \leq \|f - g\|_\infty$  (Lemma B.1),  
 993 hence  $\|\widehat{\beta}_k(G, f) - \widehat{\beta}_k(G, g)\|_1 \leq C \|f - g\|_\infty$ .*

995 *Proof of Theorem 3.1.* By Theorem 3.1, we have  $\|\widehat{\beta}_k(G, f) - \widehat{\beta}_k(G, g)\|_1 \leq C d_B(M_k^f, M_k^g)$ .

997 By Lemma B.1 (interlevel/level-set stability on the fixed clique complex),  $d_B(M_k^f, M_k^g) \leq \|f - g\|_\infty$ .  
 998 Combining the two inequalities yields the claim.  $\square$

999 **Connection to classical sublevel stability.** The inequality  $d_B \leq \|f - g\|_\infty$  is classical for sublevel  
 1000 filtrations on a fixed space (Cohen-Steiner et al., 2007). Our Lemma B.1 is the level-set (interlevel)  
 1001 analogue on the fixed clique complex, following algebraic stability for zigzag/level-set modules  
 1002 (e.g., Botnan & Lesnick, 2018). We use this interlevel version to handle windowed intervals  $[a, b]$   
 1003 appearing in Topo-Scan.

1004 **Shared thresholds.** The theorem assumes a shared grid  $\{\alpha_i\}$ . If thresholds are chosen separately (e.g.,  
 1005 per-function quantiles), a monotone reparameterization of the filtration axis induces an additional  
 1006 term proportional to the grid displacement, which can be absorbed into  $C$ .

1007 *Remark B.3* (Relation to PH invariants and stable ranks). Our stability theorem focuses on the  $\ell_1$   
 1008 robustness of the discrete Topo-Scan sequences  $\widehat{\beta}_k(G, h)$ , but these sequences implicitly encode  
 1009 familiar PH objects. For a fixed filtration function  $h$ , the map  $t \mapsto \widehat{\beta}_k^h(t)$  can be viewed as a sampled  
 1010 version of the rank invariant  $(a, b) \mapsto \text{rank } H_k((\widehat{G})_{[a, b]}^h)$  associated with the interlevel module  $M_k^h$ .  
 1011 In this sense, Topo-Scan produces a coarse, structured discretization of the same information that  
 1012 barcodes and stable vectorizations of persistence diagrams, such as persistence landscapes, silhouettes  
 1013 and persistence images (Chazal et al., 2014; Adams et al., 2017), summarize in continuous form.  
 1014 Similarly, Graph Filtration Learning (Hofer et al., 2020) can be seen as learning the filtration function  
 1015  $h$ , while our work fixes  $h$  and instead changes the representation from global barcodes to local  
 1016 interlevel sequences. A full expressivity comparison and formal information-loss bounds relative to  
 1017 complete barcodes are interesting directions for future work.

## 1019 C MORE ON TOPOFORMER

### 1021 C.1 BASE MODEL: TRANSFORMER

1023 Our TOPOFORMER model is designed for classification tasks using sequential inputs, harnessing  
 1024 transformers for efficient feature extraction. The architecture includes an embedding layer, a trans-  
 1025 former encoder, and a fully connected (FC) classification head, with regularization techniques applied  
 to mitigate overfitting.

1026 Let  $\mathbf{x} = (x_1, x_2, \dots, x_T) \in \mathbb{R}^{N \times T \times D}$  represent the input sequence, where  $N$  is number of graphs,  
 1027  $T$  is the sequence length, and  $D$  is the dimensionality of each input token. The input sequence is first  
 1028 passed through an embedding layer  $\mathbf{E} : \mathbb{R}^D \rightarrow \mathbb{R}^H$ , where  $H$  denotes the embedding dimension.  
 1029 In addition, a positional encoding matrix  $\mathbf{P} \in \mathbb{R}^{1 \times T \times H}$  is added to the embeddings to encode the  
 1030 positional information of the sequence, resulting in a sequence of embedded vectors  $\mathbf{e}_t = \mathbf{E}(x_t) + \mathbf{P}_t$   
 1031 for  $t = 1, 2, \dots, T$ , where  $\mathbf{P}_t$  is the positional encoding for position  $t$ .

1032 The sequence of embeddings is then passed through a multi-layer transformer encoder, where the  
 1033 encoder operates on the embedded sequence  $\mathbf{E}(\mathbf{x}) + \mathbf{P} \in \mathbb{R}^{T \times B \times H}$ , with  $B$  representing the  
 1034 batch size. The transformer encoder generates a new sequence of output representations  $\mathbf{z} =$   
 1035  $(z_1, z_2, \dots, z_T) \in \mathbb{R}^{T \times B \times H}$ . After processing through the encoder, the output sequence is permuted  
 1036 and reshaped to a flattened vector of size  $B \times (T \cdot H)$ , ensuring compatibility with subsequent fully  
 1037 connected layers.

1038 The flattened representation  $\mathbf{z}_{\text{flat}} \in \mathbb{R}^{B \times (T \cdot H)}$  is then passed through a batch normalization  
 1039 layer,  $\text{BN}(\mathbf{z}_{\text{flat}})$ , which normalizes the activations across the batch to stabilize the training pro-  
 1040 cess. A dropout layer  $\mathbf{D}(\cdot)$  is then applied to the normalized output to regularize the model and  
 1041 mitigate overfitting. The final classification output is obtained through a fully connected layer  
 1042  $\text{FC} : \mathbb{R}^{B \times (T \cdot H)} \rightarrow \mathbb{R}^H$ .

1043

## 1044 C.2 DUAL TRANSFORMER WITH MULTI-LAYER PERCEPTRON CLASSIFIER

1045

1046 This model combines multiple sources of input data through a hybrid architecture that integrates  
 1047 two independent base models and a multi-layer perceptron (MLP). This model is designed to handle  
 1048 diverse input modalities by leveraging the strengths of both transformers and MLPs for feature  
 1049 extraction and classification.

1050 Let  $\mathbf{X}_1 \in \mathbb{R}^{N \times T_1 \times D_1}$ ,  $\mathbf{X}_2 \in \mathbb{R}^{N \times T_2 \times D_2}$ , and  $\mathbf{X}_3 \in \mathbb{R}^{N \times L}$  represent the three distinct input graph  
 1051 encoding, where  $T_i$  denotes the sequence length,  $D_i$  the dimensionality of the inputs for each modality  
 1052 and  $L$  the dimension of fingerprints. Each input is processed through its respective component: the  
 1053 first sequence  $\mathbf{X}_1$  is passed through transformer  $\mathcal{T}_1$ , the second sequence  $\mathbf{X}_2$  through transformer  $\mathcal{T}_2$ ,  
 1054 and the third sequence  $\mathbf{X}_3$  through an MLP  $\mathcal{M}$ .

1055 The output of the first transformer  $\mathcal{T}_1$ , denoted  $\mathbf{z}_1 \in \mathbb{R}^{T_1 \times B \times H}$ , is obtained by passing  $\mathbf{X}_1$  through  
 1056 the transformer encoder. Similarly, the output of the second transformer  $\mathcal{T}_2$ , denoted  $\mathbf{z}_2 \in \mathbb{R}^{T_2 \times B \times H}$ ,  
 1057 is obtained by processing  $\mathbf{X}_2$ . Finally, the output of the MLP  $\mathcal{M}$  is denoted  $\mathbf{z}_3 \in \mathbb{R}^{B \times H}$ .

1058 The outputs  $\mathbf{z}_1$ ,  $\mathbf{z}_2$ , and  $\mathbf{z}_3$  are then combined through a learnable weighted sum. Specifically, the  
 1059 combined feature vector  $\mathbf{z}_{\text{combined}}$  is computed as:

1060

$$\mathbf{z}_{\text{combined}} = \alpha \cdot \mathbf{z}_1 + \beta \cdot \mathbf{z}_2 + (1 - \alpha - \beta) \cdot \mathbf{z}_3$$

1061

1062 where  $\alpha$  and  $\beta$  are learnable parameters that control the contribution of each modality to the final  
 1063 representation. This weighted combination allows the model to adaptively learn the most relevant  
 1064 contribution of each input sequence.

1065

1066 The combined feature vector  $\mathbf{z}_{\text{combined}}$  is then passed through a batch normalization layer  
 1067  $\text{BN}(\mathbf{z}_{\text{combined}})$  to normalize the activations, improving training stability. A final fully connected  
 1068 layer  $\text{FC}$  produces the classification output:  $\hat{y} = \text{FC}(\mathbf{z}_{\text{combined}})$  where  $\hat{y} \in \mathbb{R}^C$  represents the  
 1069 predicted class probabilities, with  $C$  being the number of possible output classes.

1070

1071

## C.3 RUNTIME ANALYSIS

1072

1073

1074 To assess the computational efficiency of our method, we report the total runtime across two key  
 1075 stages: (i) topological signature extraction using Topo-Scan (via Degree Centrality and Ollivier-Ricci  
 1076 curvature), and (ii) model training using the transformer-based classifier. Table 9 presents a detailed  
 1077 breakdown of runtimes (in minutes) for five benchmark datasets.

1078

1079

1079 As expected, Degree Centrality is extremely fast to compute and contributes negligible overhead.  
 1080 Ollivier-Ricci curvature, while more computationally intensive, remains tractable even for large  
 1081 graphs, as evidenced by reasonable runtimes on datasets such as REDDIT-5K and OGBG-MOLHIV.  
 1082 Transformer training times scale smoothly with dataset size and remain within practical limits.

1080 Overall, our method maintains scalability  
 1081 while offering strong performance, demon-  
 1082 strating the feasibility of integrating topo-  
 1083 logical signatures into deep graph models  
 1084 at scale.

1085 **TopoScan vs. PH.** We report in Table 10  
 1086 the runtime for Topo-Scan and standard PH  
 1087 on four benchmark datasets with degree  
 1088 centrality filtration (already computed), us-  
 1089 ing the same backend (pyflagser) for both  
 1090 pipelines. For Topo-Scan, we invoke the un-  
 1091 weighted flagser routine, since our method  
 1092 only requires Betti numbers on unweighted clique complexes. For PH, we use the weighted flagser  
 1093 routine, which constructs a full filtration and computes persistence diagrams.

1094 The clustering coefficient column serves as a proxy  
 1095 for graph density and hence clique complexity.  
 1096 On highly clustered graphs such as IMDB-B and  
 1097 IMDB-M (coefficients  $\approx 0.95\text{--}0.97$ ), PH is roughly  
 1098  $13\text{--}14\times$  slower than Topo-Scan, reflecting the com-  
 1099 binatorial explosion of cliques and the cost of global  
 1100 boundary-matrix reductions, whereas on the sparser  
 1101 REDDIT datasets the gap is smaller but still consis-  
 1102 tent (about  $2\times$  on REDDIT-B and  $1.5\times$  on REDDIT-  
 1103 5K). These results empirically confirm that Topo-  
 1104 Scan achieves multi-fold runtime savings over stan-  
 1105 dard PH pipelines on dense graphs while remaining uniformly more efficient across all tested datasets.

1106 **Comparison with other methods.** We  
 1107 also compare the runtimes of two PH-based  
 1108 baselines, PersLay (with degree centrality  
 1109 input) and TopoGCL (using only the topology  
 1110 derived component), against Topo-Scan  
 1111 on IMDB-B and REDDIT-B (Table 11). All  
 1112 times are reported in seconds. For Topo-  
 1113 Scan, we include both the scalar filtration  
 1114 computation and Topo-Scan feature extrac-  
 1115 tion. On IMDB-B, Topo-Scan is about 6 times faster than PersLay and around 29 times faster than  
 1116 the topological part of TopoGCL; on REDDIT-B, it remains faster than PersLay and roughly 14 times  
 1117 faster than TopoGCL. These results further support the practical efficiency of Topo-Scan compared  
 1118 with PH-based pipelines.

#### 1119 C.4 COMPARISON WITH POOLING METHODS

1120 Table 12 compares TOPOFORMER with six representative graph pooling methods designed to  
 1121 adapt GNNs to graph-level tasks. DiffPool (Ying et al., 2018) learns a soft assignment matrix that  
 1122 hierarchically clusters nodes in a differentiable, end-to-end manner. Top-KPooling with Graph  
 1123 U-Nets (Top-K) (Gao & Ji, 2019) ranks nodes using a learnable projection score and retains the top- $k$   
 1124 fraction to coarsen the graph. EigenPool (EigenGCN) (Ma et al., 2019) projects node features onto  
 1125 the leading eigenvectors of the graph Laplacian to preserve global spectral properties. SAGPool (Lee  
 1126 et al., 2019) computes attention scores through a GNN layer, pruning low-importance nodes and  
 1127 re-wiring the remaining graph. MinCutPool (Bianchi et al., 2020) casts pooling as a relaxed spectral  
 1128 clustering problem by optimizing a minimum-cut objective to form node clusters. HaarPool (Wang  
 1129 et al., 2020) applies a Haar wavelet transform to graph signals and performs pooling by selecting  
 1130 key wavelet coefficients. Our model TOPOFORMER takes a different approach by integrating  
 1131 multiscale topological filtrations with a transformer-based attention mechanism, enabling the pooling  
 1132 of substructures across scales and yielding robust higher-order graph representations. As shown in  
 1133 Table 12, TOPOFORMER consistently outperforms all baselines, achieving the best accuracy on six  
 out of seven datasets and ranking second on the remaining one.

Table 9: **Runtimes.** Total runtime (in minutes) per dataset. The second and third columns report the time to compute scalar filtration values and Topo-Scan vectorizations for degree centrality and O.Ricci curvature, respectively, and the final column shows the Transformer training time.

| Dataset     | Degree C. | O. Ricci | Transformer |
|-------------|-----------|----------|-------------|
| IMDB-B      | 0.51      | 5.04     | 3.05        |
| IMDB-M      | 0.49      | 7.62     | 4.57        |
| REDDIT-B    | 5.30      | 23.70    | 6.10        |
| REDDIT-5K   | 36.74     | 109.98   | 15.65       |
| OGBG-MOLHIV | 14.70     | 339.06   | 21.67       |

Table 10: **Runtime for PH vs. Topo-Scan.** Runtime (in seconds) per dataset for computing topo-  
 logical features using the same backend (pyflagser).

| Dataset   | Clus. Coeff. | Topo-Scan | PH     |
|-----------|--------------|-----------|--------|
| IMDB-B    | 0.947        | 9.67      | 135.48 |
| IMDB-M    | 0.969        | 17.81     | 234.30 |
| REDDIT-B  | 0.048        | 12.29     | 24.25  |
| REDDIT-5K | 0.027        | 29.51     | 43.92  |

Table 11: Runtimes (in seconds) for PH-based baselines (PersLay, TopoGCL) and Topo-Scan on IMDB-B and REDDIT-B.

| Method    | IMDB-B | REDDIT-B | Notes                          |
|-----------|--------|----------|--------------------------------|
| PersLay   | 97.02  | 454.78   | PH-based layer on degree input |
| TopoGCL   | 435.49 | 4010.08  | Only topological component     |
| Topo-Scan | 15.16  | 290.00   | Filtration + Topo-Scan         |

1134 **Table 12: Comparison with Pooling Methods.** Accuracy results of six baseline pooling methods  
 1135 and TOPOFORMER on seven graph classification benchmark datasets.

| Model      | BZR                     | COX2                    | MUTAG                   | PROTEINS                | IMDB-B                  | IMDB-M                  | REDDIT-B                |
|------------|-------------------------|-------------------------|-------------------------|-------------------------|-------------------------|-------------------------|-------------------------|
| Top-K      | 79.40 $\pm$ 1.20        | 80.30 $\pm$ 4.21        | 67.61 $\pm$ 3.36        | 69.60 $\pm$ 3.50        | 73.17 $\pm$ 4.84        | 48.80 $\pm$ 3.19        | 79.40 $\pm$ 7.40        |
| MinCutPool | 82.64 $\pm$ 5.05        | 80.07 $\pm$ 3.85        | 79.17 $\pm$ 1.64        | <b>76.52</b> $\pm$ 2.58 | 70.77 $\pm$ 4.89        | 49.00 $\pm$ 2.83        | <b>87.20</b> $\pm$ 5.00 |
| DiffPool   | 83.93 $\pm$ 4.41        | 79.66 $\pm$ 2.64        | 79.22 $\pm$ 1.02        | <b>73.63</b> $\pm$ 3.60 | 68.60 $\pm$ 3.10        | 45.70 $\pm$ 3.40        | 79.00 $\pm$ 1.10        |
| EigenGCN   | 83.05 $\pm$ 6.00        | 80.16 $\pm$ 5.80        | 79.50 $\pm$ 0.66        | 74.10 $\pm$ 3.10        | 70.40 $\pm$ 3.30        | 47.20 $\pm$ 3.00        | N/A                     |
| SAGPool    | 82.95 $\pm$ 4.91        | 79.45 $\pm$ 2.98        | 76.78 $\pm$ 2.12        | 71.86 $\pm$ 0.97        | <b>74.87</b> $\pm$ 4.09 | 49.33 $\pm$ 4.90        | 84.70 $\pm$ 4.40        |
| HaarPool   | 83.95 $\pm$ 5.68        | <b>82.61</b> $\pm$ 2.69 | <b>90.00</b> $\pm$ 3.60 | 73.23 $\pm$ 2.51        | 73.29 $\pm$ 3.40        | <b>49.98</b> $\pm$ 5.70 | N/A                     |
| TOPOFORMER | <b>92.36</b> $\pm$ 4.11 | <b>83.93</b> $\pm$ 4.03 | <b>94.68</b> $\pm$ 4.30 | <b>77.64</b> $\pm$ 3.64 | <b>78.90</b> $\pm$ 3.31 | <b>55.40</b> $\pm$ 4.78 | <b>91.50</b> $\pm$ 1.89 |

### C.5 TOPO-SCAN HYPERPARAMETERS

In the Topo-Scan algorithm, two key hyperparameters play a crucial role: the width parameter, which controls the thickness of slices, and the filtration function, which defines the hierarchical importance of nodes or edges. To determine the optimal hyperparameter settings, we conducted extensive experiments to validate their impact on model performance.

**Width Parameter Selection.** To determine the optimal width parameter  $m$ , we conducted experiments using degree centrality and Ollivier-Ricci curvature as filtration functions for the Topo-Scanner on graph classification datasets. We evaluated  $m = 2, 3$ , and  $4$ , extracting the corresponding Topo-Scanner feature vectors and using them as inputs to a transformer model. The results presented in Table 6 indicate that, for most of the datasets, the Topo-Scanner features achieve the best performance when  $m = 2$  for both filtration functions. Based on this experimental analysis, we select  $m = 2$  as the optimal parameter for our model.

**Multiple Filtrations.** Different filtration functions impose distinct hierarchical orderings on nodes (or edges), enabling our model to capture diverse topological patterns in the induced sequences. This allows the *Topo-Scan* process to effectively integrate domain-specific information. To fully leverage multiple filtrations, TOPOFORMER applies separate transformers for each filtration function and combines their outputs using a learnable attention mechanism. This mechanism dynamically assigns higher weights to the most relevant topological signatures, ensuring optimal feature selection and enhanced performance. As shown in Table 13, TOPOFORMER employing multiple functions consistently outperforms models using a single filtration function, demonstrating the advantages of multiple filtrations. This approach enhances model robustness and stability by incorporating diverse topological perspectives.

1170 **Table 13: Filtration Functions.** Performance comparison of single filtration and multiple filtrations  
 1171 with TOPOFORMER across different datasets. The best values in each column are highlighted in bold.

| Filtrations   | MUTAG                   | PROTEINS                | BZR                     | COX2                    | IMDB-B                  | IMDB-M                  | REDDIT-B                |
|---------------|-------------------------|-------------------------|-------------------------|-------------------------|-------------------------|-------------------------|-------------------------|
| Degree only   | 92.02 $\pm$ 7.24        | 77.28 $\pm$ 5.93        | 89.89 $\pm$ 3.74        | 78.36 $\pm$ 4.93        | 74.20 $\pm$ 3.36        | 51.53 $\pm$ 3.34        | 86.60 $\pm$ 2.97        |
| O. Ricci only | 89.91 $\pm$ 3.86        | 77.26 $\pm$ 4.29        | 90.60 $\pm$ 3.69        | 78.60 $\pm$ 4.79        | <b>79.10</b> $\pm$ 3.78 | <b>54.53</b> $\pm$ 3.52 | <b>91.40</b> $\pm$ 1.24 |
| HKS Only      | <b>95.32</b> $\pm$ 5.58 | <b>77.35</b> $\pm$ 2.86 | <b>90.62</b> $\pm$ 4.91 | <b>83.95</b> $\pm$ 2.99 | 76.90 $\pm$ 5.72        | 54.07 $\pm$ 2.54        | 90.05 $\pm$ 2.41        |
| Deg.+O.Ricci  | 93.01 $\pm$ 5.29        | <b>78.35</b> $\pm$ 4.22 | <b>91.12</b> $\pm$ 4.68 | 81.80 $\pm$ 5.40        | 78.80 $\pm$ 3.65        | 53.87 $\pm$ 3.52        | 90.65 $\pm$ 2.12        |
| HKS+O.Ricci   | 94.68 $\pm$ 4.30        | 77.64 $\pm$ 3.64        | 92.36 $\pm$ 4.11        | <b>83.93</b> $\pm$ 4.03 | <b>78.90</b> $\pm$ 3.31 | <b>55.40</b> $\pm$ 4.78 | <b>91.50</b> $\pm$ 1.89 |
| HKS+Degree    | <b>95.26</b> $\pm$ 3.88 | 78.08 $\pm$ 2.34        | 91.09 $\pm$ 5.53        | 83.71 $\pm$ 4.38        | 77.30 $\pm$ 2.41        | 52.60 $\pm$ 2.25        | 89.90 $\pm$ 2.35        |

### C.6 TOPOFORMER VS. PH WITH DIFFERENT VECTORIZATIONS

TOPOFORMER consistently outperforms Persistent Homology methods in both accuracy and computational efficiency. As shown in Table 14, we compare against the best PH results reported in (Cai & Wang, 2020), which evaluates 16 combinations of four filtration functions (degree, O.Ricci, Fiedler, closeness centrality) and four vectorization techniques (Sliced Wasserstein, Pervec, Filvec, SW-p) per dataset. TOPOFORMER achieves higher accuracy on all six benchmarks.

1188 Table 14: Accuracy results for TOPOFORMER (HKS) vs. Persistent Homology in graph classification  
 1189 tasks. In PH row, we report the best performance of 16 combinations with four filtration functions  
 1190 combined with four vectorizations.

|                      |  | BZR                            | COX2                           | PROTEINS                       | IMDB-B                         | IMDB-M                         |  |
|----------------------|--|--------------------------------|--------------------------------|--------------------------------|--------------------------------|--------------------------------|--|
| PH (Best of 16 comb) |  | 88.4 $\pm$ 0.6                 | <b>82.0<math>\pm</math>0.6</b> | 74.0 $\pm$ 0.4                 | 69.5 $\pm$ 0.5                 | 46.5 $\pm$ 0.3                 |  |
| TOPOFORMER           |  | <b>90.6<math>\pm</math>4.9</b> | <b>82.0<math>\pm</math>4.6</b> | <b>77.4<math>\pm</math>2.9</b> | <b>77.9<math>\pm</math>3.4</b> | <b>54.1<math>\pm</math>2.5</b> |  |

1195 Table 15: **TOPOFORMER vs. PH Performance Comparison:** Accuracy of three topological models  
 1196 under seven filtrations: Degree, Ollivier-Ricci, HKS, Betweenness centrality, Closeness centrality, Eigenvector  
 1197 centrality, and Forman-Ricci curvature. The last column reports the average accuracy improvements of our  
 1198 models PH-TR and TOPOFORMER over the classical TDA pipeline PH-MLP for the same filtration function.

| Filtration  | Model      | BZR                              | COX2                             | MUTAG                            | PROTEINS                         | IMDB-B                           | IMDB-M                           | REDDIT-B                         | Av.Imp. |
|-------------|------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|---------|
| Degree      | PH-MLP     | 82.71 $\pm$ 6.51                 | 76.44 $\pm$ 5.39                 | 84.06 $\pm$ 4.65                 | 68.37 $\pm$ 3.97                 | 65.70 $\pm$ 4.03                 | 45.07 $\pm$ 2.59                 | 89.50 $\pm$ 2.87                 | —       |
|             | PH-TR      | 86.43 $\pm$ 4.33                 | 78.15 $\pm$ 5.19                 | 86.11 $\pm$ 5.23                 | <b>77.54<math>\pm</math>2.64</b> | <b>75.00<math>\pm</math>2.11</b> | <b>50.67<math>\pm</math>3.57</b> | <b>92.30<math>\pm</math>1.77</b> | 4.91    |
|             | TOPOFORMER | <b>91.10<math>\pm</math>5.14</b> | <b>80.27<math>\pm</math>5.24</b> | <b>92.54<math>\pm</math>5.12</b> | 77.45 $\pm$ 4.02                 | 74.20 $\pm$ 5.01                 | 50.33 $\pm$ 1.52                 | 89.75 $\pm$ 2.18                 | 6.26    |
| O.Ricci     | PH-MLP     | 85.45 $\pm$ 3.36                 | 78.16 $\pm$ 5.09                 | 84.06 $\pm$ 5.21                 | 65.50 $\pm$ 4.26                 | 68.00 $\pm$ 3.55                 | 44.87 $\pm$ 3.65                 | 85.65 $\pm$ 2.62                 | —       |
|             | PH-TR      | 88.62 $\pm$ 5.40                 | 78.16 $\pm$ 5.73                 | 87.61 $\pm$ 5.70                 | 77.27 $\pm$ 5.08                 | 72.20 $\pm$ 6.24                 | 48.00 $\pm$ 4.33                 | 90.65 $\pm$ 1.08                 | 4.40    |
|             | TOPOFORMER | <b>90.38<math>\pm</math>5.50</b> | <b>80.72<math>\pm</math>6.44</b> | <b>92.54<math>\pm</math>4.47</b> | <b>77.90<math>\pm</math>3.17</b> | <b>74.70<math>\pm</math>4.95</b> | <b>51.53<math>\pm</math>3.49</b> | <b>91.90<math>\pm</math>2.73</b> | 6.85    |
| HKS         | PH-MLP     | 84.96 $\pm$ 4.42                 | 78.19 $\pm$ 4.34                 | 84.09 $\pm$ 5.72                 | 70.80 $\pm$ 4.70                 | 71.10 $\pm$ 5.28                 | 47.93 $\pm$ 3.20                 | 88.10 $\pm$ 1.67                 | —       |
|             | PH-TR      | 89.60 $\pm$ 5.84                 | 79.89 $\pm$ 4.66                 | 94.12 $\pm$ 5.42                 | 77.18 $\pm$ 3.15                 | 76.80 $\pm$ 3.97                 | 53.60 $\pm$ 3.31                 | 87.25 $\pm$ 1.95                 | 4.75    |
|             | TOPOFORMER | <b>90.62<math>\pm</math>4.91</b> | <b>83.95<math>\pm</math>2.99</b> | <b>95.32<math>\pm</math>5.58</b> | <b>77.35<math>\pm</math>2.86</b> | <b>77.90<math>\pm</math>5.72</b> | <b>54.07<math>\pm</math>2.54</b> | <b>90.05<math>\pm</math>2.41</b> | 6.30    |
| Betweenness | PH-MLP     | 84.95 $\pm$ 4.19                 | 80.99 $\pm$ 6.35                 | 89.94 $\pm$ 7.93                 | 71.61 $\pm$ 1.85                 | 68.10 $\pm$ 2.55                 | 43.80 $\pm$ 1.74                 | 79.10 $\pm$ 2.88                 | —       |
|             | PH-TR      | 85.43 $\pm$ 4.13                 | <b>81.60<math>\pm</math>6.00</b> | 90.52 $\pm$ 5.62                 | 74.13 $\pm$ 2.95                 | 69.90 $\pm$ 3.48                 | 45.40 $\pm$ 1.79                 | 84.05 $\pm$ 2.33                 | 1.79    |
|             | TOPOFORMER | <b>87.41<math>\pm</math>4.07</b> | 80.74 $\pm$ 6.15                 | <b>90.99<math>\pm</math>6.61</b> | <b>76.73<math>\pm</math>2.67</b> | <b>73.90<math>\pm</math>3.73</b> | <b>51.47<math>\pm</math>2.96</b> | <b>86.55<math>\pm</math>2.30</b> | 4.19    |
| Closeness   | PH-MLP     | 84.21 $\pm$ 2.19                 | 79.07 $\pm$ 6.56                 | 88.94 $\pm$ 7.20                 | 74.39 $\pm$ 3.19                 | 65.30 $\pm$ 4.61                 | 47.47 $\pm$ 3.93                 | 66.85 $\pm$ 2.98                 | —       |
|             | PH-TR      | <b>87.43<math>\pm</math>4.83</b> | 79.65 $\pm$ 4.98                 | 89.94 $\pm$ 6.25                 | 75.93 $\pm$ 3.32                 | 69.70 $\pm$ 4.60                 | 50.47 $\pm$ 3.72                 | 77.20 $\pm$ 3.31                 | 3.44    |
|             | TOPOFORMER | 85.13 $\pm$ 8.36                 | <b>81.17<math>\pm</math>5.28</b> | <b>90.88<math>\pm</math>6.65</b> | <b>77.64<math>\pm</math>4.36</b> | <b>73.20<math>\pm</math>2.20</b> | <b>51.07<math>\pm</math>3.02</b> | <b>86.40<math>\pm</math>2.22</b> | 5.61    |
| Eigenvector | PH-MLP     | 83.97 $\pm$ 3.46                 | 80.56 $\pm$ 6.04                 | 89.39 $\pm$ 6.64                 | 67.20 $\pm$ 5.87                 | 66.70 $\pm$ 2.61                 | 47.40 $\pm$ 3.00                 | 79.40 $\pm$ 3.21                 | —       |
|             | PH-TR      | 87.41 $\pm$ 4.21                 | 79.88 $\pm$ 6.04                 | <b>91.57<math>\pm</math>5.70</b> | 70.53 $\pm$ 4.60                 | 72.40 $\pm$ 4.48                 | 50.13 $\pm$ 3.44                 | 89.25 $\pm$ 1.40                 | 3.79    |
|             | TOPOFORMER | <b>90.59<math>\pm</math>5.63</b> | <b>82.87<math>\pm</math>3.35</b> | 90.99 $\pm$ 6.61                 | <b>77.35<math>\pm</math>2.78</b> | <b>76.10<math>\pm</math>3.63</b> | <b>51.00<math>\pm</math>2.14</b> | <b>91.85<math>\pm</math>1.43</b> | 6.59    |
| F. Ricci    | PH-MLP     | 82.46 $\pm$ 3.94                 | 80.13 $\pm$ 5.86                 | 87.81 $\pm$ 6.21                 | 73.95 $\pm$ 4.12                 | 66.60 $\pm$ 4.12                 | 45.53 $\pm$ 3.36                 | 73.70 $\pm$ 3.78                 | —       |
|             | PH-TR      | 86.18 $\pm$ 6.06                 | <b>81.97<math>\pm</math>6.62</b> | 91.99 $\pm$ 4.63                 | 76.09 $\pm$ 4.04                 | 70.80 $\pm$ 4.47                 | 50.67 $\pm$ 2.59                 | 77.20 $\pm$ 2.21                 | 3.53    |
|             | TOPOFORMER | <b>88.41<math>\pm</math>6.04</b> | 81.02 $\pm$ 6.46                 | <b>92.08<math>\pm</math>5.67</b> | <b>77.81<math>\pm</math>3.80</b> | <b>79.40<math>\pm</math>3.69</b> | <b>54.47<math>\pm</math>3.34</b> | <b>88.95<math>\pm</math>1.94</b> | 7.42    |

### C.7 TOPOFORMER VS PH PERFORMANCE

Table 15 extends our ablation (Table 5) from three to seven filtration functions and compares three topological pipelines under the same filtration function: the classical PH-MLP baseline (sublevel PH + Betti vector + MLP), our PH-TR variant (same Betti vectors but processed as sequences by a Transformer), and TOPOFORMER (Topo-Scan sequences with sliding-window interlevel filtrations). Across all seven filtrations, replacing the MLP with a Transformer already yields consistent gains: PH-TR improves over PH-MLP by roughly 2–4 accuracy points on average (see the “Av.Imp.” column), confirming that treating Betti curves as ordered sequences is beneficial even without changing the underlying filtration.

TOPOFORMER further improves on PH-TR for almost every filtration, typically adding another 1–3 points on most datasets and yielding average gains of 6–7 points over PH-MLP. The effect is especially pronounced on more challenging benchmarks such as IMDB-M and REDDIT-B, where sliding-window interlevel slices capture richer late-emerging structure than standard sublevel PH. Importantly, this pattern holds not only for the three filtrations used in the main text (degree, Ollivier–Ricci, HKS) but also for the four additional ones (betweenness, closeness, eigenvector centrality, Forman–Ricci curvature), indicating that the benefit of Topo-Scan is robust to the choice of scalar function. Together, these results support our central claim: the main performance gains come from the Topo-Scan sequential representation (and its integration with Transformers), rather than from a particular hand-picked filtration function.

### C.8 EARLY SATURATION IN PH FILTRATIONS AND TOPO-SCAN

**Goal.** We compare classical PH (sub/superlevel on a fixed clique-complex 2-skeleton) with *Topo-Scan* to show how PH frequently *early-saturates* on graphs, i.e., after a relatively small portion of the

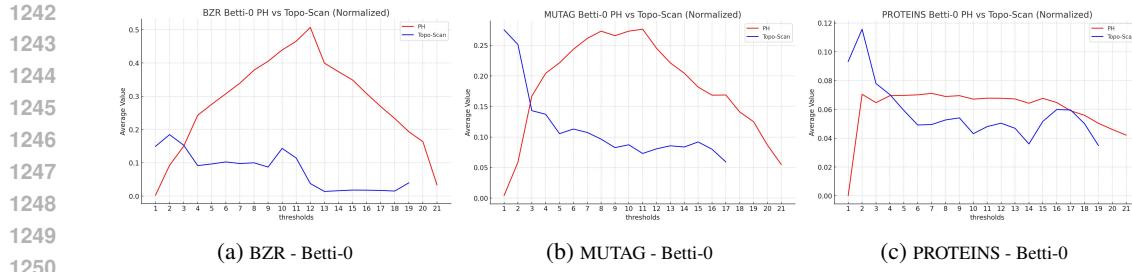


Figure 6: **PH vs. Topo-Scan.** Normalized average Betti-0 values over 20 thresholds of the degree-centrality filtration on (a) BZR, (b) MUTAG, and (c) PROTEINS. Under the classical PH pipeline, feature counts decline monotonically with increasing threshold, whereas Topo-Scan maintains elevated values at higher thresholds, revealing late-emerging topological features that PH alone misses.

threshold range, new features cease to appear, whereas Topo-Scan continues to surface structure by sliding windows over the same signal.

**Protocol.** For each dataset and filtration function  $f$  (e.g., degree or Ollivier–Ricci), we fix a common grid of  $T$  thresholds and evaluate both methods on the same clique-complex 2-skeleton (upper–star from nodes). PH: sublevel filtration evaluated at the same grid points; Betti counts are read at each threshold. Topo-Scan: window width  $m$  and stride  $s$  define  $T$  overlapping slices whose vertex sets correspond to consecutive value ranges in the same grid. Betti counts are computed per slice. To make cross-dataset plots visually comparable, we report (i) *normalized* Betti-0 curves when scales differ markedly (Fig. 6) and (ii) *unnormalized* Betti-0 when PH and Topo-Scan share similar ranges (Fig. 3). Betti-1 frequency barplots are shown to illustrate higher-order behavior (Fig. 7).

**How to read the figures.** A positive Betti-0 value at a position means additional connected components are present in that slice/threshold; persistent nonzero values toward the *right side* of the horizontal axis indicate *late-emerging* structure. For Betti-1, darker bars at higher thresholds indicate more cycles appearing later in the filtration. Because Topo-Scan slices are value–localized *ranges* rather than one-sided sublevels, they retain visibility into regions that are otherwise drowned out once early high- or low-valued nodes saturate the PH complex.

**Results on small biochemical graphs (BZR, MUTAG, PROTEINS).** Figure 6 plots normalized Betti-0 curves over 20 degree thresholds. Across all three datasets, PH curves drop quickly and remain low: after an early rise, new components rarely appear as the complex fills up. In contrast, Topo-Scan maintains elevated values deeper into the axis, indicating that as the sliding window moves, it continues to expose distinct local subgraphs in later value ranges. This pattern is precisely the late-structure retention we aim to capture.

**Results on social graphs (IMDB-B, IMDB-M).** Figure 3 shows unnormalized Betti-0 with 100 thresholds (comparable scales). Here, PH exhibits a sharp taper near the end: once the core of the

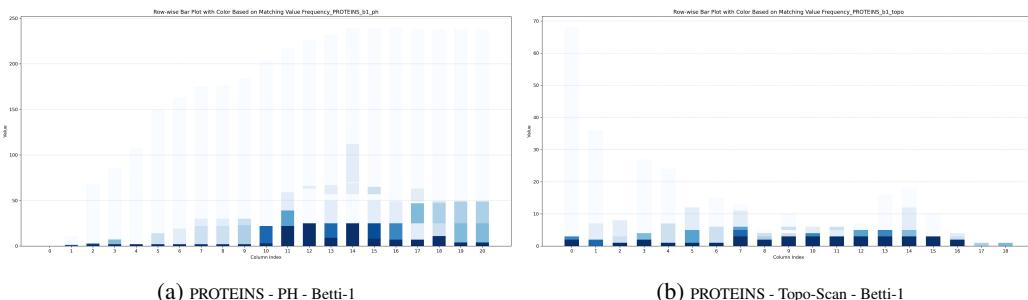


Figure 7: **PH vs. Topo-Scan.** Bar plots of Betti-1 counts at each O.Ricci filtration threshold on the PROTEINS dataset, with bar color intensity encoding the frequency of each integer value at that threshold. In (a) classical PH features rapidly taper off and plateau early, whereas in (b) Topo-Scan shows increasingly darker bars at higher thresholds, evidence of continued cycle emergence beyond PH’s saturation.

graph enters the complex, subsequent thresholds add little. Topo-Scan avoids this collapse; activity persists and often *exceeds* PH in the tail, reflecting components that are still exposed by the windowed slices even when global sublevels have already merged them away.

**Higher-order signal (Betti-1 on PROTEINS, O. Ricci).** Figure 7 provides barplots where color intensity encodes the frequency of each integer Betti-1 value per threshold. Under PH (left), bars fade and plateau early, showing few cycles after the initial growth phase. Under Topo-Scan (right), darker bars persist across later thresholds, demonstrating continued cycle emergence that PH no longer reveals once the complex has saturated.

**Why does this happen?** In sub/superlevel PH, once extreme-valued vertices enter early, the induced complex quickly fills in, so later additions create little new topology, especially on graphs where dense regions are correlated with the signal. Topo-Scan, by scanning *ranges* of values with overlap, repeatedly re-centers attention on late parts of the signal, preventing early regions from dominating the entire sequence. Importantly, this is not a claim that sublevel is intrinsically flawed; task-aligned or learned filtrations can mitigate early saturation. Our point is empirical and architectural: a fixed-budget sliding-window view preserves late signal *by design*.

**Controls and caveats.** (i) We use the same signal, grid, and complex for both methods to avoid confounding factors. (ii) Normalization is applied only for visualization when scales differ; conclusions do not depend on normalization. (iii) Sublevel and superlevel yield the same multiset of slices in reverse order; Topo-Scan’s behavior is insensitive to that choice. (iv) Window hyperparameters  $(m, s)$  trade locality for coverage; we keep them fixed across datasets in these plots for clarity.

**Takeaway.** Across biochemical and social benchmarks, PH curves commonly *flatten early*, while Topo-Scan remains *active in the tail* (Betti-0 and Betti-1), revealing late-emerging components and cycles. This supports our central design choice: turning topology into short, ordered, range-localized tokens helps retain information that standard PH pipelines often lose once the complex saturates.