# Learning Graphon Autoencoders for Generative Graph Modeling

Hongteng Xu[1, 2]    Peilin Zhao[3]    Junzhou Huang[3]    Dixin Luo[4*]

[1]Gaoling School of Artificial Intelligence, Renmin University of China
[2]Beijing Key Laboratory of Big Data Management and Analysis Methods
[3]Tencent AI Lab
[4]School of Computer Science and Technology, Beijing Institue of Technology
hongtengxu@ruc.edu.cn    dixin.luo@bit.edu.cn

June 1, 2021

## Abstract

Graphon is a nonparametric model that generates graphs with arbitrary sizes and can be induced from graphs easily. Based on this model, we propose a novel algorithmic framework called *graphon autoencoder* to build an interpretable and scalable graph generative model. This framework treats observed graphs as induced graphons in functional space and derives their latent representations by an encoder that aggregates Chebshev graphon filters. A linear graphon factorization model works as a decoder, leveraging the latent representations to reconstruct the induced graphons (and the corresponding observed graphs). We develop an efficient learning algorithm to learn the encoder and the decoder, minimizing the Wasserstein distance between the model and data distributions. This algorithm takes the KL divergence of the graph distributions conditioned on different graphons as the underlying distance and leads to a reward-augmented maximum likelihood estimation. The graphon autoencoder provides a new paradigm to represent and generate graphs, which has good generalizability and transferability.

## 1   Introduction

As a significant methodology for generative modeling, autoencoders map the data in the sample space $\mathcal{X}$ to a low-dimensional manifold embedded in a latent space $\mathcal{Z} \subset \mathbb{R}^C$. Typically, an autoencoder is specified by an encoder $f : \mathcal{X} \mapsto \mathcal{Z}$ mapping the data to latent codes in $\mathcal{Z}$, a predefined or learnable prior distribution $p_{\mathcal{Z}}$ on $\mathcal{Z}$, and a decoder $h : \mathcal{Z} \mapsto \mathcal{X}$ mapping the latent codes back to $\mathcal{X}$. By learning these modules, the autoencoder minimizes the discrepancy between the data distribution $p_{\mathcal{X}}$ and the model distribution $p_{h(\mathcal{Z})}$ [24, 52]. Compared with other generative modeling strategies like generative adversarial networks (GANs) [17] and generative flows [23], autoencoders can represent observed data explicitly in the latent space. Therefore, besides generating high-dimensional data like images [62] and texts [57], autoencoders have been widely used to learn data representations for other downstream tasks, *e.g.*, data clustering and classification.

However, most existing autoencoders are designed for the data in the same space. They are often inapplicable for complicated structured data sampled from incomparable spaces, such as a collection of arbitrarily-sized unaligned graphs (*i.e.*, the graphs have different numbers of nodes and the
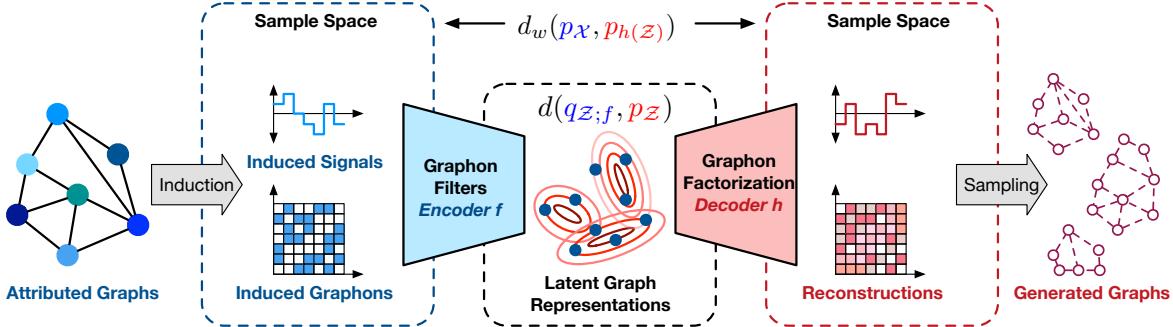
---

*Corresponding author

Figure 1: An illustration of our graphon autoencoder.

correspondence between their nodes is unknown). The variational graph autoencoder (VGAE) [25] and its variants [39, 56] obtain node-level embeddings rather than a global graph representation. Recently, some models apply attention-based pooling layers to aggregate the node embeddings as the graph representation [55, 30, 32]. However, these models often require side information to explore the clustering structure of the graphs, and they seldom consider the reconstructive and generative power of the graph representations.

To overcome the challenges above, we propose a novel **Graphon Autoencoder (GNAE)**. Leveraging the theory of graphon [31], we induce graphons (*i.e.*, two-dimensional symmetric Lebesgue measurable functions) from observed graphs and represent the node attributes associated with each graph as the signals defined on the graphons, such that the attributed graphs become the induced graphons and signals, which are in the same functional space. As illustrated in Figure 1, our GNAE essentially achieves a Wasserstein autoencoder [52] for the graphons. The encoder of our GNAE is an aggregation of Chebyshev graphon filters, which can be implemented as a graph neural network (GNN) for the induced graphons. It outputs the latent representations of the induced graphons (or, equivalently, the observed graphs). The posterior distribution of the latent representations is regularized by their prior distribution. The decoder of our GNAE is a graphon factorization model, which can reconstruct graphons from the latent representations and sample graphs with arbitrary sizes. For each induced graphon, its latent representation corresponds to the coefficients of the graphon factors in the decoder, which explicitly indicates its similarity to the factors. Therefore, our GNAE achieves an interpretable and scalable generative model for graphs.

We develop an efficient algorithm to achieve our GNAE. For each reconstructed graphon, we sample graphs and calculate their distributions conditioned on the reconstructed graphon and the input graphon, respectively. Taking the KL divergence between the conditional distributions as the underlying distance, we minimize the Wasserstein distance between the data and model distributions and learn our GNAE by a reward-augmented maximum likelihood (RAML) estimation method [38]. This algorithm avoids dense matrix multiplications and the backpropagation corresponding to the fused Gromov-Wasserstein (FGW) distance [50] between graphons, which owns low computational complexity. Experiments show that our GNAE performs well on representing and generating graphs, which have good generalizability (*i.e.*, generating graphs with various sizes but similar structures) and transferability (*i.e.*, training on a graph set and testing on others).

# 2 Proposed Model

## 2.1 From attributed graphs to graphons with signals

Mathematically, a graphon is a two-dimensional symmetric Lebesgue measurable function, denoted as $g : \Omega^2 \mapsto [0,1]$, where $\Omega$ is a measure space with a probability measure $\mu_\Omega$. Typically, we often set $\Omega = [0,1]$ and $\mu_\Omega$ as a uniform distribution on $\Omega$. Associated with the graphon, we can define a $M$-dimensional signal on it [34], which is denoted as $s : \Omega \mapsto \mathbb{R}^M$.

Denote the graphon space as $\mathcal{G}$ and the signal space as $\mathcal{S}$, respectively. For arbitrary $g_1, g_2 \in \mathcal{G}$, their $\delta_p$ distance [31] is $\delta_p(g_1, g_2) := \inf_{\phi \in \mathcal{F}_\Omega} \|g_1 - g_2^\phi\|_p$, where $\|g\|_p := (\int_{\Omega^2} |g(u,v)|^p dudv)^{\frac{1}{p}}$. Typically, we set $p = 1$ or $2$. If $\delta_p(g_1, g_2) = 0$, we say $g_1$ and $g_2$ are equivalent, denoted as $g_1 \cong g_2$. The work in [4, 15] shows that the quotient space $(\widehat{\mathcal{G}}, \delta_p)$, where $\widehat{\mathcal{G}} := \mathcal{G} \backslash \cong$, is homomorphic. $\delta_p$ is widely used in practice because of its computability. Especially, the work in [18, 61] indicates that the $\delta_p$ distance is equivalent to the order-$p$ Gromov-Wasserstein (GW) distance [33]:

**Definition 2.1.** *For arbitrary $g_1, g_2 \in \mathcal{G}$, their order-p Gromov-Wasserstein distance is*

$$d_{gw}(g_1, g_2) := \inf_{\pi \in \Pi(\mu_\Omega, \mu_\Omega)} \Big( \int_{\Omega^2 \times \Omega^2} |g_1(u, u') - g_2(v, v')|^p d\pi(u,v) d\pi(u', v') \Big)^{\frac{1}{p}}, \tag{1}$$

*where $\Pi(\mu_\Omega, \mu_\Omega) := \{\pi \geq 0 | \int_{u \in \Omega} d\pi(u,v) = \mu_\Omega, \int_{v \in \Omega} d\pi(u,v) = \mu_\Omega\}$.*

For $\mathcal{S}$, we can apply the Wasserstein distance as its metric:

**Definition 2.2.** *For arbitrary $s_1, s_2 \in \mathcal{S}$, their order-p Wasserstein distance is*

$$d_w(s_1, s_2) := \inf_{\pi \in \Pi(\mu_\Omega, \mu_\Omega)} \Big( \int_{\Omega^2} \|s_1(u) - s_2(v)\|_p^p d\pi(u,v) \Big)^{\frac{1}{p}}. \tag{2}$$

**Sampling graphs:** Graphon is a nonparametric graph generative model. We can sample graphs with arbitrary sizes from a graphon by the following steps:

$$\text{1) for } n = 1, .., N, \ v_n \sim \mu_\Omega; \ \text{2) } a_{nn'} \sim \text{Bernoulli}(g(v_n, v_{n'})); \ \text{3) } \boldsymbol{s}_n \sim \mathcal{N}(s(v_n), \sigma). \tag{3}$$

The first step is sampling $N$ nodes independently from $\mu_\Omega$. The second step generates an adjacency matrix $\boldsymbol{A} = [a_{nn'}] \in \{0,1\}^{N \times N}$, whose elements are sampled from the Bernoulli distributions determined by the graphon. When a signal is available, we can sample the attributes associated with the nodes, denoted as $\boldsymbol{S} = [\boldsymbol{s}_n] \in \mathbb{R}^{N \times M}$, from the distributions determined by the signal, *e.g.*, the Gaussian distributions in (3). For convenience, we denote $G(\boldsymbol{A}, \boldsymbol{S})$ as the sampled graph.

**Inducing graphons:** We induce a graphon and a signal from an attributed graph as follows.

**Definition 2.3** (Induced Graphon). *For a graph $G(\boldsymbol{A}, \boldsymbol{S})$, where $\boldsymbol{A} = [a_{nn'}] \in \{0,1\}^{N \times N}$ and $\boldsymbol{S} = [\boldsymbol{s}_n] \in \mathbb{R}^{N \times M}$, we can induce a graphon and its corresponding signal as two step functions:*

$$g_\mathcal{P}(v, v') = \sum_{n,n'=1}^N a_{nn'} 1_{\mathcal{P}_n}(v) 1_{\mathcal{P}_{n'}}(v'), \ \text{and} \ s_\mathcal{P}(v) = \sum_{n=1}^N \boldsymbol{s}_n 1_{\mathcal{P}_n}(v), \ \forall \ v, v' \in \Omega, \tag{4}$$

*where $\mathcal{P} = \{\mathcal{P}_n\}_{n=1}^N$ represents $N$ equitable partitions of $\Omega$, i.e., $\cup_n \mathcal{P}_n = \Omega$ and $|\mathcal{P}_n| = |\mathcal{P}_{n'}|$ for all $n \neq n'$. The indicator $1_{\mathcal{P}_n}(v) = 1$ if $v \in \mathcal{P}_n$, otherwise it equals to 0.*

Obviously, $g_\mathcal{P} \in \mathcal{G}$ and $s_\mathcal{P} \in \mathcal{S}$. The step function approximation lemma [8] shows that for the graphs sampled from a graphon $g$, the average of their induced graphons provides a consistent estimation of $g$. The estimation error reduces with the increase of the number and the size of the graphs.
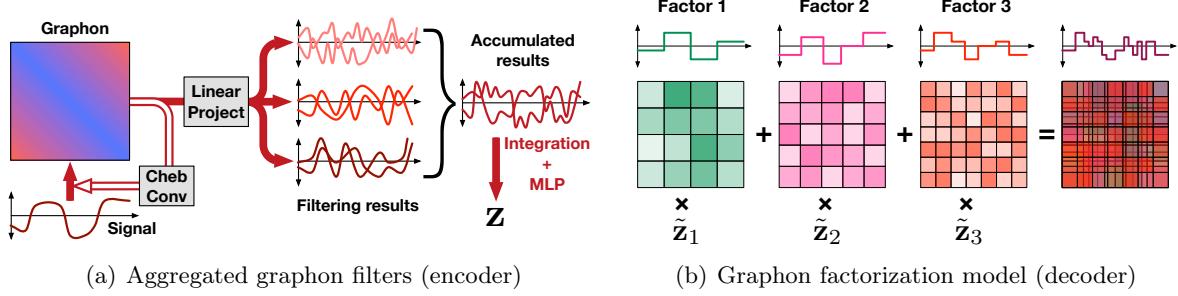
3

(a) Aggregated graphon filters (encoder)        (b) Graphon factorization model (decoder)

Figure 2: Illustrations of the encoder and the decoder of our GNAE.

## 2.2 A graphon autoencoder in functional space

The graphons and their associated signals, denoted as $\{(g, s)\}$, can be viewed as samples in a functional space $(\mathcal{X}, d_{\mathcal{X}}, \mathbb{P})$. Here, $\mathcal{X} = \mathcal{G} \times \mathcal{S}$, $d_{\mathcal{X}}$ is an underlying distance defining the discrepancy between different samples,[1] which will be introduced below and discussed in-depth in Section 3, and $\mathbb{P}$ represents the set of probability measures defined on $\mathcal{X}$.

By inducing graphons with signals, we can represent the arbitrarily-sized unaligned graphs as the samples in the same space. Accordingly, existing machine learning techniques like autoencoders become applicable. In particular, our graphon autoencoder (GNAE) can be viewed as a Wasserstein autoencoder of the graphons, which consists of an encoder $f : \mathcal{X} \mapsto \mathcal{Z}$, a decoder $h : \mathcal{Z} \mapsto \mathcal{X}$, and a learnable latent prior distribution $p_{\mathcal{Z}}$. Given attributed graphs, we obtain a set of induced graphons and associated signals, denoted as $\{(g_{\mathcal{P}}, s_{\mathcal{P}})\}$ and learn the autoencoder to minimize the order-1 Wasserstein distance between the (unknown) data distribution $p_{\mathcal{X}}$ and the model distribution $p_{h(\mathcal{Z})}$, $i.e.$, $\min d_{\mathrm{w}}(p_{\mathcal{X}}, p_{h(\mathcal{Z})})$, where $p_{\mathcal{X}}, p_{h(\mathcal{Z})} \in \mathbb{P}$. According to Theorem 1 in [52], we relax the optimization problem as

$$\min_{f,h,p_{\mathcal{Z}}} \ \mathbb{E}_{\boldsymbol{x} \sim p_{\mathcal{X}}} \mathbb{E}_{\boldsymbol{z} \sim q_{\mathcal{Z}|\mathcal{X};f}}[d_{\mathcal{X}}(\boldsymbol{x}, h(\boldsymbol{z}))] + \gamma d(q_{\mathcal{Z};f}, p_{\mathcal{Z}}), \tag{5}$$

where each $\boldsymbol{x} = (g_{\mathcal{P}}, s_{\mathcal{P}})$ represents a tuple of the graphon and signal induced from the corresponding observed graph. $d_{\mathcal{X}}(\boldsymbol{x}, h(\boldsymbol{z}))$ represents the reconstruction error of the sample $\boldsymbol{x}$. $q_{\mathcal{Z};f} = \mathbb{E}_{\boldsymbol{x} \sim p_{\mathcal{X}}}[q_{\mathcal{Z}|\boldsymbol{x};f}]$ is the expected latent posterior conditioned on different samples, which is required to be closed to the latent prior $p_{\mathcal{Z}}$ under the metric $d$. Parameter $\gamma$ achieves a trade-off between reconstruction loss and the regularizer. As shown in (5), our GNAE learns the encoder, the decoder, and the latent prior distribution jointly. These modules are implemented as follows.

**Latent prior distribution** Following the work in [62], we set the latent prior $p_{\mathcal{Z}}$ as a learnable Gaussian mixture model (GMM) and implement the regularizer in (5) as the sliced fused Gromov-Wasserstein (SFGW) distance, $i.e.$, $d(q_{\mathcal{Z};f}, p_{\mathcal{Z}}) := d_{\mathrm{sfgw}}(q_{\mathcal{Z};f}, p_{\mathcal{Z}})$. This configuration helps us to learn latent representations with clustering structures.

**Encoder** The encoder of our GNAE is designed as an aggregation of Chebyshev graphon filters. Given a graphon $g$ and its corresponding signal $s$, we achieve our encoder as follows:

$$\boldsymbol{z} = f((g, s)) = \mathrm{MLP}\Big(\int_{\Omega} \sum_{j=0}^{J} \theta_j(s^{(j)}(v))dv\Big), \ \text{where}$$

$$s^{(0)}(v) = s(v), \ s^{(1)}(v) = L_g(s^{(0)}(v)) = \int_{\Omega} g(u, v)(s^{(0)}(v) - s^{(0)}(u))du, \ \text{and} \tag{6}$$

$$s^{(j)}(v) = 2L_g(s^{(j-1)}(v)) - s^{(j-2)}(v) \ \text{for} \ j > 1.$$

---

[1]Note that the underlying distance may not be a strict metric in practice.

4

In the $j$-th step, the filtering result $s^{(j)}$ is obtained by applying a Laplacian filter to $s^{(j-1)}$ (*i.e.*, $L_g(s^{(j-1)}(v))$) and treating $s^{(j-2)}$ as an offset. $\theta^{(j)}(\cdot)$ is a linear projection, mapping each $s^{(j)}(v)$ to $\mathbb{R}^D$. We obtain the aggregated filtering result by accumulating and integrating the results of different steps. Finally, we apply a multi-layer perceptron (MLP) network to derive a $C$-dimensional latent representation. Figure 2(a) illustrates the scheme of our encoder.

For general graphons and signals, we can implement the filtering process above based on Fourier transform [6], whose complexity is high. Fortunately, for the induced graphon and signal $(g_{\mathcal{P}}, s_{\mathcal{P}})$, we can implement the graphon filters by a Chebyshev spectral graph convolutional (ChebConv) network [13]. For $j > 0$, if $\Omega = [0,1]$ and $\mathcal{P} = \{\mathcal{P}_n\}_{n=1}^N$ are equitable partitions, we have

$$\int_{\Omega} g_{\mathcal{P}}(u,v)(s_{\mathcal{P}}^{(j-1)}(v) - s_{\mathcal{P}}^{(j-1)}(u))du = \frac{1}{N}\sum_{n=1}^N (\boldsymbol{L}\boldsymbol{S}^{(j-1)})_n 1_{\mathcal{P}_n}(v), \text{ for } v \in \Omega, \tag{7}$$

where $\boldsymbol{L} = \mathrm{diag}(\boldsymbol{A1}) - \boldsymbol{A}$ is the Laplacian graph matrix, $\boldsymbol{S}^{(j-1)} = [\boldsymbol{s}_n^{(j-1)}]$ is a matrix of the $(j-1)$-th signal, and each $\boldsymbol{s}_n^{(j-1)}$ corresponds to the signal in the partition $\mathcal{P}_n$. Accordingly, $(\boldsymbol{L}\boldsymbol{S}^{(j-1)})_n$ is the $n$-th row of $\boldsymbol{L}\boldsymbol{S}^{(j-1)}$. Plugging (7) into (6), we can derive the latent representation as $\boldsymbol{z} = \mathrm{MLP}(\sum_{n=1}^N(\sum_{j=0}^J \frac{1}{N^{j+1}}\theta_j(\boldsymbol{S}^{(j)})))$, which can be implemented as a ChebConv network followed by an average pooling layer and a MLP.

**Decoder** For each $\boldsymbol{z} = [z_1,..,z_C]$ derived from the encoder, we apply a factorization model to reconstruct the corresponding graphon and signal. Specifically, the decoder consists of $C$ graphon factors, denoted as $\{(\tilde{g}_c, \tilde{s}_c)\}_{c=1}^C$. Each graphon factor corresponds to two step functions defined as (4) shows. Accordingly, the reconstructed sample can be represented as $h(\boldsymbol{z}) = (\hat{g}, \hat{s})$, where

$$\hat{g} = \sum_{c=1}^C \tilde{z}_c \tilde{g}_c, \text{ and } \hat{s} = \alpha\Big(\sum_{c=1}^C \tilde{z}_c \tilde{s}_c\Big), \text{ with } \tilde{z}_c = \mathrm{softmax}_c(\boldsymbol{z}) = \frac{\exp(z_c)}{\sum_{c'}\exp(z_{c'})}. \tag{8}$$

We reparameterize each $\tilde{g}_c$ as $\sigma(b_c)$, where $b_c(u,v) : \Omega^2 \mapsto \mathbb{R}$ is a step function with unbounded output range and $\sigma(\cdot)$ is a sigmoid function, and let the latent representation pass through a softmax layer. This setting makes $\{\tilde{g}_c\}_{c=1}^C$ and $\hat{g}$ in the space $\mathcal{G}$. The function $\alpha(\cdot)$ depends on the type of the original signal, which can be ReLU, softmax, sigmoid, *etc.* As shown in Figure 2(b), the graphon factors may have different partitions. Applying the inclusion-exclusion principle [20], we have

**Proposition 2.4.** *Suppose that* $\{\tilde{g}_c : [0,1]^2 \mapsto [0,1]\}_{c=1}^C$ *are 2D step functions, each of which has* $N_c$ *equitable partitions. Denote* $\{\mathcal{L}_c\}_{c=1}^C$ *as the sets of the landmarks indicating the partitions, where* $\mathcal{L}_c = \{\frac{1}{N_c}, ..., \frac{N_c-1}{N_c}\}$. *For the* $\hat{g}$ *derived by (8), the number of its partitions is* $|\mathcal{P}| = |\cup_{c=1}^C \mathcal{L}_c| + 1 = \sum_{\emptyset \neq \mathcal{C} \subset \{1,..,C\}} (-1)^{|\mathcal{C}|} |\cap_{c\in\mathcal{C}} \mathcal{L}_c| + 1$. *If all the* $N_c$'s *are prime numbers,* $|\mathcal{P}| = \sum_{c=1}^C |\mathcal{L}_c| + 1$.

Proposition 2.4 shows the number of the partitions of the reconstructed graphon can be much larger than that of the graphon factors, which is beneficial for the capability of our model.

# 3   Learning algorithm

## 3.1   The fused Gromov-Wasserstein distance between graphons

Besides the modules above, the key of our GNAE is the underlying distance $d_{\mathcal{X}}$. A straightforward way is implementing $d_{\mathcal{X}}$ as the fused Gromov-Wasserstein (FGW) distance [50]:

**Definition 3.1.** *For* $\boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathcal{X}$, *where* $\boldsymbol{x}_1 = (g_1, s_1)$ *and* $\boldsymbol{x}_2 = (g_2, s_2)$, *their order-p fused Gromov-Wasserstein distance, denoted as* $d_{fgw}(\boldsymbol{x}_1, \boldsymbol{x}_2)$, *is*

$$\inf_{\pi \in \Pi(\mu_\Omega, \mu_\Omega)} \Big(\int_{\Omega^2 \times \Omega^2} |g_1(u,u') - g_2(v,v')|^p d\pi(u,v)d\pi(u',v') + \int_{\Omega^2} \|s_1(u) - s_2(v)\|_p^p d\pi(u,v)\Big)^{\frac{1}{p}},$$

5

The FGW distance combines the GW distance between graphons and the Wasserstein distance between signals, enforcing them share the same optimal transport $\pi(u, v)$. It is a metric for the quotient space $\widehat{\mathcal{X}} := \mathcal{X}\backslash \cong$ when $p = 1$ and a semi-metric when $p > 1$ [50].

We can compute the FGW distance by solving an optimization problem with finite variables when dealing with the graphons and signals formulated as step functions.

**Proposition 3.2.** *Given $\boldsymbol{x}_{1,\mathcal{P}} = (g_{1,\mathcal{P}}, s_{1,\mathcal{P}})$ and $\boldsymbol{x}_{2,\mathcal{Q}} = (g_{2,\mathcal{Q}}, s_{2,\mathcal{Q}})$, where*

$$g_{1,\mathcal{P}}(v, v') = \sum\nolimits_{n,n'=1}^{N} g_{1,nn'} 1_{\mathcal{P}_n}(v) 1_{\mathcal{P}_{n'}}(v'), \ g_{2,\mathcal{Q}}(v, v') = \sum\nolimits_{m,m'=1}^{M} g_{2,mm'} 1_{\mathcal{Q}_m}(v) 1_{\mathcal{Q}_{m'}}(v'),$$

$$s_{1,\mathcal{P}}(v) = \sum\nolimits_{n=1}^{N} \boldsymbol{s}_{1,n} 1_{\mathcal{P}}(v), \ s_{2,\mathcal{Q}}(v) = \sum\nolimits_{m=1}^{M} \boldsymbol{s}_{2,m} 1_{\mathcal{Q}}(v)$$

*are step functions, we have*

$$d_{fgw}(\boldsymbol{x}_{1,\mathcal{P}}, \boldsymbol{x}_{2,\mathcal{Q}}) = \min_{\boldsymbol{T}\in\Pi(\boldsymbol{\mu}_{\mathcal{P}}, \boldsymbol{\mu}_{\mathcal{Q}})}(\langle \boldsymbol{D}_g, \boldsymbol{T}\otimes\boldsymbol{T}\rangle + \langle \boldsymbol{D}_s, \boldsymbol{T}\rangle)^{\frac{1}{p}}, \tag{9}$$

*where $\boldsymbol{D}_g = [|g_{1,nn'} - g_{2,mm'}|^p] \in \mathbb{R}^{N^2\times M^2}$, $\boldsymbol{D}_s = [\|\boldsymbol{s}_{1,n} - \boldsymbol{s}_{2,m}\|_p^p] \in \mathbb{R}^{N\times M}$, $\otimes$ represents Kronecker product, and $\Pi(\boldsymbol{\mu}_{\mathcal{P}}, \boldsymbol{\mu}_{\mathcal{Q}}) = \{\boldsymbol{T} \geq \boldsymbol{0}|\boldsymbol{T}\boldsymbol{1} = \boldsymbol{\mu}_{\mathcal{P}}, \boldsymbol{T}^\top\boldsymbol{1} = \boldsymbol{\mu}_{\mathcal{Q}}\}$ with $\boldsymbol{\mu}_{\mathcal{P}} = [\frac{|\mathcal{P}_1|}{|\Omega|}, .., \frac{|\mathcal{P}_N|}{|\Omega|}]$ and $\boldsymbol{\mu}_{\mathcal{Q}} = [\frac{|\mathcal{Q}_1|}{|\Omega|}, .., \frac{|\mathcal{Q}_M|}{|\Omega|}]$.*

Although (9) is computable, its computational complexity is so high as $\mathcal{O}(N^4)$ for the graphons with $N$ partitions (or $\mathcal{O}(N^3)$ if $p = 2$ [41]). What is worse, because the graphons reconstructed by our GNAE are non-sparse 2D step functions, some commonly-used acceleration strategies like the sliced FGW distance [62] and the sparse matrix multiplications [51, 60] become inapplicable. Additionally, the underlying distance is parameterized by the GNAE model, so we have to consider the gradient of the optimal transport matrix $\boldsymbol{T}$ to the model parameters, which is expensive on both time and memory [58].

## 3.2 The KL divergence between conditional graph distributions

Facing the issues above, we need to explore a surrogate of the FGW distance when learning our GNAE model. A competitive choice is the KL divergence between the distributions of graphs conditioned on $h(\boldsymbol{z})$ and $\boldsymbol{x}$, denoted as $d_{\mathrm{KL}}(q(G|\boldsymbol{x}), p(G|h(\boldsymbol{z})))$. Specifically, given a reconstructed sample, i.e., $h(\boldsymbol{z}) = (\hat{g}, \hat{s})$, we sample a set of attributed graphs, denoted as $\mathcal{Y} = \{G(\widehat{\boldsymbol{A}}, \widehat{\boldsymbol{S}})\}$. Each $G(\widehat{\boldsymbol{A}}, \widehat{\boldsymbol{S}})$ has $K$ nodes. According to (3), the likelihood of $G(\widehat{\boldsymbol{A}}, \widehat{\boldsymbol{S}})$ is

$$p(G|h(\boldsymbol{z})) = p(\mathcal{V})p(\widehat{\boldsymbol{A}}|\hat{g}, \mathcal{V})p(\widehat{\boldsymbol{S}}|\hat{s}, \mathcal{V}) = \frac{1}{|\Omega|^K}\prod\nolimits_{k,k'=1}^{K} p(\hat{a}_{kk'}|\hat{g}(v_k, v_{k'}))\prod\nolimits_{k=1}^{K} p(\hat{\boldsymbol{s}}_k|\hat{s}(v))$$

$$\propto \prod\nolimits_{k,k'=1}^{K} \hat{g}(v_k, v_{k'})^{\hat{a}_{kk'}}(1 - \hat{g}(v_k, v_{k'}))^{1-\hat{a}_{kk'}}\prod\nolimits_{k=1}^{K}\exp\Big(-\frac{\|\hat{\boldsymbol{s}}_k - \hat{s}(v_k)\|_2^2}{2M\sigma^2}\Big), \tag{10}$$

where $\mathcal{V} = \{v_1, .., v_K\}$ and each $v_k$ is sampled independently from $\mu_\Omega$, so that $p(\mathcal{V}) = \frac{1}{|\Omega|^K}$ when $\mu_\Omega$ is a uniform distribution.

Additionally, we leverage an exponentiated payoff distribution [38] to approximate the probability of each $G(\widehat{\boldsymbol{A}}, \widehat{\boldsymbol{S}})$ conditioned on the observed graphon $\boldsymbol{x}$:

$$q(G|\boldsymbol{x}) = \frac{\exp(r(\hat{\boldsymbol{x}}_G, \boldsymbol{x}))}{\sum_{G'\in\mathcal{Y}}\exp(r(\hat{\boldsymbol{x}}_{G'}, \boldsymbol{x}))} = \frac{\exp(-d_{\mathrm{fgw}}(\hat{\boldsymbol{x}}_G, \boldsymbol{x})/\tau)}{\sum_{G'\in\mathcal{Y}}\exp(-d_{\mathrm{fgw}}(\hat{\boldsymbol{x}}_{G'}, \boldsymbol{x})/\tau)}, \tag{11}$$

where $\hat{\boldsymbol{x}}_G$ is the graphon (and the signal) induced from the graph $G$, and $r(\hat{\boldsymbol{x}}_G, \boldsymbol{x}) = -\frac{d_{\mathrm{fgw}}(\hat{\boldsymbol{x}}_G, \boldsymbol{x})}{\tau}$ is the reward function implemented as the negative order-2 FGW distance between $\hat{\boldsymbol{x}}_G$ and $\boldsymbol{x}$. Parameter $\tau$ controls the smoothness of $q(G|\boldsymbol{x})$. In our work, we set $\tau$ adaptively as $\min_{G\in\mathcal{Y}} d_{\mathrm{fgw}}(\hat{\boldsymbol{x}}_G, \boldsymbol{x})$.

---
**Algorithm 1** Learning a GNAE by RAML
---
**input** A set of graphons and signals induced from observed attributed graphs, denoted as $\mathcal{X}$.
**output** An encoder $f$, a decoder $h$, and a latent prior $p_{\mathcal{Z}}(\boldsymbol{z}) = \frac{1}{T}\sum_t \mathcal{N}(\boldsymbol{z}; \boldsymbol{\mu}_t, \text{diag}(\boldsymbol{\sigma}_t^2))$.
 1: **for** each epoch
 2:    **for** each batch $\{\boldsymbol{x}_n\}_{n=1}^{N_b} \subset \mathcal{X}$
 3:       **for** $n = 1, .., N_b$
 4:          Samples of $q_{\mathcal{Z};f}$: $\boldsymbol{z}_n = f(\boldsymbol{x}_n)$.
 5:          Samples of $p_{\mathcal{Z}}$: $t \sim \text{Categorical}(\frac{1}{T})$, and $\boldsymbol{z}_n' \sim \mathcal{N}(\boldsymbol{\mu}_t, \text{diag}(\boldsymbol{\sigma}_t^2))$.
 6:          Sample attributed graphs $\{G_i(\widehat{\boldsymbol{A}}_i, \widehat{\boldsymbol{S}}_i)\}_{i=1}^I$ from $h(\boldsymbol{z}_n)$ and induce $\{\hat{\boldsymbol{x}}_{G_i}\}_{i=1}^I$ by (4).
 7:          Compute $\{d_{\text{fgw}}(\hat{\boldsymbol{x}}_{G_i}, \boldsymbol{x}_n)\}_{i=1}^I$ by (9) and obtain $\{q(G_i|\boldsymbol{x}_n)\}_{i=1}^I$ by (11).
 8:          Compute $\{p(G_i|h(\boldsymbol{z}_n))\}_{i=1}^I$ by (10).
 9:       Calculate $\mathcal{L} = -\sum_{n=1}^{N_b}\sum_{i=1}^I q(G_i|\boldsymbol{x}_n)\log p(G_i|h(\boldsymbol{z}_n)) + \gamma d_{\text{sfgw}}(q_{\mathcal{Z};f}, p_{\mathcal{Z}})$.
10:       Update the model by Adam optimizer [22].
---

## 3.3 Reward-augmented maximum likelihood estimation

Leveraging $d_{\text{KL}}(q(G|\boldsymbol{x}), p(G|h(\boldsymbol{z})))$, we develop an efficient learning algorithm with much lower computational complexity. Specifically, we can rewrite $d_{\text{KL}}(q(G|\boldsymbol{x}), p(G|h(\boldsymbol{z})))$ as

$$d_{\text{KL}}(q(G|\boldsymbol{x}), p(G|h(\boldsymbol{z}))) = -\mathbb{E}_{G \sim q(G|\boldsymbol{x})}[\log p(G|h(\boldsymbol{z}))] + \mathbb{E}_{G \sim q(G|\boldsymbol{x})}[\log q(G|\boldsymbol{x})], \tag{12}$$

where the second term is the entropy of the sampled graphs, which is a constant with respect to the model. Plugging (12) into (5), we learn our GNAE by

$$\min_{f,h,p_{\mathcal{Z}}} -\mathbb{E}_{\boldsymbol{x} \sim p_{\mathcal{X}}}\mathbb{E}_{\boldsymbol{z} \sim q_{\mathcal{Z}|\mathcal{X};f}}\mathbb{E}_{G \sim q(G|\boldsymbol{x})}[\log p(G|h(\boldsymbol{z}))] + \gamma d_{\text{sfgw}}(q_{\mathcal{Z};f}, p_{\mathcal{Z}}). \tag{13}$$

This optimization problem can be solved by the reward-augmented maximum likelihood (RAML) method [38]. The scheme of our learning algorithm is shown in Algorithm 1. In principle, when the sampled graph is close to the original input graph, the FGW distance between their induced graphon will be small. Accordingly, the log-likelihood of the graph (*i.e.*, $\log p(G|h(\boldsymbol{z}))$) will be assigned to a large weight (*i.e.*, $q(G|\boldsymbol{x})$), and the model will be updated to increase the likelihood.

Although our learning algorithm still involves computing FGW distances, it has obvious advantages on computational complexity compared to using $d_{\text{fgw}}(\boldsymbol{x}, h(\boldsymbol{z}))$ as underlying distance directly. Firstly, the number of each sampled graph's nodes (*i.e.*, $K$) can be much smaller than that of the reconstructed graphon's partitions (*i.e.*, $N$). Additionally, we replace dense reconstructed graphons with sparse adjacency matrices of the sampled graphs with the help of the Bernoulli sampling. Therefore, it is relatively easy to compute the $d_{\text{fgw}}(\hat{\boldsymbol{x}}_{G_i}, \boldsymbol{x}_n)$ in Line 7 of Algorithm 1 — its computational complexity is $\mathcal{O}(EK)$ and $E$ is the number of the edges of the graph inducing $\boldsymbol{x}_n$. Moreover, the gradient corresponding to the first term of $\mathcal{L}$ is $-\sum_{n=1}^{N_b}\sum_{i=1}^I q(G_i|\boldsymbol{x}_n)\nabla\log p(G_i|h(\boldsymbol{z}_n))$. Here, the $q(G_i|\boldsymbol{x}_n)$ is used as a constant, and the corresponding FGW distance is not involved in the backpropagation, which reduces the cost of time and memory greatly.

## 4 Connections to Existing Work

**Autoencoders** The principle of autoencoders is to minimize the discrepancy between the data and model distributions. The variational autoencoder (VAE) [24] and its variants [53, 57] apply the KL-divergence as the discrepancy and learns a probabilistic autoencoder via maximizing the evidence lower bound (ELBO). The Wasserstein autoencoders (WAEs) [52, 27] minimize a relaxed form of

the Wasserstein distance to learn a deterministic autoencoder. Both these two strategies lead to a learning task including a reconstruction loss of observed data and a regularizer penalizing the distance between the prior and the posterior (or the mixture of different posterior distributions) in the latent space. The prior can be a predefined normal distribution or a learnable mixture model [49, 62]. The commonly-used distances between the prior and the posterior include KL divergence, maximum mean discrepancy, GAN-based loss, FGW distance, *etc.*

**Generative graph modeling** The early graph models like the Erdős-Rényi graph [14] simulate large graphs to yield certain statistical properties but cannot capture complicated mechanisms of real-world graphs. Recently, the GNN-based graph generative models have been widely used, which can be categorized into two classes. The first class learns node-level embeddings and estimates edges based on the pairs of the embeddings [25, 26, 37, 64], which works well on link prediction [69] and conditional graph generation [66]. The second class applies various pooling layers [67, 55, 30] to obtain graph embeddings and then leverages recurrent neural networks to generate nodes and edges in an autoregressive manner [68, 47, 19, 12]. Besides the GNN-based models, the Gromov-Wasserstein factorization (GWF) model [59] reconstructs each graph as a weighted GW barycenter [41] of learnable graph factors, which achieves encouraging performance on graph clustering. Following the GWF model, the graph dictionary learning (GDL) model [54] leverages a linear factorization to reconstruct graphs, which has lower complexity than the GWF model. However, the models above seldom consider the clustering structure or the distribution of the graph embeddings they learned.

**Graphon-based graph models** Graphon is a nonparametric graph model, which has been widely used in network modeling [3, 15] and optimization [40]. To infer graphons from observed graphs, many methods have been proposed, *i.e.*, the stochastic block approximation (SBA) methods [2, 10, 8] and the low-rank approximation methods [21, 11, 63]. These methods require well-aligned graphs, which is questionable in practice. The work in [61] relaxes this requirement, learning the graphon and aligning the observed graphs alternately by solving a GW barycenter problem [41]. All the methods above are based on the weak regularity lemma [31], approximating graphons by 2D step functions. Recently, the work in [44, 43] bridges the gap between graphon-based signal processing and graph neural networks, which inspires the design of our encoder. However, existing methods either learn graphons to generate graphs or leverage graphons to process the information of nodes. None of them consider learning the distribution of graphons as we did.

**The novelties of our GNAE** To our knowledge, our graphon autoencoder makes the first attempt to build a WAE in the graphon space, which provides a new algorithmic framework for graphon distribution modeling and graph generation. Our GNAE extends the GNN-based model and the factorization model to the functional space of graphons. The encoder achieves graphon filtering, whose GNN-based implementation is a special case for induced graphons. The decoder improves the GDL model by leveraging graphon factors with different partitions. Combining these two strategies in the framework of autoencoders, our GNAE inherits their advantages and learns them with better interpretability and capability.

## 5 Experiments

### 5.1 Graph representation and classification

To demonstrate the usefulness of our GNAE model, we test it on six public graph datasets and compare it with state-of-the-art methods on graph modeling. The datasets we used can be categorized into three classes: The **MUTAG** and the **PTC-MR** in [29] contain molecules with categorical node attributes; the **PROTEIN** and the **ENZYMES** in [5] contain proteins with continuous node attributes; and the **IMDB-B** and the **IMDB-M** in [65] contain social networks without node

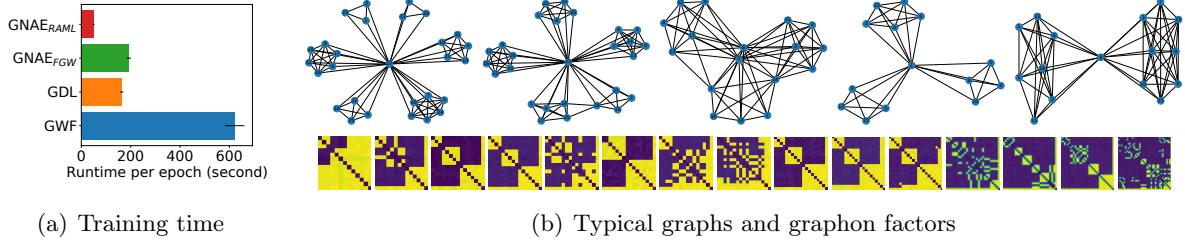(a) Training time       (b) Typical graphs and graphon factors

Figure 3: For the IMDB-B dataset: (a) Comparisons for FGW-based methods on their runtime. (b) Illustrations of typical graphs and the graphon factors learned by our $\text{GNAE}_{\text{RAML}}$.

attributes. These datasets can be downloaded from `https://chrsmrrs.github.io/datasets/` [35]. For the datasets without node attributes, we treat the local degree profiles [7] of nodes as the attributes.

The baselines include: (*i*) the kernel-based methods, *e.g.*, Random Walk Kernel (**RWK**) [16], Shortest Path Kernel (**SPK**) [5], Graphlet Kernel (**GK**) [46], Weisfeiler-Lehman Sub-tree Kernel (**WLK**) [45], Deep Graph Kernel (**DGK**) [65], Multi-Scale Laplacian Kernel (**MLGK**) [28], and Fused Gromov-Wasserstein Kernel (**FGWK**) [51]; (*ii*) the GNN-based methods, *e.g.*, **sub2vec** [1], **graph2vec** [36], and the state-of-the-art InfoGraph method [48] that uses Graph Isomorphismic Network (GIN) [64] and Differentiable Pooling (DP) [67] as its backbone model, respectively (**InfoGraph$_{\text{GIN}}$** and **InfoGraph$_{\text{DP}}$**); (*iii*) the factorization models (FMs), *e.g.*, the Gromov-Wasserstein factorization (**GWF**) [59] and the Graph Dictionary Learning (**GDL**) [54]. We reproduce the baselines either based on the code released by the authors or our own implementations and set their hyperparameters according to the released code or the corresponding references. When implementing our GNAE model, we consider two variants: applying the FGW distance directly as the underlying distance and learning the GNAE model by alternating optimization (**GNAE$_{\text{FGW}}$**), or applying the KL divergence of graph distributions as the underlying distance and learning the model by the proposed RAML (**GNAE$_{\text{RAML}}$**). For the GNAE models, the settings of their hyperparameters are given in Appendix.

We test the methods above on graph classification. For each kernel-based method, we train a kernel SVM classifier [9]. For other methods, we learn graph representations explicitly in an unsupervised way and train an SVM classifier based on the representations. The SVM classifier of each method is trained based on 10-fold cross-validation, and we use the same random seed to split data and select the most suitable SVM kernel function manually. Table 1 lists the mean and the standard deviation of the classification accuracy achieved by the methods on each dataset. We can find that the performance of our GNAE models is at least comparable to that of the state-of-the-art methods (*e.g.*, MLGK, FGWK and InfoGraph). Especially, the proposed $\text{GNAE}_{\text{RAML}}$ achieves the top-5 accuracy on four of the six datasets, which is the same as $\text{InforGraph}_{\text{GIN}}$ does. Note that for the GNN-based methods, the dimension of their graph representations is over 100. However, our GNAE models achieve competitive results based on the representations with much a lower dimension ($\leq 30$ for all the datasets). For the challenging ENZYMES dataset, our GNAE methods do not work well. A potential reason for this phenomenon is the model misspecification issue — the node attributes in this dataset are sparse and have high dynamic ranges, so the smoothed signal model we applied may not be able to describe and reconstruct such attributes well.

Our $\text{GNAE}_{\text{RAML}}$ is more efficient than its competitors that apply FGW distance (*i.e.*, GWF, GDL, and $\text{GNAE}_{\text{FGW}}$). Suppose that all the models contain $C$ graph (or graphon) factors with comparable sizes, denoted as $\mathcal{O}(N)$. Given a graph with $N$ nodes and $E$ edges, the GWF reconstructs

9

Table 1: Comparison on classification accuracy (%).

| Category | Method | MUTAG | PTC-MR | PROTEIN | ENZYMES | IMDB-B | IMDB-M | # in Top5 |
|---|---|---|---|---|---|---|---|---|
| Kernels | RWK | $83.72_{\pm1.50}$ | $57.85_{\pm1.30}$ | $73.95_{\pm0.59}$ | $28.52_{\pm1.83}$ | $50.70_{\pm0.26}$ | $34.65_{\pm0.19}$ | 0 |
| | SPK | $\mathbf{85.22}_{\pm2.43}$ | $58.24_{\pm2.44}$ | $\mathbf{74.93}_{\pm0.86}$ | $38.87_{\pm3.01}$ | $55.60_{\pm0.22}$ | $37.99_{\pm0.30}$ | 2 |
| | GK | $81.66_{\pm2.11}$ | $57.26_{\pm1.41}$ | $71.10_{\pm1.08}$ | $30.36_{\pm4.84}$ | $65.90_{\pm0.98}$ | $43.89_{\pm0.38}$ | 0 |
| | WLK | $80.72_{\pm3.00}$ | $57.97_{\pm0.49}$ | $73.01_{\pm1.09}$ | $54.69_{\pm3.27}$ | $\mathbf{72.30}_{\pm3.44}$ | $46.35_{\pm0.46}$ | 1 |
| | DGK | $\mathbf{87.44}_{\pm2.72}$ | $60.08_{\pm2.55}$ | $74.27_{\pm1.12}$ | $53.22_{\pm1.01}$ | $66.90_{\pm0.56}$ | $44.55_{\pm0.52}$ | 1 |
| | MLGK | $\mathbf{87.94}_{\pm1.61}$ | $\mathbf{62.23}_{\pm1.39}$ | $\mathbf{75.86}_{\pm0.99}$ | $61.89_{\pm1.17}$ | $66.60_{\pm0.25}$ | $41.17_{\pm0.03}$ | 3 |
| | FGWK$^*$ | $\mathbf{88.13}_{\pm4.22}$ | $\mathbf{62.98}_{\pm5.27}$ | $72.20_{\pm3.81}$ | $\mathbf{71.48}_{\pm2.96}$ | $63.50_{\pm4.01}$ | $46.27_{\pm3.85}$ | 3 |
| GNNs | sub2vec | $60.88_{\pm9.89}$ | $59.99_{\pm6.38}$ | $54.29_{\pm5.20}$ | $45.25_{\pm2.80}$ | $55.30_{\pm1.54}$ | $36.67_{\pm0.83}$ | 0 |
| | graph2vec | $83.15_{\pm9.25}$ | $60.17_{\pm6.86}$ | $72.96_{\pm1.89}$ | $\mathbf{71.65}_{\pm3.10}$ | $71.10_{\pm0.54}$ | $\mathbf{50.44}_{\pm0.87}$ | 2 |
| | InfoGraph$_{\text{GIN}}$ | $\mathbf{89.13}_{\pm1.01}$ | $61.65_{\pm1.43}$ | $\mathbf{74.88}_{\pm4.31}$ | $39.52_{\pm3.99}$ | $\mathbf{73.90}_{\pm0.87}$ | $49.29_{\pm0.53}$ | 4 |
| | InfoGraph$_{\text{DP}}{}^*$ | $84.28_{\pm3.94}$ | $\mathbf{62.26}_{\pm4.55}$ | $73.50_{\pm2.91}$ | $\mathbf{61.93}_{\pm4.64}$ | $68.50_{\pm5.07}$ | $44.79_{\pm3.33}$ | 2 |
| FMs | GWF | $78.25_{\pm3.67}$ | $\mathbf{61.87}_{\pm2.53}$ | $73.19_{\pm1.97}$ | $\mathbf{72.11}_{\pm4.00}$ | $60.90_{\pm2.68}$ | $39.97_{\pm1.35}$ | 2 |
| | GDL$^*$ | $78.18_{\pm2.37}$ | $60.32_{\pm1.35}$ | $74.29_{\pm3.60}$ | $\mathbf{71.15}_{\pm3.19}$ | $71.70_{\pm1.10}$ | $49.12_{\pm0.49}$ | 3 |
| **Ours** | GNAE$_{\text{FGW}}$ | $79.53_{\pm5.79}$ | $61.43_{\pm4.28}$ | $\mathbf{75.32}_{\pm2.88}$ | $48.00_{\pm6.36}$ | $\mathbf{72.50}_{\pm4.30}$ | $\mathbf{47.30}_{\pm1.97}$ | 3 |
| | GNAE$_{\text{RAML}}$ | $79.76_{\pm3.88}$ | $\mathbf{61.75}_{\pm6.29}$ | $\mathbf{75.78}_{\pm3.42}$ | $50.70_{\pm4.14}$ | $\mathbf{73.10}_{\pm3.75}$ | $46.67_{\pm3.33}$ | 4 |

[1] The methods marked by "$^*$" are implemented by ourselves.

[2] For each dataset, the bold numbers are the five highest accuracy (top-5 results).

it by the FGW barycenter of its factors [59], which needs to compute $C$ FGW distances iteratively. Therefore, its computational complexity is at least $\mathcal{O}(CN^3)$. Both the GNAE$_{\text{FGW}}$ and the GDL applies linear factorization models, so they only need to compute one FGW distance between the input and the reconstruction, whose complexity is $\mathcal{O}(|\mathcal{P}|^2 N)$ and $\mathcal{O}(N^3)$, respectively. Here, $\mathcal{P}$ is the partitions of the graphon reconstructed by the GNAE$_{\text{FGW}}$. Proposition 2.4 shows that $|\mathcal{P}| \geq N$, so the GNAE$_{\text{FGW}}$ is slightly slower than the GDL. Our GNAE$_{\text{RAML}}$ samples $I$ small graphs and computes $I$ FGW distances, each of which is a pair of two sparse matrices. Denote the number of nodes in each small graph as $K$. The computational complexity of our GNAE$_{\text{RAML}}$ is $\mathcal{O}(IEK)$. Because $E \ll N^2$, $K \ll N$, and we set $I = \mathcal{O}(C)$, our GNAE$_{\text{RAML}}$ owns the lowest computational complexity. Figure 3(a) shows the training time per epoch of different models on the IMDB-B dataset, which verifies our analysis above — our GNAE$_{\text{RAML}}$ is $\times 4$ faster than the GDL and GNAE$_{\text{FGW}}$ and $\times 12$ faster than the GWF. Note that because the implementation of the GDL does not support GPU computing, we test all the methods on a single core of a CPU (Core i7 2.5GHz) for fairness. Figure 3(b) visualize some typical graphs in the IMDB-B dataset and the graphon factors learned by our GNAE$_{\text{RAML}}$ (*i.e.*, $\{\tilde{g}_c\}_{c=1}^{15}$). We can find that the IMDB-B graphs are formulated as communities connected by one or two central nodes. The graphon factors we learned reflect the topological property of the graphs, which further demonstrates the rationality of our GNAE model.

## 5.2 Generalizability and transferability on social network modeling

Our GNAE can generate graphons from graph representations and sampling graphs with different sizes but similar topological structures. Again, take the IMDB-B dataset as an example. For this dataset, the average number of nodes per graph is 19.77. Given a GNAE trained on this dataset, we sample graph representations from learned prior distribution and generate graphons by the decoder of the GNAE, *i.e.*, $\hat{g} = h(\boldsymbol{z})$ with $\boldsymbol{z} \sim p_{\mathcal{Z}}$. Based on $\hat{g}$, we sample graphs with different sizes, as shown in Figure 4. We can find that the generated graphs have similar structures, each containing two communities connected by few key nodes. Note that this topological structure is typical for the real IMDB graph (as shown in Figure 3(b)). This experimental result demonstrates that our GNAE
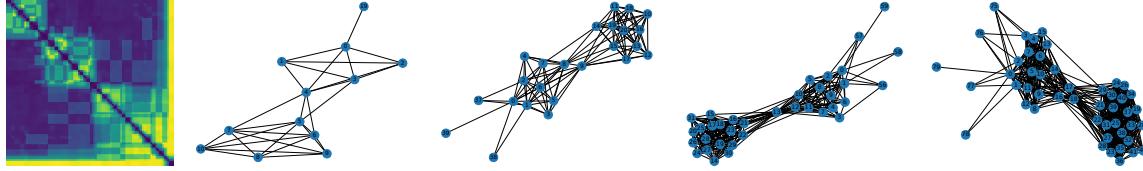
Figure 4: Illustrations of the graphs sampled from the generated graphon on the left. From left to right, the number of nodes for each graph is $20, 40, 60, 80$, respectively.

Table 2: Comparison on the classification accuracy (%) achieved by transfer learning

| Method | Training $\rightarrow$ Testing | | | |
|---|---|---|---|---|
| | IMDB-B $\rightarrow$ IMDB-B | IMDB-M $\rightarrow$ IMDB-B | IMDB-M $\rightarrow$ IMDB-M | IMDB-B $\rightarrow$ IMDB-M |
| InfoGraph$_{\text{GIN}}$ | $73.90_{\pm 0.87}$ | $66.10_{\pm 1.90}$ | $49.29_{\pm 0.53}$ | $45.29_{\pm 1.28}$ |
| GNAE$_{\text{RAML}}$ | $73.60_{\pm 3.80}$ | $70.70_{\pm 3.49}$ | $46.93_{\pm 3.14}$ | $46.20_{\pm 3.50}$ |

has the potentials as a graph generator with strong generalizability, which is especially suitable for social network modeling and simulation.

Another advantage of our GNAE is its transferability, which is seldom considered by existing work. In particular, we can train a GNAE on a dataset and use it to represent the graphs in a related but different dataset. For example, both the IMDB-B and the IMDB-M are movie collaboration datasets. Each graph in these two datasets is an ego-network of an actor/actress, which indicates his/her collaborations with other actors/actresses [65]. The IMDB-B contains 1000 ego-networks driven by two genres (*Action* and *Romance*), while the IMDB-M contains 1500 ego-networks driven by three genres (*Comedy*, *Romance* and *Sci-Fi*). Obviously, these two datasets have different structures but share some information. To demonstrate the transferability of our model, we first train a GNAE model on one dataset. Then, without any fine-tuning, we leverage the model to represent the graphs in the other dataset. Finally, we train and test an SVM classifier based on the representations and record the classification accuracy achieved by 10-fold cross-validation. Table 2 shows the performance of our GNAE and the strongest baseline InfoGraph$_{\text{GIN}}$ in the transfer learning scenarios. We can find that the performance of the InfoGraph$_{\text{GIN}}$ drops a lot when doing transfer learning. On the contrary, our GNAE shows good transferability, whose performance only degrades slightly. It captures the structural information shared by the two datasets, making the model transferable.

# 6    Conclusion and Future Work

We proposed a novel graphon autoencoder associated with an efficient learning algorithm. It is pioneering work achieving an interpretable and scalable graph generative model. Currently, the main advantages of our GNAE, *e.g.*, its generalizability and transferability, are demonstrated on social network modeling. However, as shown in Table 1, we need to improve the GNAE model for other graph types like proteins, molecules, and more complicated heterogeneous graphs and hypergraphs. Additionally, we will explore other potential substitutes for the FGW distance to further improve the efficiency of our learning algorithm.

# References

[1] Bijaya Adhikari, Yao Zhang, Naren Ramakrishnan, and B Aditya Prakash. Sub2vec: Feature learning for subgraphs. In *PAKDD*, 2018.

[2] Edo M Airoldi, Thiago B Costa, and Stanley H Chan. Stochastic blockmodel approximation of a graphon: theory and consistent estimation. In *Advances in Neural Information Processing Systems*, pages 692–700, 2013.

[3] Marco Avella-Medina, Francesca Parise, Michael Schaub, and Santiago Segarra. Centrality measures for graphons: Accounting for uncertainty in networks. *IEEE Transactions on Network Science and Engineering*, 2018.

[4] Christian Borgs, Jennifer T Chayes, László Lovász, Vera T Sós, and Katalin Vesztergombi. Convergent sequences of dense graphs I: Subgraph frequencies, metric properties and testing. *Advances in Mathematics*, 219(6):1801–1851, 2008.

[5] Karsten M Borgwardt, Cheng Soon Ong, Stefan Schönauer, SVN Vishwanathan, Alex J Smola, and Hans-Peter Kriegel. Protein function prediction via graph kernels. *Bioinformatics*, 21(suppl_1):i47–i56, 2005.

[6] Ron Bracewell. The fourier transform and its applications. *American Journal of Physics*, 34(8):712–712, 1966.

[7] Chen Cai and Yusu Wang. A simple yet effective baseline for non-attributed graph classification. *arXiv preprint arXiv:1811.03508*, 2018.

[8] Stanley Chan and Edoardo Airoldi. A consistent histogram estimator for exchangeable graph models. In *International Conference on Machine Learning*, pages 208–216, 2014.

[9] Chih-Chung Chang and Chih-Jen Lin. Libsvm: a library for support vector machines. *ACM transactions on intelligent systems and technology (TIST)*, 2(3):1–27, 2011.

[10] Antoine Channarond, Jean-Jacques Daudin, Stéphane Robin, et al. Classification and estimation in the stochastic blockmodel based on the empirical degrees. *Electronic Journal of Statistics*, 6:2574–2601, 2012.

[11] Sourav Chatterjee et al. Matrix estimation by universal singular value thresholding. *The Annals of Statistics*, 43(1):177–214, 2015.

[12] Hanjun Dai, Azade Nazi, Yujia Li, Bo Dai, and Dale Schuurmans. Scalable deep generative modeling for sparse graphs. In *International Conference on Machine Learning*, pages 2302–2312. PMLR, 2020.

[13] Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. Convolutional neural networks on graphs with fast localized spectral filtering. In *International Conference on Neural Information Processing Systems*, pages 3844–3852, 2016.

[14] Paul Erdős and Alfréd Rényi. On the evolution of random graphs. *Publ. Math. Inst. Hung. Acad. Sci*, 5(1):17–60, 1960.

[15] Shuang Gao and Peter E Caines. Graphon control of large-scale networks of linear systems. *IEEE Transactions on Automatic Control*, 2019.

[16] Thomas Gärtner, Peter Flach, and Stefan Wrobel. On graph kernels: Hardness results and efficient alternatives. In *Learning theory and kernel machines*, pages 129–143. Springer, 2003.

[17] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems*, pages 2672–2680, 2014.

[18] Svante Janson. Graphons, cut norm and distance, couplings and rearrangements. *New York journal of mathematics*, 2013.

[19] Wengong Jin, Regina Barzilay, and Tommi Jaakkola. Hierarchical generation of molecular graphs using structural motifs. In *International Conference on Machine Learning*, pages 4839–4848. PMLR, 2020.

[20] Stasys Jukna. *Extremal combinatorics: with applications in computer science*. Springer Science & Business Media, 2011.

[21] Raghunandan H Keshavan, Andrea Montanari, and Sewoong Oh. Matrix completion from a few entries. *IEEE Transactions on Information Theory*, 56(6):2980–2998, 2010.

[22] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

[23] Diederik P Kingma and Prafulla Dhariwal. Glow: generative flow with invertible $1 \times 1$ convolutions. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pages 10236–10245, 2018.

[24] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.

[25] Thomas N Kipf and Max Welling. Variational graph auto-encoders. *Stat*, 1050:21, 2016.

[26] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations*, 2017.

[27] Soheil Kolouri, Phillip E Pope, Charles E Martin, and Gustavo K Rohde. Sliced-wasserstein autoencoder: an embarrassingly simple generative model. *arXiv preprint arXiv:1804.01947*, 2018.

[28] Risi Kondor and Horace Pan. The multiscale Laplacian graph kernel. In *NeurIPS*, 2016.

[29] Nils Kriege and Petra Mutzel. Subgraph matching kernels for attributed graphs. In *ICML*, 2012.

[30] Yujia Li, Daniel Tarlow, Marc Brockschmidt, and Richard Zemel. Gated graph sequence neural networks. In *International Conference on Learning Representations*, 2016.

[31] László Lovász. *Large networks and graph limits*, volume 60. American Mathematical Soc., 2012.

[32] Giulia Luise, Alessandro Rudi, Massimiliano Pontil, and Carlo Ciliberto. Differential properties of Sinkhorn approximation for learning with Wasserstein distance. In *NeurIPS*, 2018.

[33] Facundo Mémoli. Gromov-Wasserstein distances and the metric approach to object matching. *Foundations of Computational Mathematics*, 11(4):417–487, 2011.

[34] Matthew W Morency and Geert Leus. Graphon filters: Graph signal processing in the limit. *IEEE Transactions on Signal Processing*, 69:1740–1754, 2021.

[35] Christopher Morris, Nils M. Kriege, Franka Bause, Kristian Kersting, Petra Mutzel, and Marion Neumann. Tudataset: A collection of benchmark datasets for learning with graphs. In *ICML 2020 Workshop on Graph Representation Learning and Beyond (GRL+ 2020)*, 2020.

[36] Annamalai Narayanan, Mahinthan Chandramohan, Rajasekar Venkatesan, Lihui Chen, Yang Liu, and Shantanu Jaiswal. graph2vec: Learning distributed representations of graphs. *arXiv preprint arXiv:1707.05005*, 2017.

[37] Mathias Niepert, Mohamed Ahmed, and Konstantin Kutzkov. Learning convolutional neural networks for graphs. In *International Conference on Machine Learning*, pages 2014–2023. PMLR, 2016.

[38] Mohammad Norouzi, Samy Bengio, Navdeep Jaitly, Mike Schuster, Yonghui Wu, Dale Schuurmans, et al. Reward augmented maximum likelihood for neural structured prediction. *Advances In Neural Information Processing Systems*, 29:1723–1731, 2016.

[39] Shirui Pan, Ruiqi Hu, Guodong Long, Jing Jiang, Lina Yao, and Chengqi Zhang. Adversarially regularized graph autoencoder for graph embedding. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence*, pages 2609–2615, 2018.

[40] Francesca Parise and Asuman Ozdaglar. Graphon games: A statistical framework for network games and interventions. *arXiv preprint arXiv:1802.00080*, 2018.

[41] Gabriel Peyré, Marco Cuturi, and Justin Solomon. Gromov-Wasserstein averaging of kernel and distance matrices. In *International Conference on Machine Learning*, pages 2664–2672, 2016.

[42] Kaspar Riesen and Horst Bunke. *Graph classification and clustering based on vector space embedding*. World Scientific, 2010.

[43] Luana Ruiz, Luiz Chamon, and Alejandro Ribeiro. Graphon neural networks and the transferability of graph neural networks. *Advances in Neural Information Processing Systems*, 33, 2020.

[44] Luana Ruiz, Luiz FO Chamon, and Alejandro Ribeiro. The graphon Fourier transform. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5660–5664. IEEE, 2020.

[45] Nino Shervashidze, Pascal Schweitzer, Erik Jan van Leeuwen, Kurt Mehlhorn, and Karsten M Borgwardt. Weisfeiler-Lehman graph kernels. *Journal of Machine Learning Research*, 12(Sep):2539–2561, 2011.

[46] Nino Shervashidze, SVN Vishwanathan, Tobias Petri, Kurt Mehlhorn, and Karsten Borgwardt. Efficient graphlet kernels for large graph comparison. In *Artificial Intelligence and Statistics*, pages 488–495, 2009.

[47] Chence Shi, Minkai Xu, Zhaocheng Zhu, Weinan Zhang, Ming Zhang, and Jian Tang. Graphaf: a flow-based autoregressive model for molecular graph generation. In *International Conference on Learning Representations*, 2019.

[48] Fan-Yun Sun, Jordan Hoffman, Vikas Verma, and Jian Tang. InfoGraph: Unsupervised and semi-supervised graph-level representation learning via mutual information maximization. In *International Conference on Learning Representations*, 2019.

[49] Hiroshi Takahashi, Tomoharu Iwata, Yuki Yamanaka, Masanori Yamada, and Satoshi Yagi. Variational autoencoder with implicit optimal priors. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 5066–5073, 2019.

[50] Vayer Titouan, Laetitia Chapel, Rémi Flamary, Romain Tavenard, and Nicolas Courty. Fused Gromov-Wasserstein distance for structured objects. *Algorithms*, 13(9):212, 2020.

[51] Vayer Titouan, Nicolas Courty, Romain Tavenard, and Rémi Flamary. Optimal transport for structured data with application on graphs. In *International Conference on Machine Learning*, pages 6275–6284, 2019.

[52] I Tolstikhin, O Bousquet, S Gelly, and B Schölkopf. Wasserstein auto-encoders. In *International Conference on Learning Representations*. OpenReview. net, 2018.

[53] Jakub Tomczak and Max Welling. Vae with a vampprior. In *International Conference on Artificial Intelligence and Statistics*, pages 1214–1223, 2018.

[54] Cédric Vincent-Cuaz, Titouan Vayer, Rémi Flamary, Marco Corneli, and Nicolas Courty. Online graph dictionary learning. *arXiv preprint arXiv:2102.06555*, 2021.

[55] Oriol Vinyals, Samy Bengio, and Manjunath Kudlur. Order matters: Sequence to sequence for sets. In *International Conference on Learning Representations*, 2016.

[56] Daixin Wang, Peng Cui, and Wenwu Zhu. Structural deep network embedding. In *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 1225–1234, 2016.

[57] Wenlin Wang, Zhe Gan, Hongteng Xu, Ruiyi Zhang, Guoyin Wang, Dinghan Shen, Changyou Chen, and Lawrence Carin. Topic-guided variational auto-encoder for text generation. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 166–177, 2019.

[58] Yujia Xie, Yixiu Mao, Simiao Zuo, Hongteng Xu, Xiaojing Ye, Tuo Zhao, and Hongyuan Zha. A hypergradient approach to robust regression without correspondence. In *International Conference on Learning Reresentations*, 2021.

[59] Hongteng Xu. Gromov-wasserstein factorization models for graph clustering. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 6478–6485, 2020.

[60] Hongteng Xu, Dixin Luo, and Lawrence Carin. Scalable Gromov-Wasserstein learning for graph partitioning and matching. In *Advances in Neural Information Processing Systems*, pages 3052–3062, 2019.

[61] Hongteng Xu, Dixin Luo, Lawrence Carin, and Hongyuan Zha. Learning graphons via structured gromov-wasserstein barycenters. In *AAAI*, 2021.

[62] Hongteng Xu, Dixin Luo, Ricardo Henao, Svati Shah, and Lawrence Carin. Learning autoencoders with relational regularization. In *International Conference on Machine Learning*, pages 10576–10586. PMLR, 2020.

[63] Jiaming Xu. Rates of convergence of spectral methods for graphon estimation. In *International Conference on Machine Learning*, pages 5433–5442, 2018.

[64] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? In *International Conference on Learning Representations*, 2018.

[65] Pinar Yanardag and SVN Vishwanathan. Deep graph kernels. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1365–1374, 2015.

[66] Carl Yang, Peiye Zhuang, Wenhan Shi, Alan Luu, and Pan Li. Conditional structure generation through graph variational generative adversarial nets. In *NeurIPS*, 2019.

[67] Rex Ying, Jiaxuan You, Christopher Morris, Xiang Ren, William L Hamilton, and Jure Leskovec. Hierarchical graph representation learning with differentiable pooling. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pages 4805–4815, 2018.

[68] Jiaxuan You, Bowen Liu, Rex Ying, Vijay Pande, and Jure Leskovec. Graph convolutional policy network for goal-directed molecular graph generation. In *International Conference on Neural Information Processing Systems*, pages 6412–6422, 2018.

[69] Muhan Zhang and Yixin Chen. Link prediction based on graph neural networks. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pages 5171–5181, 2018.

# 7 Appendix

## 7.1 The derivation of (9)

Suppose that $\mathcal{P} = \{\mathcal{P}_n\}_{n=1}^N$ and $\mathcal{Q} = \{\mathcal{Q}_m\}_{m=1}^M$ are two sets of partitions in $\Omega$. Given $\boldsymbol{x}_{1,\mathcal{P}} = (g_{1,\mathcal{P}}, s_{1,\mathcal{P}})$ and $\boldsymbol{x}_{2,\mathcal{Q}} = (g_{2,\mathcal{Q}}, s_{2,\mathcal{Q}})$, where $g_{1,\mathcal{P}}(v, v') = \sum_{n,n'=1}^N g_{1,nn'} 1_{\mathcal{P}_n}(v) 1_{\mathcal{P}_{n'}}(v')$, $g_{2,\mathcal{Q}}(v, v') = \sum_{m,m'=1}^M g_{2,mm'} 1_{\mathcal{Q}_m}(v) 1_{\mathcal{Q}_{m'}}(v')$, $s_{1,\mathcal{P}}(v) = \sum_{n=1}^N \boldsymbol{s}_{1,n} 1_{\mathcal{P}}(v)$, and $s_{2,\mathcal{Q}}(v) = \sum_{m=1}^M \boldsymbol{s}_{2,m} 1_{\mathcal{Q}}(v)$ are step functions, we have

$$d_{\mathrm{fgw}}(\boldsymbol{x}_{1,\mathcal{P}}, \boldsymbol{x}_{2,\mathcal{Q}})$$

$$= \inf_{\pi \in \Pi(\mu_\Omega, \mu_\Omega)} \Big( \int_{\Omega^2 \times \Omega^2} |g_{1,\mathcal{P}}(u, u') - g_{2,\mathcal{Q}}(v, v')|^p d\pi(u, v) d\pi(u', v') +$$

$$\int_{\Omega^2} \|s_{1,\mathcal{P}}(u) - s_{2,\mathcal{Q}}(v)\|_p^p d\pi(u, v) \Big)^{\frac{1}{p}}$$

$$= \inf_{\pi \in \Pi(\mu_\Omega, \mu_\Omega)} \Big( \int_{\Omega^2 \times \Omega^2} |\sum_{n,n'} g_{1,nn'} 1_{\mathcal{P}_n}(u) 1_{\mathcal{P}_{n'}}(u') -$$

$$\sum_{m,m'} g_{2,mm'} 1_{\mathcal{Q}_m}(v) 1_{\mathcal{Q}_{m'}}(v')|^p d\pi(u, v) d\pi(u', v') +$$

$$\int_{\Omega^2} \|\sum_n \boldsymbol{s}_{1,n} 1_{\mathcal{P}_m}(u) - \sum_m \boldsymbol{s}_{2,m} 1_{\mathcal{Q}_m}(v)\|_p^p d\pi(u, v) \Big)^{\frac{1}{p}}$$

$$= \inf_{\pi \in \Pi(\mu_\Omega, \mu_\Omega)} \Big( \sum_{n,n',m,m'} \int_{\mathcal{P}_n \times \mathcal{P}_{n'} \times \mathcal{Q}_m \times \mathcal{Q}_{m'}} |g_{1,nn'} - g_{2,mm'}|^p d\pi(u, v) d\pi(u', v') +$$

$$\sum_{n,m} \int_{\mathcal{P}_n \times \mathcal{Q}_m} \|\boldsymbol{s}_{1,n} - \boldsymbol{s}_{2,m}\|_p^p d\pi(u, v) \Big)^{\frac{1}{p}}$$

$$= \inf_{\pi \in \Pi(\mu_\Omega, \mu_\Omega)} \Big( \sum_{n,n',m,m'} |g_{1,nn'} - g_{2,mm'}|^p \underbrace{\int_{\mathcal{P}_n \times \mathcal{Q}_m} d\pi(u, v)}_{t_{nm}} \underbrace{\int_{\mathcal{P}_{n'} \times \mathcal{Q}_{m'}} d\pi(u', v')}_{t_{n'm'}} +$$

$$\sum_{n,m} \|\boldsymbol{s}_{1,n} - \boldsymbol{s}_{2,m}\|_p^p \underbrace{\int_{\mathcal{P}_n \times \mathcal{Q}_m} d\pi(u, v)}_{t_{nm}} \Big)^{\frac{1}{p}}$$

$$= \min_{\boldsymbol{T} \in \Pi(\boldsymbol{\mu}_\mathcal{P}, \boldsymbol{\mu}_\mathcal{Q})} (\langle \boldsymbol{D}_g, \boldsymbol{T} \otimes \boldsymbol{T} \rangle + \langle \boldsymbol{D}_s, \boldsymbol{T} \rangle)^{\frac{1}{p}},$$

where $\boldsymbol{D}_g = [|g_{1,nn'} - g_{2,mm'}|^p] \in \mathbb{R}^{N^2 \times M^2}$, $\boldsymbol{D}_s = [\|\boldsymbol{s}_{1,n} - \boldsymbol{s}_{2,m}\|_p^p] \in \mathbb{R}^{N \times M}$, $\otimes$ represents Kronecker product, $\boldsymbol{T} = [t_{nm}]$, and it is in the set $\Pi(\boldsymbol{\mu}_\mathcal{P}, \boldsymbol{\mu}_\mathcal{Q}) = \{\boldsymbol{T} \geq \boldsymbol{0} | \boldsymbol{T} \boldsymbol{1} = \boldsymbol{\mu}_\mathcal{P}, \boldsymbol{T}^\top \boldsymbol{1} = \boldsymbol{\mu}_\mathcal{Q}\}$ with $\boldsymbol{\mu}_\mathcal{P} = [\frac{|\mathcal{P}_1|}{|\Omega|}, .., \frac{|\mathcal{P}_N|}{|\Omega|}]$ and $\boldsymbol{\mu}_\mathcal{Q} = [\frac{|\mathcal{Q}_1|}{|\Omega|}, .., \frac{|\mathcal{Q}_M|}{|\Omega|}]$.

## 7.2 The computation of the FGW distance

For the FGW distance shown in (9), we set $p = 2$ as [41]. Denote $\boldsymbol{G}_1 = [g_{1,nn'}]$ and $\boldsymbol{G}_2 = [g_{2,mm'}]$ as the matrices corresponding to the step functions $g_{1,\mathcal{P}}$ and $g_{2,\mathcal{Q}}$. Applying $p = 2$, we can rewrite the objective function as

$$\langle \boldsymbol{D}(\boldsymbol{T}), \boldsymbol{T} \rangle = \langle \boldsymbol{D}_s + \boldsymbol{G}_{12} - 2\boldsymbol{G}_1 \boldsymbol{T} \boldsymbol{G}_2^\top, \boldsymbol{T} \rangle, \tag{14}$$

where $\boldsymbol{G}_{12} = (\boldsymbol{G}_1 \odot \boldsymbol{G}_1) \boldsymbol{\mu}_\mathcal{P} \boldsymbol{1}_M^\top + \boldsymbol{1}_N \boldsymbol{\mu}_\mathcal{Q}^\top (\boldsymbol{G}_2 \odot \boldsymbol{G}_2)^\top$ [41]. $\odot$ represents Hadamard product and $\boldsymbol{1}_N$ is $N$-dimensional all-one vector.

We apply the proximal gradient algorithm in [59] to compute the FGW distance, which is shown in Algorithm 2. This algorithm ensures the optimal transport matrix to converge to a stationary point. In this work, we set the number of iterations $J_1 = 20$ and the number of Sinkhorn iterations $J_2 = 5$ when computing FGW distance.

---

**Algorithm 2** $\min_{\boldsymbol{T} \in \Pi(\boldsymbol{\mu}_{\mathcal{P}}, \boldsymbol{\mu}_{\mathcal{Q}})} \langle \boldsymbol{D}_g, \boldsymbol{T} \otimes \boldsymbol{T} \rangle + \langle \boldsymbol{D}_s, \boldsymbol{T} \rangle$

---

1: Compute $\boldsymbol{G}_{12} = (\boldsymbol{G}_1 \odot \boldsymbol{G}_1)\boldsymbol{\mu}_{\mathcal{P}}\mathbf{1}_M^\top + \mathbf{1}_N \boldsymbol{\mu}_{\mathcal{Q}}^\top (\boldsymbol{G}_2 \odot \boldsymbol{G}_2)^\top$
2: Initialize $\boldsymbol{T}^{(0)} = \boldsymbol{\mu}_{\mathcal{P}}\boldsymbol{\mu}_{\mathcal{Q}}^\top$, $\boldsymbol{a} = \boldsymbol{\mu}_{\mathcal{P}}$.
3: **for** $j = 0, ..., J_1 - 1$
4: $\quad \boldsymbol{C} = \exp(-\frac{1}{\beta}(\boldsymbol{D}_s + \boldsymbol{G}_{12} - 2\boldsymbol{G}_1 \boldsymbol{T} \boldsymbol{G}_2^\top)) \odot \boldsymbol{T}^{(j)}$
5: $\quad$ Sinkhorn iteration: **for** $n = 0, .., J_2 - 1$, $\boldsymbol{b} = \frac{\boldsymbol{\mu}_{\mathcal{Q}}}{\boldsymbol{C}^\top \boldsymbol{a}}$, $\boldsymbol{a} = \frac{\boldsymbol{\mu}_{\mathcal{P}}}{\boldsymbol{C}\boldsymbol{b}}$,
6: $\quad \boldsymbol{T}^{(j+1)} = \text{diag}(\boldsymbol{a})\boldsymbol{C}\text{diag}(\boldsymbol{b})$.
7: The optimal transport $\boldsymbol{T}^* = \boldsymbol{T}^{(J)}$.

---

## 7.3   The computation of the sliced FGW distance

**Definition 7.1** (Sliced FGW Distance)*. Denote $\mathcal{S}^{M-1} = \{\boldsymbol{\theta} \in \mathbb{R}^M | \|\boldsymbol{\theta}\|_2 = 1\}$ as the $M$-dimensional hypersphere and $\mu_{\mathcal{S}^{M-1}}$ the probability measure on $\mathcal{S}^{M-1}$. For the probability measures $p_{\mathcal{Z}}$ and $q_{\mathcal{Z}}$ on the metric space $(\mathcal{Z}, d_{\mathcal{Z}})$, their sliced fused Gromov-Wasserstein distance is*

$$d_{sfgw}(p_{\mathcal{Z}}, q_{\mathcal{Z}}) = \mathbb{E}_{\boldsymbol{\theta} \sim \mu_{\mathcal{S}^{M-1}}}[d_{fgw}(R_{\boldsymbol{\theta}}\#p_{\mathcal{Z}}, R_{\boldsymbol{\theta}}\#q_{\mathcal{Z}})],$$

*where $R_{\boldsymbol{\theta}}$ as the projection on $\boldsymbol{\theta}$, where $R_{\boldsymbol{\theta}}(\boldsymbol{z}) = \langle \boldsymbol{z}, \boldsymbol{\theta} \rangle$. $R_{\boldsymbol{\theta}}\#p$ represents the distribution after the projection, and $d_{fgw}(R_{\boldsymbol{\theta}}\#p_{\mathcal{Z}}, R_{\boldsymbol{\theta}}\#q_{\mathcal{Z}})$ is the FGW distance between the $R_{\boldsymbol{\theta}}\#p_{\mathcal{Z}}$ and $R_{\boldsymbol{\theta}}\#q_{\mathcal{Z}}$ defined on the 1D metric space $(R_{\boldsymbol{\theta}}(\mathcal{Z}), d_{R_{\boldsymbol{\theta}}(\boldsymbol{Z})})$.*

In our case, we compute the order-2 sliced FGW distance between the expected posterior $q_{\mathcal{Z};f}$ and the prior $p_{\mathcal{Z}}$, i.e., $d_{\text{sfgw}}(q_{\mathcal{Z};f}, p_{\mathcal{Z}})$. The sliced FGW distance is the expectation of 1D FGW distances under different projections. We can approximate it based on the samples of the distributions and the samples of the projections. In particular, given $\{\boldsymbol{z}_{1,i}\}_{i=1}^N \sim q_{\mathcal{Z};f}$, $\{\boldsymbol{z}_{2,i}\}_{i=1}^N \sim p_{\mathcal{Z}}$, and $L$ projections $\{R_{\boldsymbol{\theta}_l}\}_{l=1}^L$, where $\boldsymbol{\theta} \sim \mu_{\mathcal{S}^{M-1}}$, the empirical sliced FGW, denoted as $\hat{d}_{\text{sfgw}}$, is

$$\hat{d}_{\text{sfgw}}(q_{\mathcal{Z};f}, p_{\mathcal{Z}})$$

$$=\frac{1}{L}\sum_{l=1}^L \hat{d}_{\text{fgw}}(R_{\boldsymbol{\theta}_l}\#q_{\mathcal{Z};f}, R_{\boldsymbol{\theta}_l}\#p_{\mathcal{Z}})$$

$$=\frac{1}{NL}\sum_{l=1}^L \min_{\sigma \in \mathcal{U}_N} \sum_{i,j=1}^N ((R_{\boldsymbol{\theta}_l}(\boldsymbol{z}_{1,i}) - R_{\boldsymbol{\theta}_l}(\boldsymbol{z}_{1,j}))^2 - (R_{\boldsymbol{\theta}_l}(\boldsymbol{z}_{2,\sigma(i)}) - R_{\boldsymbol{\theta}_l}(\boldsymbol{z}_{2,\sigma(j)}))^2)^2 +$$

$$\sum_{i=1}^N (R_{\boldsymbol{\theta}_l}(\boldsymbol{z}_{1,i}) - R_{\boldsymbol{\theta}_l}(\boldsymbol{z}_{2,\sigma(i)}))^2$$

$$=\frac{1}{NL}\sum_{l=1}^L \min_{\sigma \in \{\sigma_a, \sigma_d\}} \sum_{i,j=1}^N ((R_{\boldsymbol{\theta}_l}(\boldsymbol{z}_{1,\sigma_a(i)}) - R_{\boldsymbol{\theta}_l}(\boldsymbol{z}_{1,\sigma_a(j)}))^2 - (R_{\boldsymbol{\theta}_l}(\boldsymbol{z}_{2,\sigma(i)}) - R_{\boldsymbol{\theta}_l}(\boldsymbol{z}_{2,\sigma(j)}))^2)^2 +$$

$$\sum_{i=1}^N (R_{\boldsymbol{\theta}_l}(\boldsymbol{z}_{1,\sigma_a(i)}) - R_{\boldsymbol{\theta}_l}(\boldsymbol{z}_{2,\sigma(i)}))^2.$$

Table 3: The setting related to the type of signal

| Signal type | Signal distribution | $p(\boldsymbol{s}|s(v))$ | $\alpha(\cdot)$ |
|---|---|---|---|
| Continuous | Gaussian | $\propto \exp(-\|\boldsymbol{s} - s(v)\|_2^2/2M\sigma^2)$ | — |
| Binary | Bernoulli | $\prod_{m=1}^M s_m(v)^{\boldsymbol{s}_m}(1 - s_m(v))^{1-\boldsymbol{s}_m}$ | sigmoid |
| One hot | Categorical | $\prod_{m=1}^M s_m(v)^{\boldsymbol{s}_m}$ | softmax |

Table 4: The setting of other hyperparameters

| Dataset | MUTAG | PTC-MR | PROTEIN | ENZYMES | IMDB-B | IMDB-M |
|---|---|---|---|---|---|---|
| $J$ | 8 | 4 | 4 | 8 | 4 | 4 |
| $D$ | 30 | 30 | 30 | 50 | 30 | 30 |
| $C$ | 15 | 15 | 30 | 30 | 15 | 15 |

Here, $\sigma(\cdot) \in \mathcal{U}_N$ represents a permutation of $\{1, ..., N\}$. $\sigma_a(\cdot)$ ($\sigma_d(\cdot)$) outputs the indices of the samples in an ascending (descending) order. The second equation is based on the fact that when computing the empirical FGW distance (*i.e.*, $\hat{d}_{\text{sfgw}}$) between a pair of $N$-sample sets in 1D space, the optimal transport matrix is a permutation matrix [62]. The third equation is based on the Theorem 3.4 in [62] — the empirical FGW distance in 1D space corresponds to the distance between the samples in either identity or anti-identity order. In our work, the number of samples $N$ is equal to the batch size we set. The number of projections we used is $L = 50$.

## 7.4 The setting of hyperparameters

We set the hyperparameters of our GNAE empirically. Some hyperparameters are fixed for all the datasets: the batch size is $N_b = 50$; the learning rate is 0.005; the number of epochs is 25; the number of sampled graphs per graphon is $I = 5$; the size of the sampled graphs is $K = 10$; the weight of regularizer $\gamma = 0.1$; the number of ChebConv layers is $J = 4$; the order of the FGW distance is $p = 2$.

The other hyperparameters are specified for various datasets, *e.g.*, the number of Gaussian components in the prior distribution is equal to the number of clusters in each dataset. According to the type of signal (or equivalently, the type of node attribute), we set the distribution of signal in (3, 10) and the activation function $\alpha(\cdot)$ in (8), as shown in Table 3. The other settings are listed in Table 4.

## 7.5 Initialization of graph/graphon factors

The GWF and the GDL learn $C$ graph factors, respectively. Their factors can be explained as representative adjacency matrices. Our GNAEs learn $C$ graphon factors, each of which is explained as the step function induced from a representative graph. Both the GNAEs and the GWF allow the factors to have different sizes (numbers of partitions), so we leverage $C$ observed graphs to initialize their factors. The GDL requires the factors to have the same size, so we initialize its factors as random matrices, whose sizes are equal to the average size of the factors used in the GWF.

## 7.6 Visualizations of the learned graph representations

Besides the six datasets reported in the main paper, we consider two more datasets: the AIDS dataset [42], which contains 2000 molecules that are active or inactive to HIV virus; the PROTEIN-S dataset, which is the same with the PROTEIN dataset but applies simplified node attributes.

(a) MUTAG

(b) AIDS
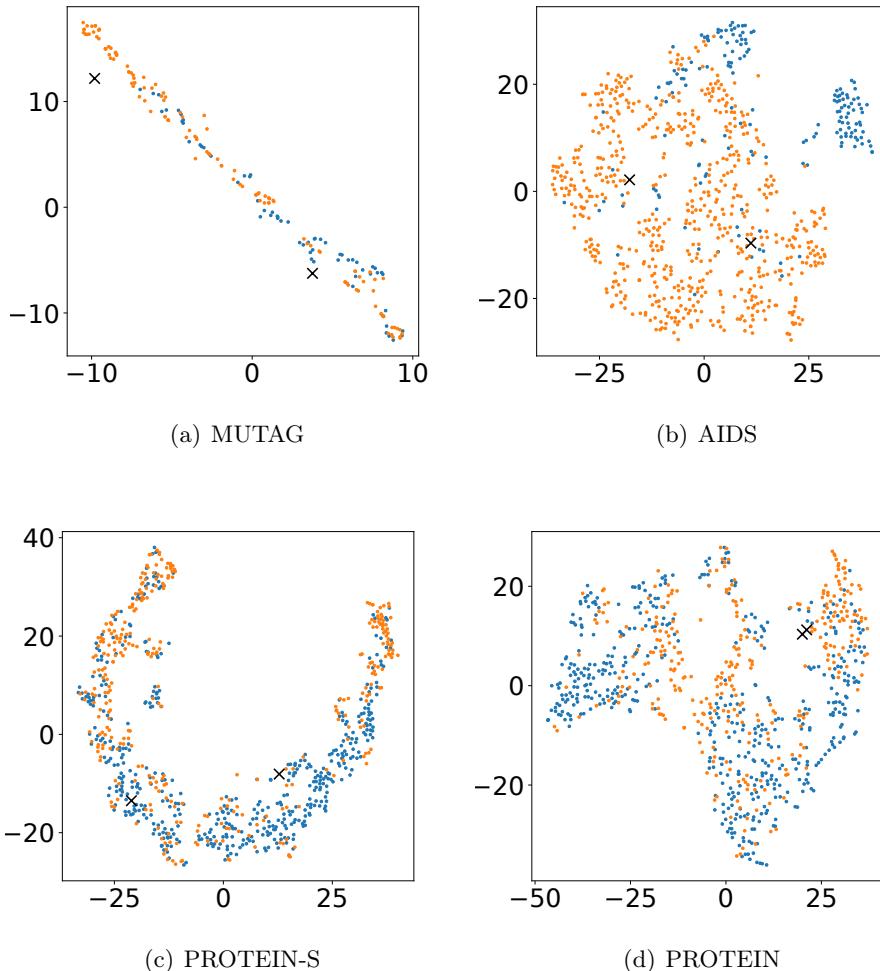
(c) PROTEIN-S

(d) PROTEIN

Figure 5: (Good cases) The t-SNE plots of the learned latent representations for some representative datasets. The colors of the points indicate the real categories of the representations. The black crosses indicate the centers of Gaussian components of the prior distribution.

Figures 5 and 6 show the t-SNE plots of the graph representations learned by our GNAE$_{\mathrm{RAML}}$. We can find that the representations indicate obvious clustering structures in some situations. However, in some challenging datasets, *e.g.*, the PTC-MR, the ENZYMES and the IMDB-M, our model does not work well. The reasons for the unsatisfying results may include:

- Data sparsity The PTC-MR just contains 344 molecules, so our model has a high risk of overfitting.

- Multi-class The graphs in the ENZYMES and the IMDB-M belong to multiple classes, which increases the difficulty of clustering.

(a) PTC-MR

(b) ENZYMES
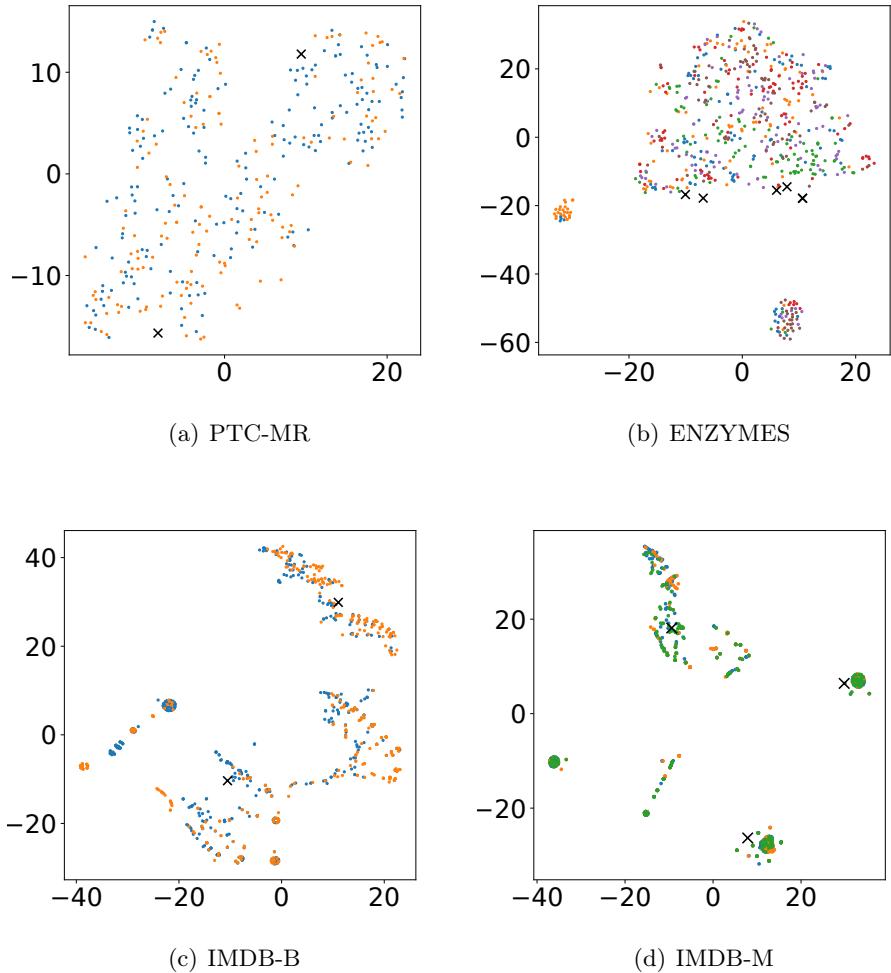
(c) IMDB-B

(d) IMDB-M

Figure 6: (Bad cases) The t-SNE plots of the learned latent representations for some representative datasets. The colors of the points indicate the real categories of the representations. The black crosses indicate the centers of Gaussian components of the prior distribution.