# Graph Mixup on Approximate Gromov-Wasserstein Geodesics

## **Abstract**

Mixup, which generates synthetic training samples on the data manifold, has been shown to be highly effective in augmenting Euclidean data. However, finding a proper data manifold for graph data is non-trivial, as graphs are non-Euclidean data in disparate spaces. Though efforts have been made, most of the existing graph mixup methods neglect the intrinsic geodesic guarantee, thereby generating inconsistent sample-label pairs. To address this issue, we propose GEOMIX to mixup graphs on the Gromov-Wasserstein (GW) geodesics. A joint space over input graphs is first defined based on the GW distance, and graphs are then transformed into the GW space through equivalence-preserving transformations. We further show that the linear interpolation of the transformed graph pairs defines a geodesic connecting the original pairs on the GW manifold, hence ensuring the consistency between generated samples and labels. An accelerated mixup algorithm on the approximate low-dimensional GW manifold is further proposed. Extensive experiments show that the proposed GEOMIX promotes the generalization and robustness of GNN models.

## 1. Introduction

In the era of big data and AI, graphs are ubiquitous in various domains carrying rich information. Graph Neural Networks (GNNs) have achieved remarkable success in enormous graph learning tasks, including graph classification (Xu et al., 2018), node classification (Kipf & Welling, 2017; Xu et al., 2022a), link prediction (Zhang & Chen, 2018; Yan et al., 2024c;b), and many more. For (semi-)supervised learning tasks, the superior performance of GNNs largely depends on the training graphs, which are often noisy and

Proceedings of the 41<sup>st</sup> International Conference on Machine Learning, Vienna, Austria. PMLR 235, 2024. Copyright 2024 by the author(s).

scarce, inevitably inducing model over-fitting (Ding et al., 2022). To address this issue, graph data augmentation has been adopted for better model generalization by generating synthetic training graphs.

Mixup (Zhang et al., 2018) has achieved great success in computer vision (Zhang et al., 2018; Verma et al., 2019) and natural language processing (Guo et al., 2019a; Guo, 2020; Chen et al., 2020) in improving model generalization and robustness. The general idea is to linearly interpolate sample pairs on the data manifold, which defines a geodesic between sample pairs in the Euclidean space. However, mixup on graph data is non-trivial due to three key challenges. First (space disparity), as different graphs lie in disparate spaces, it is a prerequisite to find a joint space of graphs, i.e., space of spaces (Sturm, 2012), for mixup. Second (non-Euclidean data), even with a joint graph space, it is still challenging to interpolate graph pairs as graphs vary in sizes and are often not well-aligned (Han et al., 2022). Third (sample–label consistency), it is essential to ensure the consistency between mixup samples and labels, i.e., similar samples should share similar labels (Guo et al., 2019b; Kim et al., 2020; Liu et al., 2023a), but it remains unknown how to define and ensure such consistency for graph data.

To address the space disparity and non-Euclidean issues, existing methods have been focusing on design practical interpolation so that mixup in the graph space can be implemented in a similar way as mixup in the Euclidean space. For example, (Wang et al., 2021; Verma et al., 2021) transform graphs into vector embeddings and (Han et al., 2022; Guo & Mao, 2021; Ling et al., 2023) transform graph into well-aligned pairs, and mixup is performed between the transformed counterparts. A fundamental assumption behind these approaches is that the original graph and the transformed graph are equivalent so that the mixup samples between transformed graphs correspond to the mixup samples between the original graphs. Nevertheless, most of the transformations are not equivalence-preserving, thereby facing the risk of generating inconsistent sample–label pairs that hurt, as opposed to benefit, model training. Therefore, we take a step further and ask:

What kind of interpolation guarantees sample—label consistency in graph mixup from the geodesic perspective?

**Contributions.** In this paper, we address the sample–label consistency issue in graph mixup and propose a novel algorithm named GEOMIX to interpolate graphs on the Gromov-Wasserstein (GW) geodesics. A comparison between the proposed GEOMIX and existing practical graph interpolation methods is shown in Figure 1. We first define samplelabel consistency based on the GW distance, and show that interpolating graphs on the GW geodesics guarantees such consistency. We then construct a GW space (Mémoli, 2011; Sturm, 2012) as a unified graph space, and employ equivalence-preserving transformations (EPT) to transform graphs to their equivalent counterparts that are well-aligned in the GW space. We theoretically prove that the linear interpolation of the transformed graph pairs defines a geodesic connecting the original graph pairs in the GW space, hence ensuring the consistency between mixup samples and labels. For faster computation, an accelerated algorithm is introduced to mixup graphs on the approximate GW geodesic in the low-dimensional GW space, achieving quadratic time complexity w.r.t. the number of nodes in the input graph. Extensive experiments on real-world graphs show that GE-OMIX improves GNN generalization and robustness, achieving up to 6.6% outperformance compared with the state-ofthe-art graph mixup methods.

The rest of the paper is organized as follows. Section 2 introduces and analyzes the background knowledge on mixup and GW geometry. Section 3 introduces our proposed GE-OMIX followed by experiments in Section 4. Related works and conclusions are given in Sections 5 and 6, respectively.

#### 2. Preliminaries

## 2.1. Notations

We use bold uppercase letters for matrices (e.g., A), bold lowercase letters for vectors (e.g., s), calligraphic letters for sets (e.g.,  $\mathcal{G}$ ), and lowercase letters for scalars (e.g.,  $\alpha$ ). The element (i,j) of a matrix A is denoted as A(i,j). The transpose of A is denoted by the superscript  $\tau$ . The simplex histogram with n bins is denoted as  $\Delta_n = \{\mu \in \mathbb{R}^n_n | \sum_{i=1}^n \mu(i) = 1\}$ . A coupling is denoted as  $\Pi(\cdot, \cdot)$ , and the inner product is denoted as  $\langle \cdot, \cdot \rangle$ . A graph is represented as a tuple  $\mathcal{G} = (A, \mu)$  where  $A \in \mathbb{R}^{n \times n}$  is the adjacency matrix, and  $\mu \in \Delta_n$  is the node weight. Without any prior knowledge on nodes, we default the node weight as uniform  $\mu = \frac{1}{n}$ . We denote the node set of a graph  $\mathcal{G}$  as  $\mathcal{V}(\mathcal{G})$ .

## 2.2. Graph Mixup

Mixup is a simple yet effective augmentation method for Euclidean data like images. Given a pair of samples  $x_1, x_2$  with their labels  $y_1, y_2$ , mixup generates synthetic samples by linearly interpolating the sample pairs and their labels:

$$\boldsymbol{x}_{(\lambda)} = (1 - \lambda)\boldsymbol{x}_1 + \lambda \boldsymbol{x}_2, \boldsymbol{y}_{(\lambda)} = (1 - \lambda)\boldsymbol{y}_1 + \lambda \boldsymbol{y}_2. \quad (1)$$

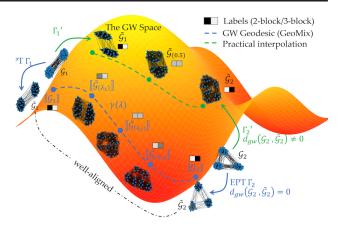


Figure 1. A comparison between GEOMIX and practical interpolation. GEOMIX (blue line) adopts equivalence-preserving transformations (EPT) to transform  $\mathcal{G}_1, \mathcal{G}_2$  into two well-aligned  $\tilde{\mathcal{G}}_1 \in [\![\mathcal{G}_1]\!], \tilde{\mathcal{G}}_2 \in [\![\mathcal{G}_2]\!]$ , and then mixup on the GW geodesics  $\gamma(\lambda) = [\![\tilde{\mathcal{G}}_{(\lambda)}]\!]$ , hence ensuring sample-label consistency. While practical interpolation (green line) often ignore the equivalence between original and transformed graphs, resulting in inconsistent sample-label pairs, e.g.,  $\tilde{\mathcal{G}}_{(0.5)}$  is a 2-block graph but is labelled as half 2-block half 3-block. Best viewed in color.

In the graph setting, in order to interpolate graphs with different sizes that are not well-aligned, most of the existing methods follow a practical interpolation approach by interpolating the well-aligned transformed graphs, that is

$$\tilde{\mathcal{G}}_{(\lambda)} = (1 - \lambda)\Gamma_1(\mathcal{G}_1) + \lambda\Gamma_2(\mathcal{G}_2), \, \tilde{\boldsymbol{y}}_{(\lambda)} = (1 - \lambda)\boldsymbol{y}_1 + \lambda\boldsymbol{y}_2,$$

where  $\Gamma_1$ ,  $\Gamma_2$  are graph transformations. For example, in G-mixup (Han et al., 2022),  $\Gamma_1$ ,  $\Gamma_2$  map original graphs to the well-aligned graphons; in M-mixup (Wang et al., 2021),  $\Gamma_1$ ,  $\Gamma_2$  are neural encoders encoding graphs into the same hidden space; in S-mixup (Ling et al., 2023),  $\Gamma_1$  corresponds to the soft alignment matrix between  $\mathcal{G}_1$  and  $\mathcal{G}_2$ .

# 2.3. The Gromov-Wasserstein Space

**Gromov–Wasserstein (GW) distance.** The GW distance is a powerful approach to measure the distance between two spaces. We formally define the GW distance in the graph context as follows. Given two graphs  $\mathcal{G}_1 = (A_1, \mu_1)$  and  $\mathcal{G}_2 = (A_2, \mu_2)$ , the *p*-GW distance between  $\mathcal{G}_1$  and  $\mathcal{G}_2$  is defined as (Mémoli, 2011):

$$d_{\text{GW}}(\mathcal{G}_1, \mathcal{G}_2) = \min_{\boldsymbol{T} \in \Pi(\boldsymbol{\mu}_1, \boldsymbol{\mu}_2)} \left( \varepsilon_{\boldsymbol{A}_1, \boldsymbol{A}_2}(\boldsymbol{T}) \right)^{1/p}, \quad (2)$$

where p is the order of the GW distance, and

$$\varepsilon_{\boldsymbol{A}_1,\boldsymbol{A}_2}(\boldsymbol{T}) = \sum_{i,j,k,l} |\boldsymbol{A}_1(i,j) - \boldsymbol{A}_2(k,l)|^p \boldsymbol{T}(i,k) \boldsymbol{T}(j,l).$$

Intuitively, the GW distance yields an optimal matching T between two graphs regarding the connectivity structure

(i.e.,  $A_1(i, j) - A_2(k, l)$ ), together with a distance measured by the sum of node pair connectivity distances weighted by T. In this paper, we adopt the 2-GW distance which has the following matrix form (Peyré et al., 2016):

$$\varepsilon_{\boldsymbol{A}_1,\boldsymbol{A}_2}\!(\boldsymbol{T}) \!=\! \mathrm{Tr}(\!\boldsymbol{A}_1^2\boldsymbol{\mu}_1\boldsymbol{\mu}_1^{\!\mathsf{T}}\!) \!+\! \mathrm{Tr}(\!\boldsymbol{\mu}_2\boldsymbol{\mu}_2^{\!\mathsf{T}}\!\boldsymbol{A}_2^{\!\mathsf{T}}\!) \!-\! 2\mathrm{Tr}(\!\boldsymbol{A}_1\boldsymbol{T}\boldsymbol{A}_2^{\!\mathsf{T}}\!\boldsymbol{T}^{\!\mathsf{T}}\!).$$

**GW geodesics.** Note that the GW distance is only a *pseudometric* but not a *metric*. To define geodesics, we employ a standard procedure (Howes, 1995) to identify an induced metric  $d_{\text{GW}}^*$ . Define an equivalence relation  $\sim$  over graphs:  $\mathcal{G}_1 \sim \mathcal{G}_2$  iff  $d_{\text{GW}}(\mathcal{G}_1, \mathcal{G}_2) = 0$ . Let  $[\![\mathcal{G}]\!] := \{\mathcal{G}': \mathcal{G}' \sim \mathcal{G}\}$  denote the equivalence class of graph  $\mathcal{G}$  w.r.t.  $\sim$ , and let  $\mathfrak{G}$  denote the space of equivalence classes w.r.t.  $\sim$ . Then, the *induced metric*  $d_{\text{GW}}^*: \mathfrak{G} \times \mathfrak{G} \to \mathbb{R}_{>0}$  is defined by

$$d_{GW}^*(\llbracket \mathcal{G}_1 \rrbracket, \llbracket \mathcal{G}_2 \rrbracket) := d_{GW}(\mathcal{G}_1, \mathcal{G}_2).$$

Note that  $d_{\mathrm{GW}}^*$  is well defined as  $d_{\mathrm{GW}}(\mathcal{G}_1, \mathcal{G}_2)$  is the same for any  $\mathcal{G}_1 \in [\![\mathcal{G}_1]\!], \mathcal{G}_2 \in [\![\mathcal{G}_2]\!]$  (Howes, 1995). Based on the notion of equivalence class, an equivalence-preserving graph transformation is defined as follows:

**Definition 2.1.** Equivalence-Preserving Transformation. Given a graph space  $\mathfrak{G}$ , a graph transformation  $\Gamma:\mathfrak{G}\to\mathfrak{G}$  is equivalence-preserving if any graph  $\mathcal{G}\in\mathfrak{G}$  is equivalent to its transformed graph  $\Gamma(\mathcal{G})$ , i.e.,  $\mathcal{G}\sim\Gamma(\mathcal{G}), \forall \mathcal{G}\in\mathfrak{G}$ .

**Definition 2.2.** Gromov–Wasserstein Geodesics.

A curve  $\gamma:[0,1]\to\mathfrak{G}$  is called a *GW geodesic* from  $\llbracket\mathcal{G}_1\rrbracket$  to  $\llbracket\mathcal{G}_2\rrbracket$  iff  $\gamma(0)=\llbracket\mathcal{G}_1\rrbracket$ ,  $\gamma(1)=\llbracket\mathcal{G}_2\rrbracket$ , and for every  $\lambda_1,\lambda_2\in[0,1]$ ,

$$d_{\mathrm{GW}}^*(\gamma(\lambda_1), \gamma(\lambda_2)) = |\lambda_1 - \lambda_2| \cdot d_{\mathrm{GW}}^*(\llbracket \mathcal{G}_1 \rrbracket, \llbracket \mathcal{G}_2 \rrbracket). \tag{3}$$

## 2.4. Geodesic Graph Mixup

It is essential to ensure the consistency between mixup samples and labels to avoid suspicious supervision signals that may mislead the model. Intuitively, if two graphs are similar to each other, we expect them to share similar labels. Based on the GW distance, we define sample—label consistency as follows.

**Definition 2.3.** *Sample–Label Consistency*.

Given two graphs  $\mathcal{G}_1, \mathcal{G}_2$  with labels  $\mathbf{y}_1, \mathbf{y}_2$ , the mixup samples  $\tilde{\mathcal{G}}_{(\lambda)}$  and labels  $\tilde{\mathbf{y}}_{(\lambda)}$  are consistent iff for every  $\lambda \in [0, 1]$ ,

$$d_{\mathrm{GW}}(\tilde{\mathcal{G}}_{(\lambda)},\mathcal{G}_1)\|\tilde{\boldsymbol{y}}_{(\lambda)}-\boldsymbol{y}_2\|=d_{\mathrm{GW}}(\tilde{\mathcal{G}}_{(\lambda)},\mathcal{G}_2)\|\tilde{\boldsymbol{y}}_{(\lambda)}-\boldsymbol{y}_1\|.$$

When the transformations ignore the equivalence-preserving property, as most of the existing graph mixup methods do, the sample-label consistency is clearly violated as  $d(\tilde{\mathcal{G}}_{(0)},\mathcal{G}_1)>0=\|\tilde{\boldsymbol{y}}_{(0)}-\boldsymbol{y}_1\|$ . The following proposition states that samples on the GW geodesics satisfy the sample-label consistency.

**Proposition 2.4.** The mixup samples  $\tilde{\mathcal{G}}_{(\lambda)}$  and labels  $\tilde{\boldsymbol{y}}_{(\lambda)}$  are consistent if  $\tilde{\mathcal{G}}_{(\lambda)}$  are on the geodesic connecting original samples  $\mathcal{G}_1, \mathcal{G}_2$ .

It is easy to verify the sample–label consistency for samples on the GW geodesics as follows:

$$d_{\mathrm{GW}}(\tilde{\mathcal{G}}_{(\lambda)},\mathcal{G}_1) = \lambda d_{\mathrm{GW}}(\mathcal{G}_1,\mathcal{G}_2), d_{\mathrm{GW}}(\tilde{\mathcal{G}}_{(\lambda)},\mathcal{G}_2) = (1-\lambda)d_{\mathrm{GW}}(\mathcal{G}_1,\mathcal{G}_2)$$
$$\|\tilde{\boldsymbol{y}}_{(\lambda)} - \boldsymbol{y}_1\| = \lambda \|\boldsymbol{y}_1 - \boldsymbol{y}_2\|, \|\tilde{\boldsymbol{y}}_{(\lambda)} - \boldsymbol{y}_2\| = (1-\lambda)\|\boldsymbol{y}_1 - \boldsymbol{y}_2\|.$$

Therefore, we study the geodesic graph mixup problem as follows.

**Definition 2.5.** *Geodesic Graph Mixup.* 

Given two graphs  $\mathcal{G}_1, \mathcal{G}_2$ , the geodesic graph mixup problem seeks for a GW geodesic  $\gamma(\lambda) = [\![\tilde{\mathcal{G}}_{(\lambda)}]\!]$  connecting  $\mathcal{G}_1, \mathcal{G}_2$  such that Eq. (3) is satisfied.

# 3. Methodology

In this section, we present and analyze our proposed GE-OMIX. To generalize the Euclidean mixup to graphs, we first propose to mixup graphs on the exact GW geodesics in Section 3.1. To avoid the high dimensionality of exact GW geodesics, an accelerated algorithm on the approximated GW geodesics is introduced in Section 3.2.

## 3.1. Mixup on Exact GW Geodesics

Most of the existing graph mixup methods generalize mixup from the practical interpolation perspectives, which may induce inconsistency between the mixup samples and labels. Instead, we derive a more principled generalization of mixup from a *geodesic* perspective. We view Euclidean mixup in Eq. (1) as the *Euclidean geodesic* connecting  $x_1$  and  $x_2$ . From this perspective, a natural generalization for graphs is the *GW geodesic* connecting two graphs. Thus, it suffices to find appropriate transformations  $\Gamma_1, \Gamma_2$  such that  $[(1 - \lambda)\Gamma_1(\mathcal{G}_1) + \lambda\Gamma_2(\mathcal{G}_2)]$  is the GW geodesic connecting two equivalence classes  $[\mathcal{G}_1], [\mathcal{G}_2]$ .

Thanks to the advancement in the optimal transport theory, the concept of geodesic has been generalized to the GW space to handle non-Euclidean data (Mémoli, 2011). According to Theorem 3.1 in (Sturm, 2012), given two graphs  $\mathcal{G}_1, \mathcal{G}_2$ , the 2-GW geodesic connecting  $[\![\mathcal{G}_1]\!], [\![\mathcal{G}_2]\!] \in \mathfrak{G}$  is of the form  $([\![\tilde{\mathcal{G}}_{(\lambda)}]\!])_{0 \leq \lambda \leq 1}$  where  $\tilde{\mathcal{G}}_{(\lambda)} = (A_{(\lambda)}, \tilde{\boldsymbol{\mu}}_{(\lambda)})$  is a graph with  $\mathcal{V}(\tilde{\mathcal{G}}_{(\lambda)}) = \mathcal{V}(\mathcal{G}_1) \times \mathcal{V}(\mathcal{G}_2)$  defined by Eq. (4).

$$\tilde{\boldsymbol{A}}_{(\lambda)} := (1 - \lambda)\boldsymbol{A}_1 \otimes \boldsymbol{1}_{n_2 \times n_2} + \lambda \boldsymbol{1}_{n_1 \times n_1} \otimes \boldsymbol{A}_2$$

$$\tilde{\boldsymbol{\mu}}_{(\lambda)} = \text{vec}(\text{OT}(\mathcal{G}_1, \mathcal{G}_2))$$
(4)

where  $\mathrm{OT}(\mathcal{G}_1,\mathcal{G}_2)$  is the optimal coupling between  $\mathcal{G}_1,\mathcal{G}_2$  given by the GW distance in Eq. (2). This theorem inspires us to consider two linear transformations  $P_1 \in \mathbb{R}^{n_1 \times n_1 n_2}, P_2 \in \mathbb{R}^{n_2 \times n_1 n_2}$  transforming  $\mathcal{G}_1, \mathcal{G}_2$  to  $\tilde{\mathcal{G}}_1 =$ 

$$(\tilde{A}_1, \tilde{\mu}_1), \tilde{\mathcal{G}}_2 = (\tilde{A}_2, \tilde{\mu}_2)$$
 as follows

$$\begin{cases} \tilde{\boldsymbol{A}}_{1} = \boldsymbol{P}_{1}^{\mathsf{T}} \boldsymbol{A}_{1} \boldsymbol{P}_{1}, \tilde{\boldsymbol{\mu}}_{1} = \text{vec}(\mathsf{OT}(\mathcal{G}_{1}, \mathcal{G}_{2})) \\ \tilde{\boldsymbol{A}}_{2} = \boldsymbol{P}_{2}^{\mathsf{T}} \boldsymbol{A}_{2} \boldsymbol{P}_{2}, \tilde{\boldsymbol{\mu}}_{2} = \text{vec}(\mathsf{OT}(\mathcal{G}_{1}, \mathcal{G}_{2})) \end{cases}$$
where 
$$\begin{cases} \boldsymbol{P}_{1} = \boldsymbol{I}_{n_{1}} \otimes \boldsymbol{1}_{1 \times n_{2}} \\ \boldsymbol{P}_{2} = \boldsymbol{1}_{1 \times n_{1}} \otimes \boldsymbol{I}_{n_{2}} \end{cases} . \tag{5}$$

Each column of  $P_1$ ,  $P_2$  is an one-hot vector mapping nodes in  $\mathcal{V}(\mathcal{G}_1)$  and  $\mathcal{V}(\mathcal{G}_2)$  to  $\mathcal{V}(\mathcal{G}_1) \times \mathcal{V}(\mathcal{G}_2)$ , and the transformation in Eq. (5) is essentially a Kronecker product between the adjacency matrix and an all-one matrix, i.e.,  $\tilde{A}_1 = A_1 \otimes \mathbf{1}_{n_2 \times n_2}, \tilde{A}_2 = \mathbf{1}_{n_1 \times n_1} \otimes A_2$ . Therefore, the transformed graphs  $\tilde{\mathcal{G}}_1$  and  $\tilde{\mathcal{G}}_2$  are in the same graph space, i.e., of the same size and well-aligned. On top of this, the mixup graph  $\tilde{\mathcal{G}}_{(\lambda)} = (\tilde{A}_{(\lambda)}, \tilde{\mu}_{(\lambda)})$  and its label  $\tilde{y}_{(\lambda)}$  can be further defined as:

$$\tilde{\mathbf{A}}_{(\lambda)} = (1 - \lambda)\tilde{\mathbf{A}}_1 + \lambda\tilde{\mathbf{A}}_2, 
\tilde{\mathbf{\mu}}_{(\lambda)} = (1 - \lambda)\tilde{\mathbf{\mu}}_1 + \lambda\tilde{\mathbf{\mu}}_2, 
\tilde{\mathbf{y}}_{(\lambda)} = (1 - \lambda)\mathbf{y}_1 + \lambda\mathbf{y}_2.$$
(6)

The following theorem shows that the mixup graphs  $\tilde{\mathcal{G}}_{(\lambda)}$  precisely defines the GW geodesic connecting  $[\![\mathcal{G}_1]\!], [\![\mathcal{G}_2]\!]$ .

**Theorem 3.1.** Given two graphs  $G_1$ ,  $G_2$ , the transformed graphs  $\tilde{G}_1$ ,  $\tilde{G}_2$  in Eq. (5) are in the equivalent class of  $G_1$ ,  $G_2$ , respectively, that is

$$\llbracket \tilde{\mathcal{G}}_1 \rrbracket = \llbracket \mathcal{G}_1 \rrbracket, \llbracket \tilde{\mathcal{G}}_2 \rrbracket = \llbracket \mathcal{G}_2 \rrbracket$$

Furthermore, the curve  $\{ [\![ \tilde{\mathcal{G}}_{(\lambda)} ]\!] \}_{0 \leq \lambda \leq 1}$  is a GW geodesic connecting  $[\![ \mathcal{G}_1 ]\!]$  and  $[\![ \mathcal{G}_2 ]\!]$ , which means that for every  $\lambda_1, \lambda_2 \in [0,1]$ ,

$$d_{\mathrm{GW}}^*([\![\tilde{\mathcal{G}}_{(\lambda_1)}]\!],[\![\tilde{\mathcal{G}}_{(\lambda_2)}]\!]) = |\lambda_1 - \lambda_2| \cdot d_{\mathrm{GW}}^*([\![\mathcal{G}_1]\!],[\![\mathcal{G}_2]\!]).$$

Proof is provided in Appendix A. This theorem, to the best of our knowledge, is the first to provide theoretical guarantee on the geodesic property for graph mixup.

#### 3.2. Accelerating Mixup on Approximate GW Geodesics

Though achieving exact GW geodesics, the high dimensionality of  $\tilde{\mathcal{G}}_{\lambda}$  prohibits efficient training. To address this issue, we propose to accelerate mixup via *approxiate* GW geodesics. Nonetheless, it is hard to define approximate GW geodesics based on the original definition of GW geodesics. Hence, we will (i) reformulate the GW geodesic as a solution to an optimization problem and (ii) define approximate GW geodesics by introducing into the optimization problem a hyperparameter that controls the degree of approximation.

**Theorem 3.2.** The tuple  $(P_1, P_2, \tilde{\mu}_{(\lambda)})$  in Eqs. (4) and (5) for exact GW geodesic is an optimal solution to the following

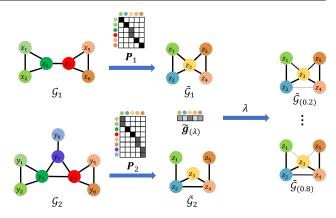


Figure 2. An example of GEOMIX transformation. Transformations  $P_1$ ,  $P_2$  map nodes  $x_i \in \mathcal{G}_1, y_j \in \mathcal{G}_2$  into common latent nodes  $z_k \in \tilde{\mathcal{G}}_1, \tilde{\mathcal{G}}_2$ , and the well-aligned  $\tilde{\mathcal{G}}_1, \tilde{\mathcal{G}}_2$  are linearly combined as the mixup graph  $\tilde{\mathcal{G}}_{(\lambda)}$ .

optimization problem:

$$\min_{\boldsymbol{P}_{1},\boldsymbol{P}_{2},\boldsymbol{g}} \left( \varepsilon_{\boldsymbol{A}_{1},\boldsymbol{A}_{2}}(\boldsymbol{P}_{1} \operatorname{diag}(\boldsymbol{g}) \boldsymbol{P}_{2}^{\mathsf{T}}) \right)^{1/2}$$
s.t. 
$$\begin{cases} \boldsymbol{P}_{1} \operatorname{diag}(\boldsymbol{g}) \in \Pi(\boldsymbol{\mu}_{1},\boldsymbol{g}) \\ \boldsymbol{P}_{2} \operatorname{diag}(\boldsymbol{g}) \in \Pi(\boldsymbol{\mu}_{2},\boldsymbol{g}) \\ \boldsymbol{g} \in \Delta_{n_{1}n_{2}} \end{cases}$$
(7)

Proof is provided in Appendix A. Theorem 3.2 reformulates the GW geodesic as a solution to an optimization problem. Intuitively, by regarding the mixup graph as the latent factor, Eq. (7) factorizes the optimal coupling T into three parts, including:  $P_1$ ,  $P_2$  transforming  $\mathcal{G}_1$ ,  $\mathcal{G}_2$  to the latent graph  $\mathcal{G}_{(\lambda)}$ , and the node weight g for  $\mathcal{G}_{(\lambda)}$ .

However, the computation in Eq. (7) is still inefficient due to the high dimensionality  $n_1n_2$ . To address this issue, we define *approximate* GW geodesics by replacing the dimensionality  $n_1n_2$  with a hyperparameter r that controls the degree of approximation.

To solve the optimization problem, we consider a change of variables:  $Q_1 = P_1 \operatorname{diag}(g)$ ,  $Q_2 = P_2 \operatorname{diag}(g)$ , based on which the optimization problem is reformulated into a low-rank GW problem (Scetbon et al., 2022):

$$\arg \min_{\mathbf{Q}_{1}, \mathbf{Q}_{2}, \mathbf{g}} \left( \varepsilon_{\mathbf{A}_{1}, \mathbf{A}_{2}}(\mathbf{Q}_{1} \operatorname{diag}(1/\mathbf{g}) \mathbf{Q}_{2}^{\mathsf{T}}) \right)^{1/2} \\
\text{s.t.} \begin{cases} \mathbf{Q}_{1} \in \Pi(\boldsymbol{\mu}_{1}, \mathbf{g}) \\ \mathbf{Q}_{2} \in \Pi(\boldsymbol{\mu}_{2}, \mathbf{g}) \\ \mathbf{g} \in \Delta_{r} \end{cases} \tag{8}$$

In general, the optimal coupling  $T \in \Pi(\mu_1, \mu_2)$  is factorized into three parts  $(Q_1, Q_2, g)$  that are closely connected to  $(P_1, P_2, \tilde{\mu}_{(\lambda)})$  for the exact GW geodesics in Eqs. (4) and (5). Specifically,  $Q_1, Q_2$  approximates the unnormalized

## Algorithm 1 GEOMIX

- 1: **Input** two graphs  $G_1 = \{A_1, \mu_1\}, G_2 = \{A_2, \mu_2\},$ mixup graph size r, mixup ratio  $\lambda$ .
- 2: Initialize  $g = \frac{\mathbf{1}_r}{r}, Q_1^{(0)} = \boldsymbol{\mu}_1 \otimes \boldsymbol{g}^{\mathsf{T}}, Q_2^{(0)} = \boldsymbol{\mu}_2 \otimes \boldsymbol{g}^{\mathsf{T}};$
- 3: **for**  $t \in \mathbb{N}^+_{\leq T}$  **do**
- Compute  $T^{(t)}, \omega^{(t)}, K_1^{(t)}, K_2^{(k)}, K_3^{(t)}$  in Eq. (9); 4:
- Solve  $\underset{\boldsymbol{Q}_{1},\boldsymbol{Q}_{2},\boldsymbol{g}}{\arg\min} \operatorname{KL}\Big((\boldsymbol{Q}_{1},\boldsymbol{Q}_{2},\boldsymbol{g}),(\boldsymbol{K}_{1}^{(t)}\!,\boldsymbol{K}_{2}^{(t)},\boldsymbol{K}_{3}^{(t)})\Big)$ by Dykstra's algorithm (Scetbon et al., 2021);

- 7: Column-normalize  $Q_1^{(T)}, Q_2^{(T)}$  as  $P_1, P_2$ ; 8: Transform graphs  $\tilde{A}_1 = P_1^{\mathsf{T}} A_1 P_1, \tilde{A}_2 = P_2^{\mathsf{T}} A_2 P_2$ ;
- 9:  $\tilde{\mathcal{G}}_{(\lambda)} = \{(1-\lambda)\tilde{\boldsymbol{A}}_1 + \lambda\tilde{\boldsymbol{A}}_2, \text{vec}(\boldsymbol{T}^{(T)})\};$
- 10:  $\tilde{\boldsymbol{y}}_{(\lambda)} = (1 \lambda)\boldsymbol{y}_1 + \lambda \boldsymbol{y}_2;$
- 11: **return** mixup graph  $\tilde{\mathcal{G}}_{(\lambda)}$  and label  $\tilde{\boldsymbol{y}}_{(\lambda)}$ .

linear transformation  $P_1$ ,  $P_2$ , and g approximates the node weight  $\tilde{\mu}_{(\lambda)}$ . In our setting, we choose  $r \leq n$  to map  $\mathcal{G}_1, \mathcal{G}_2$ to the coarsen graphs  $\tilde{\mathcal{G}}_1$  and  $\tilde{\mathcal{G}}_2$  that are well-aligned. The intuition is that graphs often share common latent meanings at the coarsen granularity (Zhang et al., 2020). An illustrative example is shown in Figure 2.

To solve the optimization problem in Eq. (8), a mirror descent scheme w.r.t. the generalized KL divergence is proposed (Scetbon et al., 2022). In general, the mirror descent scheme iteratively solves the following problem:

$$\begin{split} &\left(\boldsymbol{Q}_{1}^{(t+1)}, \boldsymbol{Q}_{2}^{(t+1)}, \boldsymbol{g}^{(t+1)}\right) = \underset{\boldsymbol{Q}_{1}, \boldsymbol{Q}_{2}, \boldsymbol{g}}{\arg\min} \operatorname{KL}\left((\boldsymbol{Q}_{1}, \boldsymbol{Q}_{2}, \boldsymbol{g}), (\boldsymbol{K}_{1}^{(t)}, \boldsymbol{K}_{2}^{(t)}, \boldsymbol{K}_{3}^{(t)})\right) \\ \text{s.t. } \boldsymbol{Q}_{1} \in \Pi(\boldsymbol{\mu}_{1}, \boldsymbol{g}), \ \boldsymbol{Q}_{2} \in \Pi(\boldsymbol{\mu}_{2}, \boldsymbol{g}), \ \boldsymbol{g} \in \Delta_{r} \\ &\left\{ \begin{aligned} &\boldsymbol{T}^{(t)} = \boldsymbol{Q}_{1}^{(t)} \operatorname{diag}(1/\boldsymbol{g}^{(k)}) \boldsymbol{Q}_{2}^{(t)^{\mathsf{T}}} \\ &\boldsymbol{\omega}^{(t)}(i) = \left(\boldsymbol{Q}_{1}^{(t)^{\mathsf{T}}} \boldsymbol{A}_{1} \boldsymbol{T}^{(t)} \boldsymbol{A}_{2} \boldsymbol{Q}_{2}^{(t)}\right)(i, i) \\ &\boldsymbol{K}_{1}^{(t)} = \exp\left(4\gamma \boldsymbol{A}_{1} \boldsymbol{T}^{(t)} \boldsymbol{A}_{2} \boldsymbol{Q}_{2}^{(t)} \operatorname{diag}(1/\boldsymbol{g}^{(t)}) + \log \boldsymbol{Q}_{1}^{(t)}\right) \\ &\boldsymbol{K}_{2}^{(t)} = \exp\left(4\gamma \boldsymbol{A}_{2} \boldsymbol{T}^{(t)^{\mathsf{T}}} \boldsymbol{A}_{1} \boldsymbol{Q}_{1}^{(t)} \operatorname{diag}(1/\boldsymbol{g}^{(t)}) + \log \boldsymbol{Q}_{2}^{(t)}\right) \\ &\boldsymbol{K}_{3}^{(t)} = \exp\left(-4\gamma \boldsymbol{\omega}^{(t)}/\boldsymbol{g}^{(t)^{2}} + \log \boldsymbol{g}^{(t)}\right) \end{aligned}$$

where  $\gamma$  is the step size for mirror descent. In general, for each step t, we first find the unconstraint solution  $K_1^{(t)}$ to  $\min_{\boldsymbol{Q}_1} \varepsilon(\boldsymbol{Q}_1 \operatorname{diag}(1/\boldsymbol{g}^{(t)} \boldsymbol{Q}_2^{(t)^{\mathsf{T}}})) + \operatorname{KL}(\boldsymbol{Q}_1, \boldsymbol{Q}_1^{(t)})$ . Similarly, we can obtain the unconstraint solutions  $m{K}_2^{(t)}, m{K}_3^{(t)}$ for  $Q_2, g$ , respectively. Then  $K_1^{(t)}, K_2^{(t)}, K_3^{(t)}$  are projected to the feasible regions  $\Pi(\mu_1, \mathbf{g}), \Pi(\mu_2, \mathbf{g}), \Delta_r$ , respectively, by minimizing the KL divergence. As proposed in (Scetbon et al., 2021), the above problem can be efficiently solve via the Dykstra's algorithm (Dykstra, 1983). Full algorithm details are provided in Appendix B. Afterwards, the transformations  $P_1, P_2$  can be obtained by the column normalization of  $Q_1, Q_2$ . The overall algorithm is summarized in Algorithm 1.

On the factorization of the optimal coupling. The key idea of transform-then-interpolate strategy is to transform a graph pair into well-aligned pairs that are easy to interpolate. We show that the factorization  $T = P_1 \operatorname{diag}(g) P_2^{\mathsf{T}}$  in Eq. (8) can be interpreted as a probabilistic alignment between  $A_1, A_2$ . Given that  $P_1, P_2$  are column-normalized matrices,  $P_1(i, u) = p_1(i|u), P_2(k, v) = p_2(k|v)$  indicate the generation probability of node  $i \in \mathcal{G}_1, k \in \mathcal{G}_2$  conditioned on node  $u, v \in \mathcal{G}_{(\lambda)}$ , respectively. The simplex g(u) = p(u)indicates the prior probability of nodes  $u \in \tilde{\mathcal{G}}_{(\lambda)}$ . Then the factorization in Eq. (7) can be written as:

$$\begin{split} & \left(\varepsilon_{\boldsymbol{A}_{1},\boldsymbol{A}_{2}}(\boldsymbol{P}_{1}\mathrm{diag}(\boldsymbol{g})\boldsymbol{P}_{2}^{\mathsf{T}})\right)^{1/2} \\ &= \sum_{i,j,k,l} \left( \left(\boldsymbol{A}_{1}(i,j) - \boldsymbol{A}_{2}(k,l)\right)^{2} \boldsymbol{T}(i,k) \boldsymbol{T}(j,l) \right)^{1/2} \\ &= \sum_{i,j,k,l} \left( \left(\boldsymbol{A}_{1}(i,j) - \boldsymbol{A}_{2}(k,l)\right)^{2} \boldsymbol{P}_{1}(i,u) \boldsymbol{g}(u) \boldsymbol{P}_{2}(k,u) \boldsymbol{P}_{1}(j,v) \boldsymbol{g}(v) \boldsymbol{P}_{2}(l,v) \right)^{1/2} \\ &= \left( \mathbb{E}_{\substack{u \sim p(u), \ v \sim p(v) \\ k \sim p_{1}(\cdot|u), \ j \sim p_{1}(\cdot|v) \\ k \sim p_{2}(\cdot|u), \ l \sim p_{2}(\cdot|v)}} \left(\boldsymbol{A}_{1}(i,j) - \boldsymbol{A}_{2}(k,l)\right)^{2} \right)^{1/2} \end{split}$$

The above equation can be interpreted as a two-step alignment. We first select a node pair (u, v) from the mixup graph  $\mathcal{G}_{(\lambda)}$  based on the prior probabilities p(u) and p(v). Then we select node pair i, j from  $\mathcal{G}_1$  based on conditional probability  $p_1(\cdot|u), p_1(\cdot|v)$ , and node pair k, l from  $\mathcal{G}_2$  based on conditional probability  $p_2(\cdot|u), p_2(\cdot|v)$ . Minimization of the above equation corresponds to finding the optimal transformations  $p_1, p_2$ , and node weight p to best align  $\mathcal{G}_1, \mathcal{G}_2$ .

Time complexity. Without loss of generality, we assume that input graphs share a comparable size that is greater than the mixup graph size r, i.e.,  $\mathcal{O}(n_1) = \mathcal{O}(n_2) = \mathcal{O}(n) >$  $\mathcal{O}(r)$ , we have the following time complexity analysis:

**Proposition 3.3.** Given K graph pairs to be mixed up, the time complexity for GEOMIX is  $\mathcal{O}(KTn^2r)$ , where T is the number of iterations in the low-rank GW algorithm.

The time complexity is quadratic w.r.t. the input graph size n and linear w.r.t. the mixup graph size r. It's worth mentioning that the exact GW geodesic requires  $\mathcal{O}(n^4)$  time complexity, while the approximate GW geodesic, as we will empirically show in Section 4, achieves a significant reduction in running time with  $\mathcal{O}(n^2r)$  time complexity at little expense of effectiveness.

# 4. Experiments

We conduct extensive experiments to evaluate the proposed GEOMIX. We first provide visualization on the mixup process and assess the geodesic property in Section 4.1. Then we evaluate the effectiveness of GEOMIX in enhancing GNN generalization and robustness in Sections 4.2 and 4.3, respectively. Further analysis is carried out in Section 4.4.

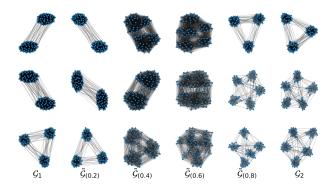


Figure 3. Visualization of the GEOMIX process. Each row corresponds to one mixup case. The leftmost and rightmost columns are input graphs, and the middle four columns are mixup graphs.

#### 4.1. Understanding the GEOMIX process

Visualize the mixup graphs. To understand the GEOMIX process, we first provide visualization of the mixup graphs in Figure 3. We use the Stochastic Block Model (SBM) to generate graphs with 2,3,5 blocks. GEOMIX is further employed to generate  $\tilde{\mathcal{G}}_{(\lambda)}$  with  $\lambda = \{0.2, 0.4, 0.6, 0.8\}$ . It can be observed that GEOMIX indeed generates a mixture of two graphs, and higher the  $\lambda$ , more similar  $\tilde{\mathcal{G}}_{(\lambda)}$  is to  $\mathcal{G}_2$ . More visualization results on real-world datasets are provided in Figure 9 in Appendix C.

**Evaluate the sample–label consistency.** We further evaluate whether the sample–label pairs generated by GE-OMIX are consistent by evaluating the correlation between  $\frac{d_{GW}(\tilde{\mathcal{G}}_{(\lambda)},\mathcal{G}_1)}{d_{GW}(\tilde{\mathcal{G}}_{(\lambda)},\mathcal{G}_2)}$  and  $\frac{\|\tilde{\boldsymbol{y}}_{(\lambda)}-\boldsymbol{y}_1\|}{\|\tilde{\boldsymbol{y}}_{(\lambda)}-\boldsymbol{y}_2\|}$ . We compare GEOMIX with two state-of-the-art graph mixup methods: G-Mixup (Han et al., 2022) and FGWMixup (Ma et al., 2023) on the IMDB-B dataset (Yanardag & Vishwanathan, 2015). As shown in Figure 4,  $\frac{d_{GW}(\tilde{\mathcal{G}}_{(\lambda)},\mathcal{G}_1)}{d_{GW}(\tilde{\mathcal{G}}_{(\lambda)},\mathcal{G}_2)}$  is closely related to  $\frac{\|\tilde{\boldsymbol{y}}_{(\lambda)}-\boldsymbol{y}_1\|}{\|\tilde{\boldsymbol{y}}_{(\lambda)}-\boldsymbol{y}_2\|}$  for GEOMIX achieving a Pearson coefficient of 0.91. For baseline methods like G-Mixup and FGWMixup, they achieve relatively lower Pearson coefficients of 0.25 and 0.71, respectively. This validates that samples generated by GEOMIX are consistent with their labels.

We also visualize the embeddings for IMDB-B and MU-TAG datasets. Specifically, we use GCN model to learn graph embeddings for both original graphs (blue) and GE-OMIX samples (orange), and adopt TSNE for visualization. As shown in Figure 5, the mixup samples form geodesics connecting the original graph pairs, which validates the geodesic property of GEOMIX. More embedding space visualization can be found in Figure 10 in Appendix C.

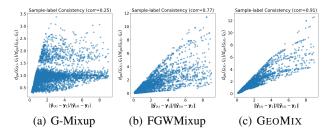


Figure 4. Sample–label consistency on IMDB-B. Each point represents a mixup sample–label pair. The Pearson coefficient between mixup samples and labels is 0.91 for GEOMIX, which is higher than those of G-Mixup (0.25) and FGWMixup (0.71).

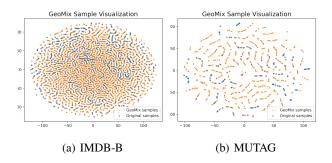


Figure 5. Embedding visualization on (a) IMDB-B and (b) MUTAG. Mixup samples (orange) reside on the geodesics connecting original samples (blue).

#### 4.2. On the Generalization of GNNs

**Experiment setup.** To assess the effectiveness of GE-OMIX in enhancing GNN generalization capability, we conduct experiments on graph classification on five real-world datasets, including PROTEINS (Borgwardt et al., 2005), MUTAG (Kriege & Mutzel, 2012), MSRC-9 (Neumann et al., 2016), IMDB-B (Yanardag & Vishwanathan, 2015), and IMDB-M (Yanardag & Vishwanathan, 2015). We split the dataset into train/test/validation set by 80%/10%/10% and use 10-fold cross validation for evaluation.

We consider two widely used GNN models, GCN (Kipf & Welling, 2017) and GIN (Xu et al., 2018), as the backbone models. We compare GEOMIX with eight state-of-the-art graph augmentation methods, including DropEdge (Rong et al., 2019), DropNode (You et al., 2020), Subgraph (You et al., 2020), M-Mixup (Wang et al., 2021), SubMix (Yoo et al., 2022), G-Mixup (Han et al., 2022), S-Mixup (Ling et al., 2023) and FGWMixup (Ma et al., 2023).

For a fair comparison, we adopt the same model architecture and hyperparameters. We use augmentation methods to generate the same number of augmented samples as the number of original samples. More experiment details can be found in Appendix D. Code and datasets are available at https:

Table 1.	Experiment result	lts on graph c	lassification.

	MODEL	PROTEINS	MUTAG	MSRC-9	IMDB-B	IMDB-M
	Vanilla	62.1±2.8	72.8±7.5	88.7±3.5	71.5±2.8	48.3±1.6
	DROPEDGE	61.6±2.9	$73.3 \pm 7.0$	89.4±4.2	71.4±2.7	47.8±1.9
	DropNode	61.1±2.6	$71.3 \pm 9.0$	89.0±4.0	71.6±3.4	47.7±1.8
	SUBGRAPH	59.6±2.6	$70.2 \pm 9.5$	83.1±3.2	71.7±3.7	45.5±2.8
GCN	M-MIXUP	61.0±2.6	$73.3 \pm 7.0$	88.2±3.6	$71.0 \pm 2.0$	48.1±1.8
$\ddot{\mathfrak{S}}$	SUBMIX	63.0±4.1	$73.5 \pm 7.1$	88.7±4.1	71.6±2.3	48.1±3.5
	G-MIXUP	63.6±3.5	$73.3 \pm 7.0$	$12.2 \pm 5.3$	$72.2 \pm 2.4$	48.0±2.5
	S-MIXUP	63.1±2.7	72.7±7.7	89.1±4.1	$71.3 \pm 2.1$	48.5±2.0
	FGWMIXUP	61.7±3.0	73.4±7.3	89.6±3.9	71.7±2.4	$47.5 \pm 2.2$
	GEOMIX	70.2±4.7	76.1±5.5	90.3±2.9	72.0±4.3	49.3±5.5
GIN	Vanilla	67.8±5.1	69.1±10.8	91.0±4.6	70.1±4.3	46.2±5.3
	DROPEDGE	66.8±4.9	71.2±7.4	88.2±5.8	$70.0 \pm 4.0$	46.5±4.0
	DropNode	66.8±4.6	71.8±6.7	91.8±4.9	71.6±3.6	44.9±5.0
	SUBGRAPH	62.3±4.9	$75.0 \pm 7.6$	85.6±10.3	68.1±4.8	45.4±4.9
	M-MIXUP	70.4±6.1	70.7±12.8	91.4±5.0	71.7±3.6	46.8±2.6
	SUBMIX	66.7±4.9	70.9±10.8	91.2±5.4	$70.9 \pm 3.1$	46.4±3.1
	G-MIXUP	69.6±6.5	$73.4 \pm 9.7$	19.5±9.1	71.1±3.4	47.0±4.1
	S-MIXUP	67.1±4.3	72.5±8.7	91.2±4.2	$70.7 \pm 3.2$	46.4±3.2
	FGWMIXUP	62.6±4.6	74.7±8.9	89.6±5.0	71.1±3.9	44.1±4.7
	GEOMIX	71.3±3.7	$80.8 \pm 6.2$	92.8±3.6	72.1±2.9	46.9±4.9

#### //github.com/zhichenz98/GeoMix-ICML24.

Results and analysis. The graph classification results are shown in Table 1. In general, our proposed GEOMIX achieves the best performance on 8 out 10 settings and the second-best performance on the rest 2 settings, with an up to 6.6% outperformance compared with the best competitor. Besides, we found that GEOMIX is more effective on datasets with fewer graphs and larger sizes (e.g., MUTAG and PROTEINS). This is because the mixup technique aims to generate synthetic samples covering the whole graph space. When dealing with larger and fewer graphs, the graph space is less covered, thus, the mixup technique achieves more significant improvements. This can be also validated by comparing the embedding visualization in Figure 5, where the MUTAG dataset is sparsely-covered while the IMDB-B dataset is more densely covered.

#### 4.3. On the Robustness of GNNs

**Experiment setup.** Besides generalization, mixup can also improve the model robustness (Zhang et al., 2021). We evaluate model robustness against two kinds of corruption, namely topology corruption and label corruption. Specifically, for topology corruption, we randomly remove/add 10% or 20% of the edges; for label corruption, we randomly change the labels for 10% or 20% of the training graphs.

**Results and analysis.** Robustness results against topology and label corruption are shown in Tables 2 and 3, respectively. In general, our proposed GEOMIX improves model robustness against both topology and label corruption, achieving the best performance in most cases. Specifi-

*Table 2.* Robustness against topology corruption.

	MODEL		10%			20%	
		IMDB-B	IMDB-M	PROTEINS	IMDB-B	IMDB-M	PROTEINS
GCN	VANILLA DROPEDGE DROPNODE SUBGRAPH GEOMIX	69.8±5.2 70.3±3.0 69.1±3.6 69.3±2.9 70.6±4.7	46.0±1.9 44.9±2.8 46.7±3.2 44.1±5.0 48.1±3.8	60.9±2.4 60.8±3.6 61.2±2.7 59.6±2.6 <b>65.9±4.7</b>	67.9±2.3 69.8±4.0 69.8±3.2 68.4±3.4 70.7±4.7	46.3±3.7 46.7±3.2 47.3±3.2 46.3±4.3 47.5±4.7	60.6±3.3 60.2±2.3 59.3±2.8 59.6±2.6 64.8±6.4
GIN	VANILLA DROPEDGE DROPNODE SUBGRAPH GEOMIX	66.4±4.4 65.7±3.8 62.8±8.8 63.9±4.9 67.1±4.6	44.6±5.1 42.8±4.6 42.9±5.5 39.4±4.8 43.2±4.9	63.2±5.5 60.0±1.9 60.1±3.6 59.6±2.6 <b>66.3±6.3</b>	67.6±3.7 65.5±5.5 63.7±5.3 66.6±5.9 67.8±3.6	41.9±3.9 40.4±5.0 40.7±3.3 41.3±5.6 42.4±5.3	59.6±2.6 59.7±2.9 59.6±2.6 59.6±2.6 61.2±6.2

cally, GEOMIX achieves up to 4.7% outperformance against topology corruption and 1.8% outperformance against label corruption compared with the best competitor. The reason for such outperformance is two-fold: (1) mixup samples that are far away from the training samples can help alleviate the model memorization of the corrupted labels (Zhang et al., 2018), and (2) the low-rank decomposition of GEOMIX naturally alleviates the topology noises, hence contributing to the robustness against topology noise.

## 4.4. On the Effect of Mixup Graph Size

We analyze how the mixup graph size affect the effectiveness and efficiency of GEOMIX. As shown in Figure 7, the performance of GEOMIX is quite stable w.r.t. different mixup graph sizes, with at most  $\pm 3\%$  fluctuation in accuracy. This indicates that (1) graphs contain redundant/noisy information and can be compressed as smaller graphs with lower r, and (2) the transformations employed by GEOMIX can capture the essential information for input graphs. Be-

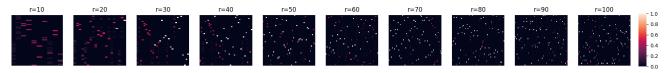


Figure 6. Transformation matrices with different mixup graph sizes r. Lighter the color, higher the value. As r increases, the transformation matrices become sparser and approximate the solutions to exact GW geodesic in Eq. (5).

Table 3. Robustness against label corruption.

	MODEL		10%		20%		
		IMDB-B	IMDB-M	PROTEINS	IMDB-B	IMDB-M	PROTEINS
GCN	VANILLA DROPEDGE DROPNODE SUBGRAPH GEOMIX	57.5±8.8 58.7±8.1 60.6±8.3 57.4±9.3 57.7±5.1	40.6±2.7 38.4±4.7 40.0±2.4 36.7±5.8 41.2±5.0	53.9±4.4 52.3±5.1 51.1±8.4 49.8±5.9 55.3±6.6	55.0±9.8 54.9±9.3 53.5±9.2 55.0±5.8 55.7±9.1	36.0±8.0 35.3±5.1 35.1±3.5 35.1±4.7 37.8±3.2	50.0±7.8 50.7±6.6 51.0±3.6 50.6±4.1 51.3±3.2
OIN	VANILLA DROPEDGE DROPNODE SUBGRAPH GEOMIX	58.0±8.5 54.7±8.8 58.4±8.6 56.9±6.7 58.0±7.5	35.7±3.6 37.8±5.8 35.7±4.9 37.7±6.7 39.1±4.3	53.7±6.1 54.6±9.2 53.4±4.6 52.2±4.1 54.7±7.4	53.2±7.7 52.1±3.9 52.8±6.4 53.6±6.8 53.8±6.6	30.3±7.5 37.2±5.6 34.7±6.3 36.7±5.2 37.9±3.2	51.1±6.1 53.0±5.8 51.9±5.6 49.9±5.8 52.2±6.6

sides, the running time of GEOMIX grows approximately linearly w.r.t. r. Compared with the exact GW geodesic with  $r=n^2$ , a small mixup graph size r helps reduce the computational cost at little cost of effectiveness.

Besides, we visualize the learned transformations  $P_1$  with different mixup graph size r. Specifically, we consider two 2-block SBM graphs with graph size n=50 and mixup size ranging from 10 to 100. Ideally, we expect the matrices to be column-wise sparse, i.e., each column to be a one-hot vector, so that it is analogous to the solution on exact GW geodesic in Eq. (5). As the results shown in Figure 6, the transformation matrices become sparser as r increases. This validates that (1) GEOMIX provides good approximation to the exact GW geodesic solution in Eq. (5), and (2) the hyperparameter r achieves a tradeoff between the degree of approximation and algorithm efficiency.

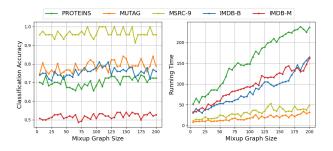


Figure 7. Hyperparameter study on mixup graph size r. Left: the classification accuracy is robust to different r. Right: the running time increase linearly w.r.t. r.

## 5. Related Work

Graph data augmentation. The superior performance of GNNs (Kipf & Welling, 2017; Xu et al., 2018; Gasteiger et al., 2018) on graph learning tasks largely depends on the quality of the training samples, which, however, are usually noisy and incomplete (Ding et al., 2022). Graph data augmentation has recently shown great power in enhancing model generalization and robustness, which can be categorized into three categories including, feature-based modification, structure-based modification and label-based modification. Feature-based methods manipulate node features by corruption (Feng et al., 2019; Yang et al., 2021), masking (You et al., 2020; 2021) and rewriting (Xu et al., 2022b), while structure-based methods perturb graph structure by node perturbation (You et al., 2020), edge perturbation (Rong et al., 2019; Xu et al., 2023a; Wei et al., 2022) and subgraph sampling (You et al., 2020). For label-based augmentation, a line of work interpolates graphs in the latent space (Wang et al., 2021; Verma et al., 2021; Navarro & Segarra, 2023), but may not be applicable to methods without latent space (e.g., kernel methods) and the mixup samples in the hidden space lack interpretability. Another line of work interpolates the aligned graph pairs (Guo & Mao, 2021; Ling et al., 2023; Kan et al., 2023), which can be regarded as a special case of our proposed GEOMIX when fixing one transformation P as the identity matrix. Some other works adopt the transplant strategy by swapping the subgraphs of two graphs. More recently, FGWMixup (Ma et al., 2023) adopts the Fused Gromov-Wasserstein barycenter as the mixup graphs, but suffers from heavy computation.

Optimal transport on graphs. OT has achieved great success in processing geometric data, such as graphs, due to its ability in capturing the intrinsic geometry (Peyré et al., 2016) and comparing irregular objects in disparate spaces (Mémoli, 2011). For example, OT has been applied in different graph-learning tasks such as graph alignment (Xu et al., 2019; Zeng et al., 2023a; 2024), graph comparison (Maretic et al., 2019; Titouan et al., 2019), graph representation learning (Kolouri et al., 2021; Vincent-Cuaz et al., 2021; Zeng et al., 2023b), and many more. Besides, it is shown that the Gromov–Wasserstein space provides a joint space for disparate spaces, and the geodesic is well-defined on the GW space (Sturm, 2012), hence providing a natural solution to

geodesic-preserving graph mixup. In this work, we utilize the OT theory from two aspects. First, the GW distance is used as the fundamental metric to define the joint graph space as well as the geodesic for mixup (Sturm, 2012). Second, the OT coupling serves as the alignment between two graphs, whose decomposition transforms two graphs into a shared space for mixup (Scetbon et al., 2022).

Graph neural networks. Graph machine learning has shown great power in various real-world applications, such as fraud detection (Ding et al., 2021; Fu et al., 2022; Warmsley et al., 2022), social recommendation (Wei et al., 2020; Jing et al., 2022b; Wei & He, 2022; Jing et al., 2024), information retrieval (Fu & He, 2021; Yoo et al., 2023), and bioinformatics (Fu & He, 2022; Xu et al., 2024). Behind these applications, GNNs are the most prominent graph machine learning models with superior performance on graph classification (Xu et al., 2018), node classification (Yan et al., 2024c; Xu et al., 2022a; Liu et al., 2023b), node clustering (Jing et al., 2021a), link prediction (Yan et al., 2022; 2024b), and many more.

The key idea behind various GNNs is to leverage node information via custom aggregation and propagation. GCN (Kipf & Welling, 2017) aggregates neighboring node information via a localized first-order approximation of spectral graph convolutions. GraphSage (Hamilton et al., 2017) handles the inductive graph learning tasks by sampling and aggregating features from local neighborhood. GAT (Velickovic et al., 2017) adopts self-attention to provide more flexibility in the aggregation stage. GIN (Xu et al., 2018) achieves better expressiveness by leveraging the WL test. APPNP (Gasteiger et al., 2018) adopts personalized pagerank to achieve better node sampling. Recent studies attempt to generalize GNNs to more challenging settings, e.g., heterophilic graphs (Yan et al., 2024a; Xu et al., 2023b; Guo et al., 2023), heterogeneous graphs (Jing et al., 2021b; 2022a), dynamic graphs (Wang et al., 2023; Fu et al., 2020; 2023; Yan et al., 2021), and many more.

# 6. Conclusion

In this paper, we address the sample—label consistency issue in graph mixup from the geodesic perspective, and propose GEOMIX to mixup graphs on the Gromov—Wasserstein geodesics. We first transform graph pairs into two well-aligned graphs through equivalence-preserving transformations, and then mixup on the exact GW geodesics. For faster computation, we parameterize the transformation by reformulating the GW geodesic as the solution to the low-rank GW problem, and an accelerated algorithm is further proposed. Extensive experiments demonstrate the effectiveness of the proposed GEOMIX in enhancing GNN model generalization and robustness.

# **Impact Statement**

This paper presents work whose goal is to advance the field of Graph Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

# Acknowledgements

This work is supported by NSF (2134079, 2316233, 2324770, and 2117902), DARPA (HR001121C0165), NIFA (2020-67021-32799), DHS (17STQAC00001-07-00), AFOSR (FA9550-24-1-0002), the C3.ai Digital Transformation Institute, and IBM-Illinois Discovery Accelerator Institute. The content of the information in this document does not necessarily reflect the position or the policy of the Government, and no official endorsement should be inferred. The U.S. Government is authorized to reproduce and distribute reprints for Government purposes notwithstanding any copyright notation here on.

#### References

- Borgwardt, K. M., Ong, C. S., Schönauer, S., Vishwanathan, S., Smola, A. J., and Kriegel, H.-P. Protein function prediction via graph kernels. *Bioinformatics*, 21(suppl\_1): i47–i56, 2005.
- Chen, J., Yang, Z., and Yang, D. Mixtext: Linguistically-informed interpolation of hidden space for semi-supervised text classification. *arXiv* preprint *arXiv*:2004.12239, 2020.
- Debnath, A. K., Lopez de Compadre, R. L., Debnath, G., Shusterman, A. J., and Hansch, C. Structure-activity relationship of mutagenic aromatic and heteroaromatic nitro compounds. correlation with molecular orbital energies and hydrophobicity. *Journal of medicinal chemistry*, 34 (2):786–797, 1991.
- Ding, K., Zhou, Q., Tong, H., and Liu, H. Few-shot network anomaly detection via cross-network meta-learning. In *Proceedings of the Web Conference 2021*, pp. 2448–2456, 2021.
- Ding, K., Xu, Z., Tong, H., and Liu, H. Data augmentation for deep graph learning: A survey. *ACM SIGKDD Explorations Newsletter*, 24(2):61–77, 2022.
- Dykstra, R. L. An algorithm for restricted least squares regression. *Journal of the American Statistical Association*, 78(384):837–842, 1983.
- Feng, F., He, X., Tang, J., and Chua, T.-S. Graph adversarial training: Dynamically regularizing based on graph structure. *IEEE Transactions on Knowledge and Data Engineering*, 33(6):2493–2504, 2019.

- Fu, D. and He, J. Sdg: A simplified and dynamic graph neural network. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 2273–2277, 2021.
- Fu, D. and He, J. Dppin: A biological repository of dynamic protein-protein interaction network data. In 2022 IEEE International Conference on Big Data (Big Data), pp. 5269–5277. IEEE, 2022.
- Fu, D., Zhou, D., and He, J. Local motif clustering on time-evolving graphs. In *Proceedings of the 26th ACM SIGKDD International conference on knowledge discovery & data mining*, pp. 390–400, 2020.
- Fu, D., Ban, Y., Tong, H., Maciejewski, R., and He, J. Disco: comprehensive and explainable disinformation detection. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*, pp. 4848–4852, 2022.
- Fu, D., Xu, Z., Tong, H., and He, J. Natural and artificial dynamics in gnns: A tutorial. In *Proceedings of the Sixteenth ACM International Conference on Web Search* and Data Mining, pp. 1252–1255, 2023.
- Gasteiger, J., Bojchevski, A., and Günnemann, S. Predict then propagate: Graph neural networks meet personalized pagerank. In *International Conference on Learning Representations*, 2018.
- Guo, H. Nonlinear mixup: Out-of-manifold data augmentation for text classification. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pp. 4044–4051, 2020.
- Guo, H. and Mao, Y. ifmixup: Towards intrusion-free graph mixup for graph classification. *arXiv e-prints*, pp. arXiv–2110, 2021.
- Guo, H., Mao, Y., and Zhang, R. Augmenting data with mixup for sentence classification: An empirical study. *arXiv preprint arXiv:1905.08941*, 2019a.
- Guo, H., Mao, Y., and Zhang, R. Mixup as locally linear out-of-manifold regularization. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pp. 3714–3722, 2019b.
- Guo, K., Cao, X., Liu, Z., and Chang, Y. Taming oversmoothing representation on heterophilic graphs. *Infor*mation Sciences, 647:119463, 2023.
- Hamilton, W., Ying, Z., and Leskovec, J. Inductive representation learning on large graphs. *Advances in neural information processing systems*, 30, 2017.

- Han, X., Jiang, Z., Liu, N., and Hu, X. G-mixup: Graph data augmentation for graph classification. In *Interna*tional Conference on Machine Learning, pp. 8230–8248. PMLR, 2022.
- Howes, N. R. Modern Analysis and Topology. Springer Science & Business Media, 1995.
- Jing, B., Park, C., and Tong, H. Hdmi: High-order deep multiplex infomax. In *Proceedings of the Web Conference* 2021, pp. 2414–2424, 2021a.
- Jing, B., Tong, H., and Zhu, Y. Network of tensor time series. In *Proceedings of the Web Conference* 2021, pp. 2425–2437, 2021b.
- Jing, B., Feng, S., Xiang, Y., Chen, X., Chen, Y., and Tong, H. X-goal: Multiplex heterogeneous graph prototypical contrastive learning. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*, pp. 894–904, 2022a.
- Jing, B., Yan, Y., Zhu, Y., and Tong, H. Coin: Co-cluster infomax for bipartite graphs. In *NeurIPS 2022 Workshop: New Frontiers in Graph Learning*, 2022b.
- Jing, B., Yan, Y., Ding, K., Park, C., Zhu, Y., Liu, H., and Tong, H. Sterling: Synergistic representation learning on bipartite graphs. In *Proceedings of the AAAI Conference* on Artificial Intelligence, volume 38, pp. 12976–12984, 2024.
- Kan, X., Li, Z., Cui, H., Yu, Y., Xu, R., Yu, S., Zhang, Z., Guo, Y., and Yang, C. R-mixup: Riemannian mixup for biological networks. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, KDD '23, pp. 1073–1085, 2023.
- Kim, J.-H., Choo, W., and Song, H. O. Puzzle mix: Exploiting saliency and local statistics for optimal mixup. In *International Conference on Machine Learning*, pp. 5275–5285. PMLR, 2020.
- Kipf, T. N. and Welling, M. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations*, 2017.
- Kolouri, S., Naderializadeh, N., Rohde, G. K., and Hoffmann, H. Wasserstein embedding for graph learning. In *International Conference on Learning Representations*, 2021.
- Kriege, N. and Mutzel, P. Subgraph matching kernels for attributed graphs. *arXiv* preprint arXiv:1206.6483, 2012.
- Ling, H., Jiang, Z., Liu, M., Ji, S., and Zou, N. Graph mixup with soft alignments. In *International Conference* on *Machine Learning*, volume 202, pp. 21335–21349. PMLR, 2023.

- Liu, Z., Li, S., Wang, G., Wu, L., Tan, C., and Li, S. Z. Harnessing hard mixed samples with decoupled regularizer. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023a.
- Liu, Z., Zeng, Z., Qiu, R., Yoo, H., Zhou, D., Xu, Z., Zhu, Y., Weldemariam, K., He, J., and Tong, H. Topological augmentation for class-imbalanced node classification. *arXiv* preprint arXiv:2308.14181, 2023b.
- Ma, X., Chu, X., Wang, Y., Lin, Y., Zhao, J., Ma, L., and Zhu, W. Fused gromov-wasserstein graph mixup for graph-level classifications. In *Thirty-seventh Conference* on Neural Information Processing Systems, 2023.
- Maretic, H. P., El Gheche, M., Chierchia, G., and Frossard, P. Got: an optimal transport framework for graph comparison. *Advances in Neural Information Processing Systems*, 32, 2019.
- Mémoli, F. Gromov-wasserstein distances and the metric approach to object matching. *Foundations of Computational Mathematics*, 11:417–487, 2011.
- Morris, C., Kriege, N. M., Bause, F., Kersting, K., Mutzel, P., and Neumann, M. Tudataset: A collection of benchmark datasets for learning with graphs. arXiv preprint arXiv:2007.08663, 2020.
- Navarro, M. and Segarra, S. Graphmad: Graph mixup for data augmentation using data-driven convex clustering. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 1–5. IEEE, 2023.
- Neumann, M., Garnett, R., Bauckhage, C., and Kersting, K. Propagation kernels: efficient graph kernels from propagated information. *Machine learning*, 102:209–245, 2016.
- Peyré, G., Cuturi, M., and Solomon, J. Gromov-wasserstein averaging of kernel and distance matrices. In *International Conference on Machine Learning*, pp. 2664–2672. PMLR, 2016.
- Rong, Y., Huang, W., Xu, T., and Huang, J. Dropedge: Towards deep graph convolutional networks on node classification. *arXiv* preprint arXiv:1907.10903, 2019.
- Scetbon, M., Cuturi, M., and Peyré, G. Low-rank sinkhorn factorization. In *International Conference on Machine Learning*, pp. 9344–9354. PMLR, 2021.
- Scetbon, M., Peyré, G., and Cuturi, M. Linear-time gromov wasserstein distances using low rank couplings and costs. In *International Conference on Machine Learning*, pp. 19347–19365. PMLR, 2022.

- Sturm, K.-T. The space of spaces: curvature bounds and gradient flows on the space of metric measure spaces. *arXiv* preprint arXiv:1208.0434, 2012.
- Titouan, V., Courty, N., Tavenard, R., and Flamary, R. Optimal transport for structured data with application on graphs. In *International Conference on Machine Learning*, pp. 6275–6284. PMLR, 2019.
- Velickovic, P., Cucurull, G., Casanova, A., Romero, A., Lio, P., Bengio, Y., et al. Graph attention networks. *stat*, 1050 (20):10–48550, 2017.
- Verma, V., Lamb, A., Beckham, C., Najafi, A., Mitliagkas, I., Lopez-Paz, D., and Bengio, Y. Manifold mixup: Better representations by interpolating hidden states. In *International conference on machine learning*, pp. 6438–6447. PMLR, 2019.
- Verma, V., Qu, M., Kawaguchi, K., Lamb, A., Bengio, Y., Kannala, J., and Tang, J. Graphmix: Improved training of gnns for semi-supervised learning. In *Proceedings of* the AAAI conference on artificial intelligence, volume 35, pp. 10024–10032, 2021.
- Vincent-Cuaz, C., Vayer, T., Flamary, R., Corneli, M., and Courty, N. Online graph dictionary learning. In *Interna*tional conference on machine learning, pp. 10564–10574. PMLR, 2021.
- Wang, D., Yan, Y., Qiu, R., Zhu, Y., Guan, K., Margenot, A., and Tong, H. Networked time series imputation via position-aware graph enhanced variational autoencoders. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 2256–2268, 2023.
- Wang, Y., Wang, W., Liang, Y., Cai, Y., and Hooi, B. Mixup for node and graph classification. In *Proceedings of the Web Conference 2021*, pp. 3663–3674, 2021.
- Warmsley, D., Waagen, A., Xu, J., Liu, Z., and Tong, H. A survey of explainable graph neural networks for cyber malware analysis. In 2022 IEEE International Conference on Big Data (Big Data), pp. 2932–2939. IEEE, 2022.
- Wei, T. and He, J. Comprehensive fair meta-learned recommender system. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 1989–1999, 2022.
- Wei, T., Wu, Z., Li, R., Hu, Z., Feng, F., He, X., Sun, Y., and Wang, W. Fast adaptation for cold-start collaborative filtering with meta-learning. In 2020 IEEE International Conference on Data Mining (ICDM), pp. 661–670. IEEE, 2020.

- Wei, T., You, Y., Chen, T., Shen, Y., He, J., and Wang, Z. Augmentations in hypergraph contrastive learning: Fabricated and generative. *Advances in neural information processing systems*, 35:1909–1922, 2022.
- Wei, T., Guo, Z., Chen, Y., and He, J. Ntk-approximating mlp fusion for efficient language model fine-tuning. In *International Conference on Machine Learning*, pp. 36821–36838. PMLR, 2023.
- Wei, T., Jin, B., Li, R., Zeng, H., Wang, Z., Sun, J., Yin, Q., Lu, H., Wang, S., He, J., et al. Towards unified multimodal personalization: Large vision-language models for generative recommendation and beyond. arXiv preprint arXiv:2403.10667, 2024.
- Xu, H., Luo, D., and Carin, L. Scalable gromov-wasserstein learning for graph partitioning and matching. Advances in Neural Information Processing Systems, 32, 2019.
- Xu, K., Hu, W., Leskovec, J., and Jegelka, S. How powerful are graph neural networks? In *International Conference on Learning Representations*, 2018.
- Xu, Z., Ding, K., Wang, Y.-X., Liu, H., and Tong, H. Generalized few-shot node classification. In 2022 IEEE International Conference on Data Mining (ICDM), pp. 608–617. IEEE, 2022a.
- Xu, Z., Du, B., and Tong, H. Graph sanitation with application to node classification. In *Proceedings of the ACM Web Conference* 2022, pp. 1136–1147, 2022b.
- Xu, Z., Chen, Y., Pan, M., Chen, H., Das, M., Yang, H., and Tong, H. Kernel ridge regression-based graph dataset distillation. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 2850–2861, 2023a.
- Xu, Z., Chen, Y., Zhou, Q., Wu, Y., Pan, M., Yang, H., and Tong, H. Node classification beyond homophily:
   Towards a general solution. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 2862–2873, 2023b.
- Xu, Z., Qiu, R., Chen, Y., Chen, H., Fan, X., Pan, M., Zeng, Z., Das, M., and Tong, H. Discrete-state continuoustime diffusion for graph generation. arXiv preprint arXiv:2405.11416, 2024.
- Yan, Y., Liu, L., Ban, Y., Jing, B., and Tong, H. Dynamic knowledge graph alignment. In *Proceedings of the AAAI* conference on artificial intelligence, volume 35, pp. 4564– 4572, 2021.
- Yan, Y., Zhou, Q., Li, J., Abdelzaher, T., and Tong, H. Dissecting cross-layer dependency inference on multi-layered inter-dependent networks. In *Proceedings of the*

- 31st ACM International Conference on Information & Knowledge Management, pp. 2341–2351, 2022.
- Yan, Y., Chen, Y., Chen, H., Xu, M., Das, M., Yang, H., and Tong, H. From trainable negative depth to edge heterophily in graphs. *Advances in Neural Information Processing Systems*, 36, 2024a.
- Yan, Y., Hu, Y., Zhou, Q., Liu, L., Zeng, Z., Chen, Y., Pan, M., Chen, H., Das, M., and Tong, H. Pacer: Network embedding from positional to structural. In *Proceedings of the ACM on Web Conference 2024*, pp. 2485–2496, 2024b.
- Yan, Y., Jing, B., Liu, L., Wang, R., Li, J., Abdelzaher, T., and Tong, H. Reconciling competing sampling strategies of network embedding. *Advances in Neural Information Processing Systems*, 36, 2024c.
- Yanardag, P. and Vishwanathan, S. Deep graph kernels. In *Proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining*, pp. 1365–1374, 2015.
- Yang, L., Zhang, L., and Yang, W. Graph adversarial self-supervised learning. Advances in Neural Information Processing Systems, 34:14887–14899, 2021.
- Yoo, H., Zeng, Z., Kang, J., Liu, Z., Zhou, D., Wang, F., Chan, E., and Tong, H. Ensuring user-side fairness in dynamic recommender systems. *arXiv* preprint *arXiv*:2308.15651, 2023.
- Yoo, J., Shim, S., and Kang, U. Model-agnostic augmentation for accurate graph classification. In *Proceedings of the ACM Web Conference* 2022, pp. 1281–1291, 2022.
- You, Y., Chen, T., Sui, Y., Chen, T., Wang, Z., and Shen, Y. Graph contrastive learning with augmentations. *Advances in neural information processing systems*, 33:5812–5823, 2020.
- You, Y., Chen, T., Shen, Y., and Wang, Z. Graph contrastive learning automated. In *International Conference on Machine Learning*, pp. 12121–12132. PMLR, 2021.
- Zeng, Z., Zhang, S., Xia, Y., and Tong, H. Parrot: Position-aware regularized optimal transport for network alignment. In *Proceedings of the ACM Web Conference 2023*, pp. 372–382, 2023a.
- Zeng, Z., Zhu, R., Xia, Y., Zeng, H., and Tong, H. Generative graph dictionary learning. In *International Conference on Machine Learning*, pp. 40749–40769. PMLR, 2023b.
- Zeng, Z., Du, B., Zhang, S., Xia, Y., Liu, Z., and Tong, H. Hierarchical multi-marginal optimal transport for network alignment. In *Proceedings of the AAAI Conference*

- *on Artificial Intelligence*, volume 38, pp. 16660–16668, 2024.
- Zhang, H., Cisse, M., Dauphin, Y. N., and Lopez-Paz, D. mixup: Beyond empirical risk minimization. In *International Conference on Learning Representations*, 2018.
- Zhang, L., Deng, Z., Kawaguchi, K., Ghorbani, A., and Zou, J. How does mixup help with robustness and generalization? In *International Conference on Learning Representations*, 2021.
- Zhang, M. and Chen, Y. Link prediction based on graph neural networks. *Advances in neural information processing systems*, 31, 2018.
- Zhang, S., Tong, H., Xia, Y., Xiong, L., and Xu, J. Nettrans: Neural cross-network transformation. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 986–996, 2020.

# **Appendix**

The content of the appendix is organized as follows:

- Appendix A includes the proofs for all the theorems in the main content.
- Appendix B provides the details of GEOMIX and discuss potential variants:
  - Appendix B.1 provides details for the Dykstra's and Low-rank GW algorithm.
  - Appendix B.2 generalizes GEOMIX to attributed graphs with node attributes.
- Appendix C provides additional experimental results, including sensitivity analysis, mixup graph visualization, and embedding space visualization.
- Appendix D details the reproducibility, including dataset description, model architecture, and experiment pipeline.
- Appendix E discusses the limitations and possible future directions of this work.

## A. Proof

**Theorem 3.1.** Given two graphs  $G_1, G_2$ , the transformed graphs  $\tilde{G}_1, \tilde{G}_2$  in Eq. (5) are in the equivalent class of  $G_1, G_2$ , respectively, that is

$$\llbracket \tilde{\mathcal{G}}_1 \rrbracket = \llbracket \mathcal{G}_1 \rrbracket, \llbracket \tilde{\mathcal{G}}_2 \rrbracket = \llbracket \mathcal{G}_2 \rrbracket$$

Furthermore, the curve  $\{ \llbracket \tilde{\mathcal{G}}_{(\lambda)} \rrbracket \}_{0 \leq \lambda \leq 1}$  is a GW geodesic connecting  $\llbracket \mathcal{G}_1 \rrbracket$  and  $\llbracket \mathcal{G}_2 \rrbracket$ , which means that for every  $\lambda_1, \lambda_2 \in [0, 1]$ ,

$$d_{\mathrm{GW}}^*(\llbracket \tilde{\mathcal{G}}_{(\lambda_1)} \rrbracket, \llbracket \tilde{\mathcal{G}}_{(\lambda_2)} \rrbracket) = |\lambda_1 - \lambda_2| \cdot d_{\mathrm{GW}}^*(\llbracket \mathcal{G}_1 \rrbracket, \llbracket \mathcal{G}_2 \rrbracket).$$

*Proof.* We first prove that  $\mathcal{G}_1 \sim \tilde{\mathcal{G}}_1$  in terms of the GW distance. According to Eq. (2), the GW distance between  $\mathcal{G}_1$  and  $\tilde{\mathcal{G}}_1$  can be written as follows:

$$d_{\text{GW}}(\mathcal{G}_1, \tilde{\mathcal{G}}_1) = \underset{\boldsymbol{T} \in \Pi(\boldsymbol{\mu}_1, \tilde{\boldsymbol{\mu}}_1)}{\arg \min} \sum_{\substack{i, j \\ (x_1, y_1), (x_2, y_2)}} \left[ \boldsymbol{A}_1(i, j) - \tilde{\boldsymbol{A}}_1((x_1, y_1), (x_2, y_2)) \right]^2 \boldsymbol{T}(i, (x_1, y_1)) \boldsymbol{T}(j, (x_2, y_2))$$

It is easy to verify that the marginal distributions satisfy  $\tilde{\mu}_1 = \tilde{P}_1 \mu_1 = \mu_1 \otimes \frac{\mathbf{1}_{n_2}}{n_2}$  and  $\tilde{\mu}_2 = \tilde{P}_2 \mu_2 = \mu_2 \otimes \frac{\mathbf{1}_{n_1}}{n_1}$ . We consider a naive coupling  $T_1 \in \Pi(\mu_1, \tilde{\mu}_1)$  as follows:

$$T_1(i, (x_1, y_1)) = \begin{cases} \frac{\mu_1(i)}{n_2}, & \text{if } i = x_1 \\ 0, & \text{else} \end{cases},$$
 (10)

which leads to a (sub)optimal solution to the GW distance, that is:

$$\begin{split} d_{\mathrm{GW}}(\mathcal{G}_{1},\tilde{\mathcal{G}}_{1}) &\leq \sum_{\stackrel{i,j}{(x_{1},y_{1}),(x_{2},y_{2})}} \left[ \boldsymbol{A}_{1}(i,j) - \tilde{\boldsymbol{A}}_{1}((x_{1},y_{1}),(x_{2},y_{2})) \right]^{2} \boldsymbol{T}_{1}(i,(x_{1},y_{1})) \boldsymbol{T}_{1}(j,(x_{2},y_{2})) \\ &= \sum_{\stackrel{i,j}{(x_{1},y_{1}),(x_{2},y_{2})}} \left[ \boldsymbol{A}_{1}(i,j) - \sum_{i,j} \boldsymbol{P}_{1}(i,(x_{1},y_{1})) \boldsymbol{A}_{1}(i,j) \boldsymbol{P}_{1}(j,(x_{2},y_{2})) \right]^{2} \boldsymbol{T}_{1}(i,(x_{1},y_{1})) \boldsymbol{T}_{1}(j,(x_{2},y_{2})) \\ &= \sum_{\stackrel{i,j}{(x_{1},y_{1}),(x_{2},y_{2})}} \left[ \boldsymbol{A}_{1}(i,j) - \boldsymbol{A}_{1}(x_{1},x_{2}) \right]^{2} \boldsymbol{T}_{1}(i,(x_{1},y_{1})) \boldsymbol{T}_{1}(j,(x_{2},y_{2})) \quad \text{due to Eq. (5)} \\ &= \sum_{i,j} \left[ \boldsymbol{A}_{1}(i,j) - \boldsymbol{A}_{1}(i,j) \right]^{2} \frac{\boldsymbol{\mu}_{1}(i)\boldsymbol{\mu}_{1}(j)}{n_{2}^{2}} \quad \text{due to Eq. (10)} \\ &= 0 \end{split}$$

Since the GW distance is non-negative, we have  $d_{\text{GW}}(\mathcal{G}_1, \tilde{\mathcal{G}}_1) = 0$ . Similarly, we can also show that  $d_{\text{GW}}(\mathcal{G}_2, \tilde{\mathcal{G}}_2) = 0$ . Therefore, we prove  $\mathcal{G}_1 \sim \tilde{\mathcal{G}}_1, \mathcal{G}_2 \sim \tilde{\mathcal{G}}_2$ . Then we show that the mixup graph  $\mathcal{G}_{(\lambda)} = (A_{(\lambda)}, \mu_{(\lambda)})$  defines a geodesic connecting  $\mathcal{G}_1$  and  $\mathcal{G}_2$  in the GW space. We can write down  $A_{(\lambda)}((x_1, y_1), (x_2, y_2))$  as:

$$\begin{split} \boldsymbol{A}_{\lambda}((x_{1},y_{1}),(x_{2},y_{2})) &= (1-\lambda) \sum_{i,j \in \mathcal{G}_{1}} \underbrace{\boldsymbol{P}_{1}(i,(x_{1},y_{1}))}_{1[i=x_{1}]} \boldsymbol{A}_{1}(i,j) \underbrace{\boldsymbol{P}_{1}(j,(x_{2},y_{2}))}_{1[j=x_{2}]} + \lambda \sum_{i,j \in \mathcal{G}_{2}} \underbrace{\boldsymbol{P}_{2}(i,(x_{1},y_{1}))}_{1[i=y_{1}]} \boldsymbol{A}_{2}(i,j) \underbrace{\boldsymbol{P}_{2}(j,(x_{2},y_{2}))}_{1[j=y_{2}]} \\ &= (1-\lambda) \boldsymbol{A}_{1}(x_{1},x_{2}) + \lambda \boldsymbol{A}_{2}(y_{1},y_{2}) \end{split}$$

Since  $\mu_{(\lambda)} = (1 - \lambda)\tilde{\mu}_1 + \lambda\tilde{\mu}_2 = \text{vecOT}(\mathcal{G}_1, \mathcal{G}_2)$  in Eqs. (5) and (6), according to Eq. (4),  $\mathcal{G}_{(\lambda)}$  corresponds to the geodesic connecting  $\mathcal{G}_1, \mathcal{G}_2$  in the GW space.

**Theorem 3.2.** The tuple  $(P_1, P_2, \tilde{\mu}_{(\lambda)})$  in Eqs. (4) and (5) for exact GW geodesic is an optimal solution to the following optimization problem:

$$\min_{\boldsymbol{P}_{1},\boldsymbol{P}_{2},\boldsymbol{g}} \left( \varepsilon_{\boldsymbol{A}_{1},\boldsymbol{A}_{2}} (\boldsymbol{P}_{1} \operatorname{diag}(\boldsymbol{g}) \boldsymbol{P}_{2}^{\mathsf{T}}) \right)^{1/2}$$
s.t. 
$$\begin{cases}
\boldsymbol{P}_{1} \operatorname{diag}(\boldsymbol{g}) \in \Pi(\boldsymbol{\mu}_{1},\boldsymbol{g}) \\
\boldsymbol{P}_{2} \operatorname{diag}(\boldsymbol{g}) \in \Pi(\boldsymbol{\mu}_{2},\boldsymbol{g}) \\
\boldsymbol{g} \in \Delta_{n_{1}n_{2}}
\end{cases} \tag{7}$$

*Proof.* We first show that  $(P_1, P_2, \tilde{\mu}_{(\lambda)})$  satisfies the constraints of the optimization problem. By definition in Eq. (4), we have  $g = \tilde{\mu}_{(\lambda)} = \text{vec}(T) \in \Delta_{n_1 n_2}$ , where  $T \in \Pi(\mu_1, \mu_2)$  is the optimal coupling between  $A_1$  and  $A_2$  in terms of the GW distance. Note that

$$\sum_{i} (\mathbf{P}_{1} \operatorname{diag}(\mathbf{g}))(x_{i}, (x_{j}, y_{k})) = \sum_{i} \mathbf{P}_{1}(x_{i}, (x_{j}, y_{k})) \mathbf{g}((x_{j}, y_{k})) = \mathbf{T}(x_{j}, y_{k})$$

$$\sum_{j,k} (\mathbf{P}_{1} \operatorname{diag}(\mathbf{g}))(x_{i}, (x_{j}, y_{k})) = \sum_{j,k} \mathbf{P}_{1}(x_{i}, (x_{j}, y_{k})) \mathbf{g}((x_{j}, y_{k})) = \sum_{k} \mathbf{T}(x_{i}, y_{k}) = \boldsymbol{\mu}_{1}(x_{i})$$

Therefore, we have  $P_1 \operatorname{diag}(g) \in \Pi(\mu_1, g)$ . Similarly, we can show  $P_2 \operatorname{diag}(g) \in \Pi(\mu_2, g)$ .

To validate that  $P_1 \operatorname{diag}(\tilde{\mu}_{(\lambda)}) P_2^{\mathsf{T}}$  is an optimal coupling between  $A_1$  and  $A_2$ , it suffices to show that  $P_1 \operatorname{diag}(\tilde{\mu}_{(\lambda)}) P_2^{\mathsf{T}} = \operatorname{vec}(T)$  as follows

$$\begin{split} \left(\boldsymbol{P}_{1} \mathrm{diag}(\mathrm{vec}(\boldsymbol{T})) \boldsymbol{P}_{2}^{\mathsf{T}}\right)(x_{i}, y_{k}) &= \sum_{j, l} \boldsymbol{P}_{1}(x_{i}, (x_{j}, y_{l})) \mathrm{vec}(\boldsymbol{T})((x_{j}, y_{l})) \boldsymbol{P}_{2}(y_{k}, (x_{j}, y_{l})) \\ &= \sum_{j, l} 1[i = j] \boldsymbol{T}(x_{j}, y_{l}) 1[k = l] \\ &= \boldsymbol{T}(x_{i}, y_{k}) \end{split}$$

Therefore, we prove that  $P_1 \operatorname{diag}(\operatorname{vec}(T)) P_2^{\mathsf{T}}$  is exactly the optimal coupling between  $A_1$  and  $A_2$ . In other words, the tuple  $(P_1, P_2, \tilde{\mu}_{(\lambda)})$  for exact GW geodesic corresponds to the solution to the given optimization problem.

## B. Algorithm

#### **B.1. Detailed Algorithm**

In this section, we present the Dykstra's algorithm (Scetbon et al., 2021) and low-rank Gromov–Wasserstein (Scetbon et al., 2022) for completeness.

**Dykstra's algorithm (Dykstra, 1983; Scetbon et al., 2021).** Given  $\boldsymbol{\xi}_1 \in \mathbb{R}^{n_1 \times r}_+, \boldsymbol{\xi}_2 \in \mathbb{R}^{n_2 \times r}_+, \boldsymbol{\xi}_3 \in \mathbb{R}^r_+$ , and a convex set  $\mathcal{C} \subset \mathbb{R}^{n_1 \times r}_+ \times \mathbb{R}^{n_2 \times r}_+ \times \mathbb{R}^r_+$ , the Dykstra's algorithm seeks for a set of the optimal variables  $(\boldsymbol{Q}_1^*, \boldsymbol{Q}_2^*, \boldsymbol{g}^*) \in \mathcal{C}$  that is closest to  $(\boldsymbol{\xi}_1, \boldsymbol{\xi}_2, \boldsymbol{\xi}_3)$  in terms of the generalized KL divergence, that is:

$$(\boldsymbol{Q}_1^*,\boldsymbol{Q}_2^*,\boldsymbol{g}^*) = \mathop{\arg\min}_{(\boldsymbol{Q}_1,\boldsymbol{Q}_2,\boldsymbol{g}) \subset \mathcal{C}} \mathrm{KL}((\boldsymbol{Q}_1,\boldsymbol{Q}_2,\boldsymbol{g}),(\boldsymbol{\xi}_1,\boldsymbol{\xi}_2,\boldsymbol{\xi}_3))$$

## Algorithm 2 Dykstra's Algorithm (LR-Dykstra) (Scetbon et al., 2021)

```
1: Input target values \boldsymbol{\xi}_1, \boldsymbol{\xi}_2, \boldsymbol{\xi}_3, marginals \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, lower bound \boldsymbol{\alpha}, error threshold \boldsymbol{\delta}

2: Initialize scaling vectors \boldsymbol{q}_1 = \boldsymbol{q}_2 = \boldsymbol{q}_3^{(1)} = \boldsymbol{q}_3^{(2)} = \tilde{\boldsymbol{v}}_1 = \tilde{\boldsymbol{v}}_2 = \mathbf{1}_r, \tilde{\boldsymbol{g}} = \boldsymbol{\xi}_3;

3: \mathbf{while} \sum_{i=1}^2 \|\boldsymbol{u}_i \odot \boldsymbol{\xi}_i \boldsymbol{v}_i - \boldsymbol{\mu}_i\|_1 \geq \boldsymbol{\delta} \, \mathbf{do}

4: Update row-normalization vectors \boldsymbol{u}_i = \boldsymbol{\mu}_i/\boldsymbol{\xi}_i \tilde{\boldsymbol{v}}_i, \forall i = \{1, 2\};

5: Stabilize \boldsymbol{g} = \max(\boldsymbol{\alpha}, \tilde{\boldsymbol{g}} \odot \boldsymbol{q}_3^{(1)}), update scale vector \boldsymbol{q}_3^{(1)} = (\tilde{\boldsymbol{g}} \odot \boldsymbol{q}_3^{(1)})/\boldsymbol{g}, and update \tilde{\boldsymbol{g}} = \boldsymbol{g};

6: Update \boldsymbol{g} = (\tilde{\boldsymbol{g}} \odot \boldsymbol{q}_3^{(2)})^{1/3} \prod_{i=1}^2 (\boldsymbol{v}_i \odot \boldsymbol{q}_i \odot \boldsymbol{\xi}_i^{\mathsf{T}} \boldsymbol{u}_i)^{1/3};

7: Update column-normalization vectors \boldsymbol{v}_i = \boldsymbol{g}/\boldsymbol{\xi}_i^{\mathsf{T}} \boldsymbol{u}_i, \forall i \in \{1, 2\};

8: Update scale vectors \boldsymbol{q}_i = (\tilde{\boldsymbol{v}}_i \odot \boldsymbol{q}_i)/\boldsymbol{v}_i, \tilde{\boldsymbol{v}}_i = \boldsymbol{v}_i, \forall i \in \{1, 2\}, \boldsymbol{q}_3^{(2)} = (\tilde{\boldsymbol{g}} \odot \boldsymbol{q}_3^{(2)})/\boldsymbol{g}, \tilde{\boldsymbol{g}} = \boldsymbol{g};

9: end while

10: Compute \boldsymbol{Q}_1 = \operatorname{diag}(\boldsymbol{u}_1)\boldsymbol{\xi}_1 \operatorname{diag}(\boldsymbol{v}_1), \boldsymbol{Q}_2 = \operatorname{diag}(\boldsymbol{u}_2)\boldsymbol{\xi}_2 \operatorname{diag}(\boldsymbol{v}_2);

11: return (\boldsymbol{Q}_1, \boldsymbol{Q}_2, \boldsymbol{g});
```

## Algorithm 3 Low-rank Gromov–Wasserstein (Scetbon et al., 2022)

```
1: Input cost matrices A_1, A_2, marginals \mu_1, \mu_2, rank r, step size \gamma, lower bound \alpha, error threshold \delta.

2: Initialize x_1 = A_1^{\odot 2} \mu_1, x_2 = A_2^{\odot 2} \mu_2, z_1 = x_1^{\odot 2}, z_2 = x_2^{\odot 2}, g \in \Delta_r, Q_1 \in \Pi(\mu_1, g), Q_2\Pi(\mu_2, g);

3: for t = 1, 2, ... do

4: \xi_1 = Q_1 \odot \exp(-4\gamma A_1 Q_1 \operatorname{diag}(1/g) Q_2^{\mathsf{T}} A_2 Q_2 \operatorname{diag}(1/g));

5: \xi_2 = Q_2 \odot \exp(-4\gamma A_1^{\mathsf{T}} Q_2 \operatorname{diag}(1/g) Q_1^{\mathsf{T}} A_1^{\mathsf{T}} Q_1 \operatorname{diag}(1/g));

6: \omega = \operatorname{diag}(-Q_1^{\mathsf{T}} A_1 Q_1 \operatorname{diag}(1/g) Q_2^{\mathsf{T}} A_2 Q_2);

7: \xi_3 = g \odot \exp(-4\gamma \omega/g^2);

8: Q_1, Q_2, g = \operatorname{LR-Dykstra}(\xi_1, \xi_2, \xi_3, \mu_1, \mu_2, \alpha, \delta);

9: end for

10: return (Q_1, Q_2, g).
```

The general idea of Dykstra's algorithm is to rescale  $\xi_1, \xi_2$  and project onto the convex hull  $\mathcal{C}$ . By iterating the above process, Dykstra's algorithm is shown to converge to the desired solution (Dykstra, 1983). The Dykstra's algorithm is shown in Algorithm 2.

Low-rank Gromov-Wasserstein (Scetbon et al., 2022). Given two adjacency matrices  $A_1 \in \mathbb{R}^{n_1 \times n_1}$ ,  $A_2 \in \mathbb{R}^{n_2 \times n_2}$ , the low-rank GW problem seeks for the optimal low-rank decomposition of the optimal coupling  $T = Q_1 \operatorname{diag}(1/g)Q_2^{\mathsf{T}}$  to the GW distance, that is ](Scetbon et al., 2022),

$$\begin{split} &(\boldsymbol{Q}_1, \boldsymbol{Q}_2, \boldsymbol{g}) = \mathop{\arg\min}_{\boldsymbol{Q}_1, \boldsymbol{Q}_2, \boldsymbol{g}} (\varepsilon_{\boldsymbol{A}_1, \boldsymbol{A}_2} (\boldsymbol{Q}_1 \mathrm{diag}(1/\boldsymbol{g}) \boldsymbol{Q}_2^{\mathsf{T}}))^{1/p} \\ \mathrm{s.t.} \; &\boldsymbol{Q}_1 \in \Pi(\boldsymbol{\mu}_1, \boldsymbol{g}), \boldsymbol{Q}_2 \in \Pi(\boldsymbol{\mu}_2, \boldsymbol{g}), \boldsymbol{g} \in \Delta_r \end{split}$$

To solve the low-rank Gromov–Wasserstein problem in Eq. (8), (Scetbon et al., 2022) provides a mirror descent scheme. The general idea is to first solve the proximal point iteration under the KL divergence regularization  $\boldsymbol{\xi}_i^{(t+1)} = \arg\min_{\boldsymbol{\xi}_i} \nabla_{\boldsymbol{\xi}_i} (\varepsilon_{\boldsymbol{A}_1, \boldsymbol{A}_2} (\boldsymbol{\xi}_1 \mathrm{diag}(1/\boldsymbol{\xi}_3) \boldsymbol{\xi}_2^{\mathsf{T}}))^{1/p} + \mathrm{KL}(\boldsymbol{\xi}_i, \boldsymbol{\xi}_i^{(t)})$  and project back to the feasible space  $\mathcal{C} = \{(\boldsymbol{Q}_1, \boldsymbol{Q}_2, \boldsymbol{g}) | \text{s.t. } \boldsymbol{Q}_1 \in \Pi(\boldsymbol{\mu}_1, \boldsymbol{g}), \boldsymbol{Q}_2 \in \Pi(\boldsymbol{\mu}_2, \boldsymbol{g}), \boldsymbol{g} \in \Delta_r\}$  using Dykstra's algorithm. We consider the most widely adopted 2-GW case in this paper, and provide the overall low-rank GW algorithm in Algorithm 3.

#### **B.2. Variants of GEOMIX**

In the main content, we focus on the plain graphs w/o node attributes. when node attributes are available, we provide the following variant of GEOMIX to incorporate node attributes in the mixup process. Specifically, when node attributes  $X_1 \in \mathbb{R}^{n_1 \times d}, X_2 \in \mathbb{R}^{n_2 \times d}$  are available, we can calculate a cross-graph cost matrix  $C \in \mathbb{R}^{n_1 \times n_2}$ , e.g., using Euclidean distance  $C(i,j) = \|X_1(i,:) - X_2(j,:)\|_2^2$ . Similarly, following the mirror descent for solution, and each subproblem can

be written as follows:

$$\begin{pmatrix} \boldsymbol{Q}_{1}^{(t+1)}, \boldsymbol{Q}_{2}^{(t+1)}, \boldsymbol{g}^{(t+1)} \end{pmatrix} = \underset{\boldsymbol{Q}_{1}, \boldsymbol{Q}_{2}, \boldsymbol{g}}{\operatorname{arg \, min}} \operatorname{KL} \left( (\boldsymbol{Q}_{1}, \boldsymbol{Q}_{2}, \boldsymbol{g}), (\boldsymbol{K}_{1}^{(t)}, \boldsymbol{K}_{2}^{(t)}, \boldsymbol{K}_{3}^{(t)}) \right)$$
s.t.  $\boldsymbol{Q}_{1} \in \Pi(\boldsymbol{\mu}_{1}, \boldsymbol{g}), \ \boldsymbol{Q}_{2} \in \Pi(\boldsymbol{\mu}_{2}, \boldsymbol{g}), \ \boldsymbol{g} \in \Delta_{r}$ 

$$\begin{cases} \boldsymbol{K}_{1}^{(t)} = \exp\left( \gamma[-\alpha \boldsymbol{C} + 4(1-\alpha)\boldsymbol{A}_{1}\boldsymbol{T}^{(t)}\boldsymbol{A}_{2}]\boldsymbol{Q}_{2}^{(t)}\operatorname{diag}(1/\boldsymbol{g}^{(t)}) + \log(\boldsymbol{Q}_{1}^{(t)}) \right) \\ \boldsymbol{K}_{2}^{(t)} = \exp\left( \gamma[-\alpha \boldsymbol{C}^{\mathsf{T}} + 4(1-\alpha)\boldsymbol{A}_{2}\boldsymbol{T}^{(t)^{\mathsf{T}}}\boldsymbol{A}_{1}]\boldsymbol{Q}_{1}^{(t)}\operatorname{diag}(1/\boldsymbol{g}^{(t)}) + \log(\boldsymbol{Q}_{2}^{(t)^{\mathsf{T}}}) \right) \\ \boldsymbol{K}_{3}^{(t)} = \exp\left( \gamma\operatorname{diag}(\boldsymbol{Q}_{1}^{(t)^{\mathsf{T}}}[-\alpha \boldsymbol{C} + 4(1-\alpha)\boldsymbol{A}_{2}\boldsymbol{T}^{(t)^{\mathsf{T}}}\boldsymbol{A}_{1}]\boldsymbol{Q}_{2}^{(t)})/\boldsymbol{g}^{(t)^{\mathsf{T}}} + \log(\boldsymbol{g}^{(t)}) \right) \\ \boldsymbol{T}^{(t)} = \boldsymbol{Q}_{1}^{(t)}\operatorname{diag}(1/\boldsymbol{g}^{(k)})\boldsymbol{Q}_{2}^{(t)^{\mathsf{T}}} \end{cases}$$

where  $\alpha \in [0,1]$  is the weight hyperparameter balancing the importance of structure and node attributes. We show that the computation for attributed and plain networks share the same big-O time complexity when adopting the Euclidean distance between node attributes as the cross-graph cost matrix, i.e.,  $C = X_1^2 \mathbf{1}_{d \times n_2} + \mathbf{1}_{n_1 \times d} X_2^{2^{\mathsf{T}}} - 2X_1 X_1^{\mathsf{T}}$ .

Compared to the computation for plain networks w/o node attributes, the computation in Eq. (11) involves three additional computations  $CQ_2^{(t)}\mathrm{diag}(1/g^{(t)}),\ C^{\mathsf{T}}Q_1^{(t)}\mathrm{diag}(1/g^{(t)}),\$ and  $Q_1^{(t)^{\mathsf{T}}}CQ_2^{(t)}/g_k^2.$  Given the couplings  $Q_1^{(t)}\in\Pi(\mu_1,g^{(t)}),Q_2^{(t)}\in\Pi(\mu_2,g^{(t)})$  and simplex  $g^{(t)}\in\Delta_r$ , we have  $\mathbf{1}_{d\times n_1}Q_1^{(t)}\mathrm{diag}(1/g^{(t)})=\mathbf{1}_{d\times n_2}Q_2^{(t)}\mathrm{diag}(1/g^{(t)})=\mathbf{1}_{d\times r}$ , we have the following time complexity analysis for calculating the additional three terms:

$$\begin{split} &C\boldsymbol{Q}_{2}^{(t)}\mathrm{diag}\left(1/\boldsymbol{g}^{(t)}\right) = \underbrace{\boldsymbol{X}_{1}^{2}\boldsymbol{1}_{d\times n_{2}}}_{\mathcal{O}(ndr)} + \underbrace{\boldsymbol{1}_{n_{1}\times d}\boldsymbol{X}_{2}^{2^{\mathsf{T}}}\boldsymbol{Q}_{2}^{(t)}\mathrm{diag}\left(1/\boldsymbol{g}^{(t)}\right)}_{\mathcal{O}(ndr)} - \underbrace{2\boldsymbol{X}_{1}\boldsymbol{X}_{2}^{\mathsf{T}}\boldsymbol{Q}_{2}^{(t)}\mathrm{diag}\left(1/\boldsymbol{g}^{(t)}\right)}_{\mathcal{O}(ndr)} \\ &C^{\mathsf{T}}\boldsymbol{Q}_{1}^{(t)}\mathrm{diag}\left(1/\boldsymbol{g}^{(t)}\right) = \underbrace{\boldsymbol{1}_{n_{2}\times d}\boldsymbol{X}_{1}^{2^{\mathsf{T}}}\boldsymbol{Q}_{1}^{(t)}\mathrm{diag}\left(1/\boldsymbol{g}^{(t)}\right)}_{\mathcal{O}(ndr)} + \underbrace{\boldsymbol{X}_{2}^{2^{\mathsf{T}}}\boldsymbol{1}_{d\times n_{2}}}_{\mathcal{O}(ndr)} - \underbrace{2\boldsymbol{X}_{2}^{\mathsf{T}}\boldsymbol{X}_{1}\boldsymbol{Q}_{1}^{(t)}\mathrm{diag}\left(1/\boldsymbol{g}^{(t)}\right)}_{\mathcal{O}(ndr)} \\ &\mathrm{diag}\left(\boldsymbol{Q}_{1}^{(t)^{\mathsf{T}}}\boldsymbol{C}\boldsymbol{Q}_{2}^{(t)}\right)/\boldsymbol{g}^{(t)^{2}} = \underbrace{\mathrm{diag}\left(\boldsymbol{Q}_{1}^{(t)^{\mathsf{T}}}\boldsymbol{X}_{1}^{2}\boldsymbol{1}_{d\times r} + \boldsymbol{1}_{r\times d}\boldsymbol{X}_{2}^{2^{\mathsf{T}}}\boldsymbol{Q}_{2}^{(t)}\right)/\boldsymbol{g}^{(t)}}_{\mathcal{O}(ndr)} - \underbrace{2\mathrm{diag}\left(\boldsymbol{Q}_{1}^{(t)^{\mathsf{T}}}\boldsymbol{X}_{1}\boldsymbol{X}_{2}^{\mathsf{T}}\boldsymbol{Q}_{2}\right)/\boldsymbol{g}^{(t)^{2}}}_{\mathcal{O}(ndr)} \end{split}$$

Therefore, the additional time complexity for node attributes is  $\mathcal{O}(KTndr)$ , which is smaller than the  $o(KTn^2r)$  of the original GEOMIX, i.e., considering node attributes does not carry additional time complexity in terms of the big-O notation.

# C. Additional experiments

Sensitivity analysis of mixup ratio  $\lambda$ . We study the effect of mixup ratio  $\lambda$  on the model performance. The mixup ratio  $\lambda$  is randomly sampled from a uniform distribution  $\mathrm{Unif}(0,\lambda_{max})$ , and we test how the model performance changes as  $\lambda_{max}$  changes from 0.05 to 0.95. Experiment results are shown in Figure 8. It is shown that GEOMIX is quite stable with different selections of  $\lambda$ .

**Mixup graph visualization.** We provide additional visualization for mixup graphs on real-world datasets MUTAG, IMDB-B, and IMDB-M in Figure 9. It is shown that  $\tilde{\mathcal{G}}_{(\lambda)}$  smoothly transforms from  $\mathcal{G}_1$  to  $\mathcal{G}_2$  when  $\lambda$  increases, indicating that our mixup samples are consistent to their mixup labels, hence avoiding misleading supervision by suspicious sample–label pairs.

Besides, for two graphs  $\mathcal{G}_1, \mathcal{G}_2$ , the transformed graph  $\Gamma_1(\mathcal{G}_1) = \tilde{\mathcal{G}}_{(0)}, \Gamma_2(\mathcal{G}_2) = \tilde{\mathcal{G}}_{(1)}$  share the same labels as  $\mathcal{G}_1, \mathcal{G}_2$ , i.e.,  $\tilde{\boldsymbol{y}}_{(0)} = \boldsymbol{y}_1, \tilde{\boldsymbol{y}}_{(1)} = \boldsymbol{y}_2$ . According to sample–label consistency, , respectively. As shown in Figure 10,  $\mathcal{G}_1, \mathcal{G}_2$  are quite similar to  $\tilde{\mathcal{G}}_{(0.0)}, \tilde{\mathcal{G}}_{(1.0)}$ , which validates that the employed transformations are equivalence-preserving and GEOMIX generate consistent sample–label pairs.

**Embedding space visualization.** We provide additional embedding space visualization for real-world datasets IMDB-M, PROTEINS, and MSRC-9 in Figure 10. These visualization results explain the experiment results in Table 1. For dataset

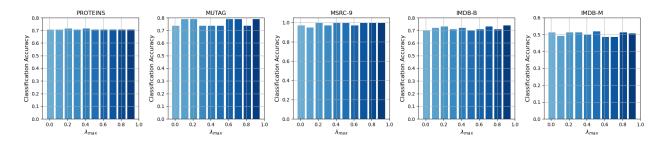


Figure 8. Sensitivity analysis on the mixup ratio. Mixup ratio  $\lambda$  is randomly sampled from Unif $(0, \lambda_{max})$ , and we test  $\lambda_{max}$  ranging from 0.05 to 0.95.

like IMDB-B and IMDB-M, almost all of the mixup methods achieve limited improvement, as the data space has already been densely covered. But for datasets like MUTAG, PROTEINS, and MSRC-9 that are relatively sparsely covered, the mixup samples can significantly improve the model performance.

# D. Reproducibility

**Dataset description.** All the real-world datasets used in the paper are from (Morris et al., 2020) and available online<sup>1</sup>. We give a brief introduction to the datasets as follows

- PROTEINS (Borgwardt et al., 2005) is a set of proteins where each graph represents the protein structure with nodes as amino acids and edges as chemical bonds. The graph labels indicate whether the proteins are enzymes or non-enzymes. The node labels correspond to atom type and the node attributes represent node chemical features.
- MUTAG (Debnath et al., 1991) is a set of chemical graphs, where each graph represents one nitro-aromatic compound with nodes as atoms and edges as chemical bonds. The binary graph labels indicate whether the compounds have mutagenicity on Salmonella typhimurium. The node labels correspond to the atom type.
- MSRC-9 (Neumann et al., 2016) is a set of semantic image networks, where each graph represents an image with nodes as superpixels and edges indicates the adjacent relationship between superpixels. The graph labels indicate the semantic meaning (e.g., building, grass, tree, etc.) of the image.
- IMDB-B and IMDB-M (Yanardag & Vishwanathan, 2015) are movie collaboration networks, where each graph represents one movie with nodes as actor/actress and edges indicating whether two actors/actresses co-appear in the same movie. The graph labels indicate the movie categories (Action/Romance for IMDB-B, and Comedy/Romance/Sci-Fi for IMDB-M).

Table 4. Dataset statistics.

DATASET	#Graphs	#Nodes	#Edges	SPARSITY	#FEATURES	#GRAPH CLASS
PROTEINS	1,113	43.31	77.79	0.04	1	2
MUTAG	188	17.93	19.79	0.06	None	2
MSRC-9	221	40.58	97.94	0.06	None	8
IMDB-B	1,000	19.77	96.53	0.25	None	2
IMDB-M	1,500	12.74	53.88	0.33	None	3

<sup>&</sup>lt;sup>1</sup>https://chrsmrrs.github.io/datasets/

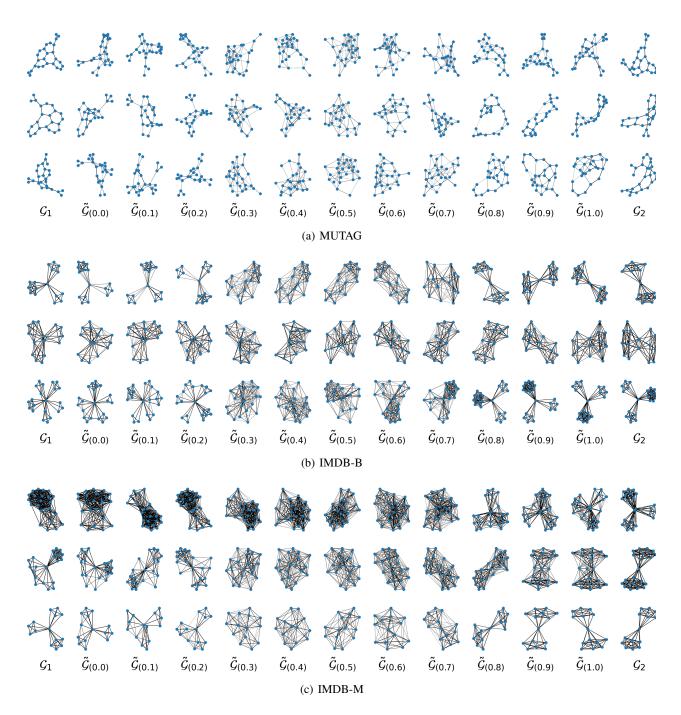


Figure 9. Visualization of mixup graphs on real-world datasets.

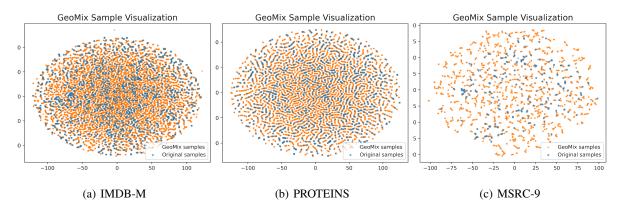


Figure 10. Visulization of the embedding space. The mixup samples lie on the geodesics connecting original samples, hence ensuring sample—label consistency.

**Model architecture.** For the GCN model (Kipf & Welling, 2017), we adopt three GCN layers with 32 hidden dimensions. We use ReLu as the activation function and global mean pooling as the readout. A dropout layer with dropout probability p = 0.5 is appended after the GCN layers, and followed by a linear layer with softmax activation to map embeddings to the classification probability.

For the GIN model (Xu et al., 2018), we adopt one 32-dimensional GIN layer with 3-layer MLP. We use ReLu as the activation function and global mean pooling as the readout. A dropout layer with dropout probability p = 0.5 is appended after the GIN layer, and followed by a linear layer with softmax activation to map embeddings to the classification probability.

**Experiment pipeline.** The proposed method is implemented in Python and all backbone models are built upon PyTorch. For model training, each model is trained for 300 epochs on the Linux platform with an Intel Xeon Gold 6240R CPU and an NVIDIA Tesla V100 SXM2 GPU.

During the experiments, for plain graphs without node attributes, we follow a standard procedure (Han et al., 2022) and use the categorized node degree as node attributes. We adopt 10-fold cross validation (train/test/validation=80%/10%/10%) to alleviate the randomness in the experiments.

## E. Future Works and Limitations.

In this paper, we design a geodesic graph mixup framework to ensure sample—label consistency for graph-level mixup. In this section, we discuss possible directions and applications (Wei et al., 2023; 2024) to explore that can further benefit and extend the current framework, including:

- Mixup sample selection. It is essential to select good mixup samples to generate informative samples that can improve model generalization and robustness. According to (Zhang et al., 2018), mixup samples are expected to be far from the training examples to avoid memorization of the training samples. On possible solution is to adopt the GW distance as the graph distance and choose samples that are far away from other samples for mixup.
- Node-level mixup. In this paper, we focus on the graph classification task. It would be of great interest to generalize the framework to the node-level tasks such as node classification. For example, following the similar idea using  $P_1$ ,  $P_2$  to transform  $G_1$ ,  $G_2$  to well-aligned  $\tilde{G}_1$ ,  $\tilde{G}_2$ , we can generate mixup node attributes as  $(1 \lambda)P_1^{\mathsf{T}}X_1 + \lambda P_2^{\mathsf{T}}X_2$  and node labels as  $(1 \lambda)P_1^{\mathsf{T}}y_1 + \lambda P_2^{\mathsf{T}}y_2$ .