

# Deep Graph Mapper: Seeing Graphs through the Neural Lens

Cristian Bodnar<sup>†,\*</sup>, Cătălina Cangea<sup>†,\*</sup> and Pietro Liò

*Department of Computer Science and Technology, University of Cambridge,  
Cambridge, United Kingdom*

Correspondence\*:

Cristian Bodnar, Cătălina Cangea  
{cb2015,ccc53}@cam.ac.uk

## 2 ABSTRACT

3 Graph summarisation has received much attention lately, with various works tackling the  
4 challenge of defining pooling operators on data regions with arbitrary structures. These contrast  
5 the grid-like ones encountered in image inputs, where techniques such as max-pooling have  
6 been enough to show empirical success. In this work, we merge the Mapper algorithm with the  
7 expressive power of graph neural networks to produce topologically-grounded graph summaries.  
8 We demonstrate the suitability of Mapper as a topological framework for graph pooling by proving  
9 that Mapper is a generalisation of pooling methods based on soft cluster assignments. Building  
10 upon this, we show how easy it is to design novel pooling algorithms that obtain competitive  
11 results with other state-of-the-art methods. Additionally, we use our method to produce GNN-aided  
12 visualisations of attributed complex networks.

13 **Keywords:** Mapper, graph neural networks, pooling, graph summarisation, graph classification

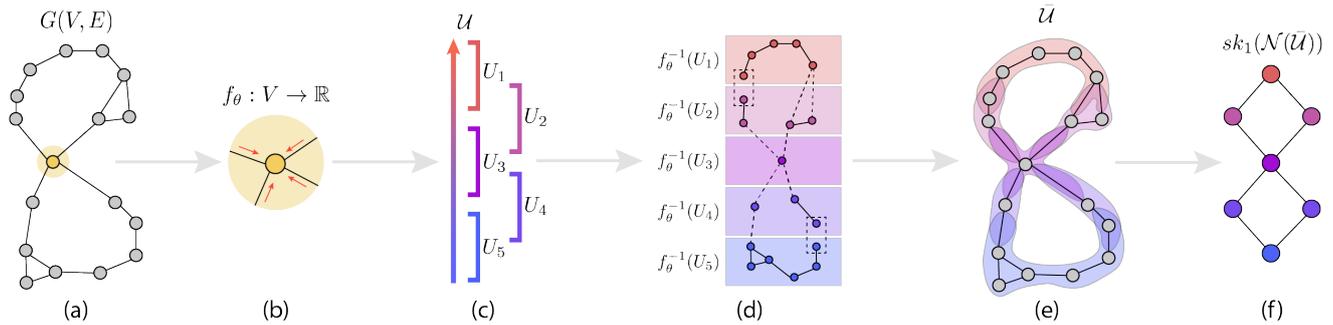
## 1 INTRODUCTION

14 The abundance of relational information in the real world and the success of deep learning techniques  
15 have brought renowned interest in learning from graph-structured data. Efforts in this direction have been  
16 primarily focused on replicating the hierarchy of convolutional filters and pooling operators, which have  
17 achieved previous success in computer vision Sperduti (1994); Goller and Kuchler (1996); Gori et al.  
18 (2005); Scarselli et al. (2009); Bruna et al. (2014); Li et al. (2015), within relational data domains. In  
19 contrast to image processing applications, where the signal is defined on a grid-like structure, designing  
20 graph coarsening (pooling) operators is a much more difficult problem, due to the arbitrary structure  
21 typically present in graphs.

22 In this work, we introduce Structural Deep Graph Mapper (SDGM)<sup>1</sup>—an adaptation of Mapper (Singh  
23 et al., 2007), an algorithm from the field of Topological Data Analysis (TDA) (Chazal and Michel, 2017),  
24 to graph domains. First, we prove that SDGM is a generalisation of pooling methods based on soft  
25 cluster assignments, which include state-of-the-art algorithms like minCUT (Bianchi et al., 2019) and

---

<sup>1</sup> Code to reproduce models and experimental results is available at <https://github.com/crisbodnar/dgm>.



**Figure 1.** A cartoon illustration of Structural Deep Graph Mapper (SDGM) where, for simplicity, a graph neural network (GNN) approximates a ‘height’ function over the nodes in the plane of the diagram. The input graph (a) is passed through the GNN, which maps the vertices of the graph to a real number (the height) (b–c). Given a cover  $\mathcal{U}$  of the image of the GNN (c), the edge-refined pull back cover  $\bar{\mathcal{U}}$  is computed (d–e). The dotted edges in (d) illustrate connections between the node clusters (structural connections), while the dotted boxes show nodes that appear in multiple clusters (semantic connections). The 1-skeleton of the nerve of the edge-refined pull back cover provides the pooled graph (f).

26 DiffPool (Ying et al., 2018). Building upon this topological perspective of graph pooling, we propose two  
 27 pooling algorithms leveraging fully-differentiable and fixed PageRank-based ‘lens’ functions, respectively.  
 28 We demonstrate that these operators achieve results competitive with other state-of-the-art pooling methods  
 29 on graph classification benchmarks. Furthermore, we show how our method offers a means to flexibly  
 30 visualise graphs and the complex data living on them through a GNN ‘lens’ function.

## 2 RELATED WORK

31 In this section, we investigate the existing work in the two broad areas that our method is part of—graph  
 32 pooling (also deemed hierarchical representation learning) and network visualisations.

33 **Graph pooling** algorithms have already been considerably explored within GNN frameworks for graph  
 34 classification. Luzhnica et al. (2019) propose a topological approach to pooling, which coarsens the graph  
 35 by aggregating its maximal cliques into new clusters. However, cliques are local topological features,  
 36 whereas our methods leverage a global perspective of the graph during pooling. Two paradigms distinguish  
 37 themselves among learnable pooling layers: Top- $k$  pooling based on a learnable ranking (Gao and Ji,  
 38 2019), and learning the cluster assignment (Ying et al., 2018) with additional entropy and link prediction  
 39 losses for more stable training (DiffPool). Following these two trends, several variants and incremental  
 40 improvements have been proposed. The Top- $k$  approach is explored in conjunction with jumping-knowledge  
 41 networks (Cangea et al., 2018), attention (Lee et al., 2019; Huang et al., 2019) and self-attention for cluster  
 42 assignment (Ranjan et al., 2019). Similarly to DiffPool, the method suggested by Bianchi et al. (2019)  
 43 uses several loss terms to enforce clusters with strongly connected nodes, similar sizes and orthogonal  
 44 assignments. A different approach is also proposed by Ma et al. (2019), who leverage spectral clustering.

45 **Graph visualisation** is a vast topic in network science. We therefore refer the reader to existing surveys,  
 46 for a complete view of the field (Nobre et al., 2019; von Landesberger et al., 2011; Beck et al., 2017), and  
 47 focus here only on methods that, similarly to ours, produce node-link-based visual summaries through  
 48 the aggregation of static graphs. Previous methods rely on grouping nodes into a set of predefined  
 49 motifs (Dunne and Shneiderman, 2013), modules (Dwyer et al., 2013) or clusters with basic topological  
 50 properties (Batagelj et al., 2010). Recent approaches have considered attribute-driven aggregation schemes  
 51 for multivariate networks. For instance, PivotGraph (Wattenberg, 2006) groups the nodes based on

52 categorical attributes, while van den Elzen and van Wijk (2014) propose a more sophisticated method  
 53 using a combination of manually-specified groupings and attribute queries. However, these mechanisms  
 54 are severely constrained by the simple types of node groupings allowed and the limited integration between  
 55 graph topology and attributes. Closest to our work, Mapper-based summaries for graphs have recently  
 56 been considered by Hajij et al. (2018). Despite the advantages provided by Mapper, their approach relies  
 57 on hand-crafted graph-theoretic ‘lenses’, such as the average geodesic distance, graph density functions  
 58 or eigenvectors of the graph Laplacian. Not only are these functions unable to fully adapt to the graph  
 59 of interest, but they are also computationally inefficient and do not take into account the attributes of the  
 60 graph.

### 3 BACKGROUND AND FORMAL PROBLEM STATEMENT

61 **Formal problem statement.** Consider a dataset whose samples are formed by a graph  $G_i = (V_i, E_i)$ , a  
 62  $d$ -dimensional signal defined over the nodes of the graph  $h_i : V \rightarrow \mathbb{R}^d$  and a label  $y_i$  associated with the  
 63 graph, where  $i \in I$ , a finite indexing set for the dataset samples. We are interested in the setting where graph  
 64 neural networks are used to classify such graphs using a sequence of (graph) convolutions and pooling  
 65 operators. While convolutional operators act like filters of the graph signal, pooling operators coarsen the  
 66 graph and reduce its spatial resolution. Unlike image processing tasks, where the inputs exhibit a regular  
 67 grid structure, graph domains pose challenges for pooling. In this work, we design topologically-inspired  
 68 pooling operators based on Mapper. As an additional contribution, we also investigate graph pooling as a  
 69 tool for the visualisation of attributed graphs.

70 We briefly review the Mapper (Singh et al., 2007) algorithm, with a focus on graph domains (Hajij et al.,  
 71 2018). We first introduce the required mathematical background.

72 **Definition 3.1.** Let  $X, Z$  be two topological spaces,  $f : X \rightarrow Z$ , a continuous function, and  $\mathcal{U} = (U_i)_{i \in I}$   
 73 a cover of  $Z$ . Then, the *pull back cover*  $f^{-1}(\mathcal{U})$  of  $X$  induced by  $(f, \mathcal{U})$  is the collection of open sets  
 74  $f^{-1}(U_i), i \in I$ , for some indexing set  $I$ . For each  $f^{-1}(U_i)$ , let  $\{C_{i,j}\}_{j \in J_i}$  be a partition of  $f^{-1}(U_i)$  indexed  
 75 by  $J_i$ . We refer to the elements of these partitions as *clusters*. The resulting collection of clusters forms  
 76 another cover of  $X$  called the *refined pull back cover*  $\mathcal{R}(f^{-1}(\mathcal{U})) = \{C_{i,j}\}_{i \in I, j \in J_i}$ .

77 **Definition 3.2.** Let  $X$  be a topological space with an open cover  $\mathcal{U} = (U_i)_{i \in I}$ . The *1-skeleton of the*  
 78 *nerve*  $\mathcal{N}(\mathcal{U})$  of  $\mathcal{U}$ , which we denote by  $sk_1(\mathcal{N}(\mathcal{U}))$ , is the graph with vertices given by  $(v_i)_{i \in I}$ , where two  
 79 vertices  $v_i, v_j$  are connected if and only if  $U_i \cap U_j \neq \emptyset$ .

80 **Mapper.** Given a topological space  $X$ , a carefully chosen *lens function*  $f : X \rightarrow Z$  and a cover  $\mathcal{U}$  of  
 81  $Z$ , Mapper produces a graph representation of the topological space by computing the 1-skeleton of the  
 82 nerve of the refined pull back cover  $sk_1(\mathcal{N}(\mathcal{R}(f^{-1}(\mathcal{U}))))$ , which we denote by  $\mathcal{M}(f, \mathcal{U})$ . We note that,  
 83 more generally, the skeleton operator might be omitted, in which case the output of the algorithm becomes  
 84 a simplicial complex. However, for the purpose of this work, we are only interested in graph outputs.  
 85 Typically, the input to the Mapper algorithm is a point cloud and the connected components are inferred  
 86 using a statistical clustering algorithm, with the help of a metric defined in the space where the points live.

87 **Mapper for Graphs.** More recently, Hajij et al. (2018) considered the case when the input topological  
 88 space  $X = G(V, E)$  is also a graph with vertices  $V$  and edge set  $E$ . In a typical point cloud setting, the  
 89 relationships between points are statistically inferred; in a graph setting, the underlying relationships are  
 90 given by the edges of the graph. The adaptation of Mapper for graphs proposed by Hajij et al. (2018) uses a  
 91 lens function  $f : V \rightarrow \mathbb{R}$  based on graph-theoretic functions and a cover  $\mathcal{U}$  formed of open intervals of the

92 real line. Additionally, the connected components  $\{C_{i,j}\}_{j \in J_i}$  are given by the vertices of the connected  
 93 components of the subgraph induced by  $f^{-1}(U_i)$ .

94 However, the graph version of Mapper described above has two main limitations. Firstly, the graph-  
 95 theoretic functions considered for  $f$  are rather limited, not taking into account the signals which are  
 96 typically defined on the graph in signal processing tasks, such as graph classification. Secondly, by using a  
 97 pull back cover only over the graph vertices, as opposed to a cover of the entire graph, the method relies  
 98 exclusively on the lens function to capture the structure of the graph and the edge-connections between the  
 99 clusters. This may end up discarding valuable structural information, as we later show in Section 7.7.

## 4 STRUCTURAL DEEP GRAPH MAPPER

100 **Structural Graph Mapper.** One of the disadvantages of the graph version of Mapper (described in the  
 101 background section) is that its output does not explicitly capture the connections between the resulting  
 102 collections of clusters. This is primarily because the lens function  $f$  is defined only over the set of vertices  $V$   
 103 and, consequently, the resulting pull-back cover only covers  $V$ . In contrast, one should aim to obtain a cover  
 104 for the graph  $G$ , which automatically includes the edges. While this could be resolved by considering a lens  
 105 function over the geometric realisation of the graph, handling only a finite set of vertices is computationally  
 106 convenient.

107 To balance these trade-offs, we add an extra step to the Mapper algorithm. Concretely, we extend  
 108 the refined pull back cover into a cover over both nodes and edges. Given the set of refined clusters  
 109  $\{C_{i,j}\}_{i \in I, j \in J_i}$ , we compute a new set of clusters  $\{C'_{i,j}\}_{i \in I, j \in J_i}$  where each cluster  $C'_{i,j}$  contains the  
 110 elements of  $C_{i,j}$  as well as all the edges incident to the vertices in  $C_{i,j}$ . We use  $\mathcal{R}_E$  (the edge-refined pull  
 111 back cover) to refer to this open cover of the graph  $G$  computed from  $f^{-1}(\mathcal{U})$ . Then, our algorithm can be  
 112 written as  $sk_1(\mathcal{N}(\mathcal{R}_E(f^{-1}(\mathcal{U}))))$  and we denote it by  $\mathcal{GM}(f, \mathcal{U})$ .

113 **Remark 1.** We note that Structural Mapper, unlike the original Mapper method, encodes two types of  
 114 relationships via the edges of the output graph. The *semantic connections* highlight a similarity between  
 115 clusters, according to the lens function (that is, two clusters have common nodes—as before), while  
 116 *structural connections* show how two clusters are connected (namely, two clusters have at least one edge in  
 117 common). This latter type of connection is the result of considering the extended cover over the edges. The  
 118 two types of connections are not mutually exclusive because two clusters might have both nodes and edges  
 119 in common.

120 We now broadly outline our proposed method, using the three main degrees of freedom of the Mapper  
 121 algorithm to guide our discussion: the lens function, the cover, and the clustering algorithm.

122 **The lens** is a function  $f : V \rightarrow \mathbb{R}^d$  over the vertices, which acts as a filter that emphasises certain features  
 123 of the graph. Typically,  $d$  is a small integer—in our case,  $d \in \{1, 2\}$ . The choice of  $f$  depends on the graph  
 124 properties that should be highlighted by the visualisation. In this work, we leverage the recent progress in  
 125 the field of graph representation learning and propose a parameterised lens function based on graph neural  
 126 networks (GNNs). We thus consider a function  $f_\theta(v) = g_\theta(V, E, \mathbf{X})_v$ , where  $g$  is a GNN with parameters  $\theta$   
 127 taking as input a graph  $G = (V, E)$  with  $n$  nodes and node features  $\mathbf{X} \in \mathbb{R}^{n \times k}$ . For visualisation purposes,  
 128 we often consider a function composition  $f_\theta(v) = (r \circ g_\theta)_v$ , where  $r : \mathbb{R}^{n \times d'} \rightarrow \mathbb{R}^{n \times d}$  is a dimensionality  
 129 reduction algorithm like  $t$ -SNE (van der Maaten and Hinton, 2008).

130 Unlike the traditional graph theoretic lens functions proposed by Hajij et al. (2018), GNNs can naturally  
 131 learn to integrate the features associated with the graph and its topology, while also scaling computationally

132 to large, complex graphs. Additionally, visualisations can be flexibly tuned for the task of interest, by  
133 adjusting the lens  $g_\theta$  through the loss function of the model.

134 **The cover  $\mathcal{U}$**  determines the resolution of the output graph. For most purposes, we leverage the usual  
135 cover choice for Mapper,  $\mathbb{R}^d$ . When  $d = 1$ , we use a set of equally-sized overlapping intervals over the real  
136 line. When  $d = 2$ , this is generalised to a grid of overlapping cells in the real plane. Using more cells will  
137 produce more detailed visualisations, while higher overlaps between the cells will increase the connectivity  
138 of the output graph. When chosen suitably, these hyperparameters are a powerful mechanism for obtaining  
139 multi-scale visualisations.

140 Another choice that we employ for designing differentiable pooling algorithms is a set of RBF kernels,  
141 where the second arguments of kernel functions are distributed over the real line. We introduce this in  
142 detail in Section 5.2.

143 **Clustering** statistically approximates the (topological) connected components of the cover sets  $U_i$ .  
144 Mapper does not require a particular type of clustering algorithm; however, when the input topological  
145 space  $X$  is a graph, a natural choice, also adopted by Hajij et al. (2018), is to take the connected components  
146 of the subgraphs induced by the vertices  $f^{-1}(U_i), i \in I$ . Therefore, in principle, there is no need to resort  
147 to statistical clustering techniques.

148 However, relying on the topological connected components introduces certain challenges when the aim  
149 is to obtain a coarsened graph. Many real-world graphs comprise thousands of connected components,  
150 which is a lower bound to the number of connected components of the graph produced by  $\mathcal{GM}$ . In the most  
151 extreme case, a graph containing only isolated nodes (namely, a point cloud) would never be coarsened by  
152 this procedure. Therefore, it is preferable to employ statistical techniques where the number of clusters  
153 can be specified. In our pooling experiments, we draw motivation from the relationship with other pooling  
154 algorithms and opt to assign all the nodes to the same cluster (which corresponds to no clustering).

155 We broadly refer to this instance of Structural Graph Mapper, with the choices described above, as  
156 Structural Deep Graph Mapper (SDGM). We summarise it step-by-step in the cartoon example in Figure 1  
157 and encourage the reader to refer to it.

## 5 STRUCTURAL GRAPH MAPPER FOR POOLING

158 We begin this section by introducing several theoretical results, which provide a connection between our  
159 version of Mapper and other graph pooling algorithms. We then use these results to show how novel pooling  
160 algorithms can be designed.

### 161 5.1 Relationship to graph pooling methods

162 An early suggestion that Mapper could be suitable for graph pooling is given by the fact that it constitutes  
163 a generalisation of binary spectral clustering, as observed by Hajij et al. (2018). This link is a strong  
164 indicator that Mapper can compute ‘useful’ clusters for pooling. We formally restate this observation below  
165 and provide a short proof.

166 **Proposition 5.1.** Let  $L$  be the Laplacian of a graph  $G(V, E)$  and  $l_2$  the eigenvector corresponding to the  
167 second lowest eigenvalue of  $L$ , also known as the Fiedler vector (Fiedler, 1973). Then, for a function  
168  $f : V \rightarrow \mathbb{R}, f(v) = l_2(v)$ , outputting the entry in the eigenvector  $l_2$  corresponding to node  $v$  and a cover  
169  $\mathcal{U} = \{(-\infty, \epsilon), (-\epsilon, +\infty)\}$ , Mapper produces a spectral bi-partition of the graph for a sufficiently small  
170 positive  $\epsilon$ .

171 PROOF. It is well known that the Fiedler vector can be used to obtain a ‘good’ bi-partition of the graph  
 172 based on the signature of the entries of the vector (i.e.  $l_2(v) > 0$  and  $l_2(v) < 0$ ) (please refer to Demmel  
 173 (1995) for a proof). Therefore, by setting  $\epsilon$  to a sufficiently small positive number  $\epsilon < \min_v |l_2(v)|$ , the  
 174 obtained pull back cover is a spectral bi-partition of the graph.

175 The result above indicates that Mapper is a generalisation of spectral clustering. As the latter is strongly  
 176 related to min-cuts (Leskovec, 2016), the proposition also links them to Mapper. We now provide a much  
 177 stronger result in that direction, showing that Structural Mapper is a generalisation of all pooling methods  
 178 based on soft-cluster assignments. *Soft cluster assignment pooling methods* use a soft cluster assignment  
 179 matrix  $\mathbf{S} \in \mathbb{R}^{N \times K}$ , where  $S_{ij}$  encodes the probability that node  $i$  belongs to cluster  $j$ ,  $N$  is the number of  
 180 nodes in the graph and  $K$  is the number of clusters. The adjacency matrix of the pooled graph is computed  
 181 via  $\mathbf{A}' = \mathbf{S}^T(\mathbf{A} + \mathbf{I})\mathbf{S}$ . Below, we prove a helpful result concerning this class of methods.

182 **Lemma 5.1.** The adjacency matrix  $\mathbf{A}' = \mathbf{S}^T(\mathbf{A} + \mathbf{I})\mathbf{S}$  defines a pooled graph, where the nodes  
 183 corresponding to clusters encoded by  $\mathbf{S}$  are connected if and only if there is a common edge (including  
 184 self-loops) between them.

185 PROOF. Let  $\mathbf{L} = \mathbf{A}\mathbf{S}$ . Then,  $A'_{ij} = \sum_k S_{ik}^T L_{kj} = 0$  if and only if  $S_{ik}^T = 0$  (node  $k$  does not belong to  
 186 cluster  $i$ ) or  $L_{kj} = 0$  (node  $k$  is not connected to any node belonging to cluster  $j$ ), for all  $k$ . Therefore,  
 187  $A'_{ij} \neq 0$  if and only if there exists a node  $k$  such that  $k$  belongs to cluster  $i$  and  $k$  is connected to a node  
 188 from cluster  $j$ . Due to the added self-loops,  $A'_{ij} \neq 0$  also holds if there is a node  $k$  belonging to both  
 189 clusters.

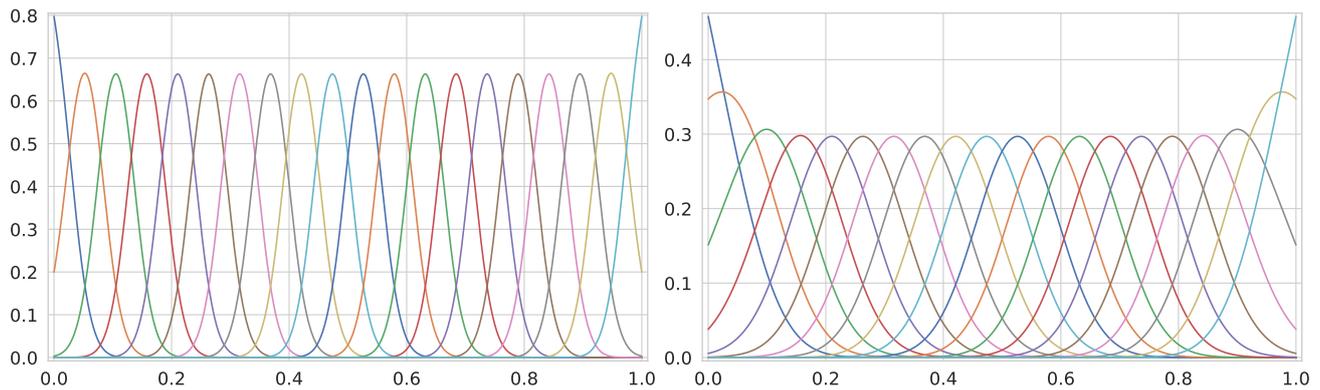
190 **Proposition 5.2.**  $\mathcal{GM}(f, \mathcal{U})$  generalises approaches based on soft-cluster assignments.

191 PROOF. Let  $s : V \rightarrow \Delta_{K-1}$  be a soft cluster assignment function that maps the vertices to the  $(K - 1)$ -  
 192 dimensional unit simplex. We denote by  $s_k(v)$  the probability that vertex  $v$  belongs to cluster  $k \leq K$  and  
 193  $\sum_k^K s_k(v) = 1$ . This function can be completely specified by a cluster assignment matrix  $\mathbf{S} \in \mathbb{R}^{N \times K}$  with  
 194  $S_{ik} = s_k(i)$ . This is the soft cluster assignment matrix computed by algorithms like minCut and DiffPool.  
 195 Let  $\mathcal{U} = \{U_i\}_{i \leq K}$  with  $U_i = \{x \in \Delta_{K-1} | x = \sum_j \lambda_j u_j, \sum_j \lambda_j = 1 \text{ and } \lambda_i > 0\}$  be an open cover of  
 196  $\Delta_{K-1}$ . Then consider an instance of  $\mathcal{GM}$  where everything is assigned to a single cluster (i.e. same as no  
 197 clustering). Clearly, there is a one-to-one correspondence between the vertices of  $\mathcal{GM}(s, \mathcal{U})$  and the soft  
 198 clusters. By Remark 1, the nodes corresponding to the clusters are connected only if the clusters share at  
 199 least one node or at least one edge. Then, by Lemma 5.1 the adjacency between the nodes of  $\mathcal{GM}(s, \mathcal{U})$   
 200 are the same as those described by  $\mathbf{A}' = \mathbf{S}^T(\mathbf{A} + \mathbf{I})\mathbf{S}$ . Thus, the two pooled graphs are isomorphic.

201 We hope that this result will enable theoreticians to study pooling operators through the topological and  
 202 statistical properties of Mapper (Carriere et al., 2018; Carrière and Oudot, 2018; Dey et al., 2017). At the  
 203 same time, we encourage practitioners to take advantage of it and design new pooling methods in terms of  
 204 a well-chosen lens function  $f$  and cover  $\mathcal{U}$  for its image. To illustrate this idea and showcase the benefits of  
 205 this new perspective over graph pooling methods, we introduce two Mapper-based operators.

## 206 5.2 Differentiable Mapper Pooling (DMP)

207 The main challenge for making pooling via Mapper differentiable is to differentiate through the pull back  
 208 computation. To address this, we replace the cover of  $n$  overlapping intervals over the real line, described  
 209 in the previous section, with a cover formed of overlapping RBF kernels  $\phi(x, x_i) = \exp(-\frac{\|x-x_i\|^2}{\delta})$ ,



**Figure 2.** Two covers of RBF kernels with different scales:  $\delta = 0.002$  and  $\delta = 0.01$ . The  $x$ -axis corresponds to the unit interval where the nodes of the graph are mapped. The  $y$ -axis represents the value of the normalised RBF kernels.

210 evaluated at  $n$  fixed locations  $x_i$ . The overlap between these kernels can be adjusted through the scale  $\delta$  of  
 211 the kernels. The soft cluster assignment matrix  $\mathbf{S}$  is given by the normalised kernel values:

$$S_{ij} = \frac{\phi(\sigma(f_\theta(\mathbf{X}_l))_i, x_j)}{\sum_{j=1}^n \phi(\sigma(f_\theta(\mathbf{X}_l))_i, x_j)}, \quad (1)$$

212 where the lens function  $f_\theta$  is a GNN layer,  $\sigma$  is a sigmoid function ensuring the outputs are in  $[0, 1]$ , and  
 213  $\mathbf{X}_l$  are the node features at layer  $l$ . Intuitively, the more closely a node is mapped to a location  $x_i$ , the  
 214 more it belongs to cluster  $i$ . By Proposition 5.2, we can compute the adjacency matrix of the pooled  
 215 graph as  $\mathbf{S}^T(\mathbf{A} + \mathbf{I})\mathbf{S}$ ; the features are given by  $\mathbf{S}^T\mathbf{X}$ . This method can also be thought as a version of  
 216 DiffPool (Ying et al., 2018), where the low-entropy constraint on the cluster assignment distribution is  
 217 topologically satisfied, since a point cannot be equally close to many other points on a line. Therefore, each  
 218 node will belong only to a few clusters if the scale  $\delta$  is appropriately set.

219 In Figure 2 we show two examples of RBF kernel covers for the output space. The scale of the kernel,  $\delta$ ,  
 220 determines the amount of overlap between the cover elements. At bigger scales, there is a higher overlap  
 221 between the clusters, as shown in the two plots. Because the line is one-dimensional, a point on the unit  
 222 interval can only be part of a small number of clusters (that is, the kernels for which the value is greater  
 223 than zero), assuming the scale  $\delta$  is not too large. Therefore, DMP can be seen as a DiffPool variant where  
 224 the low-entropy constraint on the cluster assignment is satisfied topologically, rather than by a loss function  
 225 enforcing it.

### 226 5.3 Mapper-based PageRank (MPR) Pooling

227 To evaluate the effectiveness of the differentiable pooling operator, we also consider a fixed and scalable  
 228 non-differentiable lens function  $f : V \rightarrow \mathbb{R}$  that is given by the normalised PageRank (PR) (Page et al.,  
 229 1999) of the nodes. The PageRank function assigns an importance value to each of the nodes based on their  
 230 connectivity, according to the well-known recurrence relation:

$$f(\mathbf{X})_i \triangleq \mathbf{PR}_i = \sum_{j \in N(i)} \frac{\mathbf{PR}_j}{|N(i)|}, \quad (2)$$

231 where  $N(i)$  represents the set of neighbours of the  $i$ -th node in the graph and the damping factor was  
 232 set to the typical value of  $d = 0.85$ . The resulting scores are values in  $[0, 1]$  which reflect the probability

233 of a random walk through the graph to end in a given node. Using the previously described overlapping  
 234 intervals cover  $\mathcal{U}$ , the elements of the pull back cover form a soft cluster assignment matrix  $\mathbf{S}$ :

$$S_{ij} = \frac{\mathbb{I}_{i \in f^{-1}(U_j)}}{|\{U_k | i \in f^{-1}(U_k)\}|} \quad (3)$$

235 where  $U_n$  is the  $n$ -th cover set in the cover  $\mathcal{U}$  of  $[0, 1]$ . It can be observed that the resulting clusters  
 236 contain nodes with similar PageRank scores. Intuitively, this pooling method merges the (usually few)  
 237 highly-connected nodes in the graph, at the same time clustering the (typically many) dangling nodes that  
 238 have a normalised PageRank score closer to zero. Therefore, this method favours the information attached  
 239 to the most ‘important’ nodes of the graph. The adjacency matrix of the pooled graph and the features are  
 240 computed in the same manner as for DMP.

## 241 5.4 Model

242 For the graph classification task, each example  $G$  is represented by a tuple  $(\mathbf{X}, \mathbf{A})$ , where  $\mathbf{X}$  is the node  
 243 feature matrix and  $\mathbf{A}$  is the adjacency matrix. Both our graph embedding and classification networks  
 244 consist of a sequence of graph convolutional layers (Kipf and Welling, 2016); the  $l$ -th layer operates on its  
 245 input feature matrix as follows:

$$\mathbf{X}_{l+1} = \sigma(\hat{\mathbf{D}}^{-\frac{1}{2}} \hat{\mathbf{A}} \hat{\mathbf{D}}^{-\frac{1}{2}} \mathbf{X}_l \mathbf{W}_l), \quad (4)$$

246 where  $\hat{\mathbf{A}} = \mathbf{A} + \mathbf{I}$  is the adjacency matrix with self-loops,  $\hat{\mathbf{D}}$  is the normalised node degree matrix,  $\mathbf{W}_l$  is  
 247 the weight matrix of the  $l$ -th layer and  $\sigma$  is the activation function. After  $E$  layers, the embedding network  
 248 simply outputs node features  $\mathbf{X}_{L_E}$ , which are subsequently processed by a pooling layer to coarsen the  
 249 graph. The classification network first takes as input node features of the Mapper-pooled graph<sup>2</sup>,  $\mathbf{X}_{MG}$ ,  
 250 and passes them through  $L_C$  graph convolutional layers. Following this, the network computes a graph  
 251 summary given by the feature-wise node average and applies a final linear layer which predicts the class:

$$y = \text{softmax}\left(\frac{1}{|\text{MG}|} \sum_{i=1}^{|\text{MG}|} \mathbf{X}_{L_C} \mathbf{W}_f + \mathbf{b}_f\right). \quad (5)$$

252 We note that either of these pooling operators could readily be adapted to the recently proposed message  
 253 passing simplicial neural networks (MPSNs) (Bodnar et al., 2021) as a tool for coarsening simplicial  
 254 complexes by dropping the 1-skeleton operator after computing the nerve. We leave this endeavour for  
 255 future work. The complete model details can be found in Appendix A.

## 256 5.5 Complexity

257 The topology of the output graph can be computed in  $O(V + E)$  time when using a cover over the unit  
 258 interval, as described above. The output graph can be computed via (sparse) matrix multiplication given by  
 259  $\mathbf{S}^T(\mathbf{A} + \mathbf{I})\mathbf{S}$ , to take advantage of GPU parallelism and compute the coefficients associated with the edges.

## 6 POOLING EXPERIMENTS

260 **Tasks.** We illustrate the applicability of the Mapper-GNN synthesis within a pooling framework, by  
 261 evaluating DMP and MPR in a variety of settings: social (IMDB-Binary, IMDB-Multi, Reddit-Binary,  
 262 Reddit-Multi-5k), citation networks (Collab) and chemical data (D&D, Mutag, NCI1, Proteins) (Kersting  
 263 et al., 2016).

<sup>2</sup> Note that one or more {embedding  $\rightarrow$  pooling} operations may be sequentially performed in the pipeline.

**Table 1.** Results obtained on classification benchmarks. Accuracy measures with 95% confidence intervals are reported. The highest result is **bolded** and the second highest is underlined. The first columns four are molecular graphs, while the others are social graphs. Our models perform competitively with other state of the art models.

Model	D&D	Mutag	NC11	Proteins	Collab	IMDB-B	IMDB-M	Reddit-B	Reddit-5k
DMP (Ours)	77.3 ± 3.6	84.0 ± 8.6	70.4 ± 4.2	<b>75.3 ± 3.3</b>	81.4 ± 1.2	<b>73.8 ± 4.5</b>	<b>50.9 ± 2.5</b>	86.2 ± 6.8	51.9 ± 2.1
MPR (Ours)	<b>78.2 ± 3.4</b>	80.3 ± 6.0	69.8 ± 1.8	75.2 ± 2.2	<b>81.5 ± 1.0</b>	73.4 ± 2.7	50.6 ± 2.0	<u>86.3 ± 4.8</u>	<u>52.3 ± 1.6</u>
Top- <i>k</i>	75.1 ± 2.2	82.5 ± 6.8	67.9 ± 2.3	74.8 ± 3.0	75.0 ± 1.1	69.6 ± 3.8	45.0 ± 2.8	79.4 ± 7.4	48.5 ± 1.1
minCUT	77.6 ± 3.1	82.9 ± 6.0	68.8 ± 2.1	73.5 ± 2.9	79.9 ± 0.8	70.7 ± 3.5	50.6 ± 2.1	<b>87.2 ± 5.0</b>	<b>52.9 ± 1.3</b>
DiffPool	77.9 ± 2.4	<b>94.7 ± 7.1</b>	68.1 ± 2.1	74.2 ± 0.3	81.3 ± 0.1	72.4 ± 3.1	50.3 ± 1.8	79.0 ± 1.1	50.4 ± 1.7
WL	<u>77.4 ± 2.6</u>	74.5 ± 6.5	<b>76.4 ± 2.7</b>	74.7 ± 3.2	78.5 ± 1.1	72.1 ± 3.1	<u>50.7 ± 2.9</u>	66.7 ± 10.4	49.2 ± 1.4
Flat	69.9 ± 2.2	71.8 ± 4.3	65.5 ± 1.7	70.2 ± 2.6	80.9 ± 1.4	<u>73.6 ± 4.2</u>	48.5 ± 2.4	70.0 ± 10.8	49.5 ± 1.7
avg-MLP	63.7 ± 1.4	69.1 ± 5.8	55.7 ± 2.8	61.8 ± 1.7	74.8 ± 1.3	71.5 ± 2.9	49.5 ± 2.2	53.6 ± 6.2	45.9 ± 1.6

264 **Experimental setup.** We adopt a 10-fold cross-validation approach to evaluating the graph classification  
 265 performance of DMP, MPR and other competitive state-of-the-art methods. The random seed was set to  
 266 zero for all experiments (with respect to dataset splitting, shuffling and parameter initialisation), in order to  
 267 ensure a fair comparison across architectures. All models were trained on a single Titan Xp GPU, using the  
 268 Adam optimiser (Kingma and Ba, 2014) with early stopping on the validation set, for a maximum of 30  
 269 epochs. We report the classification accuracy using 95% confidence intervals calculated for a population  
 270 size of 10 (the number of folds).

271 **Baselines.** We compare the performance of DMP and MPR to two other pooling methods that we  
 272 identify mathematical connections with: minCUT (Bianchi et al., 2019) and DiffPool (Ying et al., 2018).  
 273 Additionally, we include Graph U-Net (Gao and Ji, 2019) in our evaluation, as it has been shown to yield  
 274 competitive results while performing pooling from the perspective of a learnable node ranking; we denote  
 275 this approach by Top-*k* in the remainder of this section. The non-pooling baselines evaluated are the  
 276 WL kernel (Shervashidze et al., 2011), a ‘flat’ model (2 MP steps and global average pooling) and an  
 277 average-readout linear classifier.

278 We optimise both DMP and MPR with respect to the cover cardinality  $n$ , the cover overlap ( $\delta$  for DMP,  
 279 overlap percentage  $g$  for MPR), learning rate and hidden size. The Top-*k* architecture is evaluated using  
 280 the code provided in the official repository, where separate configurations are defined for each of the  
 281 benchmarks. The minCUT architecture is represented by the sequence of operations described by Bianchi  
 282 et al. (2019): *MP(32)-pooling-MP(32)-pooling-MP(32)-GlobalAvgPool*, followed by a linear softmax  
 283 classifier. The *MP(32)* block represents a message-passing operation performed by a graph convolutional  
 284 layer with 32 hidden units:

$$\mathbf{X}^{(t+1)} = \text{ReLU}(\tilde{\mathbf{A}}\mathbf{X}^{(t)}\mathbf{W}_m + \mathbf{X}^{(t)}\mathbf{W}_s), \quad (6)$$

285 where  $\tilde{\mathbf{A}} = \mathbf{D}^{-\frac{1}{2}}\mathbf{A}\mathbf{D}^{-\frac{1}{2}}$  is the symmetrically-normalised adjacency matrix and  $\mathbf{W}_m, \mathbf{W}_s$  are learnable  
 286 weight matrices representing the message passing and skip-connection operations within the layer.  
 287 The DiffPool model follows the same sequence of steps. Full details of the model architectures and  
 288 hyperparameters can be found in Appendix A.

289 **Evaluation Procedure** The best procedure for evaluating GNN pooling layers remains a matter of debate  
 290 in the graph machine learning community. One may consider a fixed GNN architecture with a different  
 291 pooling layer for each baseline; alternatively, the whole architecture can be optimised for each type of  
 292 pooling layer. The first option, more akin to the typical procedure for evaluating pooling layers in CNNs on  
 293 image domains, is used in papers like minCUT (Bianchi et al., 2019). The second option is more particular

294 to GNNs and it is employed, for instance, by DiffPool (Ying et al., 2018). In this work, we choose the latter  
295 option for our evaluation.

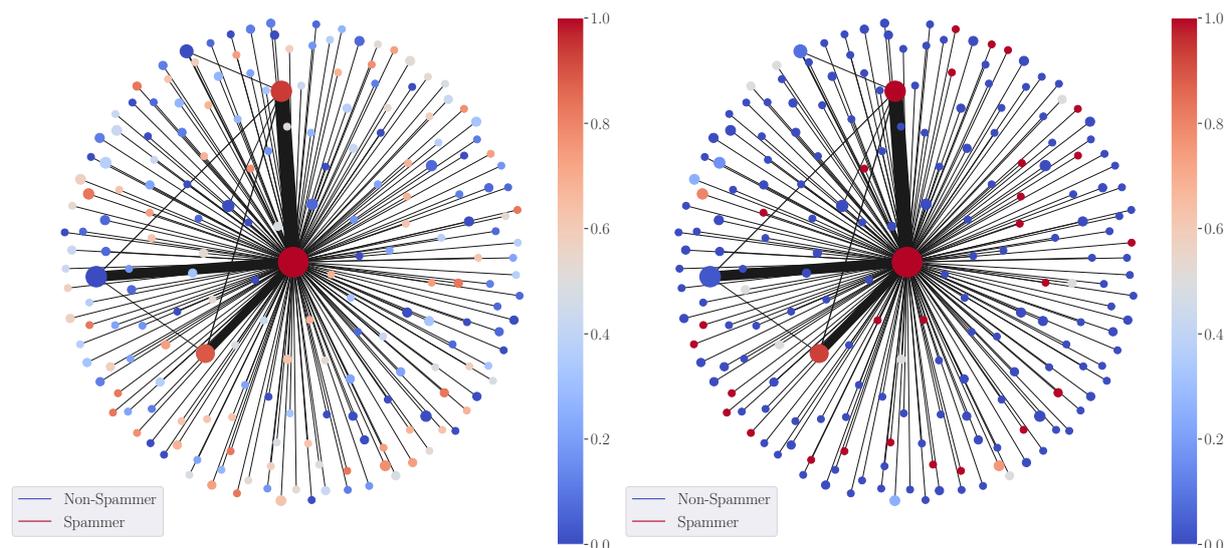
296 We argue that for non-Euclidean domains, such as graph ones, the relationships between the nodes of  
297 the pooled graph and the ones of the input graph are semantically different from one pooling method to  
298 another. This is because pooling layers have different behaviours and may interact in various ways with the  
299 interleaved convolutional layers. Therefore, evaluating the same architecture with only the pooling layer(s)  
300 swapped is restrictive and might hide the benefits of certain operators. For example, Top- $k$  pooling (one of  
301 our baselines) simply drops nodes from the input graph, instead of computing a smaller number of clusters  
302 from all nodes. Assume we fix the pooled graph to have only one node. Then Top- $k$  would only select one  
303 node from the original graph. In contrast, DiffPool would combine the information from the entire graph  
304 in a single node. DiffPool would thus have access to additional information with respect to Top- $k$ , so it  
305 would be unfair to conclude that one model is better than the other in such a setting. These differences  
306 implicitly affect the features of the output graph at that layer, which in turn affect the next pooling layer, as  
307 its computation depends on the features. This can have a cascading effect on the overall performance of  
308 the model. One might also argue that this procedure makes the evaluated models more homogeneous and,  
309 therefore, easier to compare. While this is true, the conclusions one can draw from such a comparison are  
310 much more limited because they are restricted to the particular architecture that was chosen.

311 For this reason, we have either run models with hyperparameters as previously reported by the authors,  
312 or optimised them ourselves end-to-end, where applicable. The best-performing configurations were:

- 313 • MPR—learning rate  $5e^{-4}$ , hidden sizes  $\{128, 128\}$  (except for  $\{64, 64\}$  on IMDB-Binary and  $\{32, 32\}$   
314 on IMDB-Multi), interval overlap 25% on Proteins, Reddit-Binary, Mutag, IMDB-Multi and 10%  
315 otherwise, batch size 32 (except for 128 on Proteins) and:
  - 316 • D&D, Collab, Reddit-Binary, Reddit-Multi-5K: cover sizes  $\{20, 5\}$ ;
  - 317 • Proteins, NCI1: cover sizes  $\{8, 2\}$ ;
  - 318 • Mutag, IMDB-Binary, IMDB-Multi: cover sizes  $\{4, 1\}$ ;
- 319 • DMP—learning rate  $5e^{-4}$ , hidden sizes  $\{128, 128\}$ ,  $\delta = 1/(\text{cluster\_size})^2$  and:
  - 320 • Proteins: cover sizes  $\{8, 2\}$ , batch size 128;
  - 321 • Others: cover sizes  $\{20, 5\}$ , batch size 32;
- 322 • Top- $k$ —specific dataset configurations, as provided in the official GitHub repository<sup>3</sup>;
- 323 • minCUT—learning rate  $1e^{-3}$ , same architecture as reported by the authors in the original work (Bianchi  
324 et al., 2019);
- 325 • DiffPool—learning rate  $1e^{-3}$ , hidden size 32, two pooling steps, pooling ratio  $r = 0.1$  for D&D,  
326 Proteins, Collab and Reddit-Binary and  $r = 0.25$  for Mutag, NCI1, IMDB-Binary, IMDB-Multi and  
327 Reddit-Multi-5K, global average mean readout layer, with the exception of Collab and Reddit-Binary,  
328 where the hidden size was 128;
- 329 • Flat: hidden size 32.

330 **Pooling Results.** The graph classification performance obtained by these models is reported in Table 1. We  
331 reveal that MPR ranks either first or second on all social datasets, or achieves accuracy scores within 0.5%  
332 of the best-performing model. This result confirms that PageRank-based pooling exploits the power-law  
333 distributions in this domain. The performance of DMP is similar on social data and generally higher on

<sup>3</sup> [https://github.com/HongyangGao/Graph-U-Nets/blob/48aa171b16964a2466fcea4cb06fc940d649294/run\\_GUNet.sh](https://github.com/HongyangGao/Graph-U-Nets/blob/48aa171b16964a2466fcea4cb06fc940d649294/run_GUNet.sh)



**Figure 3.** SDGM visualisation using as a lens function the the GNN-predicted probability of a node in the network to be Spam. The left plot is coloured with the average predicted spam probability in each cluster, whereas the right plot is coloured by the proportion of true spammers in each node.

334 molecular graphs. We attribute this to the fact that all nodes in molecular graphs tend to have a similar  
 335 PageRank score—MPR is therefore likely to assign all nodes to one cluster, effectively performing a  
 336 readout. In this domain, DMP performs particularly well on Mutag, where it is second-best and improves  
 337 by 3.7% over MPR, showing the benefits of having a differentiable lens in challenging data settings.  
 338 Overall, MPR achieves the best accuracy on 2 datasets (D&D, Collab) and the next best result on 3 more  
 339 (Proteins, Reddit-Binary and Reddit-Multi-5k). DMP improves on MPR by less than 1% on NC11, Proteins,  
 340 IDMB-Binary and IMDB-Multi, showing the perhaps surprising strength of the simple, fixed-lens pooling  
 341 MPR operator.

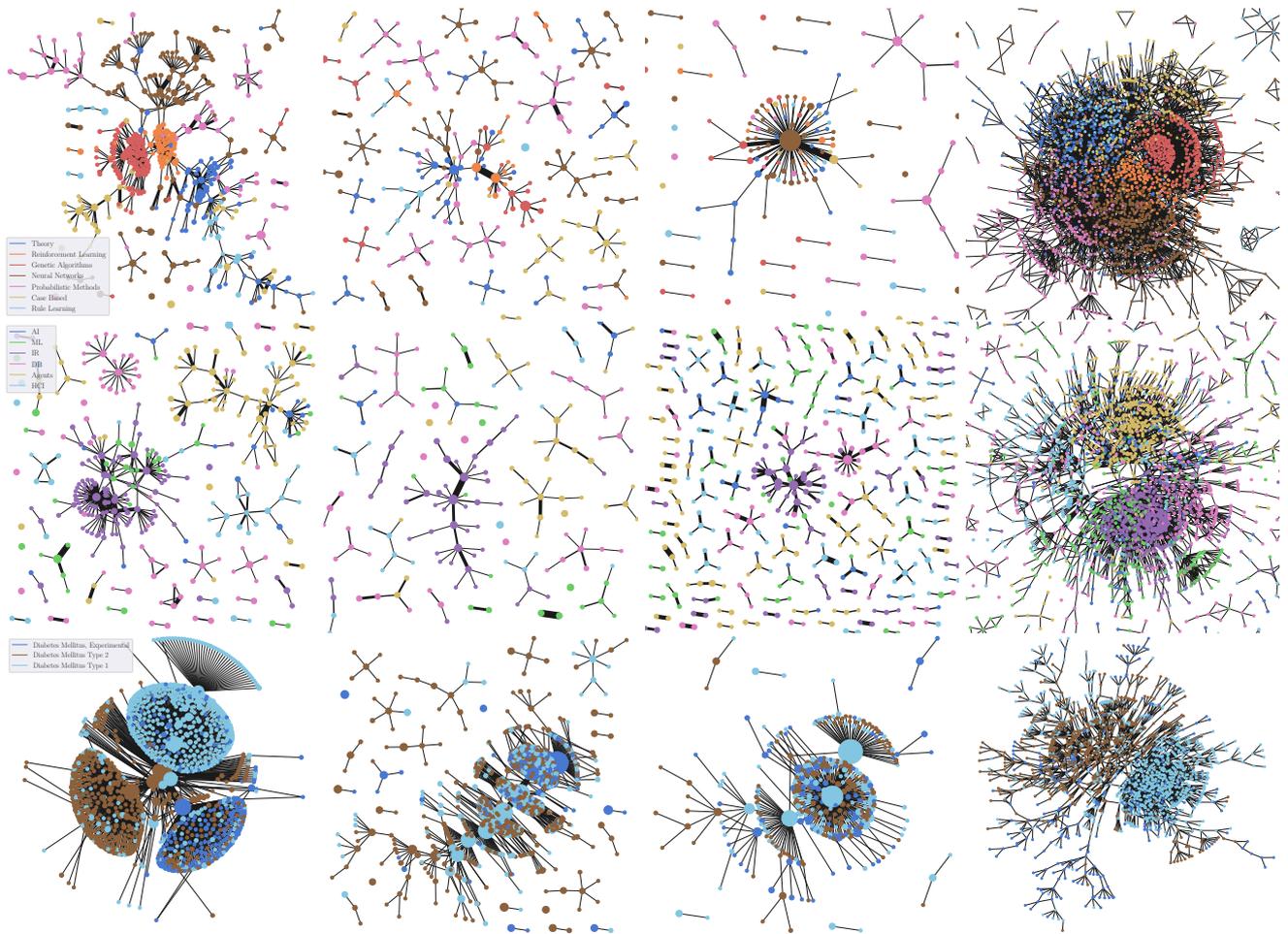
## 7 MAPPER FOR VISUALISATIONS

342 Graph pooling methods and summarised graph visualisations methods can be seen as two sides of the same  
 343 coin, since both aim to condense the information in the graph. We now turn our attention to the latter.

### 344 7.1 Visualisations in Supervised Learning

345 The first application of DGM is in a supervised learning context, where  $f_\theta$  is trained via a cross entropy  
 346 loss function to classify the nodes of the graph. When the classification is binary,  $f_\theta : V \rightarrow [0, 1]$  outputs  
 347 the probability that a node belongs to the positive class. This probability acts directly as the parameterisation  
 348 of the graph nodes. An example is shown in Figure 3 (left) for a synthetic dataset a network formed of  
 349 spammers and non-spammers. Spammers are highly connected to many other nodes in the network, whereas  
 350 non-spammers generally have fewer neighbours. For the lens function, we use a Graph Convolutional  
 351 Network (GCN) (Kipf and Welling, 2016) with four layers (with 32, 64, 128, 128 hidden units) and ReLU  
 352 activations trained to classify the nodes of the graph. For the spammer graph, the lens is given by the  
 353 predicted spam probability of each node and the cover consists of 10 intervals over  $[0, 1]$ , with 10% overlap.

354 Through the central cluster node, the SDGM visualisation correctly shows how spammers occupy an  
 355 essential place in the network, while non-spammers tend to form many smaller disconnected communities.  
 356 When labels are available, we also produce visualisations augmented with ground-truth information. These  
 357 visualisations can provide a label-driven understanding of the graph. For instance, in Figure 3 (right) we  
 358 colour each node of the SDGM visualisation according to the most frequent class in the corresponding



**Figure 4.** Qualitative comparison between SDGM (first column), Mapper with an RBF graph density function (Hajij et al., 2018) (second), and Mapper with a PageRank function (Hajij et al., 2018) (third). The Graphviz visualisation of the graph cores (fourth column) are added for reference. The rows show plots for Cora, CiteSeer, and PubMed, respectively. The graphs are coloured based on the most frequent class in each cluster to aid the comparison. SDGM with unsupervised lens implicitly makes all dataset classes appear in the visualisation more clearly separated. This does not happen in the baseline visualisations, which mainly focus on the class with the highest number of nodes from each graph.

359 cluster. This second visualisation, augmented with the ground-truth information, can also be used to  
 360 compare with the model predictions.

## 361 7.2 Visualisation in Unsupervised Learning

362 The second application corresponds to an unsupervised learning scenario, where the challenge is obtaining  
 363 a parameterisation of the graph in the absence of labels. This is the typical use case for unsupervised  
 364 graph representation learning models (Chami et al., 2020). The approach we follow is to train a model  
 365 to learn node embeddings in  $\mathbb{R}^{d'}$  (in our experiments,  $d' = 512$ ), which can be reduced, as before, to a  
 366 low-dimensional space via a dimensionality reduction method  $r$ . Unsupervised visualisations can be found  
 367 in the qualitative evaluation in Section 7.3.

## 368 7.3 Qualitative Evaluation

369 In this section, we qualitatively compare SDGM against the two best-performing graph theoretic lens  
 370 functions proposed by Hajij et al. (2018), on the Cora and CiteSeer (Sen et al., 2008) and PubMed (Yang

371 et al., 2016) citation networks. Namely, we compare against a PageRank (Page et al., 1999) lens function  
372 and a graph density function  $f(v) = \sum_{u \in V} \exp(-D(u, v)/\delta)$ , where  $D$  is the distance matrix of the graph.  
373 For SDGM, we use a composition of an unsupervised Deep Graph Infomax (DGI) (Veličković et al.,  
374 2018) model  $g_\theta : V \rightarrow \mathbb{R}^{512}$  and a dimensionality reduction function  $r : \mathbb{R}^{512} \rightarrow \mathbb{R}^2$  based on  $t$ -SNE. To  
375 aid the comparison, we mark the nodes with the colour of the most frequent class in the corresponding  
376 cluster. Additionally, we include a Graphviz (Gansner and North, 2000) plot of the full graph. We carefully  
377 fine-tuned the covers for each combination of model and graph.

378 As depicted by Figure 4, SDGM successfully summarises many of the properties of the graphs that are  
379 also reflected by full graph visualisations. For instance, on Cora, Genetic Algorithms (in dark orange) are  
380 shown to be primarily connected to Reinforcement Learning (orange). At the same time, related classes  
381 that largely overlap in the full visualisation—Probabilistic Methods and Neural Networks (NNs) on Cora  
382 or Information Retrieval (IR) and ML on CiteSeer—are connected in the SDGM plot. In contrast, the  
383 baselines do not have the same level of granularity and fail to capture many such properties. Both PageRank  
384 and the graph density function tend to focus on the classes with the highest number of nodes, such as the  
385 IR class on CiteSeer or the NNs class on Cora, while largely de-emphasising other classes.

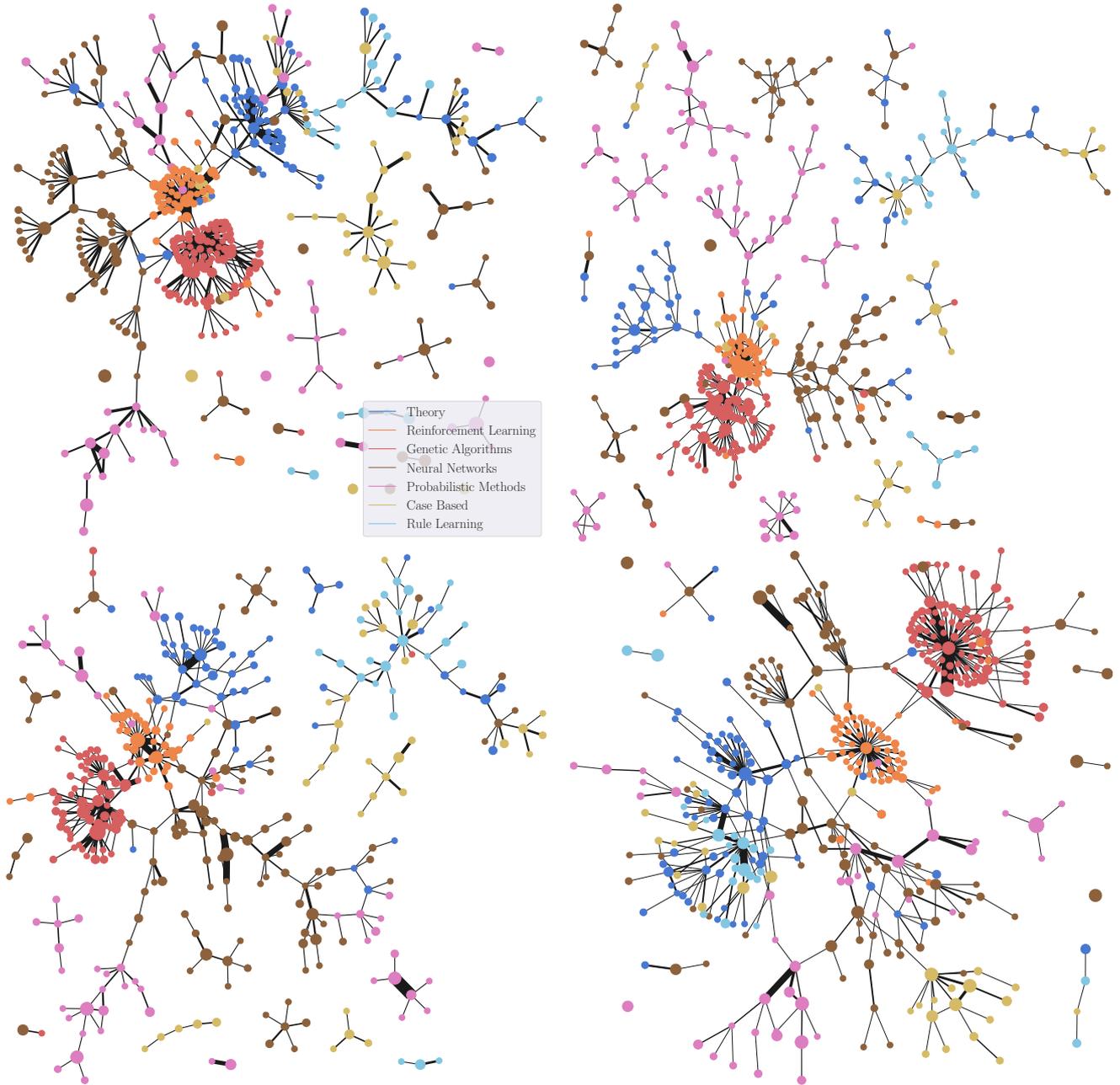
386 **Limitations.** The proposed visualisations also present certain limitations. In an unsupervised learning  
387 setting, in the absence of any labels or attributes for colouring the graph, the nodes have to be coloured  
388 based on a colourmap associated with the abstract embedding space, thus affecting the interpretability  
389 of the visualisations. In contrast, even though the graph theoretic lens functions produce lower quality  
390 visualisations, their semantics are clearly understood mathematically. This is, however, a drawback shared  
391 even by some of the most widely used data visualisation methods, such as  $t$ -SNE or UMAP (McInnes et al.,  
392 2018). In what follows, we present additional visualisations and ablation studies.

#### 393 7.4 Ablation study for dimensionality reduction

394 We study how the choice of the dimensionality reduction method for the unsupervised visualisations  
395 affects the output. To test this, we consider the following dimensionality reduction methods:  $t$ -SNE (van der  
396 Maaten and Hinton, 2008), UMAP (McInnes et al., 2018), IsoMap (Tenenbaum et al., 2000) and PCA.  
397 We use the same model as in Section 7.2 and 81 2D cells for the cover of all models. The overlap was  
398 set after fine-tuning to 0.2 for  $t$ -SNE and UMAP, and to 0.1 for the other two models. Figure 5 displays  
399 the four visualisations. As expected,  $t$ -SNE and UMAP produce more visually-pleasing outputs, due to  
400 their superior ability to capture variation in the GNN embedding space. However, the features highlighted  
401 by all visualisations are largely similar, generally indicating the same binary relations between clusters.  
402 This demonstrates that the GNN embedding space is robust to the choice of the dimensionality reduction  
403 method.

#### 404 7.5 Ablation for the unsupervised lens

405 To better understand the impact of GNNs on improving the quality of the Mapper visualisations, we  
406 perform an ablation study on the type of unsupervised lens functions used within Mapper. The first model  
407 we consider is simply the identity function taking as input only graph features. The second model is  
408 a randomly initialised DGI model. Despite the apparent simplicity of a randomly initialised model, it  
409 was shown that such a method produces reasonably good embeddings, often outperforming other more  
410 sophisticated baselines (Veličković et al., 2018). Finally, we use our trained DGI model from Section 7.2.  
411 For all models, we perform a  $t$ -SNE reduction of their embedding space to obtain a 2D output space and  
412 use 81 overlapping cells that cover this space. An overlap of 0.2 is used across all models.

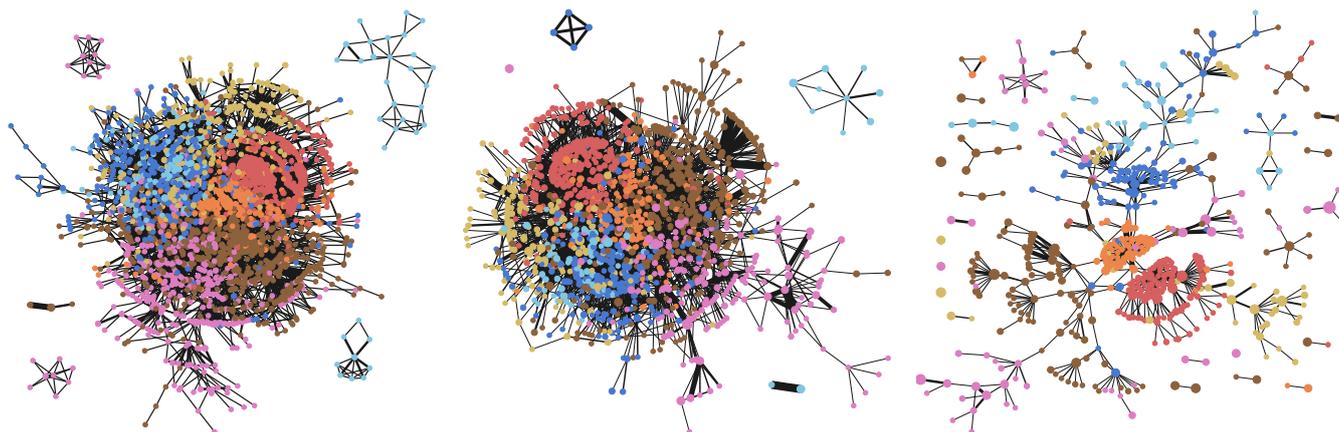


**Figure 5.** Ablation for dimensionality reduction methods; left–right, top–bottom:  $t$ -SNE, PCA, Isomap, UMAP. While  $t$ -SNE and UMAP produce slightly better visualisations, the graph features displayed by the visualisations are roughly consistent across all of the four dimensionality reduction techniques.

413 The three resulting visualisations are depicted in Figure 6. The identity model and the untrained DGI  
 414 model do not manage to exploit the dataset structure and neither does particularly well. In contrast, the  
 415 trained DGI model emphasises all the classes in the visualisation, together with their main interactions.

## 416 7.6 Hierarchical visualisations

417 One of the most powerful features of Mapper is the ability to produce multi-resolution visualisations  
 418 through the flexibility offered by the cover hyperparameters. As described in Section 4, having a higher  
 419 number of cells covering the output space results in more granular visualisations containing more nodes,  
 420 while a higher overlap between these cells results in increased connectivity. We highlight these trade-offs in



**Figure 6.** Ablation for different types of unsupervised lenses (identity, untrained DGI, trained DGI). The trained DGI model significantly improves the quality of the visualisations.

421 Figure 7, where we visualise the Cora citation network using 9 combinations of cells and overlaps. These  
 422 kinds of hierarchical visualisations can help one identify the persistent features of the graph. For instance,  
 423 when inspecting the plots that use  $n = 64$  cells, the connections between the light blue class and the  
 424 yellow class persist for all 3 degrees of overlap, which indicates that this is a persistent feature of the graph.  
 425 In contrast, the connection between the red and orange classes is relatively reduced ( $g = 0.25$ ) or none  
 426 ( $g = 0.1$ ) for low values of overlap, but it clearly appears at  $g = 0.35$  in the top-right corner, suggesting  
 427 that the semantic similarity between the two classes is very scale-sensitive (that is, less persistent).

#### 428 **7.7 The importance of capturing structural information**

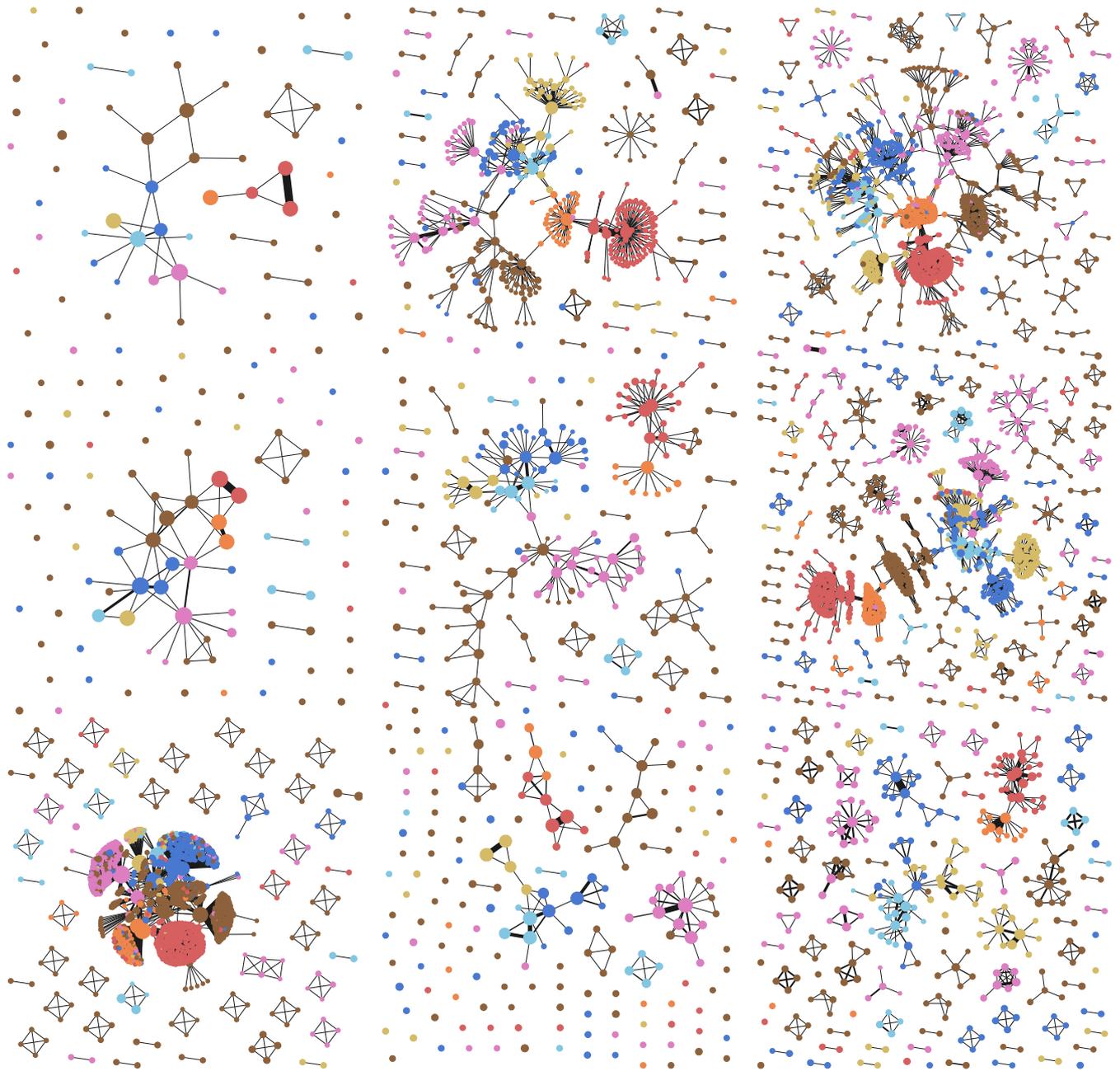
429 In this section, we revisit the synthetic spammer dataset to illustrate the importance of capturing structural  
 430 information via the edge-refined pull back cover operator. To that end, we compare SDGM with a version  
 431 using the usual refined pull back cover as in Hajij et al. (2018), while using the same lens function for  
 432 both (a GCN classifier). We refer to the latter as DGM. The visualisations produced by the two models  
 433 are included in Figure 8. We note that while both models capture the large cluster of spammers at the  
 434 center of the network and the smaller communities of non-spammers, DGM does not capture the structural  
 435 relationships between spammers and non spammers since it encodes only semantic relations.

## 8 CONCLUSION

436 We have introduced Deep Graph Mapper, a topologically-grounded method for producing informative  
 437 graph visualisations with the help of GNNs. We have shown these visualisations are not only useful for  
 438 understanding various graph properties, but can also aid in visually identifying classification mistakes.  
 439 Additionally, we have proved that Mapper is a generalisation of soft cluster assignment methods, effectively  
 440 providing a bridge between graph pooling and the TDA literature. Based on this connection, we have  
 441 proposed two Mapper-based pooling operators: a simple one that scores nodes using PageRank and a  
 442 differentiable one that uses RBF kernels to simulate the cover. Our experiments show that both layers yield  
 443 architectures competitive with several state-of-the-art methods on graph classification benchmarks.

## ACKNOWLEDGEMENT

444 CC is supported by EPSRC NRAG/465 NERC CDT Dream. PL is funded by EPSRC. We are grateful to  
 445 both reviewers for their constructive feedback and useful iterations on the manuscript. We would like to  
 446 thank Petar Veličković, Ben Day, Felix Opolka, Simeon Spasov, Alessandro Di Stefano, Duo Wang, Jacob

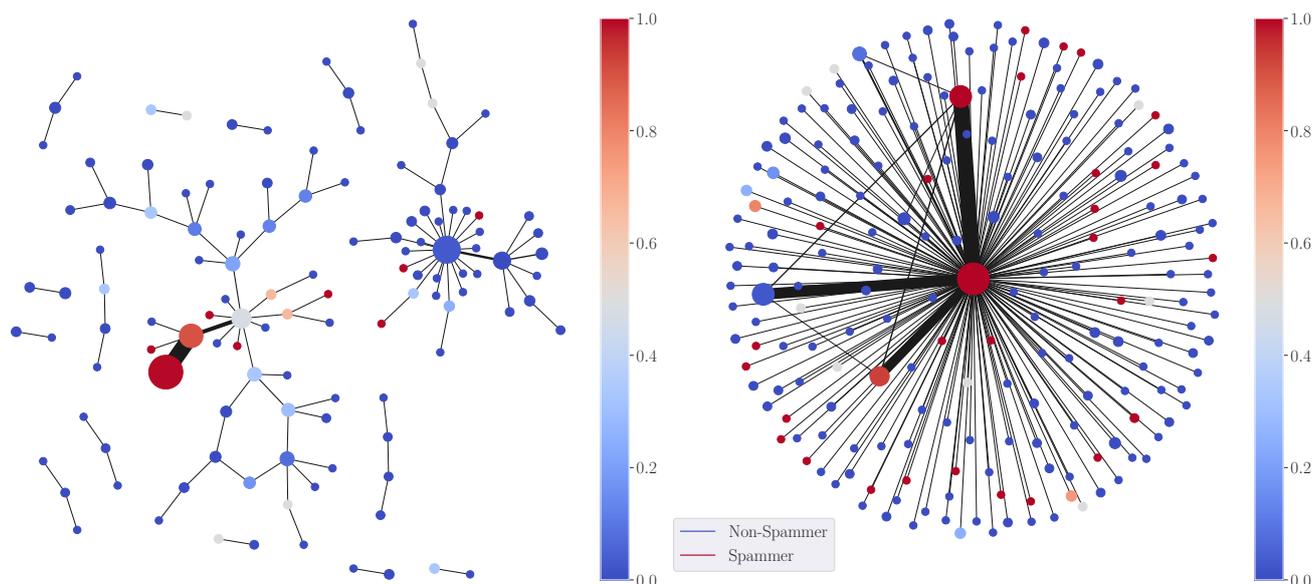


**Figure 7.** Hierarchical visualisations of the Cora citation network using various number of cover cells and degrees of overlap. Rows (top–bottom) have a different overlap ( $g$ ) between intervals:  $g = 0.1$ ,  $g = 0.25$ ,  $g = 0.35$ ; columns (left–right):  $n = 16$ ,  $n = 64$ ,  $n = 256$ .

447 Deasy, Ramon Viñas, Alex Dumitru and Teodora Reu for their constructive comments. We are also grateful  
 448 to Teo Stoleru for helping with the diagrams.

## REFERENCES

449 Batagelj, V., Didimo, W., Liotta, G., Palladino, P., and Patrignani, M. (2010). Visual analysis of large  
 450 graphs using (x,y)-clustering and hybrid visualizations. In *2010 IEEE Pacific Visualization Symposium*  
 451 (*PacificVis*). 209–216



**Figure 8.** DGM (left) vs SDGM (right) visualisation of the sythetic spammer datasets. DGM does not capture important relational information between spammers and non-spammers

- 452 Beck, F., Burch, M., Diehl, S., and Weiskopf, D. (2017). A taxonomy and survey of dynamic graph  
 453 visualization. *Computer Graphics Forum* 36, 133–159. doi:10.1111/cgf.12791
- 454 Bianchi, F. M., Grattarola, D., and Alippi, C. (2019). Mincut Pooling in Graph Neural Networks. *arXiv*  
 455 *preprint arXiv:1907.00481*
- 456 Bodnar, C., Frasca, F., Wang, Y. G., Otter, N., Montúfar, G., Liò, P., et al. (2021). Weisfeiler and lehman  
 457 go topological: Message passing simplicial networks. *arXiv preprint arXiv:2103.03212*
- 458 Bruna, J., Zaremba, W., Szlam, A., and Lecun, Y. (2014). Spectral networks and locally connected  
 459 networks on graphs. In *ICLR*
- 460 Cangea, C., Veličković, P., Jovanović, N., Kipf, T., and Liò, P. (2018). Towards Sparse Hierarchical Graph  
 461 Classifiers. *arXiv preprint arXiv:1811.01287*
- 462 Carriere, M., Michel, B., and Oudot, S. (2018). Statistical Analysis and Parameter Selection for Mapper.  
 463 *The Journal of Machine Learning Research* 19, 478–516
- 464 Carrière, M. and Oudot, S. (2018). Structure and stability of the one-dimensional mapper. *Foundations of*  
 465 *Computational Mathematics* 18, 1333–1396
- 466 Chami, I., Abu-El-Haija, S., Perozzi, B., Ré, C., and Murphy, K. (2020). Machine learning on graphs: A  
 467 model and comprehensive taxonomy. *ArXiv abs/2005.03675*
- 468 Chazal, F. and Michel, B. (2017). An introduction to Topological Data Analysis: fundamental and practical  
 469 aspects for data scientists. *arXiv preprint arXiv:1710.04019*
- 470 [Dataset] Demmel, J. (1995). UC Berkeley CS267 - Lecture 20: Partitioning Graphs without Coordinate  
 471 Information II
- 472 Dey, T. K., Mémoli, F., and Wang, Y. (2017). Topological analysis of nerves, reeb spaces, mappers, and  
 473 multiscale mappers. In *Symposium on Computational Geometry*
- 474 Dunne, C. and Shneiderman, B. (2013). Motif Simplification: Improving Network Visualization Readability  
 475 with Fan, Connector, and Clique Glyphs. In *Proceedings of the SIGCHI Conference on Human Factors in*  
 476 *Computing Systems* (New York, NY, USA: Association for Computing Machinery), CHI '13, 3247–3256.  
 477 doi:10.1145/2470654.2466444

- 478 Dwyer, T., Henry Riche, N., Marriott, K., and Mears, C. (2013). Edge Compression Techniques for  
479 Visualization of Dense Directed Graphs. *IEEE Transactions on Visualization and Computer Graphics*  
480 19, 2596–2605. doi:10.1109/TVCG.2013.151
- 481 Fiedler, M. (1973). Algebraic Connectivity of Graphs. *Czechoslovak mathematical journal* 23, 298–305
- 482 Gansner, E. R. and North, S. C. (2000). An open graph visualization system and its applications to software  
483 engineering. *Software: practice and experience* 30, 1203–1233
- 484 Gao, H. and Ji, S. (2019). Graph U-Nets. In *International Conference on Machine Learning*. 2083–2092
- 485 Goller, C. and Kuchler, A. (1996). Learning task-dependent distributed representations by backpropagation  
486 through structure. In *ICNN*
- 487 Gori, M., Monfardini, G., and Scarselli, F. (2005). A new model for learning in graph domains. In *ICNN*
- 488 [Dataset] Hajij, M., Rosen, P., and Wang, B. (2018). Mapper on Graphs for Network Visualization
- 489 Huang, J., Li, Z., Li, N., Liu, S., and Li, G. (2019). AttPool: Towards Hierarchical Feature Representation  
490 in Graph Convolutional Networks via Attention Mechanism. In *Proceedings of the IEEE International*  
491 *Conference on Computer Vision*. 6480–6489
- 492 [Dataset] Kersting, K., Kriege, N. M., Morris, C., Mutzel, P., and Neumann, M. (2016). Benchmark Data  
493 Sets for Graph Kernels
- 494 Kingma, D. P. and Ba, J. (2014). Adam: A Method for Stochastic Optimization. *arXiv preprint*  
495 *arXiv:1412.6980*
- 496 Kipf, T. N. and Welling, M. (2016). Semi-Supervised Classification with Graph Convolutional Networks.  
497 *arXiv preprint arXiv:1609.02907*
- 498 Lee, J., Lee, I., and Kang, J. (2019). Self-Attention Graph Pooling. In *International Conference on*  
499 *Machine Learning*. 3734–3743
- 500 [Dataset] Leskovec, J. (2016). CS224W: Social and Information Network Analysis - Graph Clustering
- 501 Li, Y., Tarlow, D., Brockschmidt, M., and Zemel, R. (2015). Gated graph sequence neural networks.  
502 *arXiv:1511.05493*
- 503 Luzhnica, E., Day, B., and Lio, P. (2019). Clique pooling for graph classification. *arXiv preprint*  
504 *arXiv:1904.00374*
- 505 Ma, Y., Wang, S., Aggarwal, C. C., and Tang, J. (2019). Graph Convolutional Networks with EigenPooling.  
506 In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data*  
507 *Mining*. 723–731
- 508 McInnes, L., Healy, J., and Melville, J. (2018). UMAP: Uniform Manifold Approximation and Projection  
509 for Dimension Reduction. *ArXiv e-prints*
- 510 Nobre, C., Streit, M., Meyer, M., and Lex, A. (2019). The state of the art in visualizing multivariate  
511 networks. *Computer Graphics Forum (EuroVis '19)* 38, 807–832. doi:10.1111/cgf.13728
- 512 Page, L., Brin, S., Motwani, R., and Winograd, T. (1999). *The PageRank Citation Ranking: Bringing*  
513 *Order to the Web*. Tech. rep., Stanford InfoLab
- 514 Ranjan, E., Sanyal, S., and Talukdar, P. P. (2019). ASAP: Adaptive Structure Aware Pooling for Learning  
515 Hierarchical Graph Representations. *arXiv preprint arXiv:1911.07979*
- 516 Scarselli, F., Gori, M., Tsoi, A. C., Hagenbuchner, M., and Monfardini, G. (2009). Computational  
517 capabilities of graph neural networks. *IEEE Transactions on Neural Networks* 20, 81–102. doi:10.1109/  
518 TNN.2008.2005141
- 519 Sen, P., Namata, G., Bilgic, M., Getoor, L., Galligher, B., and Eliassi-Rad, T. (2008). Collective  
520 Classification in Network Data. *AI magazine* 29, 93–93
- 521 Shervashidze, N., Schweitzer, P., Van Leeuwen, E. J., Mehlhorn, K., and Borgwardt, K. M. (2011).  
522 Weisfeiler-lehman graph kernels. *Journal of Machine Learning Research* 12, 2539–2561

- 523 Singh, G., Mémoli, F., and Carlsson, G. E. (2007). Topological methods for the analysis of high dimensional  
524 data sets and 3d object recognition. *SPBG* 91, 100
- 525 Sperduti, A. (1994). Encoding labeled graphs by labeling raam. In *NIPS*
- 526 Tenenbaum, J. B., Silva, V. d., and Langford, J. C. (2000). A global geometric framework for nonlinear  
527 dimensionality reduction. *Science* 290, 2319–2323. doi:10.1126/science.290.5500.2319
- 528 van den Elzen, S. and van Wijk, J. J. (2014). Multivariate network exploration and presentation: From  
529 detail to overview via selections and aggregations. *IEEE Transactions on Visualization and Computer*  
530 *Graphics* 20, 2310–2319
- 531 van der Maaten, L. and Hinton, G. (2008). Visualizing Data using t-SNE. *Journal of Machine Learning*  
532 *Research* 9, 2579–2605
- 533 Veličković, P., Fedus, W., Hamilton, W. L., Liò, P., Bengio, Y., and Hjelm, R. D. (2018). Deep Graph  
534 Infomax. *arXiv preprint arXiv:1809.10341*
- 535 von Landesberger, T., Kuijper, A., Schreck, T., Kohlhammer, J., van Wijk, J., Fekete, J.-D., et al. (2011).  
536 Visual analysis of large graphs: State-of-the-art and future research challenges. *Computer Graphics*  
537 *Forum* 30, 1719–1749. doi:10.1111/j.1467-8659.2011.01898.x
- 538 Wattenberg, M. (2006). Visual exploration of multivariate graphs. In *Proceedings of the SIGCHI Conference*  
539 *on Human Factors in Computing Systems* (New York, NY, USA: Association for Computing Machinery),  
540 CHI '06, 811–819. doi:10.1145/1124772.1124891
- 541 Yang, Z., Cohen, W. W., and Salakhutdinov, R. (2016). Revisiting semi-supervised learning with graph  
542 embeddings. In *ICML*
- 543 Ying, Z., You, J., Morris, C., Ren, X., Hamilton, W., and Leskovec, J. (2018). Hierarchical Graph  
544 Representation Learning with Differentiable Pooling. In *Advances in Neural Information Processing*  
545 *Systems*. 4800–4810



## A MODEL ARCHITECTURE AND HYPERPARAMETERS

546 We additionally performed a hyperparameter search for DiffPool on hidden sizes 32, 64, 128 and for DGM,  
547 over the following sets of possible values:

- 548 • all datasets: cover sizes  $\{[40, 10], [20, 5]\}$ , interval overlap  $\{10\%, 25\%\}$ ;
- 549 • D&D: learning rate  $\{5e^{-4}, 1e^{-3}\}$ ;
- 550 • Proteins: learning rate  $\{2e^{-4}, 5e^{-4}, 1e^{-3}\}$ , cover sizes  $\{[24, 6], [16, 4], [12, 3], [8, 2]\}$ , hidden sizes  
551  $\{64, 128\}$ .