# Differentiable Graph Module (DGM) for Graph Convolutional Networks

Anees Kazi*, Luca Cosmo*, Seyed-Ahmad Ahmadi, Nassir Navab, and Michael M. Bronstein

**Abstract**—Graph deep learning has recently emerged as a powerful ML concept allowing to generalize successful deep neural architectures to non-Euclidean structured data. Such methods have shown promising results on a broad spectrum of applications ranging from social science, biomedicine, and particle physics to computer vision, graphics, and chemistry. One of the limitations of the majority of current graph neural network architectures is that they are often restricted to the transductive setting and rely on the assumption that the underlying graph is *known* and *fixed*. Often, this assumption is not true since the graph may be noisy, or partially and even completely unknown. In such cases, it would be helpful to infer the graph directly from the data, especially in inductive settings where some nodes were not present in the graph at training time. Furthermore, learning a graph may become an end in itself, as the inferred structure may provide complementary insights next to the downstream task. In this paper, we introduce Differentiable Graph Module (DGM), a learnable function that predicts edge probabilities in the graph which are optimal for the downstream task. DGM can be combined with convolutional graph neural network layers and trained in an end-to-end fashion. We provide an extensive evaluation of applications from the domains of healthcare (disease prediction), brain imaging (age prediction), computer graphics (3D point cloud segmentation), and computer vision (zero-shot learning). We show that our model provides a significant improvement over baselines both in transductive and inductive settings and achieves state-of-the-art results.

**Index Terms**—Graph convolution, Graph learning, Disease prediction

✦

## 1 INTRODUCTION

G EOMETRIC deep learning (GDL) is a novel emerging branch of deep learning attempting to generalize deep neural networks to non-Euclidean structured data such as graphs and manifolds [1], [2], [3]. Graphs, being general abstract descriptions of relation and interaction systems, are ubiquitous in different branches of science. Graph-based learning models have been successfully applied in social sciences [4], [5], computer vision and graphics [6], [7], [8], physical [9], [10], [11], [12], as well as medical and biological sciences [13], [14], [15], [16], [17], [18], [19].

Graph Neural Networks (GNNs) are a popular approach for learning on graphs. While dating back to at least [20], it is mainly the recent progress that has made GNNs a useful and popular tool. Today's wide variety of GNN architectures includes spectral [21] and spectral-like [22], [23], [24], [25] methods, local charting [7], and attention [26], [27], [28], [29]. Battaglia et al. [3] showed that most GNNs can be formulated in terms of message passing [11]. Recent efforts have been devoted to scalable [30], [31], [32], [33] and deep [34], [35], [36], [37] GNN architectures, and theoretical analyses of their expressive power [38], [39], [40], [41].

A notable drawback of most GNN architectures is the assumption that the underlying graph is *given* and *fixed*. Given this graph, convolution-like operations typically amount to modifying the node-wise features. Architectures like message passing neural networks [11] or primal-dual convolutions [29] also allow to update the edge features, but the graph *topology* is always kept the same. This often happens to be a limiting assumption. In many problems, the data can be assumed to have some underlying graph structure, however, the graph itself might not be explicitly given [42], a setting we refer to as *latent graph*. In medical and healthcare applications, for example in patient population models, the graph between patients may be noisy or partially and even completely unknown, and one is thus interested in inferring it from the data. Such latent graphs can capture the actual topology of structured data, aiding towards efficient processing and analysis while simultaneously learning the relationship between the behavioral influence among the nodes. In fact, sometimes the graph may be even more important than the downstream task, as it conveys some interpretability of the model and may provide a means for knowledge discovery, e.g. in form of patient relations in a population model. Being able to learn a latent graph is especially important in case of inductive settings, where some nodes might be present in the graph at testing but not training time.

Graph topology inference is a long-standing problem and has recently been addressed using signal processing techniques [43], [44]. We provide a detailed literature survey for different types of graph learning techniques in the next sections. In our recent work [45], we proposed a graph learning model targeted towards node classification as the primary task. The main methodological contribution is the

• A. Kazi and N. Navab are with the department of Computer Aided Medical Procedures and Augmented Reality at Technical University of Munich. N. Navab is also with the Whiting School of Engineering, Johns Hopkins University, Baltimore, USA;
• L. Cosmo is with the Ca' Foscari University of Venice, Italy, and with the Faculty of Informatics, University of Lugano, Switzerland;
• S.-A. Ahmadi is with NVIDIA, Munich, Germany;
• M. M. Bronstein is with the department of Computer Science at the University of Oxford, UK and with Twitter, UK

* A. Kazi and L. Cosmo contributed equally to this work. Part of this work was done during A. Kazi's research visit at Imperial College London. Corresponding author: A. Kazi
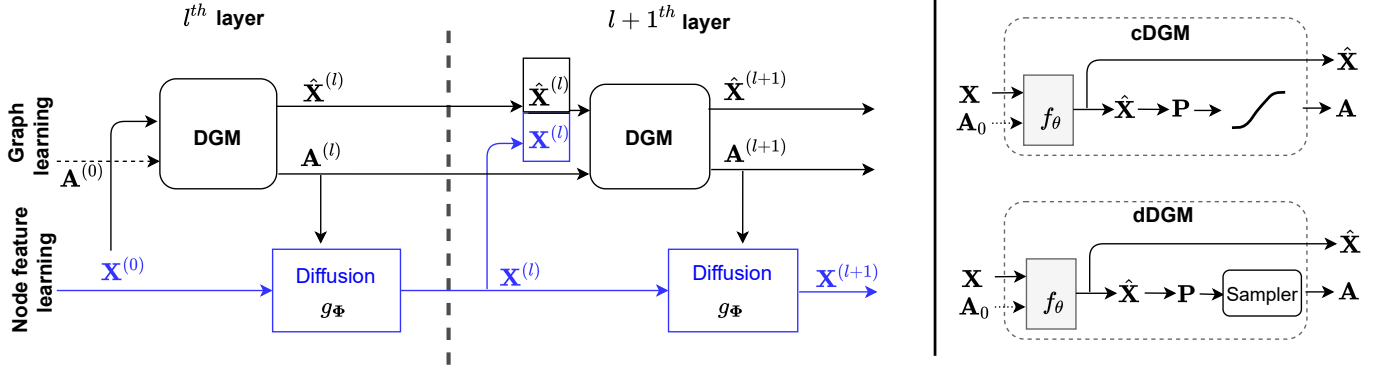
Fig. 1: *Left:* Two-layered architecture including Differentiable Graph Module (DGM) that learns the graph, and Diffusion Module that uses the graph convolutional filters. *Right:* Details of DGM in its two variants, cDGM and dDGM.

latent-graph learning block, which is capable of learning the probabilistic graph in form of a weighted adjacency matrix for the optimal classification outcome. The main limitation of this approach [45] was to model the graph in quasi-probabilistic and thus fully connected manner, since the model lacked the capability to exploit a possible sparsity of the graph. In this work, we propose a more general Differentiable Graph Module (DGM) for latent graph prediction. We propose both a continuous, generalizing [45], and a discrete sampling strategy that overcomes the dense graph limitations and allows addressing bigger graphs. Further, our model explicitly models graph sparsity, resulting in a faster computation and lower memory consumption.

*Main contributions*

- We propose an end-to-end pipeline for simultaneous learning of the downstream task, e.g node classification, along with an optimal underlying latent graph.
- We propose DGM with two sampling strategies, a continuous version that improves and generalizes our previously proposed method [45], and a novel discrete version that can efficiently learn larger graphs.
- We show proof of concept experiments, on three benchmark graph datasets for node and graph classification.
- We provide extensive ablation studies for both our model and discrete sampling strategy. Our evaluation demonstrates applications from several domains, i.e. healthcare and brain imaging (disease and age prediction), computer graphics (3D point cloud segmentation), and computer vision (zero-shot learning). Our model shows significant improvement over baselines and achieves state-of-the-art results.

The Pytorch implementation of DGM can be found at the following Github repository: https://github.com/lcosmo/DGM_pytorch.

## 2 STATE OF THE ART

Recently, graph learning has come to focus due to its applicability in various domains [46]. In this section, we provide a detailed state of the art on graph learning, subdividing it into three main categories.

**Attention-based graph learning**: In the machine learning literature, several models dealing with latent graphs have recently been proposed. In this category, we include works that deal with attention to the nodes and edges, since edge-based attention can be seen as a latent structure learning process. Kipf et al. [47] used a variational autoencoder, in which the latent code represents the interaction graph underlying a physical system, and the reconstruction is based on graph neural networks. Further, graph attention networks (GAT) [26] proposed an additional attention mechanism that operates on the edges of a fixed graph to learn the importance of each neighbor. Gated Attention Network (GAAN) [48] proposes a self-attention mechanism that computes an additional attention score for each attention head. RNN-based attention models have been proposed in [49], [50], [51]. The common drawback of these methods is that they need an input graph. Further, attention is computed directly on edge features, requiring the graph to be sparse in order to make the approach computationally feasible, especially for bigger graphs.

**Models dealing with dynamic graph**: These methods do not assume a fixed graph a-priori, but build them dynamically during training. In their seminal work, Wang et al. [8] proposed Dynamic Graph CNNs (DGCNN) for the analysis of point clouds, where a k-nearest neighbor (kNN) graph is constructed on the fly in the feature space of the neural network. This idea is further extended in PGC-DGCNN [52], which aims at increasing the contributions of distant neighbors using shortest paths connections. In both of these methods [8], [52], the latent space in which the graph is constructed is not optimized for the sought latent structure. Instead, they sample local neighborhoods of feature representations that are used for the downstream task. Importantly, the kNN operation used to build the graph is not directly differentiable.

**Models dealing with graph learning**: These methods focus on learning the graph directly. We subdivide this category into two groups.

The first group operates in the spectral domain. Zhan et al. [53] proposed to learn the latent structure by constructing multiple Laplacians and combining them with learnable weights. Similarly, Li et al. [54] proposed a spectral graph convolutional method, in which a residual Laplacian computed on the feature output from each layer and the input

Laplacian are updated after each layer. Huang et al. [55] proposed another version of spectral filters that parametrize the Laplacian instead of the coefficient of the filter. All the spectral approach based methods need an initial graph and suffer from the limitation of transductive models, i.e. that the validation and/or test set need to be embedded into the training graph.

The second group operates in the spatial domain. Jiang et al. [56] proposed a model where graph learning and graph convolution are integrated in a unified network architecture. Franceschi et al. [57] formulated graph learning as a bi-level optimization problem, by modeling the graph as a hyper-parameter and optimizing it with a separate loss. A different strategy was adopted by Yang et al. [58], who proposed a method for giving scores to the nodes during sub sampling, but without learning the graph. Similarly, [59] proposed DiffPool, a differentiable graph pooling module to generate hierarchical representations of graphs. The method learns a differentiable soft cluster assignment for nodes at each layer mapping nodes to a set of clusters. Each cluster then becomes a node for the next GNN layer. Norcliffe et al. [60] specifically tailored their method for visual question answering. They construct a joint embedding of word and image features and apply top-k sampling on the embedding distances. Their method is a two-step graph learning pipeline, and requires a heavily customized optimization techniques. Yu et al. [61] proposed to learn a fully connected graph in an end-to-end pipeline, sparsifying the graph according to a threshold in the graph embedding space. While effective, this strategy does not provide any means to control the graph topology, and importantly, it does not prevent the graph structure from degenerating e.g. into many disconnected sub-graphs. More recently, Cosmo et al. [62] proposed a graph convolution operator that directly learns both the connectivity and features of a set of graphs (i.e. structural masks) to be used as local filters on the input graph.

Compared to the methods above, we propose a technique for latent graph discovery which learns the graph end-to-end, which is inductive, scalable to large numbers of nodes, and which gives the network designer a means to control the graph topology.

## 3 BACKGROUND

Given a set of $N$ data points of dimension $d$, denoted $\mathbf{X} \in \mathbb{R}^{N \times d}$, a common problem in machine learning is to produce a representation that is aware of the underlying structure of the data. Such structure can be represented as a (weighted) graph $\mathcal{G} = (\mathcal{V}, \mathbf{A})$ where $\mathcal{V} = \{1, \dots, n\}$ is the vertex set and $\mathbf{A} = (a_{ij})$ is a (weighted) adjacency matrix. We use $\mathbf{A}$ to define the edges of the graph $\mathcal{E} = \{(i, j) : a_{ij} > 0, i, j \in \mathcal{V}\}$; the edge weight $a_{ij} \geq 0$ represents the affinity between points $\mathbf{x}_i$ and $\mathbf{x}_j$.

### 3.0.1 Graph neural networks

Assuming this structure is provided together with the data, we have a node-attributed graph $\mathcal{G} = (\mathcal{V}, \mathbf{A}, \mathbf{X})$, on which a graph neural network (GNN) can be applied. GNN attempts to find an *embedding* $\mathbf{Z} = g_{\Theta}(\mathbf{X}, \mathbf{A})$ by doing message passing [3], [11]

$$\mathbf{z}_i = \sum_{j \in \mathcal{N}_i} h_{\Theta}(\mathbf{x}_i, \mathbf{x}_j, a_{ij}) \tag{1}$$

in a local neighborhood $\mathcal{N}_i = \{j : (i, j) \in \mathcal{E}\}$ of the node. Here $h_{\Theta}$ denotes a learnable function shared across nodes, whose parameters $\Theta$ are chosen to minimize a downstream loss. Equation (1) is also referred to as *edge convolution* (EC) [8] due to its generalization of the classical convolution operation on grids. A particular case of (1) with node-wise linear transformation $h_{\Theta} = a_{ij}\Theta\mathbf{x}_j$ by matrix $\Theta$, is called *graph convolution* (GC) [22], [23]. The node embeddings can be used for node-wise classification, or pooled for graph-wise classification tasks.

### 3.0.2 Latent graphs

We are interested in the setting when the underlying graph is *unknown*, and thus needs to be learned. Learning the graph serves two purposes: First, it is used to represent the structure of the data. Second, it is used as the support for graph-based convolutions to obtain the embeddings of the data points.

The main obstacle for including the graph construction as a part of the deep learning pipeline is that it is a discrete structure and as such non-differentiable. To avoid this problem, in DGCNN [8], the graph convolutional filters and the layer activations are optimised towards the downstream classification and segmentation tasks. The graph, however, is constructed ad-hoc as a kNN graph on the activations after each layer, without a dedicated loss. As such, the graph is built dynamically but not learned, and the underlying latent graph of the domain is not recovered. Our method, described in the next section, aims at addressing this issue.

## 4 METHOD

### 4.1 Architecture

We propose a general technique for learning the graph based on the output features of each layer, along with the optimization of the network parameters during the training. Our architecture comprises two main blocks, the **Differentiable Graph Module (DGM)** and the **Diffusion Module**. The interaction between these two modules is shown in Figure 1 and further detailed in Algorithm 1.

### 4.1.1 Differentiable Graph Module:

The DGM is tasked with building the (weighted) graph representing the input space. It takes the feature matrix $\mathbf{X} \in \mathbb{R}^{N \times d}$ as input and yields a graph $\mathcal{G}$ as output. DGM can also take an initial graph $\mathcal{G}_0$ as input, in case an expert or domain knowledge exists allowing to define the initial edge weights between the nodes. Importantly, such an initial graph is optional, making DGM usable also in cases where domain knowledge is absent or where it is not desirable to impose a graph prior. Since the node set $\mathcal{V}$ is fixed, the two graphs can be represented by their adjacency matrices, $\mathbf{A}_0$ and $\mathbf{A}$.

Figure 2 depicts the inner functioning of the Differentiable Graph Module. Input features $\mathbf{X} \in \mathbb{R}^{N \times d}$ are first transformed into *auxiliary features* $\hat{\mathbf{X}} = f_{\Theta}(\mathbf{X}) \in \mathbb{R}^{N \times \hat{d}}$ by means of a parametric function $f_{\Theta}$, typically reducing the input dimension ($\hat{d} \ll d$). If the initial graph $\mathcal{G}_0$ is provided,
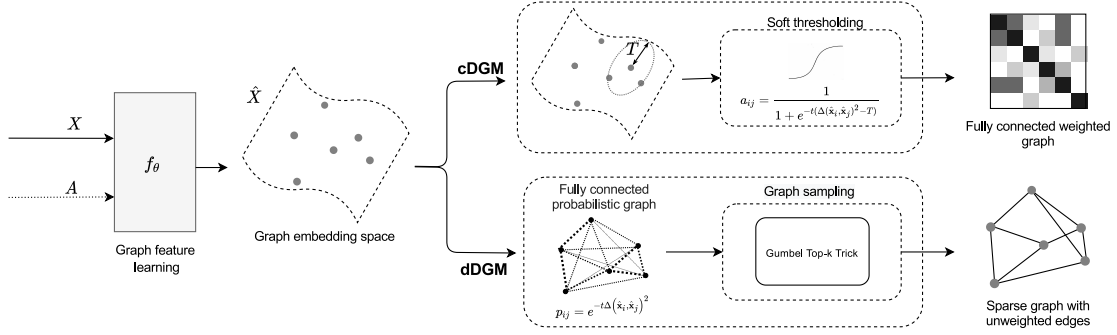
Fig. 2: Differentiable Graph Module (DGM) described in detail. After learning the lower dimensional embedding space the diagram shows proposed sampling techniques. Upper and lower part shows the continuous (cDGM) and discrete (dDGM) sampling variants.

we can use $f_\Theta$ of the general form (1), where new features $\hat{\mathbf{X}}$ are computed by edge- or graph-convolutions on $\mathcal{G}_0$. Otherwise, $f_\Theta$ is applied to each node feature independently, acting row-wise on the matrix $\mathbf{X}$.

Second, the auxiliary features $\hat{\mathbf{X}}$ are used for graph construction. We define the edge probabilities as:

$$p_{ij}(\mathbf{X}; \boldsymbol{\Theta}, t) = e^{-t\Delta(\hat{\mathbf{x}}_i, \hat{\mathbf{x}}_j)^2} = e^{-t\Delta(f_\Theta(\mathbf{x}_i), f_\Theta(\mathbf{x}_j))^2}, \quad (2)$$

where $t$ is a learnable parameter and $\Delta(\cdot, \cdot)$ is the distance between two nodes in the graph embedding space. In the experimental session, we investigate the use of two different metric definitions, namely the Euclidean and the Hyperbolic metrics [63], [64].

**Continuous sampling:** A straightforward way to derive a graph $\mathcal{G}$ is to transform the probability matrix $\mathbf{P}(\mathbf{X}; \boldsymbol{\Theta}, t)$ into a weighted adjacency matrix using the sigmoid function $a_{ij} = 1/(1 + p_{ij}e^{tT})$. Here, as depicted in figure 2, $T$ can be interpreted as a threshold on the squared distances $\Delta(\hat{\mathbf{x}}_i - \hat{\mathbf{x}}_j)^2$ defined on the embedding space. Indeed, elements $a_{ij}$ can be rewritten as $a_{ij} = \frac{1}{1+e^{-t(\Delta(\hat{\mathbf{x}}_i, \hat{\mathbf{x}}_j)^2 - T)}}$. In this way, the graph is represented by the adjacency matrix $\mathbf{A}(\mathbf{X}; \boldsymbol{\Theta}, t, T)$ parametrized through $\boldsymbol{\Theta}, t$ and the additional parameter $T$, and is differentiable w.r.t. these parameters.

This sampling strategy generalizes the method proposed in [45]. We refer on this paper to this variant as *continuous DGM (cDGM)*.

**Discrete sampling:** One of the shortcomings of cDGM is that it can produce a dense adjacency matrix, i.e. a fully connected graph in which many edges have near-zero weight. As an efficient alternative, we propose to construct a sparse $k$-degree graph by using the *Gumbel-Top-k trick* [65] to sample edges from the probability $\mathbf{P}(\mathbf{X}; \boldsymbol{\Theta}, t)$. Such sampling can be regarded as a stochastic relaxation of the kNN rule. For each node $i$, we extract $k$ edges $(i, j_{i,1}), \dots, (i, j_{i,k})$ as the first $k$ elements of

$$\text{argsort}(\log(\mathbf{p}_i) - \log(-\log(\mathbf{q}))), \quad (3)$$

where $\mathbf{q} \in \mathbb{R}^N$ is uniform i.i.d. in the interval $[0, 1]$. Samples extracted this way follow the categorical distribution $p_{ij}/\sum_r p_{ir}$ [65]. We define the edge set of the sparse graph $\mathcal{G}$ constructed this way as

$$\mathcal{E}(\mathbf{X}; \boldsymbol{\Theta}, t) = \{(i, j_{i,1}), \dots, (i, j_{i,k}) : i = 1, \dots, N\} \quad (4)$$

and represent it by the unweighted adjacency matrix $\mathbf{A}(\mathbf{X}; \boldsymbol{\Theta}, t)$. The key advantage is that this matrix is sparse, which results in a lower computational and memory complexity of the diffusion operation. We refer to this variant of our architecture as *discrete DGM (dDGM)*.

Note that, since the dDGM graph sampling scheme is stochastic, the prediction of the network at inference time is not deterministic. We can actually take advantage of this, and implement a consensus scheme. In our experiments, we run the classification for 8 times and select the maximum of the cumulative soft predictions as the predicted class.

### 4.1.2 Diffusion Module:

This module takes the graph affinity $\mathbf{A}$ produced by the DGM and the features $\mathbf{X}^{(l)}$ as inputs, and yields a new set of features $\mathbf{X}^{(l+1)} = g_\Phi(\mathbf{X}^{(l)})$ as output as shown in Fig 1. Here, $g_\Phi$ represents a general function of the form (1); in our experiments, it is either edge- or graph-convolution on $\mathbf{A}$.

## 4.2 Classification model (Combined model):

We use a multi-layer network, with layers numbered as $l = 1, \dots, L$. Each layer comprises a **DGM** and **Diffusion Module**, as shown in Figure 1. The $l$th layer of the architecture produces the output as:

$$\hat{\mathbf{X}}^{(l+1)} = f_\Theta^{(l+1)}([\mathbf{X}^{(l)} \mid \hat{\mathbf{X}}^{(l)}], \mathbf{A}^{(l)}) \quad (5)$$

$$\mathbf{A}^{(l+1)} \sim \mathbf{P}^{(l)}(\hat{\mathbf{X}}^{(l+1)}) \quad (6)$$

$$\mathbf{X}^{(l+1)} = g_\Phi(\mathbf{A}^{(l+1)}, \mathbf{X}^{(l)}) \quad (7)$$

We assume $\mathbf{X}^{(0)} = \mathbf{X}$ and unless some initial knowledge of the structure of the data is available, $\mathbf{A}^{(0)} = \mathbf{I}$ (i.e., the initial graph is $\mathcal{G}^{(0)} = (\mathcal{V}, \emptyset)$) and $f_\Theta^{(0)}$ is a node-wise function (MLP). In case of classification task, the final node features $\mathbf{X}^{(L)}$ of the last layer $L$ can then be given as input to a MLP to obtain the final node predictions. Such an end-to-end model can be defined as $\hat{Y} = \mathbb{M}(X^0, A^0)$ where, $\hat{Y}$ is predicted label vector.

Note that DGCNN can be obtained as a particular setting of our model with $f_\Theta = \text{id}$ in the **DGM** module and using edge convolution in the **Diffusion** module. Here id stands for the identity mapping.

### 4.2.1  Loss function and Training

To train end-to-end our cDGM model we optimize directly over the cross-entropy loss of the final classification task. However, the sampling scheme we adopt in dDGM does not allow the gradient of the downstream classification loss function to flow through the graph prediction branch of our network, as it involves only graph features $\hat{\mathbf{X}}$.

To allow its optimization, we create a compound loss that rewards edges involved in a correct classification and penalizes edges that led to misclassification.

Let $\mathbf{y}$ denote the vector of node-wise labels predicted by our model at step $(t)$ and $\tilde{\mathbf{y}}$ the groundtruth labels.

We define the reward function $\delta(y_i, \tilde{y}_i) = \mathbb{E}((a_i)) - a_i$ as the difference between the average accuracy of the $i$th sample and the current success value $a_i = 1$ if $y_i = \tilde{y}_i$ and 0 otherwise. We then derive the graph loss as:

$$L_{\text{graph}}(\boldsymbol{\Theta}^{(1)}, \dots, \boldsymbol{\Theta}^{(L)}) = \sum_{\substack{i=1\dots N \\ l=1\dots L \\ j:(i,j)\in\mathcal{E}^{(l)}}} \delta(y_i, \tilde{y}_i) \log p_{ij}^{(l)}(\boldsymbol{\Theta}^{(l)})$$

(8)

whose gradient approximates the gradient of the expectation $\mathbb{E}_{(\mathcal{G}^{(1)},\dots,\mathcal{G}^{(L)})\sim(\mathbf{P}(\boldsymbol{\Theta}^{(1)}),\dots,\mathbf{P}(\boldsymbol{\Theta}^{(L)}))} \sum_i \delta(y_i, \tilde{y}_i)$ with respect to the parameters of the graphs in all the layers.

We estimate $\mathbb{E}(\mathbf{a}_i)$ with two different strategies depending on the application. When the node set remains fixed for each sample and during training, we adopt an exponential moving average strategy on the accuracy values during the training process. In this case, $\mathbb{E}(a_i)^{(t+1)} = \alpha\mathbb{E}(a_i)^{(t)} + (1 - \alpha)a_i$, with $\alpha = 0.9$ in all our experiments. If the node set changes during training or between different samples, we estimate $\mathbb{E}(\mathbf{a}_i)$ individually for each training step $(t)$ evaluating the model multiple times on different sampled graphs.

## 5  EXPERIMENTS AND RESULTS

We start the experimental session by showing ablation study results on popular benchmark datasets from [66] (CiteSeer, Cora, Pubmed). We then validate our approach against state-of-the-art methods on the domains of healthcare applications (disease and age prediction) and show experiments on computer graphics (3D point cloud segmentation), and computer vision (zero-shot learning).

### 5.1  Benchmark graph based datasets

In this section we show the ablation study with respect to different hyper-parameters and architectural choices of our method to better analyze and evaluate the discrete sampling strategy proposed in this work (dDGM). We make use of three popular citation graph datasets [66] for this task, namely Cora, PuMed and CiteSeer. Each of these datasets consists of a single graph in which nodes correspond to documents and edges to citation links. The goal is to predict the category of each document transductively, i.e. the category of some documents is not known during training. Further details on these datasets are given in Table 1.

Since in these datasets the input graph is given, we take advantage of it by using a Graph Convolutional Layer (GCN) [23] as the graph embedding function $f_\Theta$. Note that, even if the DGM module uses the input graph, the Diffusion

---

**Algorithm 1:** Algorithm for the full model that produces $\hat{\mathbf{Y}}$ i.e, the predicted labels which are then used by the loss function for training the model. With $(a_{ij})$ we refer to the matrix $\mathbf{A}$ composed by elements $a_{ij}$, while $\mathbf{p_i}$ refers to the $i$th row of $P$. $f_\Theta$ and $g_\Phi$ are generic function with learnable parameters $\Theta$ and $\Phi$ respectively, MLP is the Multi Layer Perceptron used for classification with parameters $\Psi$.

**Data:** $\mathbf{X}, \mathbf{A}$
**Result:** $\hat{\mathbf{Y}}$

1 **begin**
2    $\mathbf{X}^{(0)} \leftarrow \mathbf{X}$
3    $\mathbf{A}^{(0)} \leftarrow \mathbf{A}$
4    $\hat{\mathbf{X}}^{(0)} \leftarrow \emptyset$
5    **for** $x \in 0 \dots L-1$ **do**
6      $\hat{\mathbf{X}}^{(l+1)}, \mathbf{A}^{(l+1)} \leftarrow$
     DGM$([\mathbf{X}^{(l)} \,|\, \hat{\mathbf{X}}^{(l)}], \mathbf{A}^{(l)}, f_\Theta^{(l)}, t^{(l)}, T^{(l)})$
7      $\mathbf{X}^{(l+1)} \leftarrow g_\Phi^{(l)}(\mathbf{A}^{(l+1)}, \mathbf{X}^{(l)})$
8    $\hat{\mathbf{Y}} \leftarrow MLP_\Psi(\mathbf{X}^{(L)})$

9

10 **Function** DGM $(\mathbf{X}, \mathbf{A}, f_\Theta, t, T)$:
11    $\hat{\mathbf{X}} \leftarrow f_\Theta(\mathbf{X})$
12    $p_{ij} \leftarrow e^{-t\Delta(\hat{x}_i, \hat{x}_j)^2}$
13    **Switch** *sampling* :
14      **case** *continuous* **do**
15        **for** $i, j \in 1 \dots N$ **do**
16          $a_{ij} \leftarrow 1/(1 + p_{ij}e^{tT})$
17      **case** *discrete* **do**
18        **for** $i \in 1 \dots N$ **do**
19          $\mathbf{q} \sim U(0,1)$
20          $\mathbf{j}_{\{\mathbf{k}\}} = \text{argtopK}(\log(\mathbf{p_i}) - \log(-\log(\mathbf{q})))$
21          $a_{ij} = \begin{cases} 1, & j \in \mathbf{j}_{\{\mathbf{k}\}} \\ 0, & \text{otherwise} \end{cases}$

22    **return** $\hat{\mathbf{X}}, (a_{ij})$

---

module will have access only to the predicted graph (in Figure 1 the input and the predicted graphs are represented by $\mathbf{A}^{(0)}$ and $\mathbf{A}^{(l)}$ respectively).

**Architecture details:** The base architecture for this ablation study is composed by one DGM layer where the graph is given as input to the associated diffusion layer. Its output is passed through a final MLP of size $(8, 8, c)$, which classifies each node into one of $c$ classes. The DGM layer function $f$ is composed by two GCN layers with output feature dimensions of 16 and $d$ respectively. We set $d = 4$ as the graph embedding space dimension. The diffusion function $g$ is composed of three GCN layers with output sizes 32, 16 and 8, respectively. We sample the kNN graph with a number of $k = 5$ neighbors. Before being processed by the DGM layer, input features are transformed by a perceptron (a linear transformation followed by a ReLU activation) of output size equal to 32. We train each model 10 times for 100 epochs each, using Adam optimizer with

TABLE 1: Dataset description of Cora, Citeceer and Pubmed

| Dataset | nodes | edges | features | classes |
|---------|-------|-------|----------|---------|
| Cora | 2,708 | 5,429 | 1,433 | 7 |
| CiteSeer | 3,327 | 4,372 | 3,703 | 6 |
| Pubmed | 19,717 | 44,338 | 500 | 3 |

TABLE 2: Performance of our method with different numbers of $k$ nearest neighbors when sampling the graph. We report the mean and standard deviation (in parenthesis) of the accuracy on 10 runs.

| $k$ | Cora | PubMed | CiteSeer |
|-----|------|--------|----------|
| 1 | 72.80 (2.71e+0) | **88.60 (1.16e+0)** | 67.00 (3.55e+0) |
| 3 | 83.60 (9.91e-1) | **88.60 (8.85e-1)** | 73.20 (1.49e+0) |
| 5 | **84.60 (8.52e-1)** | 87.60 (7.51e-1) | **74.80 (9.24e-1)** |
| 10 | 81.80 (1.94e+0) | 85.80 (1.71e+0) | 72.00 (1.06e+0) |
| 20 | 70.60 (4.05e+0) | 81.60 (1.38e+0) | 64.80 (3.46e+0) |

a learning-rate of 1e-2, and report the mean and standard deviation of accuracy (in percent).

### 5.1.1 Latent graph sparsity

In this experiment we consider the sparsity of the graph induced by the choice of the parameter $k$ while building the kNN graph. The motivation behind this experiment is to evaluate the graph learning module at different sparsity levels and find the generic range of sparsity in terms of 'k'. As shown in Table 2 the value of this parameter can have a significant impact on the performance of the network. In all the considered datasets a good range for $k$ is between 3 and 5. This means that the relations within the latent network can be captured by a highly sparse graph, which in turn leads to a high efficiency of our discrete sampling strategy.

### 5.1.2 Graph embedding space geometry

Recent literature in neural networks [67] suggest that the geometry of the latent space in which data samples are embedded plays an important role in capturing the underlying relations. A $d$-dimensional hyperbolic space is a homogeneous space with constant negative curvature. For instance, hyperbolic spaces have been shown to better capture hierarchical structures compared to the standard Euclidean space [68]. With this experiment, we aim at investigating whether or not embedding the graph nodes in a hyperbolic space is beneficial for node classification.

We model it using the Poincaré ball [69], where the distance between two points inside the unitary open ball, $x, y \in \mathcal{B}^d$, is computed as:

$$d(\boldsymbol{x}, \boldsymbol{y}) = \cosh^{-1}\left(1 + 2\frac{\|\boldsymbol{x} - \boldsymbol{y}\|_2^2}{\left(1 - \|\boldsymbol{x}\|_2^2\right)\left(1 - \|\boldsymbol{y}\|_2^2\right)}\right), \quad (9)$$

where $\|.\|_2^2$ is the standard euclidean squared norm. To constrain the graph embedding function to predict values within the unit ball we apply a non linearity which clips the norm of the latent vectors to $1 - \epsilon$.

In case of the Euclidean latent space geometry, we simply define the distance as $d(x, y) = \|x - y\|_2$.

Table 3 shows the performance of models that learned a graph under the assumption of a hyperbolic vs Euclidean latent space. For dDGM, as can be seen, the model with Euclidean embedding performs consistently better on the

TABLE 3: Performance of our method with euclidean and hyperbolic space geometries. We report the mean and standard deviation (within parenthesis) of the accuracy on 10 runs.

| | Cora | PubMed | CiteSeer |
|---|------|--------|----------|
| Euclidean | **84.60 (8.52e-01)** | **87.60 (7.51e-01)** | **74.80 (9.24e-01)** |
| Hyperbolic | 84.40 (1.70e+00) | 86.60 (9.52e-01) | 74.60 (7.63e-01) |

TABLE 4: Performance of our method with different graph embedding space dimensions $d$. We report the mean and standard deviation (within parenthesis) of the accuracy on 10 runs.

| $d$ | Cora | PubMed | CiteSeer |
|-----|------|--------|----------|
| 2 | 80.20 (5.22e-01) | 87.80 (5.31e-01) | 71.40 (1.08e+00) |
| 4 | **84.60 (8.52e-01)** | 87.60 (7.51e-01) | **74.80 (9.24e-01)** |
| 6 | **84.60 (5.12e-01)** | 88.20 (6.01e-01) | 74.00 (6.76e-01) |
| 8 | 84.40 (1.56e+00) | **88.40 (7.74e-01)** | 73.20 (1.20e+00) |

three investigated datasets. Further, the standard deviation shows that Euclidean space is more stable compared to the hyperbolic, indicating that the considered datasets do not draw any advantage from a hierarchical representation of node relations.

### 5.1.3 Graph embedding space dimension

The graph embedding dimension directly determines the subset of graphs that the DGM module is able to represent. A higher dimension allows to represent more complex relations between the nodes, while a smaller embedding space has a regularization effect that prevents overfitting. The accuracy for different graph embedding space dimensions is shown in table 4. For the smaller datasets, Cora and CiteSeer, the best accuracy is obtained with a dimension of 4, while PubMed requires a higher embedding dimension of 8. The optimal value for the embedding space dimension thus depends on the complexity of the dataset.

### 5.1.4 Number of DGM Layers

In this experiment, we investigate the impact of using the proposed DGM module rather than just performing graph convolutions on the input graph, and whether there is any benefit in stacking multiple DGM layers, i.e learning a dynamic graph. In table 5, we report the accuracy with respect to the number of DGM layers. Here, a value of "0" denotes that we do not use DGM and perform the convolution directly on the input graph, corresponding to a conventional GCN based model. In practice, when using $L > 0$ DGM modules, the diffusion module of the intermediate layers is composed by just one convolution layer, while the final diffusion layer will have $4 - L$ convolution layers. This keeps the total number of convolution layers on the diffusion part equal to three.

Our results indicate that using a single graph, i.e a single DGM, is sufficient for increasing the classification accuracy w.r.t. the input graph in all datasets. Only PubMed benefits from a second DGM layer, again implying that it exhibits more complex relations between its nodes.

We remark that our method does not require an a-priori graph as input. If given, however, our results indicate that the graph embedding function $f$ can potentially exploit

TABLE 5: Performance of our method with a different number of DGM layers. Zero layers means that DGM is not used. We report the mean and standard deviation (within parenthesis) of the accuracy on 10 runs.

|   | Cora | PubMed | CiteSeer |
|---|---|---|---|
| 0 | 80.90 (1.19e+00) | 86.40 (1.37e+00) | 71.90 (1.50e+00) |
| 1 | **84.70 (1.40e+00)** | 87.60 (7.94e-01) | **74.70 (1.18e+00)** |
| 2 | 83.90 (1.77e+00) | **88.50 (5.59e-01)** | 74.50 (1.36e+00) |
| 3 | 82.40 (1.26e+00). | 88.40 (5.37e-01) | 72.90 (1.77e+00) |

TABLE 6: Performance of our method with different graph embedding functions $f$. We report the mean and standard deviation (within parenthesis) of the accuracy on 10 runs.

|   | Cora | PubMed | CiteSeer |
|---|---|---|---|
| gcn | 84.70 (1.40e-01) | **87.60 (7.94e-01)** | 74.70 (1.18e+00) |
| gat | **84.90 (1.47e+00)** | 86.80 (9.14e-01) | **74.80 (1.17e+00)** |
| mlp | 63.30 (2.86e+00) | 87.50 (8.00e-01) | 64.90 (3.18e+00) |
| id | 61.80 (4.11e+00) | 84.20 (1.65e-01) | 59.90 (3.56e+00) |

complementary information in its embedding, e.g. domain knowledge given by experts.

### 5.1.5 Graph embedding function $f$

We further show analysis of different graph embedding functions $f$ within DGM module. Table 6 shows the classification accuracy and its standard deviation using four different embedding functions $f$ on the 3 datasets. GCN [23] and GAT [26] functions make use of both node features and the input graph connectivity ($A_0$ in Figure 1). On the other hand, MLP applies nodewise a multilayer perceptron directly on the input features, discarding the input graph connectivity. To keep the number of parameters comparable, the MLP is composed by a linear layer of the same input and output dimensions of the convolution operations followed by a ReLU activation. Finally, id represents the identity mapping, which builds a knn graph directly on the input feature space without any transformation.

Results show that using the input connectivity graph in the DGM module is essential to obtain reasonable performance, especially for Cora and CiteSeer datasets. Indeed, discarding the input graph connectivity by building *de novo* the graph based on input features similarities (knn), or not using any connectivity information at all (MLP), leads to a huge drop ($10 - 20\%$) in accuracy. This means that, in these datasets, the graph structure itself conceive important information about the class each node belongs to.

### 5.1.6 Diffusion function $g$

Here, we check the behaviour of the end-to-end pipeline with respect to different choices for the diffusion function $g$. We use three different convolutional layers with different characteristics: GCN [23] performs a simple diffusion operation over the node features; GAT [26] further learns a per-edge weight based on the attention mechanism; finally, EdgeConv operates over edge features to compute the output node feature. As can be seen in table 7a, GCN is more stable and generally obtains the best performance, while using EdgeConv results in a comparably large drop of accuracy, especially on smaller datasets. A potential explanation is that the graph predicted by DGM is already

TABLE 7: (a) Performance of our method with different diffusion functions $g$ on the graph predicted by DGM. (b) Same diffusion functions applied directly on the input graph (without DGM). We report the mean and standard deviation (within parenthesis) of the accuracy on 10 runs.

(a) with DGM

|   | Cora | PubMed | CiteSeer |
|---|---|---|---|
| gcn | **84.60 (8.52e-01)** | 87.60 (7.51e-01) | **74.80 (9.24e-01)** |
| gat | 81.30 (2.71e+00) | **87.80 (8.41e-01)** | 74.00 (1.32e+00) |
| edgeconv | 71.60 (4.20e+00) | 87.40 (5.74e-01) | 52.20 (2.39e+00) |

(b) without DGM

|   | Cora | PubMed | CiteSeer |
|---|---|---|---|
| gcn | 80.90 (1.19e+00) | **86.70 (1.37e+00)** | **71.90 (1.50e+00)** |
| gat | **81.10 (1.79e+00)** | 84.60 (1.26e+00) | 71.20 (1.66e+00) |
| edgeconv | 67.00 (3.94e+00) | 85.00 (1.54e+00) | 49.00 (4.52e+00) |

optimal for the task at hand, and the edge weighting performed by GAT and intrinsically by EdgeConv does not give any contribution. Moreover, the graph sampled by DGM is different every time it gets sampled, and the resulting changes in the loss landscape may hinder convergence. In consequence, learning edge-related properties becomes a more challenging task, which can lead to a lower accuracy of GAT and EdgeConv.

Further, we investigate the advantage given by using our DGM with respect to standard graph convolution neural networks. Table 7a shows the experiments with DGM i.e. performing the diffusion on the graph being learned, while table 7b without DGM, i.e. performing the diffusion directly on the graph given as input. As can be seen, irrespective of the specific diffusion function $g$, in most of the cases the model with DGM performs better than using just the input graph. Specifically, edgeconv shows worst performance as a diffusion model. The model with DGM and edgeconv performs up to $3.4\%$ better than without DGM. This experiment clearly shows the ability of our DGM to predict a graph more suitable for the task at hand, even when a predefined structure is already given as input to the classification model. Note that, both architectures with and without DGM have exactly the same number of parameters involved in the diffusion step and final classification MLP ($\approx 50K$ for the base architecture in Cora), while DGM needs to optimize some few extra parameters ($\leq 150$ in the considered experiment) dedicated to the graph prediction.

### 5.1.7 Time and memory requirements.

The key advantage of the proposed kNN discrete graph sampling strategy (dDGM) is the possibility to obtain a sparse graph, which makes graph convolutional operations more efficient in computational complexity and memory requirements.

In Figure 3, we report both the time required for each training iteration and the allocated memory for three versions of our model depending on input number of nodes $n$. In this experiment, node features are randomly generated and the input graph is not considered, and the DGM embedding function $f$ is thus replaced with an MLP. cDGM refers to the continuous sampling strategy generalizing [45]. It is the most demanding implementation, since graph convolution is performed through dense matrix multiplication.
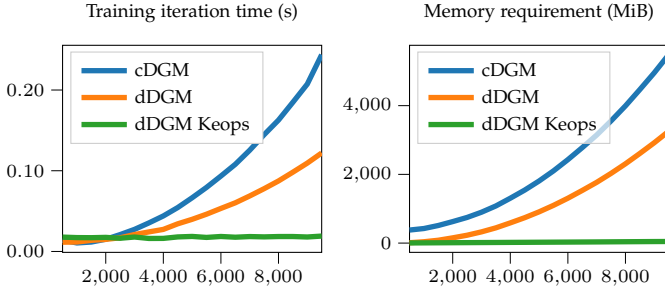
Fig. 3: Time per iteration and memory requirement for training three different implementations of DGM: continuous graph sampling (cDGM), 5-nn discrete graph sampling implemented with classical Pytorch tensor operations (dDGM) and using the Keops library (dDGM Keops).

With the kNN discrete graph sampling strategy (dDGM), graph convolutions operate on a sparse graph representation, resulting in lower memory consumption and faster computation, particularly on larger graphs with $n \gtrsim 2000$. Still, the pairwise distance in the graph embedding space is performed through standard PyTorch operations, resulting in a quadratic time and memory complexity. In order to overcome this quadratic complexity as well, we provide a further, symbolic implementation using the KeOps library [70]. Using KeOps, the kNN sampling computes the pairwise distances on-the-fly instead of storing them in memory, resulting in a linear memory complexity and a much faster computation.

## 5.2 Application to Healthcare and Brain imaging

In this section, we show experiments on two medical datasets for disease and age prediction. The main motivation behind the age prediction task is that the difference between brain age and chronicle age is an important biomarker for identifying certain neurological disorders [71].

**Dataset description**: We use two datasets: *Tadpole* [72] comprises multimodal data of 564 patients, each represented by a 354-dimensional representation derived from imaging (MRI, fMRI, PET) and non-imaging (demographics and genotypes) features. The task is to classify each subject into three classes: 'Normal Control', 'Alzheimer's Disease' and 'Mild Cognitive Impairment'. *UK Biobank* [73] comprises multimodal data of 14,503 individuals, each represented by a 440 dimensional feature vector derived from brain MRI and fMRI imaging. The task is to classify the age group of the patient (50-59,60-69, 70-79, and 80-89). We perform the two classification tasks both in a transductive and inductive manner, where in the former setting all the

TABLE 8: Accuracy on Tadpole transductive task with different architectural choices for our cDGM and dDGM models. The notation $f+g$ refers to the choice of the DGM and Diffusion modules (GC: graph convolution, EC: edge convolution; MLP: multilayer perceptron)

.

|      | MLP + GC | GC + GC | MLP + EC | GC + EC |
|------|----------|---------|----------|---------|
| cDGM | 92.42±3.82 | 90.68±4.58 | 92.29±4.18 | 91.78±3.21 |
| dDGM | 93.47±3.82 | **94.09±1.81** | 93.27±3.20 | **94.14±2.12** |

TABLE 9: Classification accuracy in % on Tadpole in the transductive setting. The top three performance scores are highlighted in color as: First, Second, **Third**.

| Method | Accuracy |
|--------|----------|
| Linear classifier | 70.22±6.32 |
| Multi-GCN [74] | 76.06±0.72 |
| Spectral-GCN [14] | 81.00±6.40 |
| InceptionGCN [16] | 84.11±4.50 |
| DGCNN [8] | 84.59±4.33 |
| LDS [57] | **87.06±3.67** |
| **cDGM** | 92.91±2.50 |
| **dDGM** | 94.14±2.12 |

TABLE 10: Classification accuracy in % for disease and age prediction tasks in the *transductive* and *inductive* settings on the Tadpole (left) and UK Biobank datasets (right). †Does not support inductive setting.

| Method | TADPOLE | | UK Biobank | |
|--------|---------|---------|---------|---------|
|        | Transduct. | Inductive | Transduct. | Inductive |
| DGCNN | 84.59±4.33 | **82.99±4.91** | 58.35±0.91 | **51.84±8.16** |
| LDS | **87.06±3.67** | † | OOM | † |
| **cDGM** | 92.91±2.50 | 91.85±2.62 | 61.32±1.51 | 55.91±3.49 |
| **dDGM** | 94.10±2.12 | 92.17±3.65 | 63.22±1.12 | 57.34±5.32 |

nodes are given during training but the labels of the test nodes are withheld, while in the latter setting, test nodes are completely removed from the population during training, and reintroduced only during testing.

**Architecture**: For medical applications, the architecture backbone is composed by two diffusion layers with output size of 16 and a final linear layer of size 16 for classification. Each layer is composed by a DGM block and one convolutional layer for diffusion. Since an input graph is not provided with the dataset, we always use a MLP $(d, 16, 8)$ as the the first DGM graph embedding function. We study the impact of using different operators for the second DGM block and the diffusion convolutions in Table 8, showing accuracy performance on the transductive setting for Tadpole. Following these results, we employ Graph Convolution for the second DGM block and Edge Convolution for the diffusion blocks.

**Comparison**: Previous methods [14], [16], [74], [75] used GNNs with hand-crafted patient population graphs built from non-imaging meta-features like the age and sex of patients. Our method allows to learn the graph directly from the input patients features, avoiding any handcrafted construction. We use the following baselines: simple linear

TABLE 11: Scalability: training and test iteration times for different number of nodes.

| Training iteration | | | |
|--------|--------|--------|--------|
| Method | $n =564$ | 5k | 10k |
| DGCNN | 6.99ms | 28.2ms | 104ms |
| LDS [57] | 1.84s | >30m | >30m |
| **cDGM** | 7.35ms | 47.8ms | 211ms |
| **dDGM** | 8.29ms | 37.0ms | 141ms |

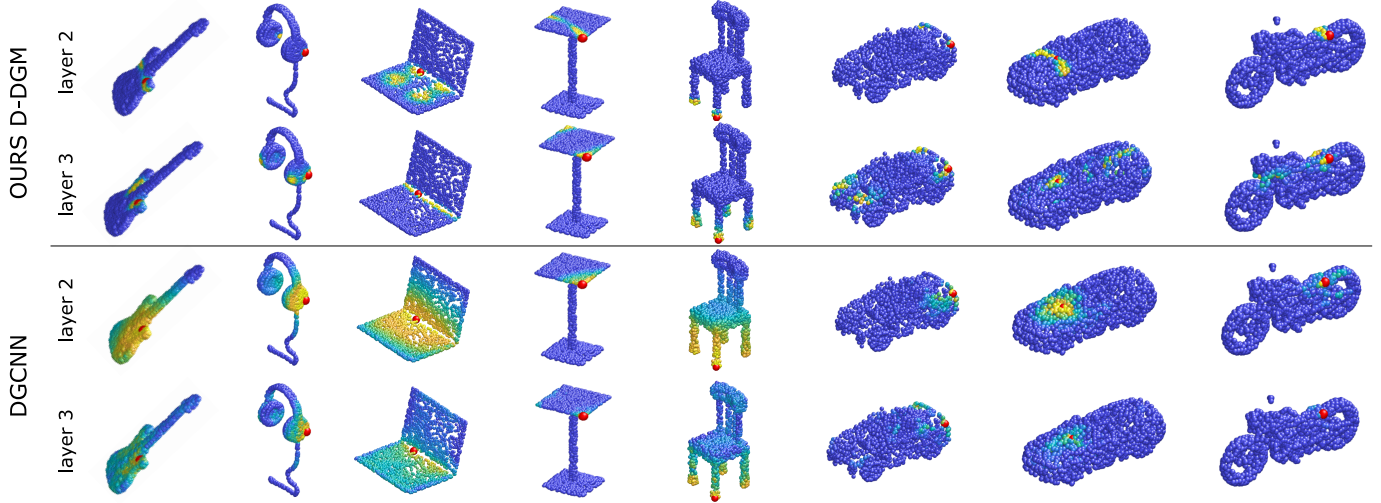| Test iteration | | | |
|--------|--------|--------|--------|
| Method | $n =564$ | 5k | 10k |
| DGCNN | 4.60ms | 25.2ms | 102ms |
| LDS [57] | 1.84s | >30m | >30m |
| **cDGM** | 3.02ms | 15.1ms | 51ms |
| **dDGM** | 3.97ms | 24.6ms | 104ms |

Fig. 4: Comparison between our dDGM and DGCNN sampling on the feature space in the last two convolutional layers of the network. In dDGM the colormap encodes the probability of each point to be connected to the red point. For DGCNN we plot the exponential of the negative Euclidean distance on feature space.

classifier as a non-graph method; Multi-GCN [74], Spectral-GCN [14], InceptionGCN [16] as graph methods with hand-crafted graph; and DGCNN [8] and LDS [57] as methods that learn the graph. We note that LDS does not support inductive learning.

Results are reported in Tables 9 and 10 using 10-fold cross validation. Our model significantly outperforms the state-of-the-art on all the tasks. Table 11 reports the training and testing times of our model for different dataset size, which shows that it is on par with DGCNN and about three orders of magnitude faster than LDS.

### 5.3 Application to computer vision and Computer graphics

We also show results of our method on two further domains and applications, namely point cloud segmentation and zero shot learning.

**Dataset description:** We perform the segmentation of 3D point clouds following [8] with the same data split. We use 16,881 point clouds from ShapeNet [76]. Each point cloud is sampled at 2048 points, and each point is annotated with one of the 50 part category labels. We mimic the DGCNN architecture used in [8], replacing their graph kNN sampling scheme by our DGM with a feature depth of 16. We keep the remainder of the network architecture the same, including the value of $k = 20$ and training parameters. During inference, given the stochastic nature of our graph, we repeat the classification of each point for 8 times and choose the $argmax$ of the cumulative soft predictions.

Table 12 reports the mean Intersection-over-Union (mIoU) values calculated by averaging the IoUs of all testing shapes. Our approach allows to increase the performance over the original kNN sampling scheme on almost all shape classes, which is considered very hard. Figure 4 depicts the sampling probabilities of some points (denoted by red) on different shapes at the last two layers of the network. We can notice that the probability of connecting two points is not related to the point feature space which is used for part

TABLE 12: Comparison of mIoU(%) score in ShapeNet part segmentation task. It can be observed in most of the classes dDGM performs better. Since the dataset size is quite large, the overall performance is significantly better.

|            | # Shapes | DGCNN | dDGM |
|------------|----------|-------|------|
| Airplane   | 2690     | 84.0  | **84.1** |
| Bag        | 76       | **83.4** | 82.5 |
| Cap        | 55       | **86.7** | 84.6 |
| Car        | 898      | 77.8  | **77.9** |
| Chair      | 3758     | 90.6  | **91.3** |
| Earphone   | 69       | 74.7  | **79.0** |
| Guitar     | 787      | 91.2  | **92.5** |
| Knife      | 392      | 87.5  | **87.7** |
| Lamp       | 1547     | 82.8  | **83.7** |
| Laptop     | 451      | 95.7  | **96.5** |
| Motorbike  | 202      | 66.3  | **66.8** |
| Mug        | 184      | 94.9  | **95.1** |
| Pistol     | 283      | 81.1  | **83.1** |
| Rocket     | 66       | **63.5** | 62.3 |
| Skateboard | 152      | 74.5  | **77.8** |
| Table      | 5271     | **82.6** | 82.2 |
| MEAN       |          | 85.2  | **85.6** |

classification, but it rather retains some spatial information and seems to be inspecting symmetries of the shape.

#### 5.3.1 Zero-shot learning (ZSL) in computer vision

In ZSL, the task is to learn classifiers for the unseen classes and solve the classification problem for samples belonging to unseen classes based on training data of only seen classes. The most popular approach is to train a network to predict a vector representation for a class starting from some implicit prior knowledge, i.e. semantic embedding [77]. Recent works showed that using additional explicit relations between classes in the form of knowledge graphs can help to significantly improve the learning of classifier for the unknown classes and hence the classification accuracy for unseen data samples. Notably, their method is explicitly devised for using a knowledge graph.

Formally, let $\mathbf{X} \in \mathbb{R}^{N \times S}$ be the semantic embeddings (i.e. word vectors) associated with each category class. The
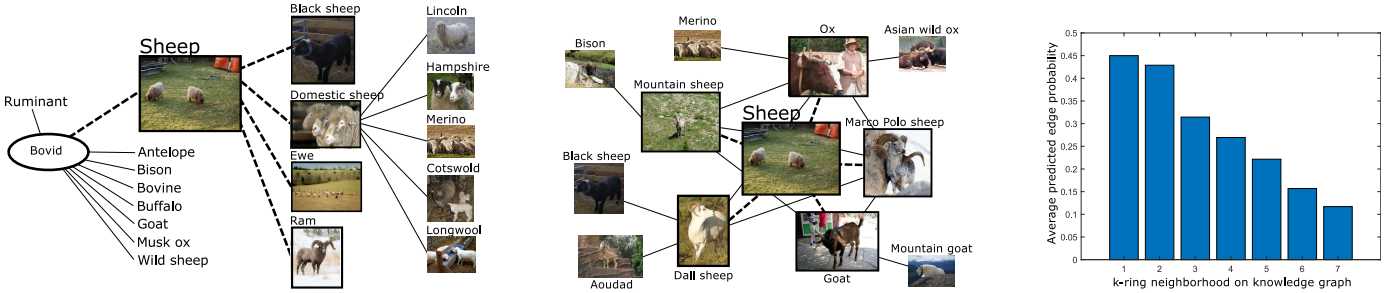
Fig. 5: Example of the 2-ring neighborhood of the "sheep" category on the knowledge graph (left) and on our predicted graph (center) sampled considering the 5 most probable edges. On the right the average predicted probability of edges belonging to the k-ring neighborhood (AwA2 test categories). Higher probabilities corresponding to nearest neighbors suggest that the predicted graph structure is loosely related to the knowledge graph.

Zero-Shot task loss is defined as the summation over all the $M < N$ training classes of $\sum_{i=1}^{M} \|\mathbf{w}_i - \tilde{\mathbf{w}}_i\|_2^2$, where $\mathbf{w}_i$ and $\tilde{\mathbf{w}}_i$ are the predicted and ground-truth vector representation of the $i$th class, respectively. Note that even though we deal with a regression problem in ZSL, it is straightforward to adapt it to deal with our graph loss defined in equation 8, considering $\arg\min_j \|\mathbf{w}_i - \mathbf{w}_j\|_2$ as the predicted category for sample $\mathbf{x}_i$.

Mimicking [78], our model consists of two graph convolution layers with hidden and output layer of dimension 2048 and 2049, paired with two DGM layers of dimension 16 for graph representation and $k = 3$. Each layer is composed by the following convolution on graphs:

$$\mathbf{X}^{(l+1)} = \sigma\left((\mathbf{D}^{(l)})^{-1}\mathbf{A}^{(l)}\mathbf{X}^{(l)}\mathbf{\Theta}^{(l)}\right) \qquad (10)$$

where $\sigma(\cdot)$ is a *LeakyReLU* non linearity, $\mathbf{\Theta}^{(l)}$ are the learned weights and $(\mathbf{D}^{(l)})^{-1}\mathbf{A}^{(l)}$, with $d_{ii}^{(l)} = \sum_j a_{ij}^{(l)}$ is the non-symmetric normalization of the adjacency matrix $\mathbf{A}^{(l)}$. Essentially, our model is the same as SGCN [78], where we replace the input knowledge graph by our DGM module for learning $\mathbf{A}$.

Following [78], we use weights of the last fully connected layer of a ResNet-50 [79] (pre-trained on ImageNet 2012 dataset [80]) as our target vector representation $\tilde{\mathbf{y}}_i \in \mathbb{R}^{2049}$. Input semantic features $\mathbf{x}_i \in \mathbb{R}^{300}$ are extracted with the GloVe text model [81] (trained on Wikipedia dataset).

We train our model on the 21K ImageNet dataset classes, where we have the semantic embedding for all classes as input, but only the first 1K have a corresponding ground-truth vector representation. The model is trained for 5000 iterations on a randomly subsampled set of 7K categories containing all the 1K training samples.

Testing is performed on the AWA2 dataset, composed of 37,322 images belonging to 50 different animal classes. We used the test split proposed in [82], comprising images from 10 classes not present in the first 1K training samples of ImageNet. The Top-1 accuracy for dDGM is 74.7%, which is reater than that of GCNZ [82] (70.5%), but lower than DGP (77.3%). Note that compared to the latter two methods, our dDGM approach does not make use of the knowledge graph. As shown in Figure 5, the knowledge graph indeed seems to be a good graph representation for the zero-shot task. Our predicted graph shows some similarity to it, however, fails to capture its hierarchical structure. We leave

it for future work to impose additional constraints on the graph structure, e.g. encouraging a more tree-like structure.

## 6 DISCUSSION AND CONCLUSION

In this paper, we tackled the challenge of graph learning in convolutional graph neural networks. We have proposed a novel Differentiable Graph Module (DGM): it predicts a probabilistic graph, from which a discrete graph can be sampled, which can then be used in any graph convolutional operator towards the downstream task. Furthermore, we devised a weighted loss that models the optimization over edge probabilities.

Our DGM is generic and adaptable to any graph convolution based method. We prove this by using our method to solve a wide variety of tasks, starting from applications in healthcare (disease prediction), brain imaging (age and gender prediction), computer graphics (3D point cloud segmentation) and computer vision (zero-shot learning). These applications feature both multimodal datasets and inductive inference setups.

There are several avenues for future work in our proposed method. For example, despite being computationally more lightweight than existing approaches (e.g. [83]), our method still has a quadratic complexity with respect to the number of input nodes, as it requires the computation of all pairwise distances. Restricting the computation of probabilities in a neighborhood of the node and using tree-based algorithms could help in reducing the complexity to $\mathcal{O}(n \log n)$. Further, our choice of sampling $k$ neighbors does not consider the heterogeneity of the graph in terms of the degree distribution of nodes. Other sampling schemes (e.g. threshold-based sampling [83]) could be investigated. It would be also interesting to consider previous knowledge about the graph, e.g. by imposing a node degree distribution, or providing an initial input graph to be optimized for a specific task.

# References

[1] M. M. Bronstein, J. Bruna, Y. LeCun, A. Szlam, and P. Vandergheynst, "Geometric deep learning: going beyond euclidean data," *IEEE Signal Process. Mag.*, vol. 34, no. 4, pp. 18–42, 2017.

[2] W. L. Hamilton, R. Ying, and J. Leskovec, "Representation learning on graphs: Methods and applications," *arXiv:1709.05584*, 2017.

[3] P. W. Battaglia *et al.*, "Relational inductive biases, deep learning, and graph networks," *arXiv:1806.01261*, 2018.

[4] M. Zhang and Y. Chen, "Link prediction based on graph neural networks," in *Proc. NeurIPS*, 2018.

[5] S. Qi, W. Wang, B. Jia, J. Shen, and S.-C. Zhu, "Learning human-object interactions by graph parsing neural networks," in *ECCV*, 2018, pp. 401–417.

[6] X. Qi, R. Liao, J. Jia, S. Fidler, and R. Urtasun, "3d graph neural networks for rgbd semantic segmentation," in *ICCV*, 2017, pp. 5199–5208.

[7] F. Monti, D. Boscaini, J. Masci, E. Rodola, J. Svoboda, and M. M. Bronstein, "Geometric deep learning on graphs and manifolds using mixture model cnns," in *Proc. CVPR*, 2017.

[8] Y. Wang, Y. Sun, Z. Liu, S. E. Sarma, M. M. Bronstein, and J. M. Solomon, "Dynamic graph cnn for learning on point clouds," *ACM TOG*, vol. 38, no. 5, p. 146, 2019.

[9] N. Choma *et al.*, "Graph neural networks for icecube signal classification," in *Proc. ICMLA*, 2018.

[10] D. K. Duvenaud, D. Maclaurin, J. Iparraguirre, R. Bombarell, T. Hirzel, A. Aspuru-Guzik, and R. P. Adams, "Convolutional networks on graphs for learning molecular fingerprints," in *NeurIPS*, 2015, pp. 2224–2232.

[11] J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals, and G. E. Dahl, "Neural message passing for quantum chemistry," in *International conference on machine learning.* PMLR, 2017, pp. 1263–1272.

[12] Y. Li, O. Vinyals, C. Dyer, R. Pascanu, and P. Battaglia, "Learning deep generative models of graphs," *arXiv preprint arXiv:1803.03324*, 2018.

[13] S. Parisot, S. I. Ktena, E. Ferrante, M. Lee, R. Guerrero, B. Glocker, and D. Rueckert, "Disease prediction using graph convolutional networks: Application to autism spectrum disorder and alzheimer's disease," *Med Image Anal*, vol. 48, pp. 117–130, 2018.

[14] S. Parisot, S. I. Ktena, E. Ferrante, M. Lee, R. G. Moreno, B. Glocker, and D. Rueckert, "Spectral graph convolutions for population-based disease prediction," in *MICCAI.* Springer, 2017, pp. 177–185.

[15] C. Mellema, A. Treacher, K. Nguyen, and A. Montillo, "Multiple deep learning architectures achieve superior performance diagnosing autism spectrum disorder using features previously extracted from structural and functional mri," in *2019 IEEE ISBI).* IEEE, 2019, pp. 1891–1895.

[16] A. Kazi, S. Shekarforoush, S. A. Krishna, H. Burwinkel, G. Vivar, K. Kortüm, S.-A. Ahmadi, S. Albarqouni, and N. Navab, "Inceptiongcn: Receptive field aware graph convolutional network for disease prediction," in *IPMI.* Springer, 2019.

[17] M. Zitnik, M. Agrawal, and J. Leskovec, "Modeling polypharmacy side effects with graph convolutional networks," *Bioinformatics*, vol. 34, no. 13, pp. i457–i466, 2018.

[18] M. Zitnik, F. Nguyen, B. Wang, J. Leskovec, A. Goldenberg, and M. M. Hoffman, "Machine learning for integrating data in biology and medicine: Principles, practice, and opportunities," *Information Fusion*, vol. 50, pp. 71–91, 2019.

[19] P. Gainza, F. Sverrisson, F. Monti, E. Rodola, M. M. Bronstein, and B. E. Correia, "Deciphering interaction fingerprints from protein molecular surfaces using geometric deep learning," *Nature Methods*, vol. 17, pp. 184–192, 2019.

[20] F. Scarselli, M. Gori, A. C. Tsoi, M. Hagenbuchner, and G. Monfardini, "The graph neural network model," *IEEE Trans. Neural Netw*, vol. 20, no. 1, pp. 61–80, 2008.

[21] J. Bruna, W. Zaremba, A. Szlam, and Y. LeCun, "Spectral networks and locally connected networks on graphs," *arXiv preprint arXiv:1312.6203*, 2013.

[22] M. Defferrard, X. Bresson, and P. Vandergheynst, "Convolutional neural networks on graphs with fast localized spectral filtering," in *NeurIPS*, 2016, pp. 3844–3852.

[23] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," *arXiv preprint arXiv:1609.02907*, 2016.

[24] R. Levie, F. Monti, X. Bresson, and M. M. Bronstein, "Cayleynets: Graph convolutional neural networks with complex rational spectral filters," *IEEE Trans. Signal Process.*, vol. 67, no. 1, pp. 97–109, 2018.

[25] F. M. Bianchi, D. Grattarola, C. Alippi, and L. Livi, "Graph neural networks with convolutional arma filters," *arXiv:1901.01343*, 2019.

[26] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Lio, and Y. Bengio, "Graph attention networks," *arXiv:1710.10903*, 2017.

[27] R. Kondor, "N-body networks: a covariant hierarchical neural network architecture for learning atomic potentials," *arXiv:1803.01588*, 2018.

[28] J. Bruna and X. Li, "Community detection with graph neural networks," *Stat*, vol. 1050, p. 27, 2017.

[29] F. Monti *et al.*, "Dual-primal graph convolutional networks," *arXiv:1806.00770*, 2018.

[30] W. Hamilton, Z. Ying, and J. Leskovec, "Inductive representation learning on large graphs," in *Proc. NIPS*, 2017.

[31] R. Ying, R. He, K. Chen, P. Eksombatchai, W. L. Hamilton, and J. Leskovec, "Graph convolutional neural networks for web-scale recommender systems," in *KDD*, 2018.

[32] W.-L. Chiang, X. Liu, S. Si, Y. Li, S. Bengio, and C.-J. Hsieh, "Cluster-gcn: An efficient algorithm for training deep and large graph convolutional networks," in *KDD*, 2019.

[33] H. Zeng, H. Zhou, A. Srivastava, R. Kannan, and V. K. Prasanna, "Graphsaint: Graph sampling based inductive learning method," *arXiv:1907.04931*, 2019.

[34] Y. Rong, W. Huang, T. Xu, and J. Huang, "Dropedge: Towards deep graph convolutional networks on node classification," *arXiv preprint arXiv:1907.10903*, 2019.

[35] L. Zhao and L. Akoglu, "Pairnorm: Tackling oversmoothing in gnns," in *Proc. ICLR*, 2020.

[36] G. Li, M. Muller, A. Thabet, and B. Ghanem, "Deepgcns: Can gcns go as deep as cnns?" in *Proc. ICCV*, 2019.

[37] S. Gong, M. Bahri, M. M. Bronstein, and S. Zafeiriou, "Geometrically principled connections in graph neural networks," in *Proc. CVPR*, 2020.

[38] K. Xu, W. Hu, J. Leskovec, and S. Jegelka, "How powerful are graph neural networks?" in *International Conference on Learning Representations (ICLR)*, 2019.

[39] H. Maron, H. Ben-Hamu, N. Shamir, and Y. Lipman, "Invariant and equivariant graph networks," in *International Conference on Learning Representations, ICLR*, 2019.

[40] N. Keriven and G. Peyré, "Universal invariant and equivariant graph neural networks," in *Advances in Neural Information Processing Systems*, 2019, pp. 7090–7099.

[41] A. Loukas, "What graph neural networks cannot learn: depth vs width," in *International Conference on Learning Representations (ICLR)*, 2020.

[42] W. Liu, J. Wang, and S.-F. Chang, "Robust and scalable graph-based semisupervised learning," *Proc. IEEE*, vol. 100, no. 9, pp. 2624–2638, 2012.

[43] X. Dong, D. Thanou, M. Rabbat, and P. Frossard, "Learning graphs from data: A signal representation perspective," *IEEE Signal Processing Magazine*, vol. 36, no. 3, pp. 44–63, 2019.

[44] G. Mateos, S. Segarra, A. G. Marques, and A. Ribeiro, "Connecting the dots: Identifying network structure via graph signal processing," *IEEE Signal Processing Magazine*, vol. 36, no. 3, pp. 16–43, 2019.

[45] C. et al., "Latent-graph learning for disease prediction," in *MICCAI.* Springer, 2020.

[46] X. Dong, D. Thanou, M. Rabbat, and P. Frossard, "Learning graphs from data: A signal representation perspective," *IEEE Signal Processing Magazine*, vol. 36, no. 3, pp. 44–63, 2019.

[47] T. Kipf, E. Fetaya, K.-C. Wang, M. Welling, and R. Zemel, "Neural relational inference for interacting systems," *arXiv preprint arXiv:1802.04687*, 2018.

[48] J. Zhang, X. Shi, J. Xie, H. Ma, I. King, and D.-Y. Yeung, "Gaan: Gated attention networks for learning on large and spatiotemporal graphs," *arXiv preprint arXiv:1803.07294*, 2018.

[49] J. B. Lee, R. Rossi, and X. Kong, "Graph classification using structural attention," in *Proceedings of the 24th ACM SIGKDD*

*International Conference on Knowledge Discovery & Data Mining*, 2018, pp. 1666–1674.

[50] S. Abu-El-Haija, B. Perozzi, R. Al-Rfou, and A. Alemi, "Watch your step: Learning node embeddings via graph attention," *arXiv preprint arXiv:1710.09599*, 2017.

[51] Z. Liu, C. Chen, L. Li, J. Zhou, X. Li, L. Song, and Y. Qi, "Geniepath: Graph neural networks with adaptive receptive paths," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, no. 01, 2019, pp. 4424–4431.

[52] D. V. Tran, N. Navarin, and A. Sperduti, "On filter size in graph convolutional networks," in *2018 IEEE Symposium Series on Computational Intelligence (SSCI)*. IEEE, 2018, pp. 1534–1541.

[53] K. Zhan, X. Chang, J. Guan, L. Chen, Z. Ma, and Y. Yang, "Adaptive structure discovery for multimedia analysis using multiple features," *IEEE transactions on cybernetics*, vol. 49, no. 5, pp. 1826–1834, 2018.

[54] R. Li, S. Wang, F. Zhu, and J. Huang, "Adaptive graph convolutional neural networks," in *AAAI*, 2018.

[55] W. Huang, T. Zhang, Y. Rong, and J. Huang, "Adaptive sampling towards fast graph representation learning," in *NeurIPS*, 2018, pp. 4558–4567.

[56] B. Jiang, Z. Zhang, D. Lin, J. Tang, and B. Luo, "Semi-supervised learning with graph learning-convolutional networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 11 313–11 320.

[57] L. Franceschi, M. Niepert, M. Pontil, and X. He, "Learning discrete structures for graph neural networks," in *International conference on machine learning*. PMLR, 2019, pp. 1972–1982.

[58] Y. et al., "Modeling point clouds with self-attention and gumbel subset sampling," in *CVPR*, 2019.

[59] Z. Ying, J. You, C. Morris, X. Ren, W. Hamilton, and J. Leskovec, "Hierarchical graph representation learning with differentiable pooling," *Advances in neural information processing systems*, vol. 31, 2018.

[60] N. et al., "Learning conditioned graph structures for interpretable visual question answering," *arXiv preprint arXiv:1806.07243*, 2018.

[61] Y. et al., "Iterative deep graph learning for graph neural networks:better and robust node embeddings," *NeurIPS*, 2020.

[62] L. Cosmo, G. Minello, M. Bronstein, E. Rodolà, L. Rossi, and A. Torsello, "Graph kernel neural networks," *arXiv preprint arXiv:2112.07436*, 2021.

[63] D. Krioukov, F. Papadopoulos, M. Kitsak, A. Vahdat, and M. Boguná, "Hyperbolic geometry of complex networks," *Physical Review E*, vol. 82, no. 3, p. 036106, 2010.

[64] M. Nickel and D. Kiela, "Poincaré embeddings for learning hierarchical representations," in *Advances in neural information processing systems*, 2017, pp. 6338–6347.

[65] W. Kool, H. van Hoof, and M. Welling, "Stochastic beams and where to find them: The gumbel-top-k trick for sampling sequences without replacement," *CoRR*, vol. abs/1903.06059, 2019. [Online]. Available: http://arxiv.org/abs/1903.06059

[66] P. Sen, G. Namata, M. Bilgic, L. Getoor, B. Galligher, and T. Eliassi-Rad, "Collective classification in network data," *AI Magazine*, vol. 29, no. 3, p. 93, Sep. 2008.

[67] A. L. Smith, D. M. Asta, and C. A. Calder, "The geometry of continuous latent space models for network data," *Statistical science: a review journal of the Institute of Mathematical Statistics*, vol. 34, no. 3, p. 428, 2019.

[68] I. Chami, R. Ying, C. Ré, and J. Leskovec, "Hyperbolic graph convolutional neural networks," *Advances in neural information processing systems*, vol. 32, p. 4869, 2019.

[69] E. Beltrami, *Teoria fondamentale degli spazii di curvatura costante: memoria*. F. Zanetti, 1868.

[70] B. Charlier, J. Feydy, J. A. Glaunès, F.-D. Collin, and G. Durif, "Kernel operations on the gpu, with autodiff, without memory overflows," *arXiv preprint arXiv:2004.11127*, 2020.

[71] X. Niu, F. Zhang, J. Kounios, and H. Liang, "Improved prediction of brain age using multimodal neuroimaging data," *Human brain mapping*, vol. 41, no. 6, pp. 1626–1643, 2020.

[72] R. V. Marinescu, N. P. Oxtoby, A. L. Young, E. E. Bron, A. W. Toga, M. W. Weiner, F. Barkhof, N. C. Fox, S. Klein, D. C. Alexander *et al.*, "Tadpole challenge: Prediction of longitudinal evolution in alzheimer's disease," *arXiv preprint arXiv:1805.03909*, 2018.

[73] K. L. Miller, F. Alfaro-Almagro, N. K. Bangerter, D. L. Thomas, E. Yacoub, J. Xu, A. J. Bartsch, S. Jbabdi, S. N. Sotiropoulos, J. L. Andersson *et al.*, "Multimodal population brain imaging in the uk biobank prospective epidemiological study," *Nature neuroscience*, vol. 19, no. 11, p. 1523, 2016.

[74] A. Kazi, S. Shekarforoush, K. Kortuem, S. Albarqouni, N. Navab *et al.*, "Self-attention equipped graph convolutions for disease prediction," in *ISBI*. IEEE, 2019, pp. 1896–1899.

[75] A. Kazi, S. Shekarforoush, S. A. Krishna, H. Burwinkel, G. Vivar, B. Wiestler, K. Kortüm, S.-A. Ahmadi, S. Albarqouni, and N. Navab, "Graph convolution based attention model for personalized disease prediction," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2019, pp. 122–130.

[76] L. Yi, V. G. Kim, D. Ceylan, I.-C. Shen, M. Yan, H. Su, C. Lu, Q. Huang, A. Sheffer, and L. Guibas, "A scalable active framework for region annotation in 3d shape collections," *ACM TOG*, vol. 35, no. 6, pp. 1–12, 2016.

[77] Y. Xian, C. H. Lampert, B. Schiele, and Z. Akata, "Zero-shot learning-a comprehensive evaluation of the good, the bad and the ugly," *IEEE PAMI*, 2018.

[78] M. Kampffmeyer, Y. Chen, X. Liang, H. Wang, Y. Zhang, and E. P. Xing, "Rethinking knowledge graph propagation for zero-shot learning," in *CVPR*, 2019, pp. 11 487–11 496.

[79] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016, pp. 770–778.

[80] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *CVPR*. Ieee, 2009.

[81] J. Pennington, R. Socher, and C. D. Manning, "Glove: Global vectors for word representation," in *EMNLP*, 2014, pp. 1532–1543.

[82] X. Wang, Y. Ye, and A. Gupta, "Zero-shot recognition via semantic embeddings and knowledge graphs," in *CVPR*, 2018, pp. 6857–6866.

[83] S. Jang, S. Moon, and J. Lee, "Brain signal classification via learning connectivity structure," *CoRR*, vol. abs/1905.11678, 2019. [Online]. Available: http://arxiv.org/abs/1905.11678