

# Gromov-Wasserstein Learning for Structured Data Modeling

Hongteng Xu

Gaoling School of Artificial Intelligence, Renmin University of China  
Beijing Key Laboratory of Big Data Management and Analysis Methods

February 23, 2022



中國人民大學  
RENMIN UNIVERSITY OF CHINA

高瓴人工智能学院  
Gaoling School of Artificial Intelligence

# Outline

## Part 1 Introduction of Gromov-Wasserstein Distance

- ▶ Preliminary and basic concepts
- ▶ Connections to structured data modeling and analysis

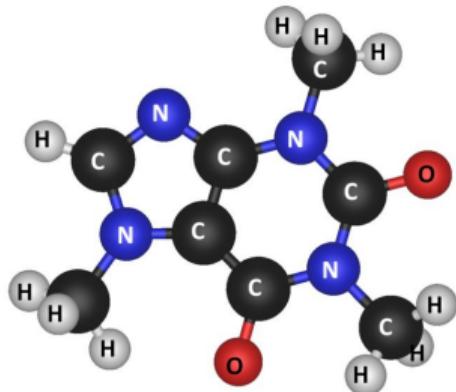
## Part 2 The Computation of Gromov-Wasserstein Distance

- ▶ Typical optimization algorithms
- ▶ The variants of Gromov-Wasserstein distance

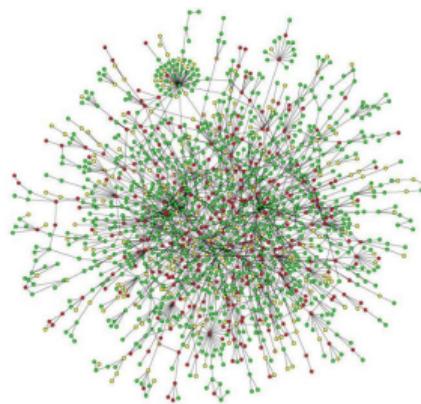
## Part 3 GW-based Models and Their ML Applications

- ▶ Generative modeling
- ▶ Graph representation
- ▶ Graph generation

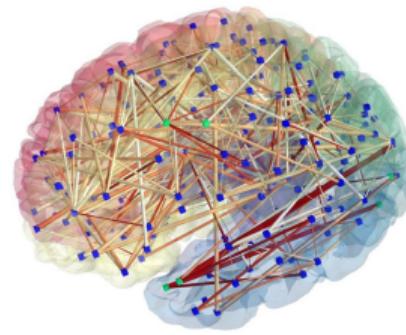
# Real-world Structured Data



Molecule

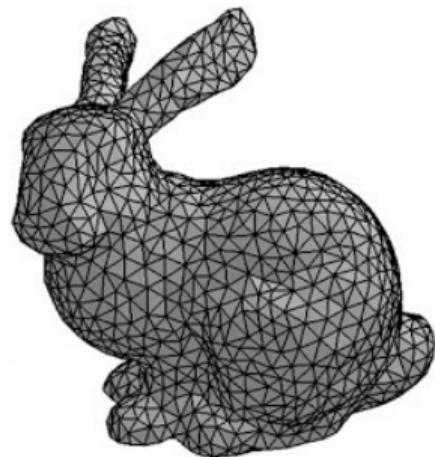


PPI Network

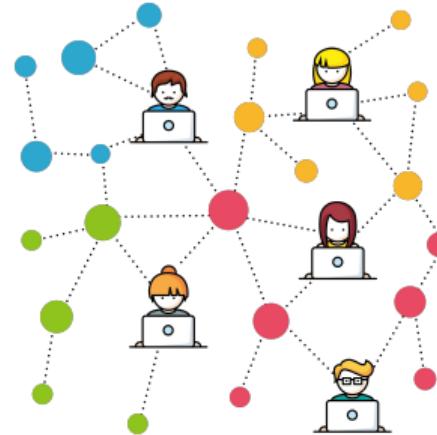


Brain Connectivity

# Real-world Structured Data



3D mesh



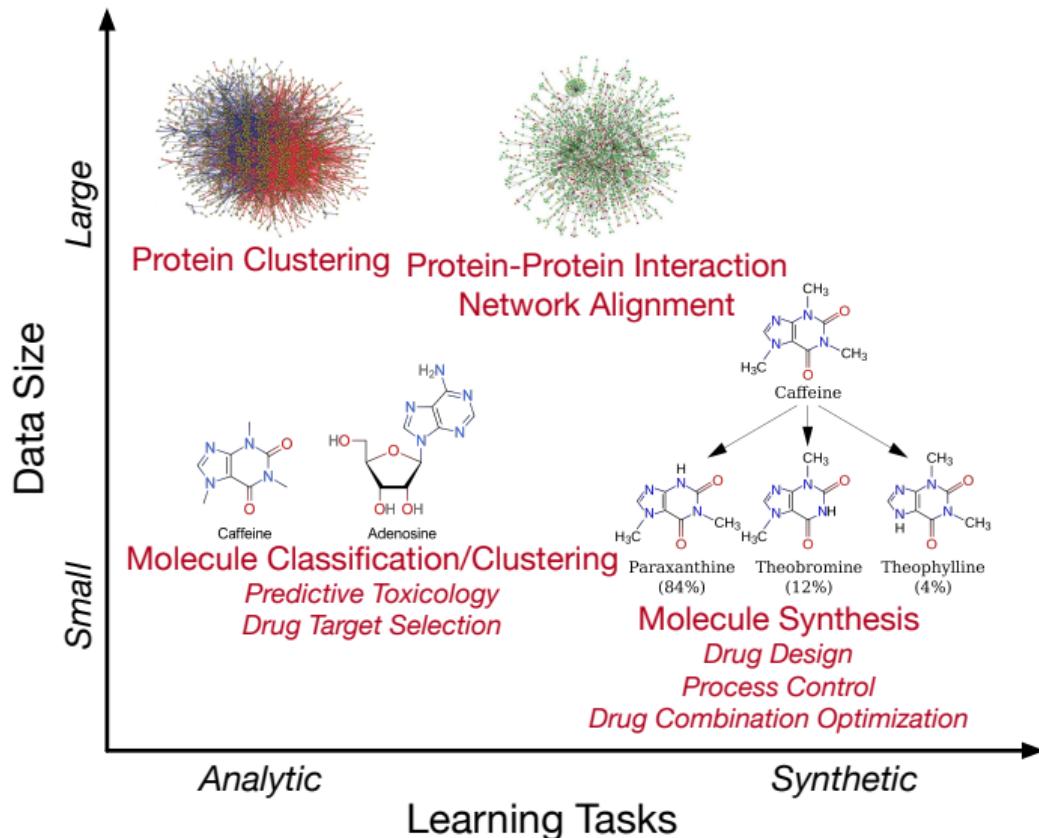
Social Network

≡ Google Scholar  🔍

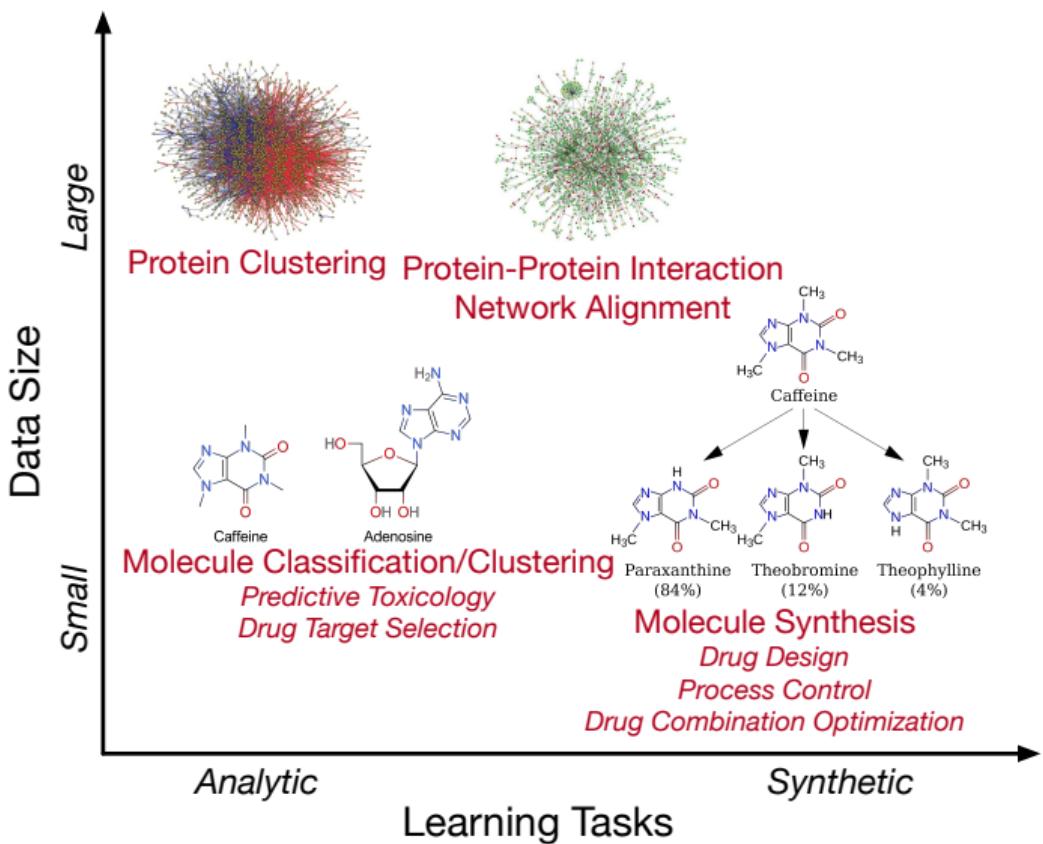
Articles About 4,700,000 results (0.13 sec)

Extract and predict relations among entities effectively and efficiently.

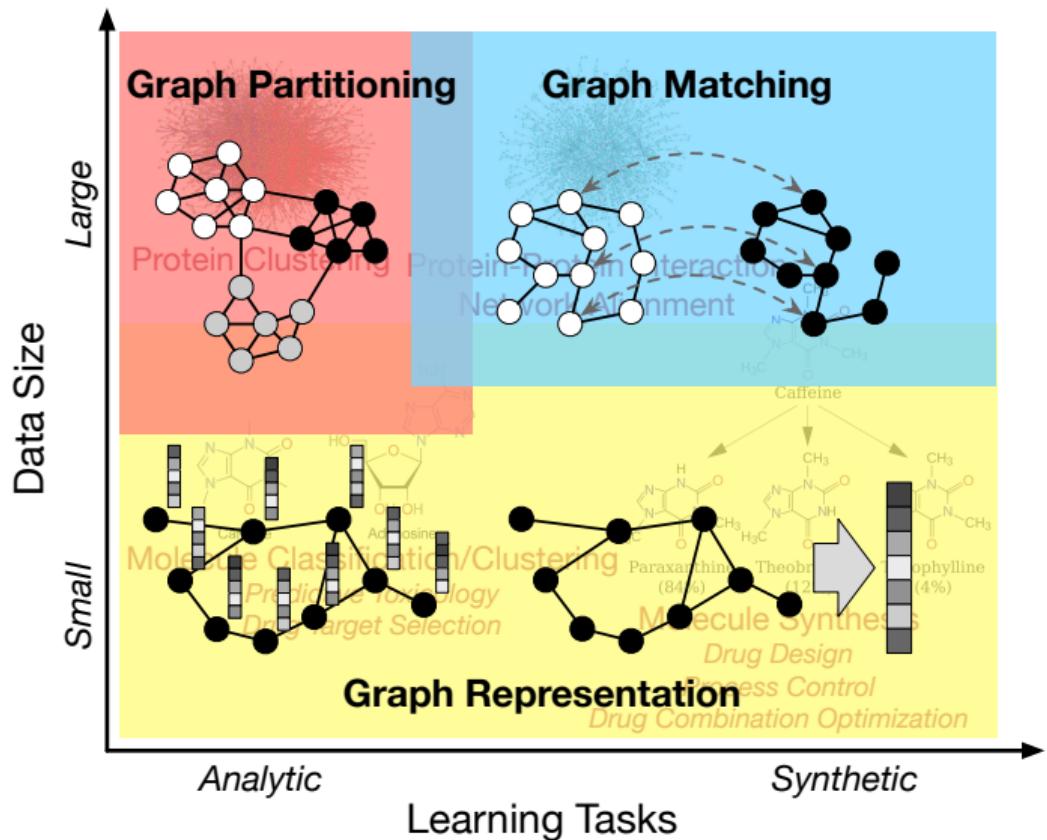
# Typical (Biochemical) Problems of Structured Data



# Typical (Biochemical) Problems of Structured Data



# Typical (Biochemical) Problems of Structured Data

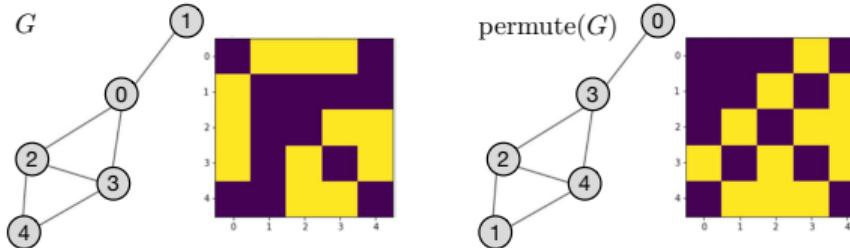


## The Challenges of The Problems

- ▶ NP-completeness
  - ▶ Approximation algorithms with high stability and scalability

# The Challenges of The Problems

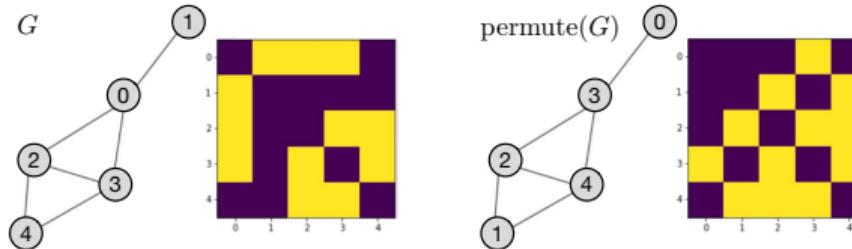
- ▶ NP-completeness
  - ▶ Approximation algorithms with high stability and scalability
- ▶ Permutation invariance



- ▶ A permutation-invariant metric  $d$ :  $d(G_X, G_Y) = d(G_X, \text{permute}(G_Y))$
- ▶ A permutation-invariant representation model  $f$ :  $f(G) = f(\text{permute}(G))$

# The Challenges of The Problems

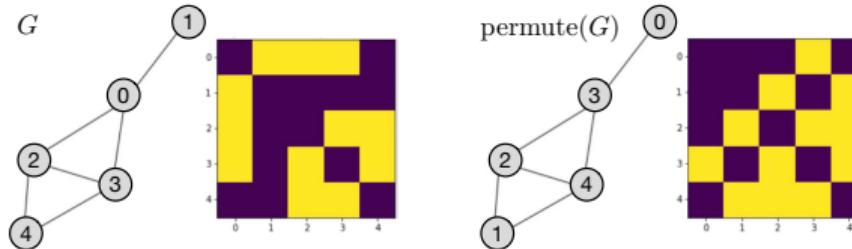
- ▶ NP-completeness
  - ▶ Approximation algorithms with high stability and scalability
- ▶ Permutation invariance



- ▶ A permutation-invariant metric  $d$ :  $d(G_X, G_Y) = d(G_X, \text{permute}(G_Y))$
- ▶ A permutation-invariant representation model  $f$ :  $f(G) = f(\text{permute}(G))$
- ▶ (Often) No labels
  - ▶ Unsupervised or semi-supervised learning

# The Challenges of The Problems

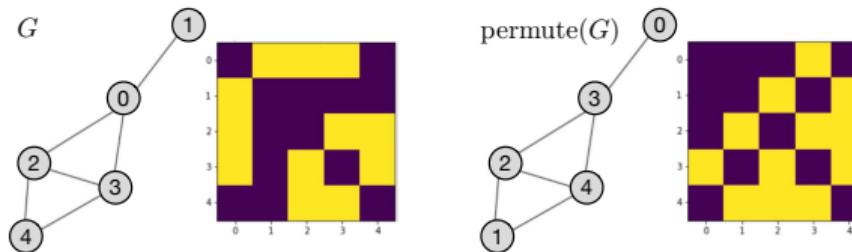
- ▶ NP-completeness
  - ▶ Approximation algorithms with high stability and scalability
- ▶ Permutation invariance



- ▶ A permutation-invariant metric  $d$ :  $d(G_X, G_Y) = d(G_X, \text{permute}(G_Y))$
- ▶ A permutation-invariant representation model  $f$ :  $f(G) = f(\text{permute}(G))$
- ▶ (Often) No labels
  - ▶ Unsupervised or semi-supervised learning
- ▶ Structured data generation
  - ▶ Generative model with generalizability and transferability

# The Challenges of The Problems

- ▶ NP-completeness
  - ▶ Approximation algorithms with high stability and scalability
- ▶ Permutation invariance



- ▶ A permutation-invariant metric  $d$ :  $d(G_X, G_Y) = d(G_X, \text{permute}(G_Y))$
- ▶ A permutation-invariant representation model  $f$ :  $f(G) = f(\text{permute}(G))$
- ▶ (Often) No labels
  - ▶ Unsupervised or semi-supervised learning
- ▶ Structured data generation
  - ▶ Generative model with generalizability and transferability

**Gromov-Wasserstein Learning (GWL): An optimal transport-based machine learning framework to solve the problems.**

# What Is GWL and Where Is It From?

**Landmarks in History (A very personal and non-theoretical viewpoint)**

## What Is GWL and Where Is It From?

### Landmarks in History (A very personal and non-theoretical viewpoint)

- ▶ 1781, Monge-form of optimal transport problem (Gaspard Monge)

# What Is GWL and Where Is It From?

## Landmarks in History (A very personal and non-theoretical viewpoint)

- ▶ **1781, Monge-form** of optimal transport problem (Gaspard Monge)
- ▶ **1939, Kantorovich-form** of optimal transport problem (Wasserstein distance) and its dual problem (Leonid Kantorovich)

# What Is GWL and Where Is It From?

## Landmarks in History (A very personal and non-theoretical viewpoint)

- ▶ **1781, Monge-form** of optimal transport problem (Gaspard Monge)
- ▶ **1939, Kantorovich-form** of optimal transport problem (Wasserstein distance) and its dual problem (Leonid Kantorovich)
- ▶ **1967, Sinkhorn-Knopp algorithm** for doubly-stochastic matrix (Richard Sinkhorn & Paul Knopp)

# What Is GWL and Where Is It From?

## Landmarks in History (A very personal and non-theoretical viewpoint)

- ▶ **1781, Monge-form** of optimal transport problem (Gaspard Monge)
- ▶ **1939, Kantorovich-form** of optimal transport problem (Wasserstein distance) and its dual problem (Leonid Kantorovich)
- ▶ **1967, Sinkhorn-Knopp algorithm** for doubly-stochastic matrix (Richard Sinkhorn & Paul Knopp)
- ▶ **1981, Gromovization of metric**, e.g., Gromov-Hausdorff (Mikhail Gromov)

# What Is GWL and Where Is It From?

## Landmarks in History (A very personal and non-theoretical viewpoint)

- ▶ **1781, Monge-form** of optimal transport problem (Gaspard Monge)
- ▶ **1939, Kantorovich-form** of optimal transport problem (Wasserstein distance) and its dual problem (Leonid Kantorovich)
- ▶ **1967, Sinkhorn-Knopp algorithm** for doubly-stochastic matrix (Richard Sinkhorn & Paul Knopp)
- ▶ **1981, Gromovization of metric**, e.g., Gromov-Hausdorff (Mikhail Gromov)
- ▶ **2006 - 2007, Gromov-Wasserstein distance** for geometric problems like shape matching (Karl-Theodor Sturm, Facundo Mémoli)

# What Is GWL and Where Is It From?

## Landmarks in History (A very personal and non-theoretical viewpoint)

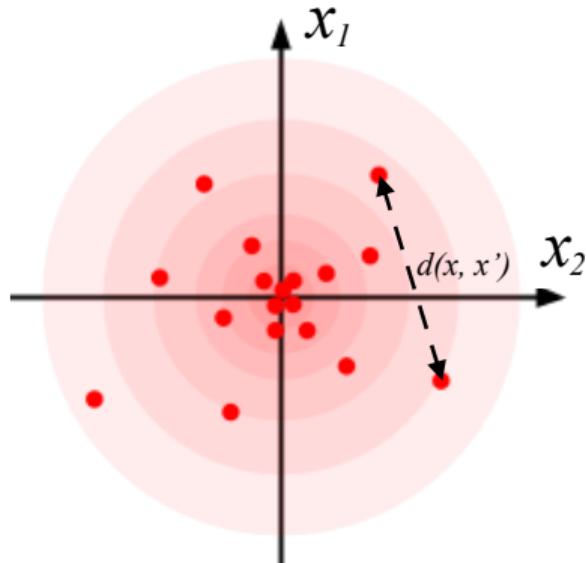
- ▶ **1781, Monge-form** of optimal transport problem (Gaspard Monge)
- ▶ **1939, Kantorovich-form** of optimal transport problem (Wasserstein distance) and its dual problem (Leonid Kantorovich)
- ▶ **1967, Sinkhorn-Knopp algorithm** for doubly-stochastic matrix (Richard Sinkhorn & Paul Knopp)
- ▶ **1981, Gromovization of metric**, e.g., Gromov-Hausdorff (Mikhail Gromov)
- ▶ **2006 - 2007, Gromov-Wasserstein distance** for geometric problems like shape matching (Karl-Theodor Sturm, Facundo Mémoli)
- ▶ **2013 - Present, Progress on computational optimal transport**
  - ▶ Entropic optimal transport and Sinkhorn scaling (2013, Marco Cuturi)
  - ▶ Variants of OT problems and GW distance, and Python OT package (Gabriel Peyré, Justin Solomon, Rémi Flamary, Nicolas Courty, ...)

# What Is GWL and Where Is It From?

## Landmarks in History (A very personal and non-theoretical viewpoint)

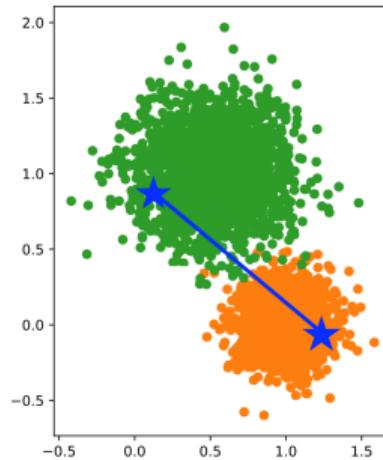
- ▶ **1781, Monge-form** of optimal transport problem (Gaspard Monge)
- ▶ **1939, Kantorovich-form** of optimal transport problem (Wasserstein distance) and its dual problem (Leonid Kantorovich)
- ▶ **1967, Sinkhorn-Knopp algorithm** for doubly-stochastic matrix (Richard Sinkhorn & Paul Knopp)
- ▶ **1981, Gromovization of metric**, e.g., Gromov-Hausdorff (Mikhail Gromov)
- ▶ **2006 - 2007, Gromov-Wasserstein distance** for geometric problems like shape matching (Karl-Theodor Sturm, Facundo Mémoli)
- ▶ **2013 - Present, Progress on computational optimal transport**
  - ▶ Entropic optimal transport and Sinkhorn scaling (2013, Marco Cuturi)
  - ▶ Variants of OT problems and GW distance, and Python OT package (Gabriel Peyré, Justin Solomon, Rémi Flamary, Nicolas Courty, ...)
- ▶ **2016 - Present, Blooming in machine learning and data science**
  - ▶ Wasserstein GAN (2017, Arjovsky, Martin, Soumith Chintala, and Léon Bottou)
  - ▶ Gromov-Wasserstein learning for graph modeling and analysis (2019, Hongteng Xu, Titouan Vayer, Makoto Yamada, Tam Le, Samir Chowdhury, Tom Needham, ...)

## Metric-Measure Space (MM-Space)



- ▶  $\mathcal{X}_{d,\mu} := (\mathcal{X}, d, \mu)$ : A metric-measure space, where  $x \in \mathcal{X}$  is a sample in the space.
- ▶  $d$ : A distance metric of samples (e.g., Euclidean distance).
- ▶  $\mathbb{P}$ : A space of (probabilistic) measures defined on  $\mathcal{X}$ .
- ▶  $\mu \in \mathbb{P}$ : a measure on  $\mathcal{X}$ .

## Optimal Transport between the Measures in the Same Space

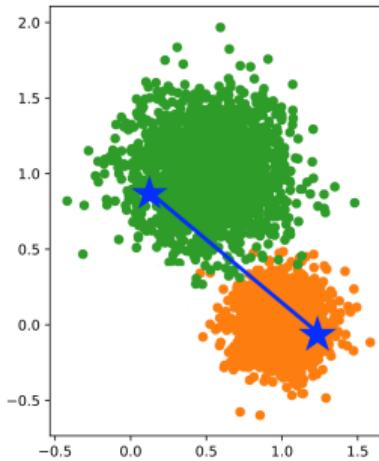


Kantorovich-form of optimal transport problem:

$$d_w(\mu, \eta) := \inf_{\pi \in \Pi(\mu, \eta)} \mathbb{E}_{x,y \sim \pi} [d(x, y)] = \inf_{\pi \in \Pi(\mu, \eta)} \iint_{(x,y) \in \mathcal{X}^2} d(x, y) \pi(x, y) dx dy, \quad (1)$$

where  $\Pi(\mu, \eta) = \{\pi > 0 \mid \int_x \pi(x, y) dx = \eta(y), \int_y \pi(x, y) dy = \mu(x)\}$ , the optimum  $\pi^*$  is called optimal transport (plan).

## Optimal Transport between the Measures in the Same Space

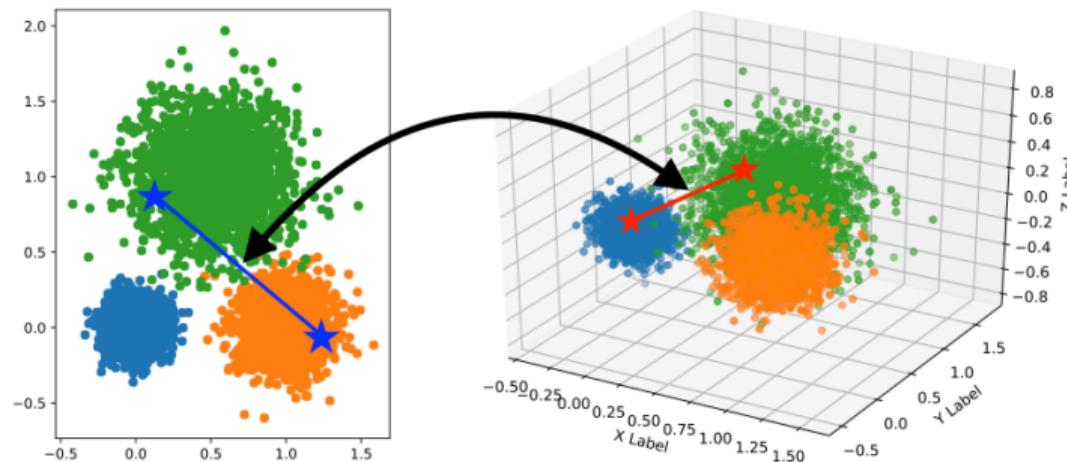


Given finite numbers of samples and their empirical distributions, i.e.,  $\{x_m\}_{m=1}^M$ ,  $\{y_n\}_{n=1}^N$ ,  $\mu \in \Delta^{N-1}$  and  $\eta \in \Delta^{M-1}$ , we have:

$$\hat{d}_w(\mu, \eta) = \min_{T \in \Pi(\mu, \eta)} \langle D, T \rangle = \min_{T \in \Pi(\mu, \eta)} \sum_{m=1}^M \sum_{n=1}^N d(x_m, y_n) T_{mn}, \quad (2)$$

where  $T = [T_{mn}]$ ,  $\Pi(\mu, \eta) = \{T > 0 | T\mathbf{1}_M = \mu, T^\top \mathbf{1}_N = \eta\}$ , the optimum  $T^*$  is called optimal transport matrix.

# Gromov-Wasserstein Distance for MM-Spaces

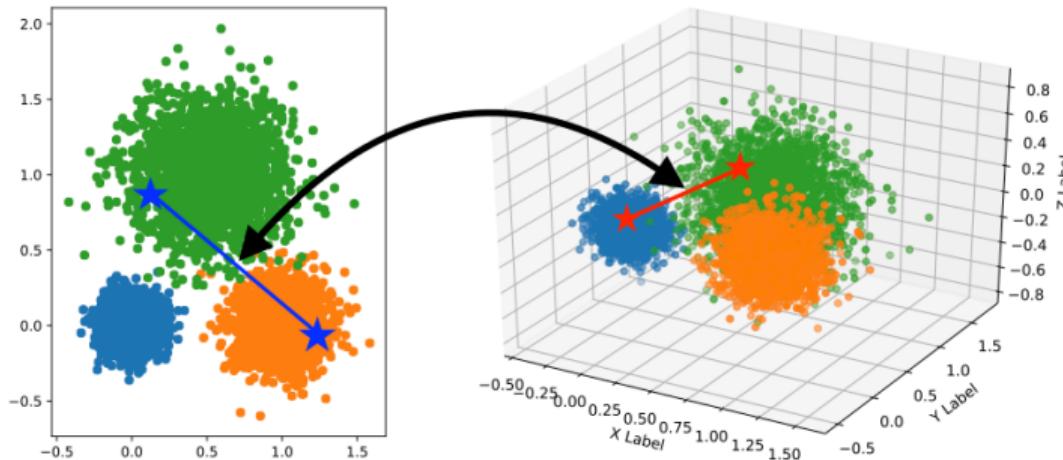


The Gromov-Wasserstein distance between  $\mathcal{X}_{d_X, \mu_X}$  and  $\mathcal{Y}_{d_Y, \mu_Y}$ :

$$\begin{aligned} d_{gw}(\mathcal{X}_{d_X, \mu_X}, \mathcal{Y}_{d_Y, \mu_Y}) &:= \inf_{\pi \in \Pi(\mu_X, \mu_Y)} \mathbb{E}_{x, x', y, y' \sim \pi \times \pi}[r(x, x', y, y')] \\ &= \inf_{\pi \in \Pi(\mu_X, \mu_Y)} \iint_{(x, x') \in \mathcal{X}^2} \iint_{(y, y') \in \mathcal{Y}^2} r(x, x', y, y') \pi(x, y) \pi(x', y') dx dx' dy dy', \end{aligned} \tag{3}$$

where  $r(x, x', y, y') = d(d_X(x, x'), d_Y(y, y'))$  is called **relational distance**.

# Gromov-Wasserstein Distance for Structured Data (Point Clouds)

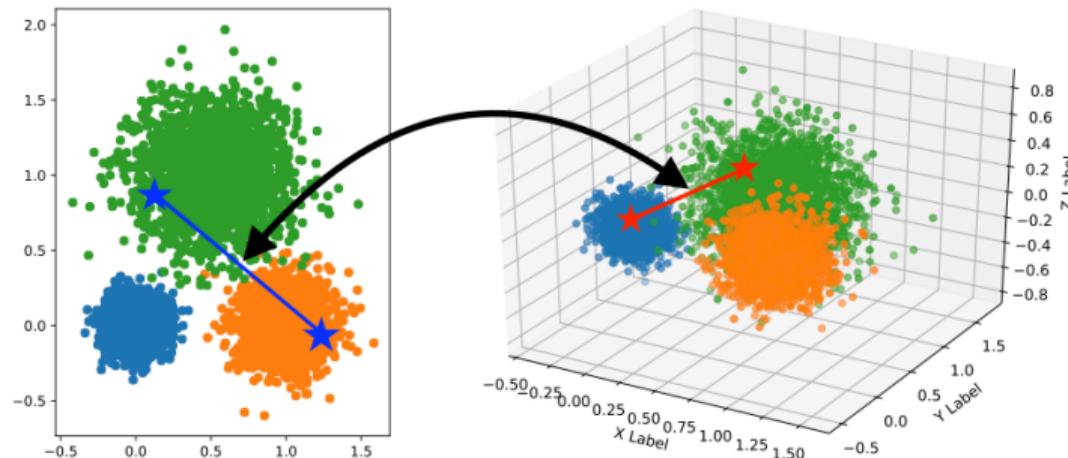


Given finite numbers of samples and their empirical distributions, i.e.,  $\mathbf{X} = \{x_m\}_{m=1}^M$ ,  $\mathbf{Y} = \{y_n\}_{n=1}^N$ ,  $\mu_X \in \Delta^{N-1}$  and  $\mu_Y \in \Delta^{M-1}$ , we have:

$$\begin{aligned}\hat{d}_{gw}(\mathbf{X}, \mathbf{Y}) &= \min_{\mathbf{T} \in \Pi(\mu_X, \mu_Y)} \langle \mathbf{R}, \mathbf{T} \otimes \mathbf{T} \rangle \\ &= \min_{\mathbf{T} \in \Pi(\mu_X, \mu_Y)} \sum_{m,m'=1}^M \sum_{n,n'=1}^N r(x_m, x_{m'}, y_n, y_{n'}) T_{mn} T_{m'n'},\end{aligned}\tag{4}$$

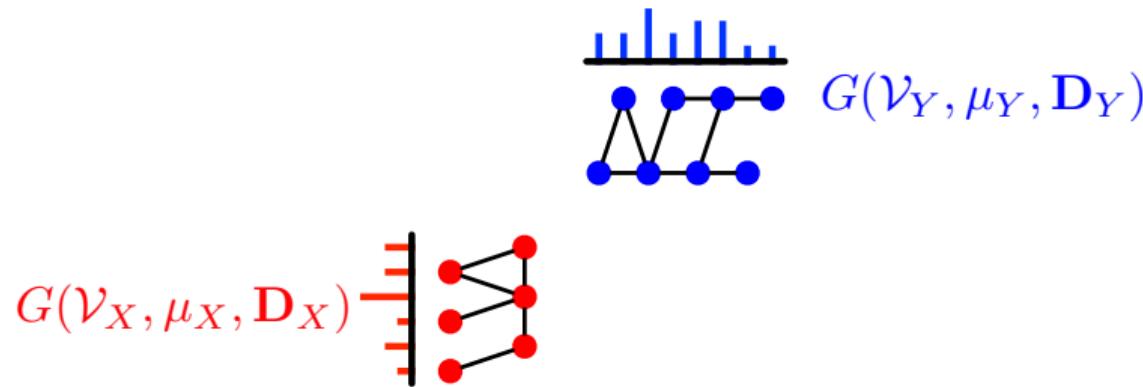
where  $\mathbf{R} = [r(x_m, x_{m'}, y_n, y_{n'})] \in \mathbb{R}^{M^2 \times N^2}$ .

# Gromov-Wasserstein Distance for Structured Data (Point Clouds)



- ▶ Commonly-used relational distance:  
 $r(x_m, x_{m'}, y_n, y_{n'}) = |d_X(x_m, x_{m'}) - d_Y(y_n, y_{n'})|^2$ , where  $d_X$  and  $d_Y$  are often Euclidean.
- ▶ Useful properties:
  - ▶ Translation-, rotation-, and permutation-invariance

# Gromov-Wasserstein Distance for Structured Data (Graphs)

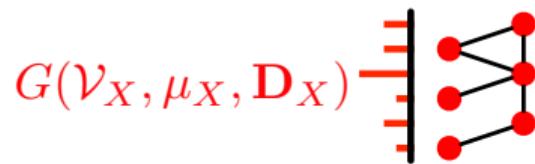
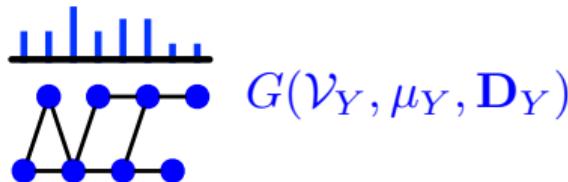


From a mm-space  $\mathcal{X}_{d,\mu}$  to a measure graph  $G_{\mathbf{D},\mu} := G(\mathcal{V}, \mu, \mathbf{D})$  [Chowdhury, et al., 2019]

- ▶  $\mathcal{V}$ : A node set
- ▶  $\mu$ : a predefined distribution of nodes
- ▶  $\mathbf{D} = [d_{ii'}]$ : a matrix recording the relations among nodes.

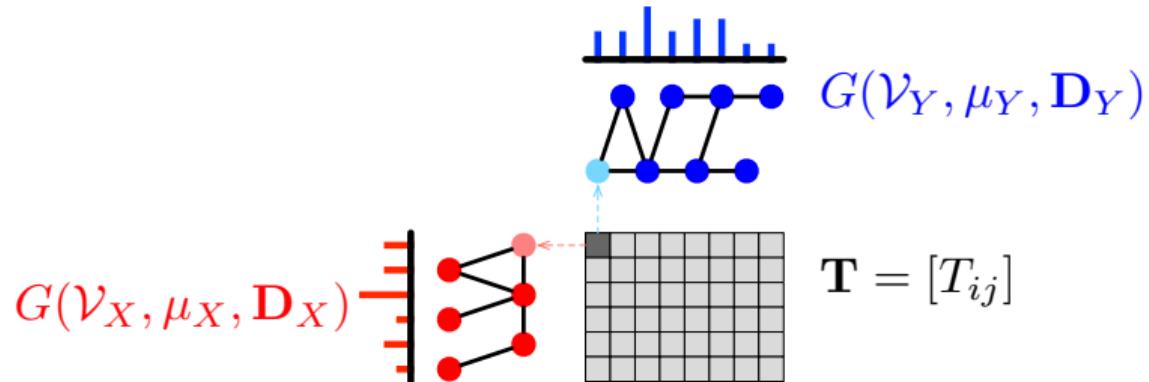
[Chowdhury, et al., 2019] Chowdhury, S., & Mémoli, F. The gromovwasserstein distance between networks and stable network invariants. Information and Inference: A Journal of the IMA, 2019.

# Gromov-Wasserstein Distance for Structured Data (Graphs)



- ▶ Options of  $\mu$ :
  - ▶ Uniform distribution:  $\mu = \frac{1}{|\mathcal{V}|} \mathbf{1}_{|\mathcal{V}|}$ .
  - ▶ Degree distribution [Xu, et al. NeurIPS 2019]:  $\mu = \frac{1}{\mathbf{1}^T \mathbf{p}} \mathbf{p}$ ,  $p_i = (\deg(v_i) + a)^b$ ,  $a, b \geq 0$ .
- ▶ Options of  $\mathbf{D} = [d_{ii'}]$ :
  - ▶ Adjacency matrix  $\mathbf{A} = [a_{ii'}]$  or Laplacian matrix  $\mathbf{L} = \text{diag}(\mathbf{A}\mathbf{1}) - \mathbf{A}$ .
  - ▶ Distance matrix  $\mathbf{D} = [d_{ii'}]$ , where  $d_{ii'}$  is the shortest path distance.
  - ▶ Kernel matrix  $\mathbf{K}$ , e.g., the  $t$ -th order heat kernel  $\mathbf{K}^t = \exp(-t\mathbf{L})$ .

# Gromov-Wasserstein Distance for Structured Data (Graphs)

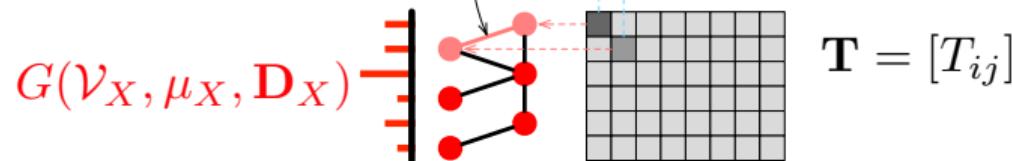


- ▶  $\mathbf{T} = [T_{ij}]$ : a joint distribution of nodes
- ▶  $\mathbf{T} \in \Pi(\mathbf{u}_X, \mathbf{u}_Y) = \{\mathbf{T} \geq \mathbf{0} \mid \mathbf{T}\mathbf{1} = \mathbf{u}_X, \mathbf{T}^\top\mathbf{1} = \mathbf{u}_Y\}$
- ▶ The pair of nodes ( $i \in \mathcal{V}_X, j \in \mathcal{V}_Y$ )  $\sim \mathbf{T}$ .

# Gromov-Wasserstein Distance for Structured Data (Graphs)

Relational Distance

$$r(i, i', j, j') = |d_{ii'}^X - d_{jj'}^Y|^2$$

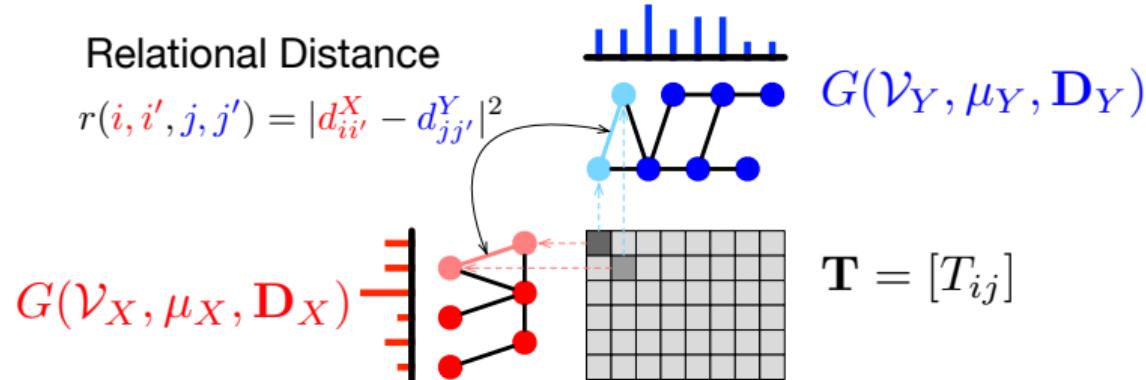


$$G(\mathcal{V}_Y, \mu_Y, \mathbf{D}_Y)$$

$$\mathbf{T} = [T_{ij}]$$

- ▶ The pair of relations  $(d_{ii'}^X, d_{jj'}^Y) \sim \mathbf{T} \times \mathbf{T}$ .
- ▶  $r(i, i', j, j')$ : the relational distance between node pairs.

# Gromov-Wasserstein Distance for Structured Data (Graphs)



- The pair of relations  $(d_{ii'}^X, d_{jj'}^Y) \sim \mathbf{T} \times \mathbf{T}$ .
- $r(i, i', j, j')$ : the relational distance between node pairs.

The GWD is **the minimum expectation of the relational distance**:

$$d_{gw}(G_X, G_Y) := \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \boldsymbol{\mu}_Y)} \mathbb{E}_{(i, i', j, j') \sim \mathbf{T} \times \mathbf{T}} [r(i, i', j, j')] = \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \boldsymbol{\mu}_Y)} \langle \mathbf{R}, \mathbf{T} \otimes \mathbf{T} \rangle. \quad (5)$$

## Simplification Under Specific Relational Distance

- ▶ Commonly used relational distance:
  - ▶ Euclidean distance (MSE):  $r(i, i', j, j') = |d_{ii'}^X - d_{jj'}^Y|^2$
  - ▶ KL-divergence:  $r(i, i', j, j') = d_{ii'}^X \log \frac{d_{ii'}^X}{d_{jj'}^Y} - d_{ii'}^X + d_{jj'}^Y.$
  - ▶ Cross-entropy loss:  $r(i, i', j, j') = -d_{ii'}^X \log d_{jj'}^Y - (1 - d_{ii'}^X) \log(1 - d_{jj'}^Y).$

# Simplification Under Specific Relational Distance

- ▶ Commonly used relational distance:
  - ▶ Euclidean distance (MSE):  $r(i, i', j, j') = |d_{ii'}^X - d_{jj'}^Y|^2$
  - ▶ KL-divergence:  $r(i, i', j, j') = d_{ii'}^X \log \frac{d_{ii'}^X}{d_{jj'}^Y} - d_{ii'}^X + d_{jj'}^Y$ .
  - ▶ Cross-entropy loss:  $r(i, i', j, j') = -d_{ii'}^X \log d_{jj'}^Y - (1 - d_{ii'}^X) \log(1 - d_{jj'}^Y)$ .
- ▶ A special form  $r(i, i', j, j') = f_1(d_{ii'}^X) + f_2(d_{jj'}^Y) - h_1(d_{ii'}^X)h_2(d_{jj'}^Y)$ :
  - ▶ Euclidean distance:  
 $f_1(d_{ii'}^X) = (d_{ii'}^X)^2$ ,  $f_2(d_{jj'}^Y) = (d_{jj'}^Y)^2$ ,  $h_1(d_{ii'}^X) = d_{ii'}^X$ , and  $h_2(d_{jj'}^Y) = 2d_{jj'}^Y$ .

# Simplification Under Specific Relational Distance

- ▶ Commonly used relational distance:
  - ▶ Euclidean distance (MSE):  $r(i, i', j, j') = |d_{ii'}^X - d_{jj'}^Y|^2$
  - ▶ KL-divergence:  $r(i, i', j, j') = d_{ii'}^X \log \frac{d_{ii'}^X}{d_{jj'}^Y} - d_{ii'}^X + d_{jj'}^Y$ .
  - ▶ Cross-entropy loss:  $r(i, i', j, j') = -d_{ii'}^X \log d_{jj'}^Y - (1 - d_{ii'}^X) \log(1 - d_{jj'}^Y)$ .
- ▶ A special form  $r(i, i', j, j') = f_1(d_{ii'}^X) + f_2(d_{jj'}^Y) - h_1(d_{ii'}^X)h_2(d_{jj'}^Y)$ :
  - ▶ Euclidean distance:  
 $f_1(d_{ii'}^X) = (d_{ii'}^X)^2$ ,  $f_2(d_{jj'}^Y) = (d_{jj'}^Y)^2$ ,  $h_1(d_{ii'}^X) = d_{ii'}^X$ , and  $h_2(d_{jj'}^Y) = 2d_{jj'}^Y$ .
  - ▶ KL-divergence:  
 $f_1(d_{ii'}^X) = d_{ii'}^X(\log d_{ii'}^X - 1)$ ,  $f_2(d_{jj'}^Y) = d_{jj'}^Y$ ,  $h_1(d_{ii'}^X) = d_{ii'}^X$ , and  $h_2(d_{jj'}^Y) = \log d_{jj'}^Y$ .

# Simplification Under Specific Relational Distance

- ▶ Commonly used relational distance:
  - ▶ Euclidean distance (MSE):  $r(i, i', j, j') = |d_{ii'}^X - d_{jj'}^Y|^2$
  - ▶ KL-divergence:  $r(i, i', j, j') = d_{ii'}^X \log \frac{d_{ii'}^X}{d_{jj'}^Y} - d_{ii'}^X + d_{jj'}^Y$ .
  - ▶ Cross-entropy loss:  $r(i, i', j, j') = -d_{ii'}^X \log d_{jj'}^Y - (1 - d_{ii'}^X) \log(1 - d_{jj'}^Y)$ .
- ▶ A special form  $r(i, i', j, j') = f_1(d_{ii'}^X) + f_2(d_{jj'}^Y) - h_1(d_{ii'}^X)h_2(d_{jj'}^Y)$ :
  - ▶ Euclidean distance:  
 $f_1(d_{ii'}^X) = (d_{ii'}^X)^2$ ,  $f_2(d_{jj'}^Y) = (d_{jj'}^Y)^2$ ,  $h_1(d_{ii'}^X) = d_{ii'}^X$ , and  $h_2(d_{jj'}^Y) = 2d_{jj'}^Y$ .
  - ▶ KL-divergence:  
 $f_1(d_{ii'}^X) = d_{ii'}^X(\log d_{ii'}^X - 1)$ ,  $f_2(d_{jj'}^Y) = d_{jj'}^Y$ ,  $h_1(d_{ii'}^X) = d_{ii'}^X$ , and  $h_2(d_{jj'}^Y) = \log d_{jj'}^Y$ .
  - ▶ Cross-entropy loss:  
 $f_1(d_{ii'}^X) = 0$ ,  $f_2(d_{jj'}^Y) = -\log(1 - d_{jj'}^Y)$ ,  $h_1(d_{ii'}^X) = d_{ii'}^X$ , and  $h_2(d_{jj'}^Y) = \log \frac{d_{jj'}^Y}{1-d_{jj'}^Y}$ .

## Simplification Under Specific Relational Distance

- ▶ Proposition 1 in [Peyré, et al., ICML 2016]

$$\langle \mathbf{R}, \mathbf{T} \otimes \mathbf{T} \rangle = \langle f_1(\mathbf{D}_X) \boldsymbol{\mu}_X \mathbf{1}_M^\top + \mathbf{1}_N \boldsymbol{\mu}_Y^\top f_2(\mathbf{D}_Y^\top) - h_1(\mathbf{D}_X) \mathbf{T} h_2(\mathbf{D}_Y^\top), \mathbf{T} \rangle. \quad (6)$$

## Simplification Under Specific Relational Distance

- ▶ Proposition 1 in [Peyré, et al., ICML 2016]

$$\langle \mathbf{R}, \mathbf{T} \otimes \mathbf{T} \rangle = \langle f_1(\mathbf{D}_X) \boldsymbol{\mu}_X \mathbf{1}_M^\top + \mathbf{1}_N \boldsymbol{\mu}_Y^\top f_2(\mathbf{D}_Y^\top) - h_1(\mathbf{D}_X) \mathbf{T} h_2(\mathbf{D}_Y^\top), \mathbf{T} \rangle. \quad (6)$$

- ▶ Furthermore, for  $\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \boldsymbol{\mu}_Y)$ , we have

$$\begin{aligned} & \langle f_1(\mathbf{D}_X) \boldsymbol{\mu}_X \mathbf{1}_M^\top + \mathbf{1}_N \boldsymbol{\mu}_Y^\top f_2(\mathbf{D}_Y^\top), \mathbf{T} \rangle \\ &= \text{tr}(f_1(\mathbf{D}_X) \boldsymbol{\mu}_X \underbrace{\mathbf{1}_M^\top \mathbf{T}^\top}_{\boldsymbol{\mu}_X^\top}) + \text{tr}(\underbrace{\mathbf{T}^\top \mathbf{1}_N}_{\boldsymbol{\mu}_Y^\top} \boldsymbol{\mu}_Y^\top f_2(\mathbf{D}_Y^\top)). \end{aligned} \quad (7)$$

## Simplification Under Specific Relational Distance

- ▶ Proposition 1 in [Peyré, et al., ICML 2016]

$$\langle \mathbf{R}, \mathbf{T} \otimes \mathbf{T} \rangle = \langle f_1(\mathbf{D}_X) \boldsymbol{\mu}_X \mathbf{1}_M^\top + \mathbf{1}_N \boldsymbol{\mu}_Y^\top f_2(\mathbf{D}_Y^\top) - h_1(\mathbf{D}_X) \mathbf{T} h_2(\mathbf{D}_Y^\top), \mathbf{T} \rangle. \quad (6)$$

- ▶ Furthermore, for  $\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \boldsymbol{\mu}_Y)$ , we have

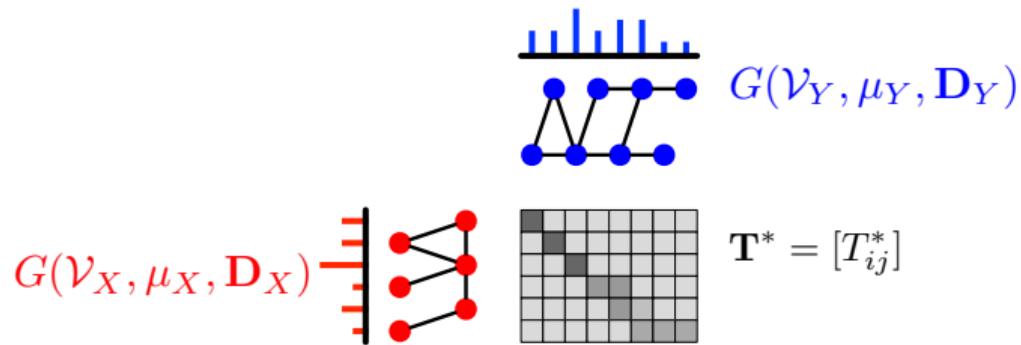
$$\begin{aligned} & \langle f_1(\mathbf{D}_X) \boldsymbol{\mu}_X \mathbf{1}_M^\top + \mathbf{1}_N \boldsymbol{\mu}_Y^\top f_2(\mathbf{D}_Y^\top), \mathbf{T} \rangle \\ &= \text{tr}(f_1(\mathbf{D}_X) \boldsymbol{\mu}_X \underbrace{\mathbf{1}_M^\top \mathbf{T}^\top}_{\boldsymbol{\mu}_X^\top}) + \text{tr}(\underbrace{\mathbf{T}^\top \mathbf{1}_N}_{\boldsymbol{\mu}_Y^\top} \boldsymbol{\mu}_Y^\top f_2(\mathbf{D}_Y^\top)). \end{aligned} \quad (7)$$

- ▶ Accordingly, we can simplify the computation of GWD from  $\mathcal{O}(M^2 N^2)$  to  $\mathcal{O}(M^2 N + N^2 M)$  (or  $\mathcal{O}(N^3)$  when  $M \approx N$ ):

$$\begin{aligned} \mathbf{T}^* &= \arg \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \boldsymbol{\mu}_Y)} -\langle h_1(\mathbf{D}_X) \mathbf{T} h_2(\mathbf{D}_Y^\top), \mathbf{T} \rangle. \\ \xrightarrow{\text{Euclidean } r} \mathbf{T}^* &= \arg \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \boldsymbol{\mu}_Y)} -\langle \mathbf{D}_X \mathbf{T} \mathbf{D}_Y^\top, \mathbf{T} \rangle. \end{aligned} \quad (8)$$

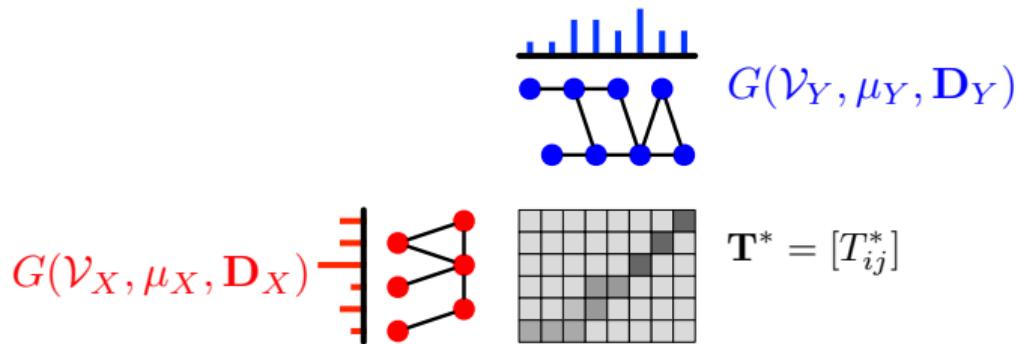
[Peyré, et al., ICML 2016] Peyré, G., Cuturi, M., & Solomon, J. Gromov-Wasserstein averaging of kernel and distance matrices. ICML, 2016.

## Advantages of GWD



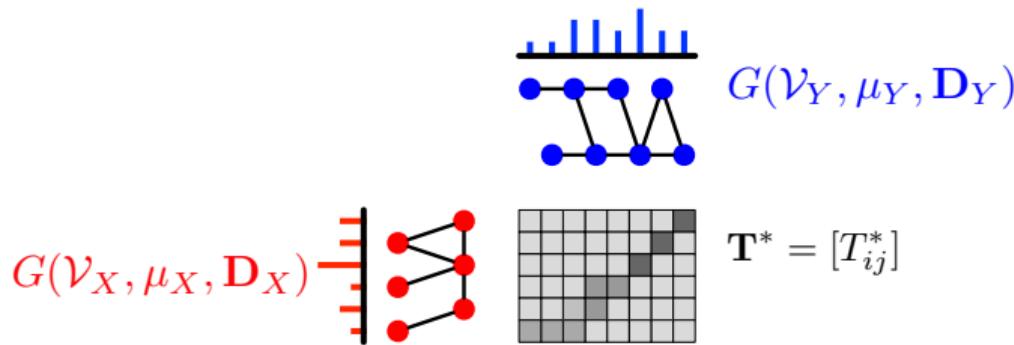
- ▶ The optimal joint distribution  $\mathbf{T}^*$  (or called “**optimal transport**” matrix) indicates the correspondence between the two graphs/point clouds.

## Advantages of GWD



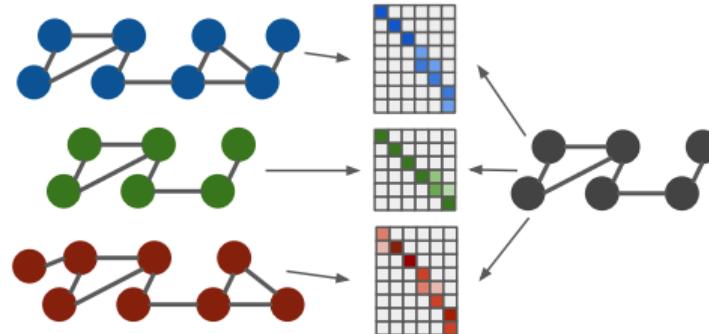
- ▶ The optimal joint distribution  $\mathbf{T}^*$  (or called “**optimal transport**” matrix) indicates the correspondence between the two graphs.
- ▶ A (pseudo) metric with symmetry and following invariance properties
  - ▶ **Translation+Rotation-invariance for point clouds:**  
 $d_{gw}(\mathbf{X}, \mathbf{Y}) = d_{gw}(\mathbf{X}, \text{Rot+Trans}(\mathbf{Y}))$
  - ▶ **Permutation-invariance for graphs:**  
 $d_{gw}(G_X, G_Y) = d_{gw}(G_X, \text{Permute}(G_Y))$

## Advantages of GWD



- ▶ The optimal transport  $\mathbf{T}^*$  is a joint distribution indicating the correspondence between the two graphs.
- ▶ A (pseudo) metric with symmetry and following invariance properties
  - ▶ **Translation+Rotation-invariance for point clouds:**  
 $d_{gw}(\mathbf{X}, \mathbf{Y}) = d_{gw}(\mathbf{X}, \text{Rot+Trans}(\mathbf{Y}))$
  - ▶ **Permutation-invariance for graphs:**  
 $d_{gw}(G_X, G_Y) = d_{gw}(G_X, \text{Permute}(G_Y))$
- ▶ Applicable to the graphs with different sizes, i.e.,  $|\mathcal{V}_X| \neq |\mathcal{V}_Y|$ .

## From GW Distance to a (Weighted) GW Barycenter

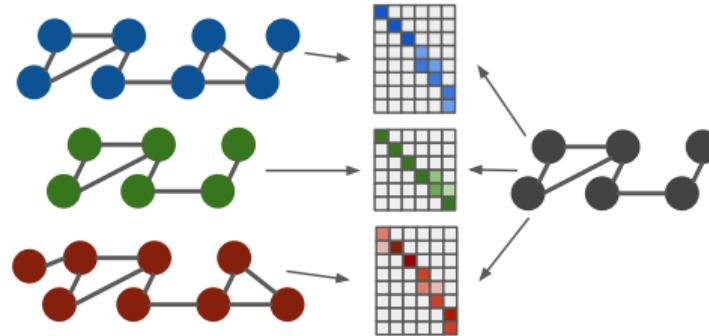


Given  $\{G_k\}_{k=1}^K$ ,  $K \geq 2$ , their (weighted) GW barycenter [Peyré, et al., ICML 2016] is

$$\underbrace{\bar{G}(\bar{\mathcal{V}}, \bar{\mu}, \mathbf{B}^*)}_{\text{Barycenter graph}}, \quad \underbrace{\{\mathbf{T}_k^*\}_{k=1}^K}_{\text{OT matrices}} := \arg \min_G \sum_{k=1}^K \lambda_k d_{gw}(G_k, G) \quad (9)$$

$$= \arg \min_{\mathbf{B}} \min_{\{\mathbf{T}_k \in \Pi(\mu_k, \bar{\mu})\}_{k=1}^K} - \sum_{k=1}^K \lambda_k \langle \mathbf{D}_k \mathbf{T}_k \mathbf{B}^\top, \mathbf{T}_k \rangle$$

## From GW Distance to a (Weighted) GW Barycenter



Given  $\{G_k\}_{k=1}^K$ ,  $K \geq 2$ , their (weighted) GW barycenter [Peyré, et al., ICML 2016] is

$$\underbrace{\bar{G}(\bar{\mathcal{V}}, \bar{\mu}, \bar{\mathbf{B}}^*)}_{\text{Barycenter graph}}, \quad \underbrace{\{\mathbf{T}_k^*\}_{k=1}^K}_{\text{OT matrices}} := \arg \min_G \sum_{k=1}^K \lambda_k d_{gw}(G_k, G) \quad (9)$$

$$= \arg \min_{\mathbf{B}} \min_{\{\mathbf{T}_k \in \Pi(\mu_k, \bar{\mu})\}_{k=1}^K} - \sum_{k=1}^K \lambda_k \langle \mathbf{D}_k \mathbf{T}_k \mathbf{B}^\top, \mathbf{T}_k \rangle$$

**Permutation-invariance:** If  $\bar{G}$  is a GW barycenter of  $\{G_k\}_{k=1}^K$ , then  $\text{permute}(\bar{G})$  is a valid GW barycenter as well.

## Implementation of GW Barycenter

- ▶ Approximation of  $\bar{\mu}$  [Xu, et al., NeurIPS 2019]:

$$\bar{\mu} = \sum_{k=1}^K \lambda_k \text{interpolate}_{|\bar{\mathcal{V}}|}(\text{sort}(\mu_k)), \quad (10)$$

- ▶  $\text{sort}(\cdot)$  sorts the elements of the input vector in descending order.
- ▶  $\text{interpolate}_{|\bar{\mathcal{V}}|}(\cdot)$  samples  $|\bar{\mathcal{V}}|$  values from the input vector via a certain interpolation method (e.g., bilinear or cubic interpolation).

# Implementation of GW Barycenter

- ▶ Approximation of  $\bar{\mu}$  [Xu, et al., NeurIPS 2019]:

$$\bar{\mu} = \sum_{k=1}^K \lambda_k \text{interpolate}_{|\bar{\mathcal{V}}|}(\text{sort}(\mu_k)), \quad (10)$$

- ▶  $\text{sort}(\cdot)$  sorts the elements of the input vector in descending order.
- ▶  $\text{interpolate}_{|\bar{\mathcal{V}}|}(\cdot)$  samples  $|\bar{\mathcal{V}}|$  values from the input vector via a certain interpolation method (e.g., bilinear or cubic interpolation).
- ▶ Alternating optimization strategy:
  - ▶ Obtain  $T_k = \arg \min_{T \in \Pi(\mu_k, \bar{\mu})} -\langle D_k T B^\top, T \rangle$  for  $k = 1, \dots, K$ .
  - ▶ Update barycenter in a closed form (the first-order optimality condition):

$$B^* = \frac{1}{\bar{\mu} \bar{\mu}^\top} \sum_{k=1}^K \lambda_k T_k^\top D_k T_k. \quad (11)$$

[Xu, et al. NeurIPS, 2019] Xu, H., Luo, D., & Carin, L. Scalable Gromov-Wasserstein learning for graph partitioning and matching. NeurIPS, 2019.

## Straightforward Applications

### **Graph matching and partitioning**

**Apply the Gromov-Wasserstein distance/barycenter as objective functions**

# Graph Matching via Calculating GW Distance

Quadratic assignment problem (QAP):

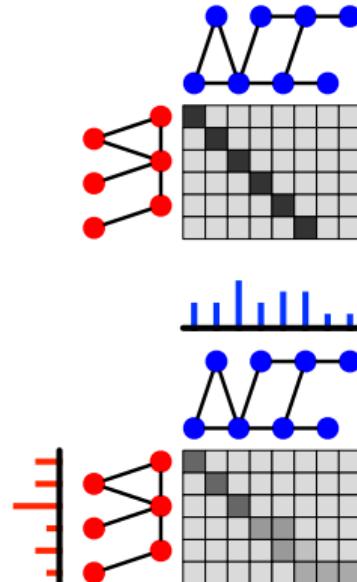
$$\max_{\mathbf{P} \in \mathcal{P}} \langle \mathbf{D}_X \mathbf{P} \mathbf{D}_Y^\top, \mathbf{P} \rangle,$$

$$\mathcal{P} = \{\mathbf{P} \in \{0, 1\}^{|\mathcal{V}_X| \times |\mathcal{V}_Y|} \mid \mathbf{P}\mathbf{1} = \mathbf{1}, \mathbf{P}^\top \mathbf{1} \leq \mathbf{1}\}.$$

Gromov-Wasserstein distance (with Euclidean relational distance):

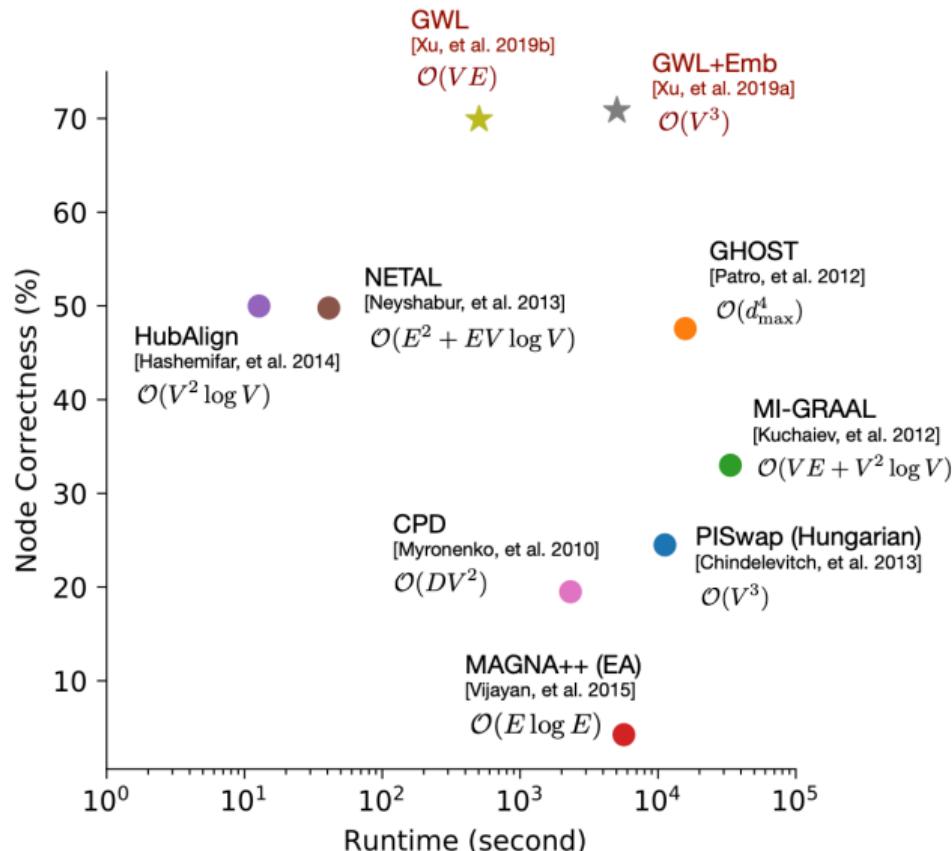
$$\min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \boldsymbol{\mu}_Y)} -\langle \mathbf{D}_X \mathbf{T} \mathbf{D}_Y^\top, \mathbf{T} \rangle,$$

$$\Pi(\boldsymbol{\mu}_X, \boldsymbol{\mu}_Y) = \{\mathbf{T} \geq \mathbf{0} \mid \mathbf{T}\mathbf{1} = \boldsymbol{\mu}_X, \mathbf{T}^\top \mathbf{1} = \boldsymbol{\mu}_Y\}$$



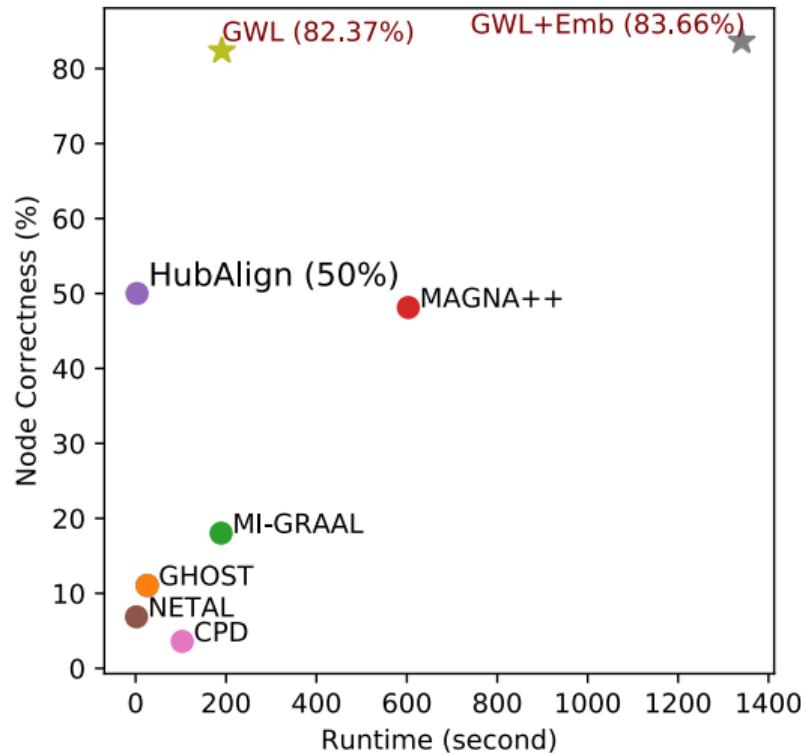
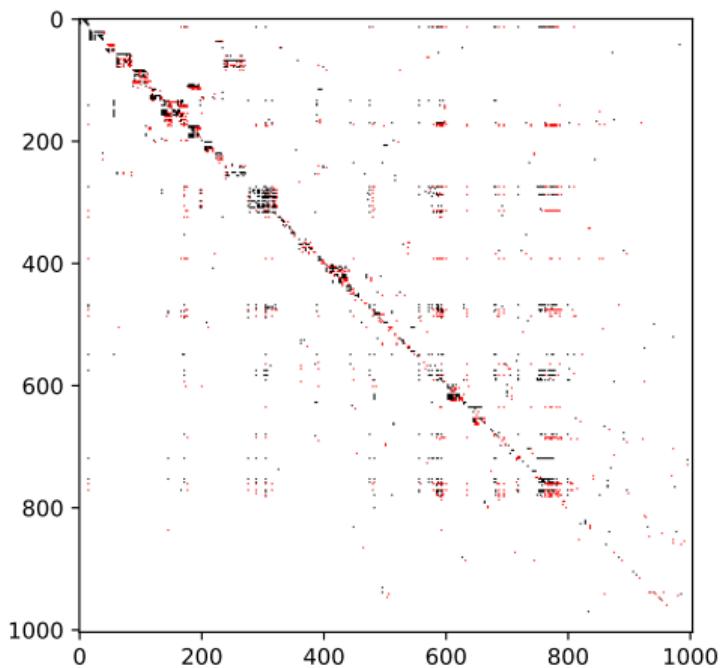
- ▶ Conditional distribution:  $\mathbf{P}^* = \frac{\mathbf{T}^*}{\boldsymbol{\mu}_X \mathbf{1}^\top}$ .
- ▶ For each node  $i \in \mathcal{V}_X$ ,  $j^* = \arg \max_j P_{ij}^*$ .

# Matching Synthetic Graphs



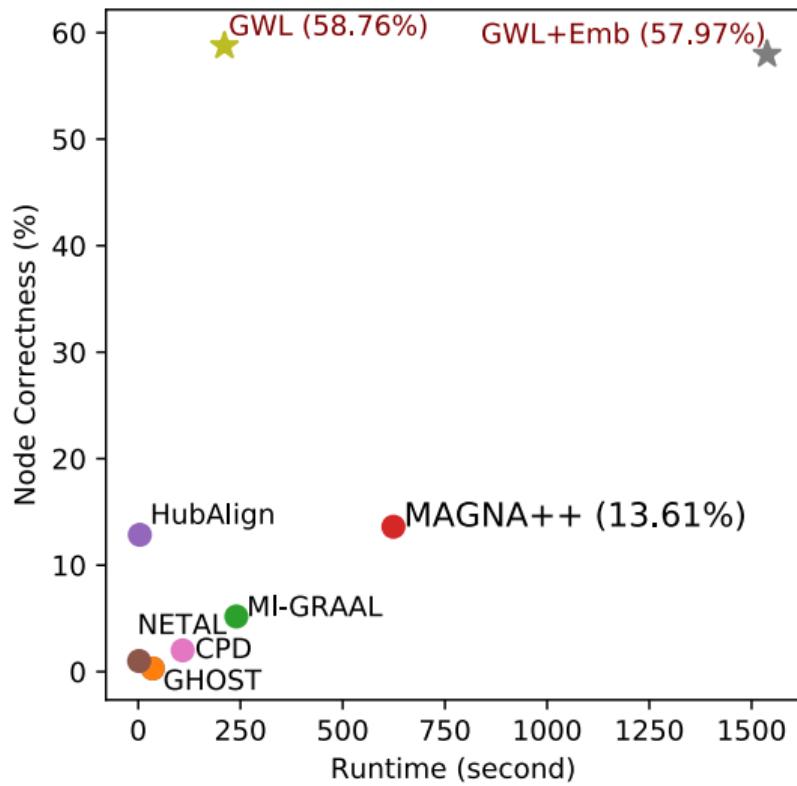
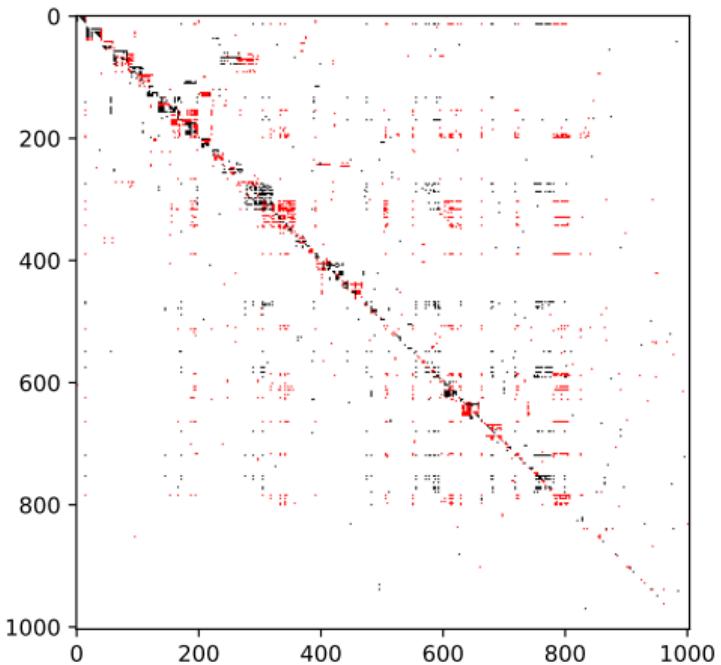
# Real-world PPI Network Alignment

Yeast PPI  $\leftrightarrow$  Yeast PPI + 5% LC edges



# Real-world PPI Network Alignment

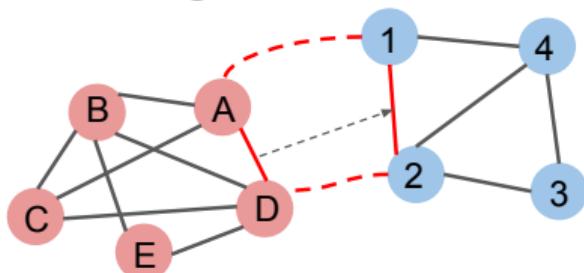
Yeast PPI  $\leftrightarrow$  Yeast PPI + 25% LC edges



# Real-world PPI Network Alignment

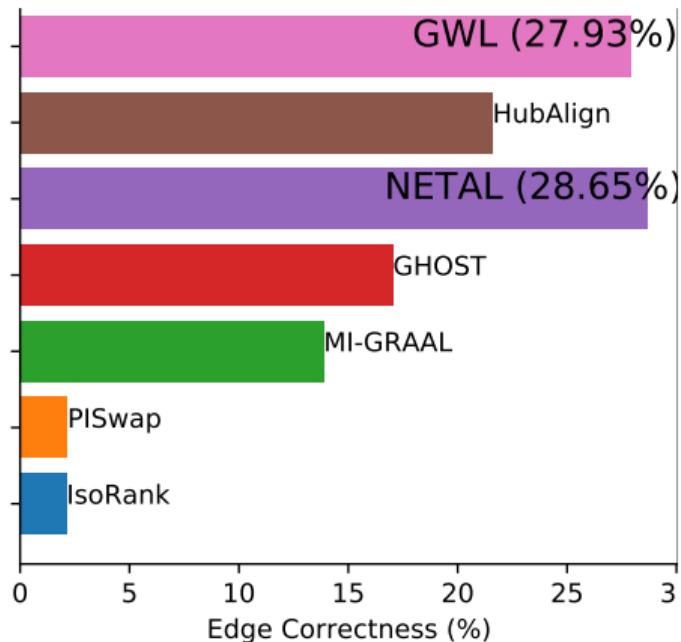
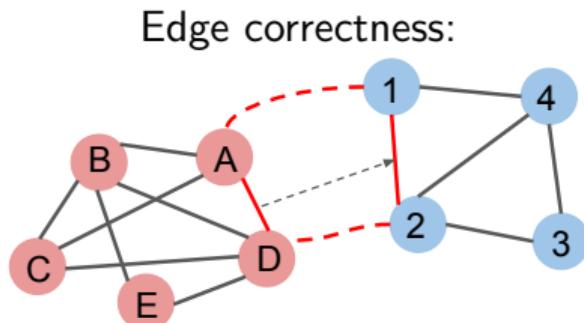
Matching the PPI network of yeast (2,340 proteins) to that of human (9,141 proteins)

Edge correctness:



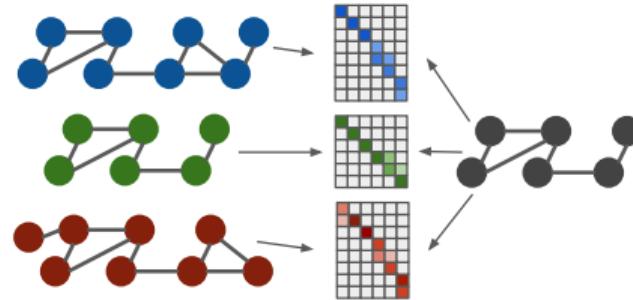
# Real-world PPI Network Alignment

Matching the PPI network of yeast (2,340 proteins) to that of human (9,141 proteins)



# Multi-Graph Matching via Calculating GW Barycenters

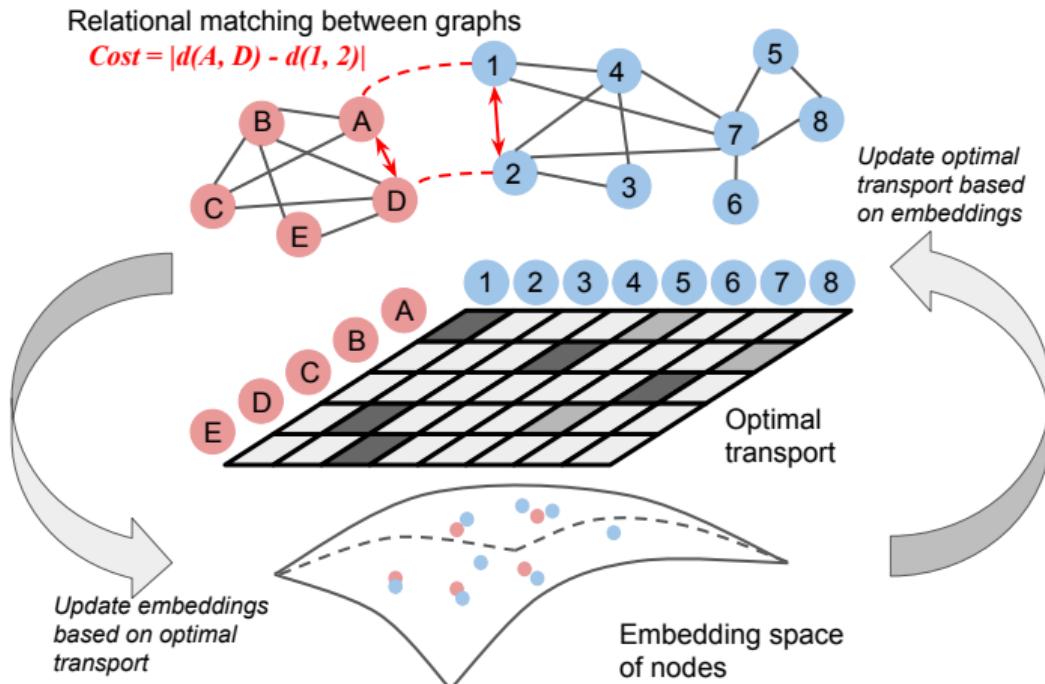
$$\bar{G}, \{\mathbf{T}_k^*\}_{k=1}^K = \arg \min \sum_{k=1}^K \lambda_k d_{gw}(G_k, \bar{G}),$$



For each node  $v$  in the barycenter graph  $\bar{G}$ , its correspondences in each graph, i.e.,  $\{v_k\}_{k=1}^K$ , where  $v_k = \arg \max_{v'} T_k^*(v', v)$  are matched with each other [Xu, et al. NeurIPS, 2019].

Method	3 graphs		4 graphs		5 graphs		6 graphs	
	NC@1	NC@all	NC@1	NC@all	NC@1	NC@all	NC@1	NC@all
MultiAlign	62.97	45.19	—	—	—	—	—	—
GWL	<b>63.84</b>	<b>46.22</b>	<b>68.73</b>	<b>39.14</b>	71.61	31.57	76.49	28.39

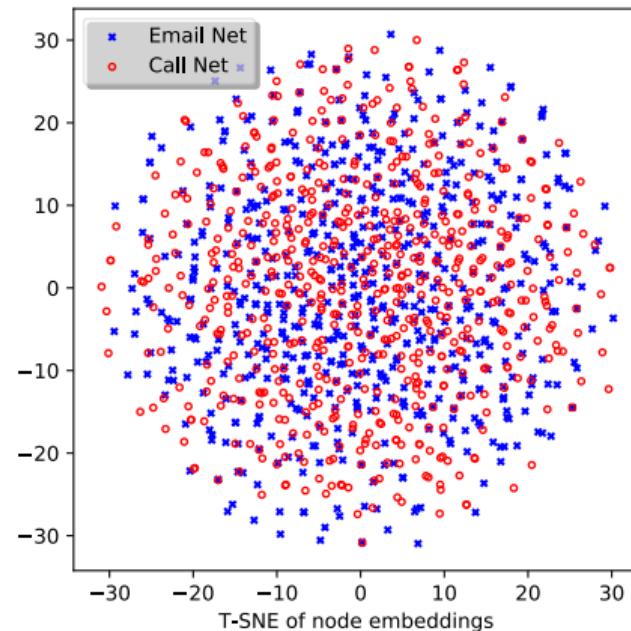
# Joint Graph Matching and Node Embedding



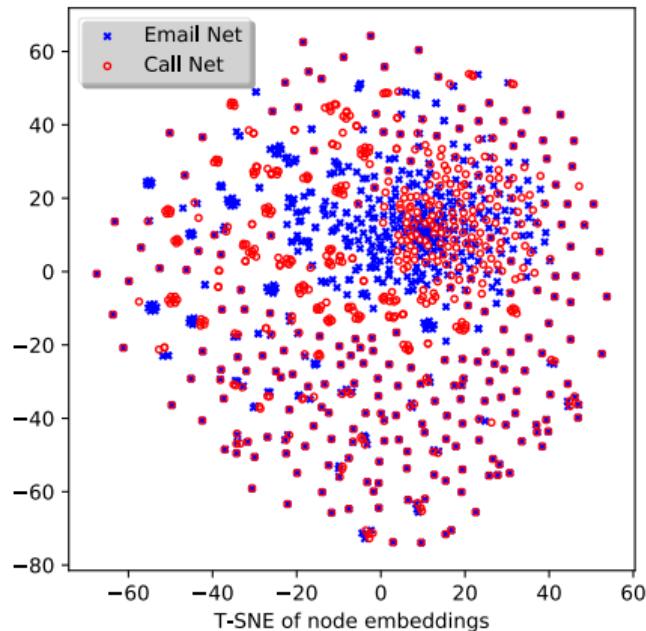
[Xu, et al. ICML 2019] Xu, H., Luo, D., Zha, H., & Carin, L. Gromov-wasserstein learning for graph matching and node embedding. ICML, 2019.

# Joint Graph Matching and Node Embedding

Binding accounts across communication networks



(a) Epoch 0

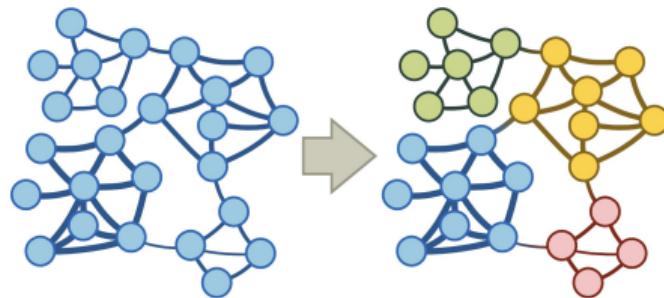


(b) Epoch 30

# Graph Partitioning via Calculating GW Distance

## Modularity maximization principle

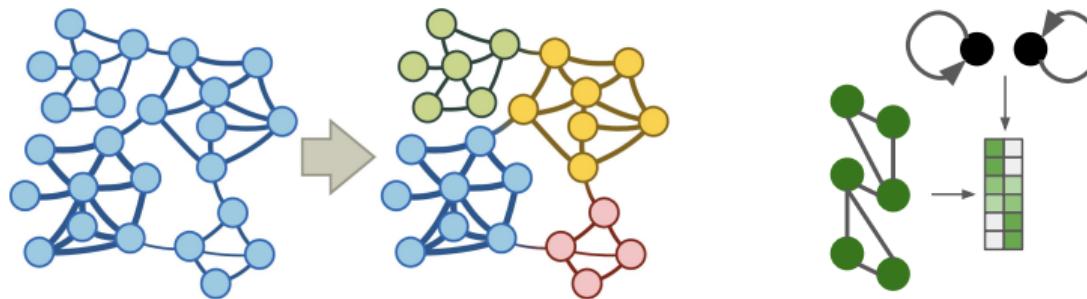
- ▶ Dense internal edges + sparse external edges.



# Graph Partitioning via Calculating GW Distance

## Modularity maximization principle

- Dense internal edges + sparse external edges.



A GWD-based solution [Xu, et al. NeurIPS, 2019]:

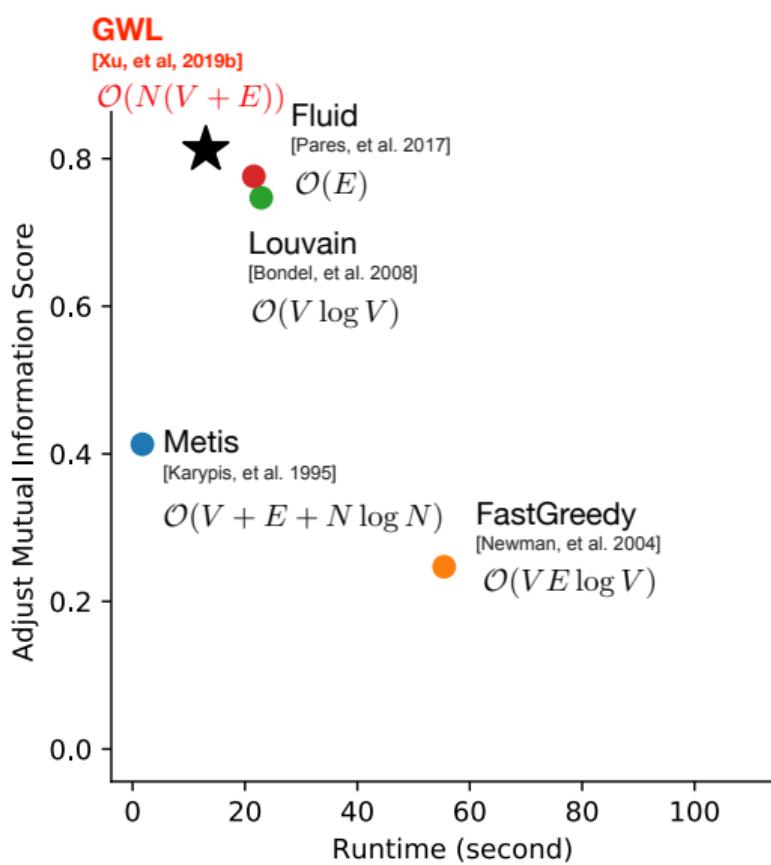
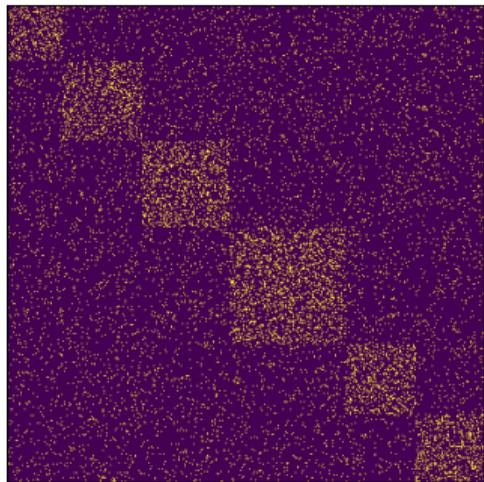
- $T^* \in \mathbb{R}^{|\mathcal{V}| \times N} \leftarrow d_{gw}(G, G_{iso})$
- $G_{iso}(\mathcal{V}_{iso}, \mu \in \Delta^{N-1}, \text{diag}(\mu))$
- For each node  $i \in G$ , its cluster is  $j^* = \arg \max_j T_{ij}^*$ .

# Partitioning Synthetic Graphs

$V = 4,000$

$p_{\text{within}} = 0.2$

$p_{\text{across}} = 0.05$

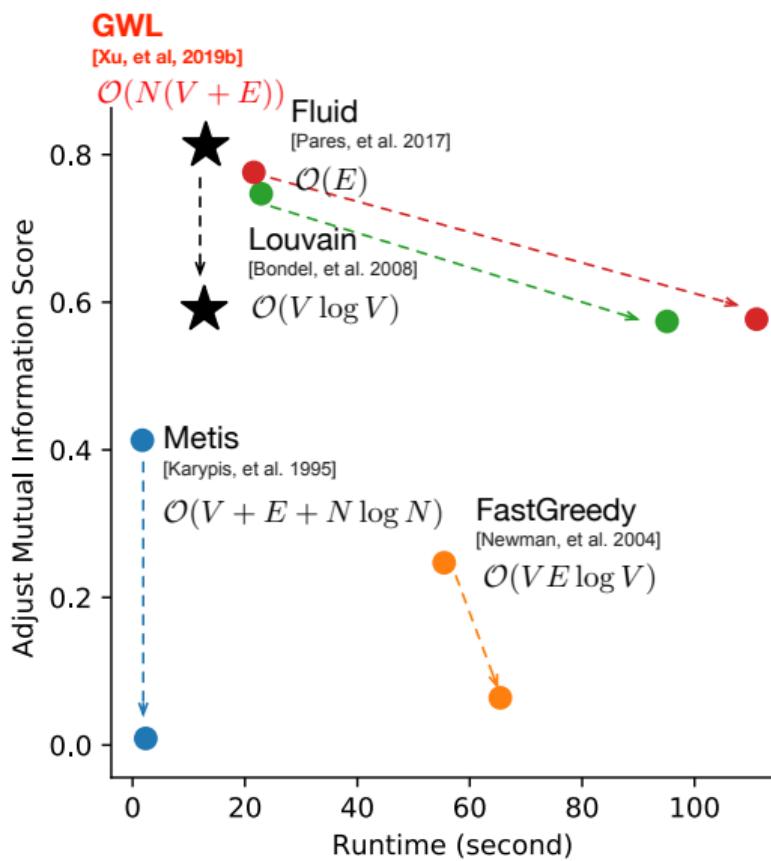
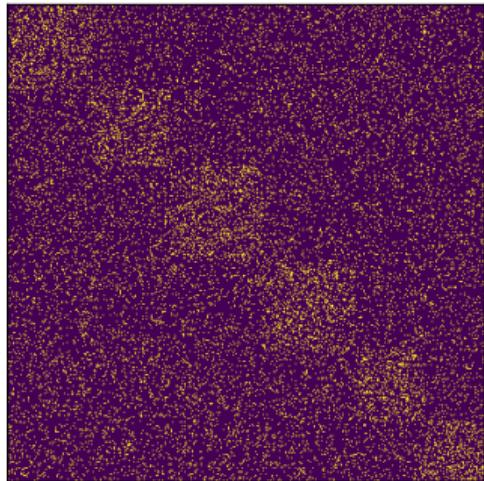


# Partitioning Synthetic Graphs

$V = 4,000$

$p_{\text{within}} = 0.2$

$p_{\text{across}} = 0.1$



## Graph Partitioning via Calculating (Spectral) GW Distance

- ▶ The work in [Chowdhury, et al. AISTATS, 2021] applies  $t$ -th order heat kernels, rather than adjacency matrices, when implementing  $d_{gw}(G, G_{iso})$ , which improves both partitioning accuracy and efficiency.

Table 2: Node Correctness, Mean  $\pm$  St. Dev. (Time).

Dataset	Adj	Spec <sup>10</sup>
Proteins	.68 $\pm$ .22 (31.9)	<b>.78</b> $\pm$ .22 (5.1)
Enzymes	.70 $\pm$ .18 (8.9)	<b>.79</b> $\pm$ .17 (1.4)
Reddit	.29 $\pm$ .21 (3941.7)	<b>.50</b> $\pm$ .11 (206.1)
Collab	<b>.50</b> $\pm$ .27 (4.3)	<b>.50</b> $\pm$ .27 (5.6)

[Chowdhury, et al. AISTATS, 2021] Chowdhury, S., & Needham, T. Generalized spectral clustering via Gromov-Wasserstein learning. AISTATS, 2021.

## Graph Partitioning via Calculating (Spectral) GW Distance

- ▶ Assumptions of a graph  $G$ :
  - ▶ The multiplicity of the smallest positive eigenvalue of the graph Laplacian  $\mathbf{L}$  is 1.
  - ▶ Uniform node probability distribution  $\boldsymbol{\mu} = \frac{1}{|\mathcal{V}|} \mathbf{1}_{|\mathcal{V}|}$ .

## Graph Partitioning via Calculating (Spectral) GW Distance

- ▶ Assumptions of a graph  $G$ :
  - ▶ The multiplicity of the smallest positive eigenvalue of the graph Laplacian  $\mathbf{L}$  is 1.
  - ▶ Uniform node probability distribution  $\boldsymbol{\mu} = \frac{1}{|\mathcal{V}|} \mathbf{1}_{|\mathcal{V}|}$ .
- ▶ **Fiedler vector of  $G$ :** The eigenvector corresponding to the smallest positive eigenvalue of graph Laplacian  $\mathbf{L}$ .

## Graph Partitioning via Calculating (Spectral) GW Distance

- ▶ Assumptions of a graph  $G$ :
  - ▶ The multiplicity of the smallest positive eigenvalue of the graph Laplacian  $\mathbf{L}$  is 1.
  - ▶ Uniform node probability distribution  $\boldsymbol{\mu} = \frac{1}{|\mathcal{V}|}\mathbf{1}_{|\mathcal{V}|}$ .
- ▶ **Fiedler vector of  $G$ :** The eigenvector corresponding to the smallest positive eigenvalue of graph Laplacian  $\mathbf{L}$ .
- ▶ **Fiedler partition of  $G$ :** The nodes of  $G$  are partitioned according to the sign of their entry in the Fiedler vector.

# Graph Partitioning via Calculating (Spectral) GW Distance

- ▶ Assumptions of a graph  $G$ :
  - ▶ The multiplicity of the smallest positive eigenvalue of the graph Laplacian  $\mathbf{L}$  is 1.
  - ▶ Uniform node probability distribution  $\boldsymbol{\mu} = \frac{1}{|\mathcal{V}|}\mathbf{1}_{|\mathcal{V}|}$ .
- ▶ **Fiedler vector of  $G$ :** The eigenvector corresponding to the smallest positive eigenvalue of graph Laplacian  $\mathbf{L}$ .
- ▶ **Fiedler partition of  $G$ :** The nodes of  $G$  are partitioned according to the sign of their entry in the Fiedler vector.
- ▶ Connect optimal transport with classic spectral clustering:
  - ▶ **Theorem 3** in [Chowdhury, et al. AISTATS, 2021]: *The 2-way partition of such  $G$  derived from the spectral version of  $d_{gw}(G, G_{iso})$  agrees with the **Fiedler partitioning** when  $t$  is sufficiently large.*

## Summary

- ▶ The GW distance works as a metric of mm-spaces, which can be used as a (pseudo) metric of point clouds / graphs.
- ▶ The GW barycenter provides a way to obtain the average of point clouds / graphs under GW distance.
- ▶ The optimal transport associated with the GW distance achieves a soft assignment of graph nodes.
- ▶ Based on GW distance, we actually achieves a unified algorithmic framework for graph matching, partitioning, and node embedding.
- ▶ The GW-based framework has the potential to many real-world problems, e.g., PPI network alignment, community detection.

## Summary

- ▶ The GW distance works as a metric of mm-spaces, which can be used as a (pseudo) metric of point clouds / graphs.
- ▶ The GW barycenter provides a way to obtain the average of point clouds / graphs under GW distance.
- ▶ The optimal transport associated with the GW distance achieves a soft assignment of graph nodes.
- ▶ Based on GW distance, we actually achieves a unified algorithmic framework for graph matching, partitioning, and node embedding.
- ▶ The GW-based framework has the potential to many real-world problems, e.g., PPI network alignment, community detection.

Next,

- ▶ **How to compute/approximate GW distance efficiently?**
- ▶ **How to extend its applications based on its variants?**

5-min break for Q & A

# Outline

## Part 1 Introduction of Gromov-Wasserstein Distance

- ▶ Preliminary and basic concepts
- ▶ Connections to structured data modeling and analysis

## Part 2 The Computation of Gromov-Wasserstein Distance

- ▶ Typical optimization algorithms
- ▶ The variants of Gromov-Wasserstein distance

## Part 3 GW-based Models and Their ML Applications

- ▶ Generative modeling
- ▶ Graph representation
- ▶ Graph generation

## The Optimization Problem Corresponding to GW Distance

As aforementioned,  $d_{gw}(G_X, G_Y)$  leads to a **non-convex non-smooth** optimization problem:

$$\min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \boldsymbol{\mu}_Y)} -\langle \mathbf{D}_X \mathbf{T} \mathbf{D}_Y^\top, \mathbf{T} \rangle \quad (12)$$

- ▶ Iterative optimization is required.
- ▶ Low complexity, guaranteed convergence, stable approximation.

## Typical Method: Entropic Regularization + Projected Gradient Descent

**Sinkhorn Distance (Entropic OT Problem)** [Cuturi, NeurIPS 2013]

$$\min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}, \boldsymbol{\eta})} \langle \mathbf{D}, \mathbf{T} \rangle - \gamma \underbrace{\mathsf{H}(\mathbf{T})}_{\text{Entropy}},$$
$$\mathsf{H}(\mathbf{T}) = -\langle \log \mathbf{T} - \mathbf{1}, \mathbf{T} \rangle.$$

## Typical Method: Entropic Regularization + Projected Gradient Descent

**Sinkhorn Distance (Entropic OT Problem)** [Cuturi, NeurIPS 2013]

$$\min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}, \boldsymbol{\eta})} \langle \mathbf{D}, \mathbf{T} \rangle - \gamma \underbrace{\mathsf{H}(\mathbf{T})}_{\text{Entropy}}, \quad (13)$$
$$\mathsf{H}(\mathbf{T}) = -\langle \log \mathbf{T} - \mathbf{1}, \mathbf{T} \rangle.$$

**Entropic Gromov-Wasserstein Distance** [Peyré, et al., ICML 2016]

$$\min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \boldsymbol{\mu}_Y)} \underbrace{-\langle \mathbf{D}_X \mathbf{T} \mathbf{D}_Y^\top, \mathbf{T} \rangle}_{\text{GW term}} - \gamma \underbrace{\mathsf{H}(\mathbf{T})}_{\text{Entropy}}. \quad (14)$$

[Cuturi, NeurIPS 2013] Cuturi, M. Sinkhorn distances: Lightspeed computation of optimal transport. NeurIPS, 2013.

## Typical Method: Entropic Regularization + Projected Gradient Descent

In each iteration

- ▶ Fix one optimal transport and obtain an entropic OT problem:

$$\mathbf{T}^{(m+1)} = \arg \min_{\mathbf{T} \in \Pi(\mu_X, \mu_Y)} -\underbrace{\langle \mathbf{D}_X \mathbf{T}^{(m)} \mathbf{D}_Y^\top, \mathbf{T} \rangle}_{\text{Constant}} - \gamma \mathsf{H}(\mathbf{T}) \quad (15)$$

## Typical Method: Entropic Regularization + Projected Gradient Descent

In each iteration

- ▶ Fix one optimal transport and obtain an entropic OT problem:

$$\mathbf{T}^{(m+1)} = \arg \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \boldsymbol{\mu}_Y)} -\underbrace{\langle \mathbf{D}_X \mathbf{T}^{(m)} \mathbf{D}_Y^\top, \mathbf{T} \rangle}_{\text{Constant}} - \gamma \mathsf{H}(\mathbf{T}) \quad (15)$$

- ▶ Sinkhorn-Knopp algorithm: Define  $\Phi = \exp(-\frac{\mathbf{D}_X \mathbf{T}^{(m)} \mathbf{D}_Y^\top}{\gamma})$ ,  $\mathbf{a} = \boldsymbol{\mu}_X$ , and
  1. Repeat  $\mathbf{b} = \frac{\boldsymbol{\mu}_Y}{\Phi^\top \mathbf{a}}$  and  $\mathbf{a} = \frac{\boldsymbol{\mu}_X}{\Phi \mathbf{b}}$  until convergence
  2.  $\mathbf{T}^{(m+1)} = \text{diag}(\mathbf{a}) \Phi \text{diag}(\mathbf{b})$ .

## Typical Method: Proximal Gradient

### Motivation:

- ▶ The entropic regularizer might be unnecessary for the GW distance.

# Typical Method: Proximal Gradient

## Motivation:

- ▶ The entropic regularizer might be unnecessary for the GW distance.

## Principle:

- ▶ In each iteration, consider the penalty between the optimal transport and its previous approximation [Xu, et al., ICML 2019]

$$\min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \boldsymbol{\mu}_Y)} -\langle \mathbf{D}_X \mathbf{T}^{(m)} \mathbf{D}_Y^\top, \mathbf{T} \rangle + \gamma \underbrace{\text{KL}(\mathbf{T} \parallel \mathbf{T}^{(m)})}_{\text{Proximal term}}$$

# Typical Method: Proximal Gradient

## Motivation:

- ▶ The entropic regularizer might be unnecessary for the GW distance.

## Principle:

- ▶ In each iteration, consider the penalty between the optimal transport and its previous approximation [Xu, et al., ICML 2019]

$$\begin{aligned} & \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \boldsymbol{\mu}_Y)} -\langle \mathbf{D}_X \mathbf{T}^{(m)} \mathbf{D}_Y^\top, \mathbf{T} \rangle + \gamma \underbrace{\text{KL}(\mathbf{T} \| \mathbf{T}^{(m)})}_{\text{Proximal term}} \\ \Rightarrow & \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \boldsymbol{\mu}_Y)} -\underbrace{\langle \mathbf{D}_X \mathbf{T}^{(m)} \mathbf{D}_Y^\top + \gamma \log \mathbf{T}^{(m)}, \mathbf{T} \rangle}_{\text{Constant}} - \gamma \mathcal{H}(\mathbf{T}). \end{aligned} \tag{16}$$

# Typical Method: Proximal Gradient

## Motivation:

- ▶ The entropic regularizer might be unnecessary for the GW distance.

## Principle:

- ▶ In each iteration, consider the penalty between the optimal transport and its previous approximation [Xu, et al., ICML 2019]

$$\begin{aligned} & \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \boldsymbol{\mu}_Y)} -\langle \mathbf{D}_X \mathbf{T}^{(m)} \mathbf{D}_Y^\top, \mathbf{T} \rangle + \gamma \underbrace{\text{KL}(\mathbf{T} \| \mathbf{T}^{(m)})}_{\text{Proximal term}} \\ & \Rightarrow \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \boldsymbol{\mu}_Y)} -\underbrace{\langle \mathbf{D}_X \mathbf{T}^{(m)} \mathbf{D}_Y^\top + \gamma \log \mathbf{T}^{(m)}, \mathbf{T} \rangle}_{\text{Constant}} - \gamma \mathcal{H}(\mathbf{T}). \end{aligned} \tag{16}$$

- ▶ Similar to the proximal point method of OT problem in [Xie, et al., UAI 2020]  
[Xie, et al., UAI 2020] Xie, Y., Wang, X., Wang, R., & Zha, H. A fast proximal point method for computing exact Wasserstein distance. UAI 2020.

## Typical Method: Proximal Gradient

$$\min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \boldsymbol{\mu}_Y)} -\langle \mathbf{D}_X \mathbf{T} \mathbf{D}_Y^\top, \mathbf{T} \rangle. \text{ (Without any modification or relaxation)}$$

## Typical Method: Proximal Gradient

$\min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \boldsymbol{\mu}_Y)} -\langle \mathbf{D}_X \mathbf{T} \mathbf{D}_Y^\top, \mathbf{T} \rangle.$  (Without any modification or relaxation)

- ▶ Initialize  $\mathbf{T}^{(0)} = \boldsymbol{\mu}_X \boldsymbol{\mu}_Y^\top$  and in each iteration

$$\min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \boldsymbol{\mu}_Y)} -\langle \mathbf{D}_X \mathbf{T}^{(m)} \mathbf{D}_Y^\top + \gamma \log \mathbf{T}^{(m)}, \mathbf{T} \rangle - \gamma \mathsf{H}(\mathbf{T}). \quad (17)$$

## Typical Method: Proximal Gradient

$\min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \boldsymbol{\mu}_Y)} -\langle \mathbf{D}_X \mathbf{T} \mathbf{D}_Y^\top, \mathbf{T} \rangle$ . (Without any modification or relaxation)

- ▶ Initialize  $\mathbf{T}^{(0)} = \boldsymbol{\mu}_X \boldsymbol{\mu}_Y^\top$  and in each iteration

$$\min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \boldsymbol{\mu}_Y)} -\langle \mathbf{D}_X \mathbf{T}^{(m)} \mathbf{D}_Y^\top + \gamma \log \mathbf{T}^{(m)}, \mathbf{T} \rangle - \gamma \mathsf{H}(\mathbf{T}). \quad (17)$$

1. Compute  $\Phi = \exp\left(-\frac{\mathbf{D}_X \mathbf{T}^{(m)} \mathbf{D}_Y^\top}{\gamma}\right) \odot \mathbf{T}^{(m)}$ ,  $\mathbf{a} = \boldsymbol{\mu}_X$ .
2. Sinkhorn-Knopp Iterations: Repeat  $\mathbf{b} = \frac{\boldsymbol{\mu}_Y}{\Phi^\top \mathbf{a}}$  and  $\mathbf{a} = \frac{\boldsymbol{\mu}_X}{\Phi \mathbf{b}}$  until convergence.
3.  $\mathbf{T}^{(m+1)} = \text{diag}(\mathbf{a}) \Phi \text{diag}(\mathbf{b})$ .

## Typical Method: Conditional Gradient (CG)

**ProxGrad** is a special case of the conditional gradient (**CG**) used in [Titouan, et al. ICML 2019]

- ▶ Initialization:  $\mathbf{T}^0 = \boldsymbol{\mu}_X \boldsymbol{\mu}_Y^\top$ .
- ▶ In each iteration:
  - 1 Compute the gradient w.r.t.  $\mathbf{T}^{(m)}$ :  $\nabla_{\mathbf{T}^{(m)}} L = -\mathbf{D}_X \mathbf{T}^{(m)} \mathbf{D}_Y^\top$ .
  - 2 Solve an OT problem:  $\tilde{\mathbf{T}}^{(m+1)} = \arg \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \boldsymbol{\mu}_Y)} \langle \nabla_{\mathbf{T}^{(m)}} L, \mathbf{T} \rangle$ .
  - 3 Line-search to check whether the objective function decreases given  $\tilde{\mathbf{T}}^{(m+1)}$  and determine a momentum  $\tau \in [0, 1]$ .
  - 4  $\mathbf{T}^{(m+1)} = (1 - \tau) \mathbf{T}^{(m)} + \tau \tilde{\mathbf{T}}^{(m+1)}$ .

## Typical Method: Conditional Gradient (CG)

**ProxGrad** is a special case of the conditional gradient (**CG**) used in [Titouan, et al. ICML 2019]

- ▶ Initialization:  $\mathbf{T}^0 = \boldsymbol{\mu}_X \boldsymbol{\mu}_Y^\top$ .
- ▶ In each iteration:
  - 1 Compute the gradient w.r.t.  $\mathbf{T}^{(m)}$ :  $\nabla_{\mathbf{T}^{(m)}} L = -\mathbf{D}_X \mathbf{T}^{(m)} \mathbf{D}_Y^\top$ .
  - 2 Solve an OT problem:  $\tilde{\mathbf{T}}^{(m+1)} = \arg \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \boldsymbol{\mu}_Y)} \langle \nabla_{\mathbf{T}^{(m)}} L, \mathbf{T} \rangle$ .
  - 3 Line-search to check whether the objective function decreases given  $\tilde{\mathbf{T}}^{(m+1)}$  and determine a momentum  $\tau \in [0, 1]$ .
  - 4  $\mathbf{T}^{(m+1)} = (1 - \tau) \mathbf{T}^{(m)} + \tau \tilde{\mathbf{T}}^{(m+1)}$ .
- ▶ The Step 1+2 corresponds to solve Eq. (5) (The sub-problem of ProxGrad).
- ▶ Without line-search, the ProxGrad always sets  $\tau = 1$ .

[Titouan, et al. ICML 2019] Titouan, V., Courty, N., Tavenard, R., & Flamary, R. Optimal transport for structured data with application on graphs. ICML, 2019.

## Typical Method: Bregman ADMM

### Motivation:

- ▶ Linear Programming is expensive, while Sinkhorn-Knopp algorithm often suffers from numerical instability issue.
- ▶ Logarithmic stabilization [Chizat, et al., 2018] is not suitable to GW distance.

## Typical Method: Bregman ADMM

### Motivation:

- ▶ Linear Programming is expensive, while Sinkhorn-Knopp algorithm often suffers from numerical instability issue.
- ▶ Logarithmic stabilization [Chizat, et al., 2018] is not suitable to GW distance.

### Principle:

- ▶ Bregman ADMM for OT problem [Wang, et al. NeurIPS 2014]:

$$\min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}, \cdot), \mathbf{S} \in \Pi(\cdot, \boldsymbol{\eta}), \mathbf{T} = \mathbf{S}} \langle \mathbf{D}, \mathbf{T} \rangle$$

# Typical Method: Bregman ADMM

## Motivation:

- ▶ Linear Programming is expensive, while Sinkhorn-Knopp algorithm often suffers from numerical instability issue.
- ▶ Logarithmic stabilization [Chizat, et al., 2018] is not suitable to GW distance.

## Principle:

- ▶ Bregman ADMM for OT problem [Wang, et al. NeurIPS 2014]:

$$\begin{aligned} & \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}, \cdot), \mathbf{S} \in \Pi(\cdot, \boldsymbol{\eta}), \mathbf{T} = \mathbf{S}} \langle \mathbf{D}, \mathbf{T} \rangle \\ & \Rightarrow \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}, \cdot), \mathbf{S} \in \Pi(\cdot, \boldsymbol{\eta}), \mathbf{Z}} \langle \mathbf{D}, \mathbf{T} \rangle + \underbrace{\langle \mathbf{Z}, \mathbf{T} - \mathbf{S} \rangle}_{\text{Augmented Lagrangian}} + \overbrace{\gamma \mathcal{B}_\phi(\mathbf{T}, \mathbf{S})}^{\text{Bregman Div.}} \end{aligned} \quad (18)$$

[Chizat, et al., 2018] Chizat, L., Peyré, G., Schmitzer, B., & Vialard, F. X. Scaling algorithms for unbalanced optimal transport problems. Mathematics of Computation, 2018.

[Wang, et al. NeurIPS 2014] Wang, H., & Banerjee, A. Bregman alternating direction method of multipliers. NeurIPS 2014.

## Typical Method: Bregman ADMM

Bregman Divergence: Given a differentiable and strictly convex function  $\phi$ ,

$$B_\phi(x, y) = \phi(x) - \phi(y) - \langle \nabla \phi(y), x - y \rangle. \quad (19)$$

## Typical Method: Bregman ADMM

Bregman Divergence: Given a differentiable and strictly convex function  $\phi$ ,

$$B_\phi(x, y) = \phi(x) - \phi(y) - \langle \nabla \phi(y), x - y \rangle. \quad (19)$$

Commonly-used Bregman divergence:

- ▶  $\phi(x) = \frac{1}{2}x^2$ : Euclidean distance  $B_\phi(x, y) = \frac{1}{2}\|x - y\|^2$ .
- ▶  $\phi(x) = x \log x - x$ : KL-divergence  $B_\phi(x, y) = \text{KL}(x\|y) = x \log \frac{x}{y} - x + y$ .

## Typical Method: Bregman ADMM

Bregman Divergence: Given a differentiable and strictly convex function  $\phi$ ,

$$B_\phi(x, y) = \phi(x) - \phi(y) - \langle \nabla \phi(y), x - y \rangle. \quad (19)$$

Commonly-used Bregman divergence:

- ▶  $\phi(x) = \frac{1}{2}x^2$ : Euclidean distance  $B_\phi(x, y) = \frac{1}{2}\|x - y\|^2$ .
- ▶  $\phi(x) = x \log x - x$ : KL-divergence  $B_\phi(x, y) = \text{KL}(x\|y) = x \log \frac{x}{y} - x + y$ .

Apply the Bregman ADMM to compute GW distance [Xu, AAAI 2020]:

$$\begin{aligned} & \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \cdot), \mathbf{S} \in \Pi(\cdot, \boldsymbol{\mu}_Y), \mathbf{T} = \mathbf{S}} -\langle \mathbf{D}_X \mathbf{S} \mathbf{D}_Y^\top, \mathbf{T} \rangle \\ \Rightarrow & \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \cdot), \mathbf{S} \in \Pi(\cdot, \boldsymbol{\mu}_Y), \mathbf{Z}} -\langle \mathbf{D}_X \mathbf{S} \mathbf{D}_Y^\top, \mathbf{T} \rangle + \langle \mathbf{Z}, \mathbf{T} - \mathbf{S} \rangle + \gamma \mathcal{B}(\mathbf{T}, \mathbf{S}) \end{aligned} \quad (20)$$

[Xu, AAAI 2020] Xu, H. Gromov-Wasserstein factorization models for graph clustering. AAAI 2020.

## Typical Method: Bregman ADMM (B-ADMM)

Initialize  $\mathbf{Z}^{(0)} = \mathbf{0}$ ,  $\mathbf{T}^{(0)} = \mathbf{S}^{(0)} = \boldsymbol{\mu}_X \boldsymbol{\mu}_Y^\top$ :

## Typical Method: Bregman ADMM (B-ADMM)

Initialize  $\mathbf{Z}^{(0)} = \mathbf{0}$ ,  $\mathbf{T}^{(0)} = \mathbf{S}^{(0)} = \boldsymbol{\mu}_X \boldsymbol{\mu}_Y^\top$ :

$$\begin{aligned}\mathbf{T}^{(m+1)} &= \arg \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \cdot)} -\langle \mathbf{D}_X \mathbf{S}^{(m)} \mathbf{D}_Y^\top, \mathbf{T} \rangle + \langle \mathbf{Z}^{(m)}, \mathbf{T} \rangle + \gamma \text{KL}(\mathbf{T} \| \mathbf{S}^{(m)}) \\ \xrightarrow{\text{First-order optimality}} \mathbf{T}^{(m+1)} &= \exp \left( \frac{\gamma \log \mathbf{S}^{(m)} + \mathbf{D}_X \mathbf{S}^{(m)} \mathbf{D}_Y^\top - \mathbf{Z}^{(m)}}{\gamma} \right) \quad (21) \\ \xrightarrow{\text{Project to } \Pi(\boldsymbol{\mu}^{(k)}, \cdot)} \mathbf{T}^{(m+1)} &= \text{diag}(\boldsymbol{\mu}_X) \sigma_{row} \left( \frac{\gamma \log \mathbf{S}^{(m)} + \mathbf{D}_X \mathbf{S}^{(m)} \mathbf{D}_Y^\top - \mathbf{Z}^{(m)}}{\gamma} \right)\end{aligned}$$

## Typical Method: Bregman ADMM (B-ADMM)

Initialize  $\mathbf{Z}^{(0)} = \mathbf{0}$ ,  $\mathbf{T}^{(0)} = \mathbf{S}^{(0)} = \boldsymbol{\mu}_X \boldsymbol{\mu}_Y^\top$ :

$$\begin{aligned} \mathbf{T}^{(m+1)} &= \arg \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \cdot)} -\langle \mathbf{D}_X \mathbf{S}^{(m)} \mathbf{D}_Y^\top, \mathbf{T} \rangle + \langle \mathbf{Z}^{(m)}, \mathbf{T} \rangle + \gamma \text{KL}(\mathbf{T} \| \mathbf{S}^{(m)}) \\ \xrightarrow{\text{First-order optimality}} \mathbf{T}^{(m+1)} &= \exp \left( \frac{\gamma \log \mathbf{S}^{(m)} + \mathbf{D}_X \mathbf{S}^{(m)} \mathbf{D}_Y^\top - \mathbf{Z}^{(m)}}{\gamma} \right) \quad (21) \\ \xrightarrow{\text{Project to } \Pi(\boldsymbol{\mu}^{(k)}, \cdot)} \mathbf{T}^{(m+1)} &= \text{diag}(\boldsymbol{\mu}_X) \sigma_{row} \left( \frac{\gamma \log \mathbf{S}^{(m)} + \mathbf{D}_X \mathbf{S}^{(m)} \mathbf{D}_Y^\top - \mathbf{Z}^{(m)}}{\gamma} \right) \end{aligned}$$

$$\begin{aligned} \mathbf{S}^{(m+1)} &= \arg \min_{\mathbf{S} \in \Pi(\cdot, \boldsymbol{\mu}_Y)} -\langle \mathbf{D}_X^\top \mathbf{T}^{(m+1)} \mathbf{D}_Y, \mathbf{S} \rangle - \langle \mathbf{Z}^{(m)}, \mathbf{S} \rangle + \gamma \text{KL}(\mathbf{S} \| \mathbf{T}^{(m+1)}) \\ &= \sigma_{col} \left( \frac{\gamma \log \mathbf{T}^{(m+1)} + \mathbf{D}_X^\top \mathbf{T}^{(m+1)} \mathbf{D}_Y + \mathbf{Z}^{(m)}}{\gamma} \right) \text{diag}(\boldsymbol{\mu}_Y) \quad (22) \end{aligned}$$

## Typical Method: Bregman ADMM (B-ADMM)

Initialize  $\mathbf{Z}^{(0)} = \mathbf{0}$ ,  $\mathbf{T}^{(0)} = \mathbf{S}^{(0)} = \boldsymbol{\mu}_X \boldsymbol{\mu}_Y^\top$ :

$$\begin{aligned} \mathbf{T}^{(m+1)} &= \arg \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \cdot)} -\langle \mathbf{D}_X \mathbf{S}^{(m)} \mathbf{D}_Y^\top, \mathbf{T} \rangle + \langle \mathbf{Z}^{(m)}, \mathbf{T} \rangle + \gamma \text{KL}(\mathbf{T} \| \mathbf{S}^{(m)}) \\ \xrightarrow{\text{First-order optimality}} \mathbf{T}^{(m+1)} &= \exp \left( \frac{\gamma \log \mathbf{S}^{(m)} + \mathbf{D}_X \mathbf{S}^{(m)} \mathbf{D}_Y^\top - \mathbf{Z}^{(m)}}{\gamma} \right) \end{aligned} \quad (21)$$

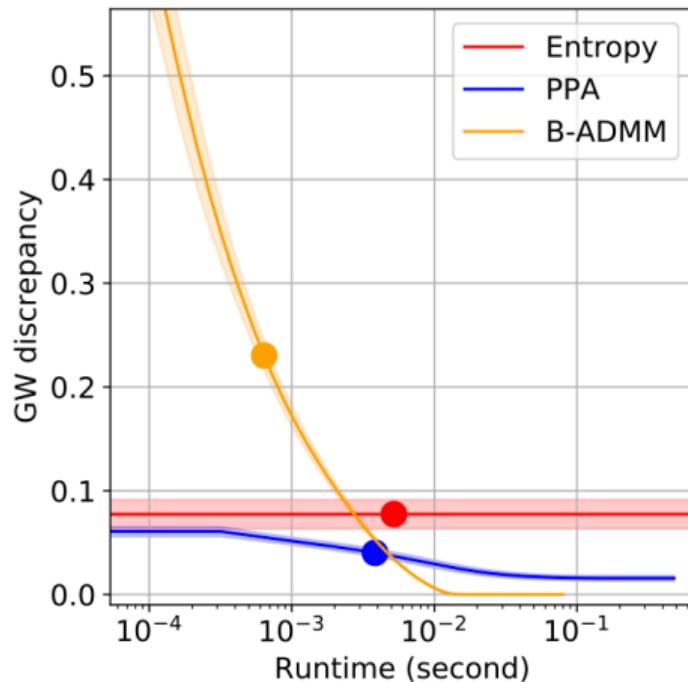
$$\xrightarrow{\text{Project to } \Pi(\boldsymbol{\mu}^{(k)}, \cdot)} \mathbf{T}^{(m+1)} = \text{diag}(\boldsymbol{\mu}_X) \sigma_{row} \left( \frac{\gamma \log \mathbf{S}^{(m)} + \mathbf{D}_X \mathbf{S}^{(m)} \mathbf{D}_Y^\top - \mathbf{Z}^{(m)}}{\gamma} \right)$$

$$\begin{aligned} \mathbf{S}^{(m+1)} &= \arg \min_{\mathbf{S} \in \Pi(\cdot, \boldsymbol{\mu}_Y)} -\langle \mathbf{D}_X^\top \mathbf{T}^{(m+1)} \mathbf{D}_Y, \mathbf{S} \rangle - \langle \mathbf{Z}^{(m)}, \mathbf{S} \rangle + \gamma \text{KL}(\mathbf{S} \| \mathbf{T}^{(m+1)}) \\ &= \sigma_{col} \left( \frac{\gamma \log \mathbf{T}^{(m+1)} + \mathbf{D}_X^\top \mathbf{T}^{(m+1)} \mathbf{D}_Y + \mathbf{Z}^{(m)}}{\gamma} \right) \text{diag}(\boldsymbol{\mu}_Y) \end{aligned} \quad (22)$$

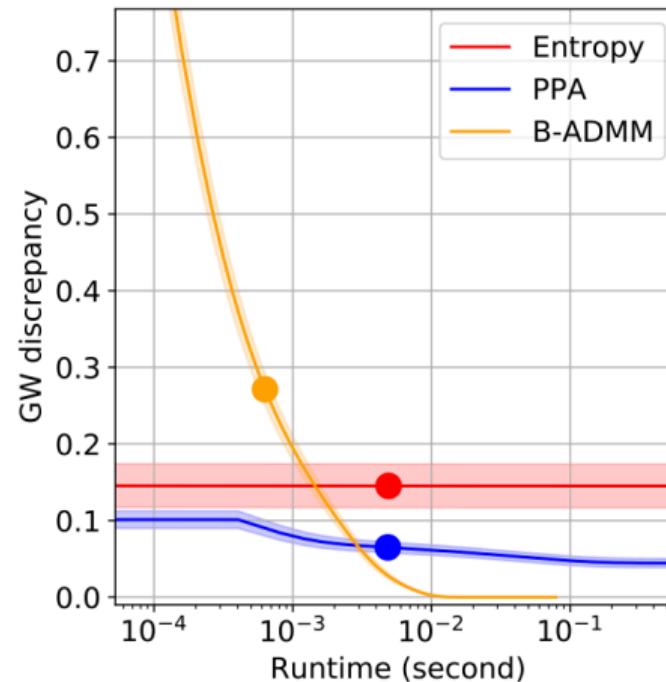
$$\mathbf{Z}^{(m+1)} = \mathbf{Z}^{(m)} + \gamma (\mathbf{T}^{(m+1)} - \mathbf{S}^{(m+1)}). \quad (23)$$

# Comparisons on Weighted Graphs

Undirected weighted graphs

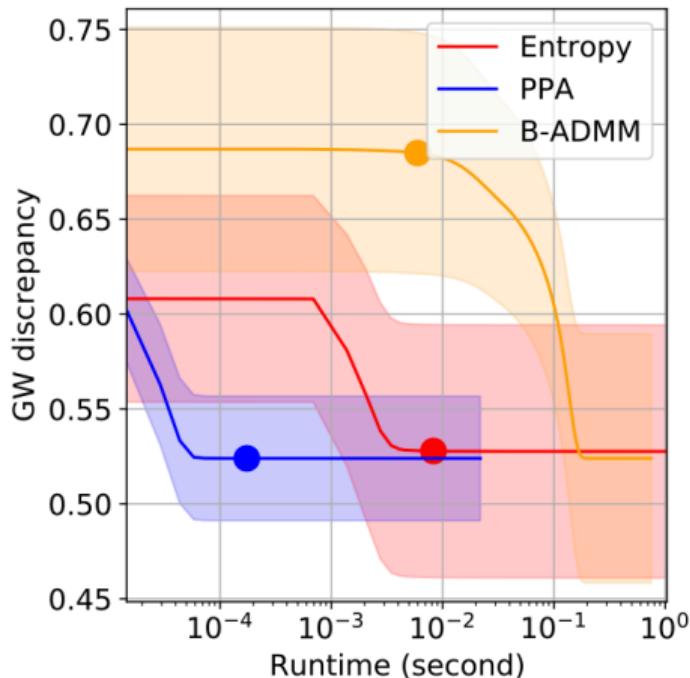


Directed weighted graphs

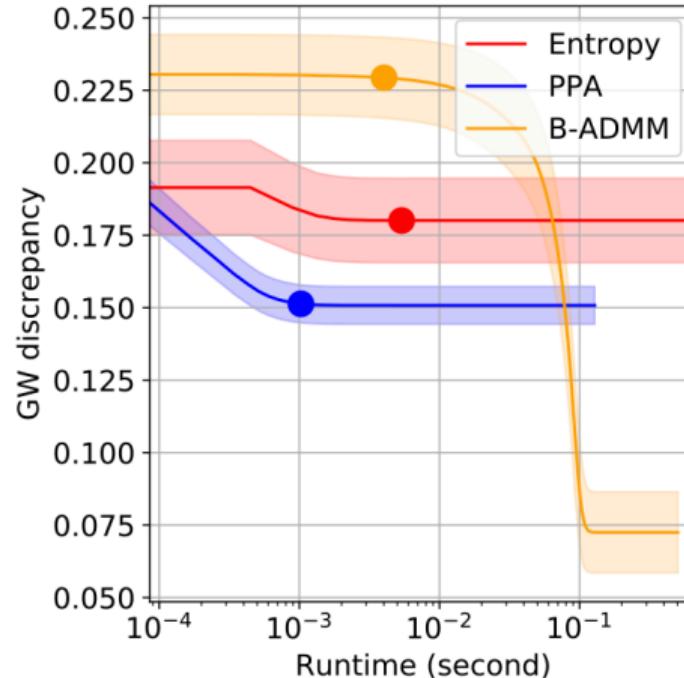


# Comparisons on Unweighted Graphs

Undirected unweighted graph



Directed unweighted graph



# Convergence and Complexity

## Convergence

- ▶  $\lim_{m \rightarrow \infty} \mathbf{T}^{(m)}$  is a stationary point [Xu, et al., ICML 2019].
- ▶ Convergence: Both Sinkhorn-based and Bregman-ADMM methods own  $\mathcal{O}(1/\epsilon^2)$ .

# Convergence and Complexity

## Convergence

- ▶  $\lim_{m \rightarrow \infty} \mathbf{T}^{(m)}$  is a stationary point [Xu, et al., ICML 2019].
- ▶ Convergence: Both Sinkhorn-based and Bregman-ADMM methods own  $\mathcal{O}(1/\epsilon^2)$ .

## Computational complexity per iteration: $D_X T D_Y^\top$

- ▶  $\mathcal{O}(V^3)$  for dense distance/kernel matrix ( $V \times V$ )

# Convergence and Complexity

## Convergence

- ▶  $\lim_{m \rightarrow \infty} \mathbf{T}^{(m)}$  is a stationary point [Xu, et al., ICML 2019].
- ▶ Convergence: Both Sinkhorn-based and Bregman-ADMM methods own  $\mathcal{O}(1/\epsilon^2)$ .

## Computational complexity per iteration: $D_X T D_Y^\top$

- ▶  $\mathcal{O}(V^3)$  for dense distance/kernel matrix ( $V \times V$ )
- ▶  $\mathcal{O}(VE)$  for the graph with  $V$  nodes and  $E$  edges.

# Convergence and Complexity

## Convergence

- ▶  $\lim_{m \rightarrow \infty} \mathbf{T}^{(m)}$  is a stationary point [Xu, et al., ICML 2019].
- ▶ Convergence: Both Sinkhorn-based and Bregman-ADMM methods own  $\mathcal{O}(1/\epsilon^2)$ .

## Computational complexity per iteration: $D_X T D_Y^\top$

- ▶  $\mathcal{O}(V^3)$  for dense distance/kernel matrix ( $V \times V$ )
- ▶  $\mathcal{O}(VE)$  for the graph with  $V$  nodes and  $E$  edges.
- ▶ **Can we further accelerate the process?**

# Acceleration Strategy 1: Divide and Conquer

## Motivation:

- ▶ When  $V = |\mathcal{V}_X| \gg |\mathcal{V}_Y| = N$  (graph partitioning),  $\mathcal{O}(N(E + V))$ .

# Acceleration Strategy 1: Divide and Conquer

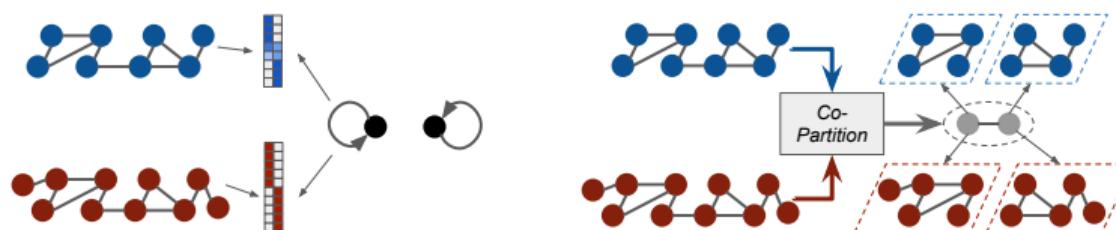
## Motivation:

- When  $V = |\mathcal{V}_X| \gg |\mathcal{V}_Y| = N$  (graph partitioning),  $\mathcal{O}(N(E + V))$ .

## Co-partition two graphs:

$$\bar{\mathbf{D}}, \mathbf{T}_X, \mathbf{T}_Y = \arg \min \frac{|\mathcal{V}_X|}{|\mathcal{V}_X| + |\mathcal{V}_Y|} d_{gw}(B, G_X) + \frac{|\mathcal{V}_Y|}{|\mathcal{V}_X| + |\mathcal{V}_Y|} d_{gw}(B, G_Y)$$

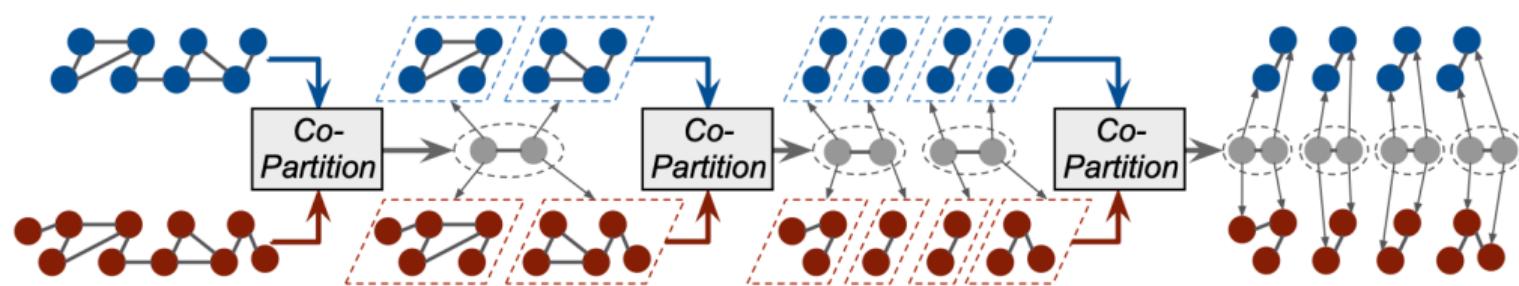
Initialize the barycenter graph by a disconnected graph.



Computational complexity:  $\mathcal{O}(2(V + E))$

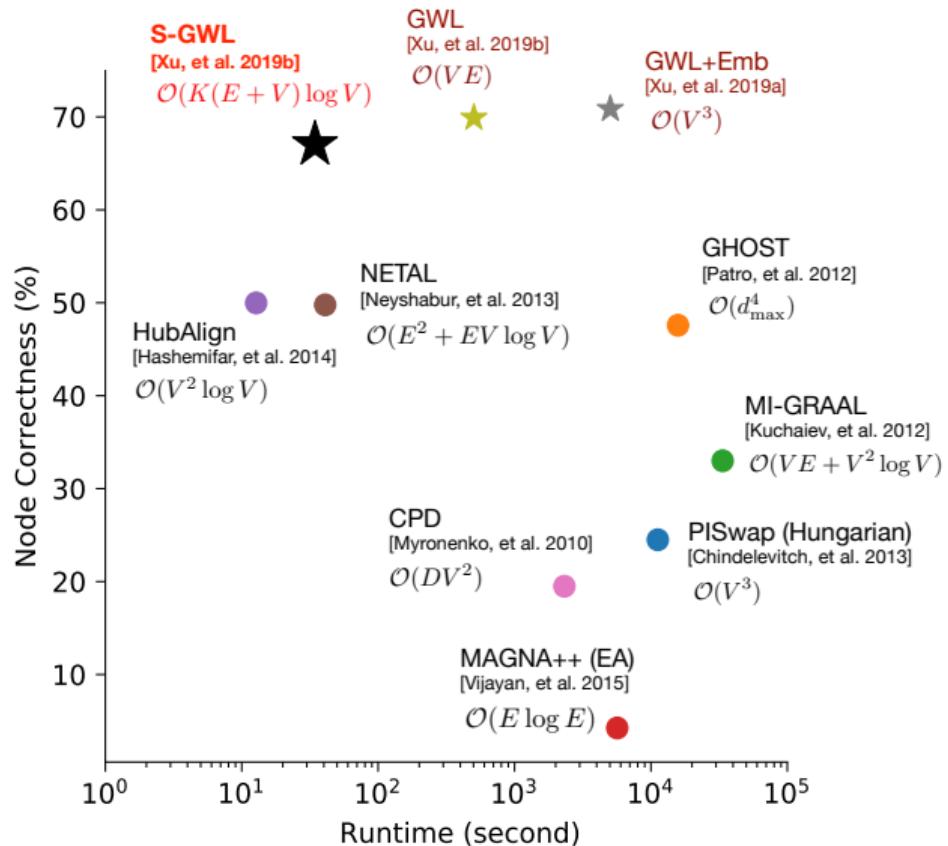
# Acceleration Strategy 1: Divide and Conquer

A “Divide and Conquer” strategy based on recursive co-partitioning



Computational complexity:  $\mathcal{O}((E + V) \log V)$  (Proposition 3.1 in [Xu, et al., NeurIPS 2019])

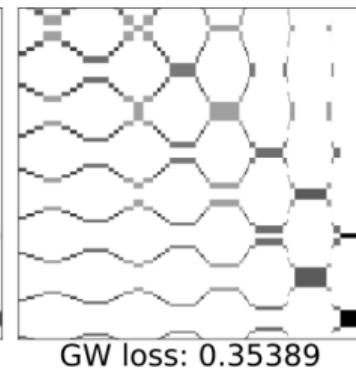
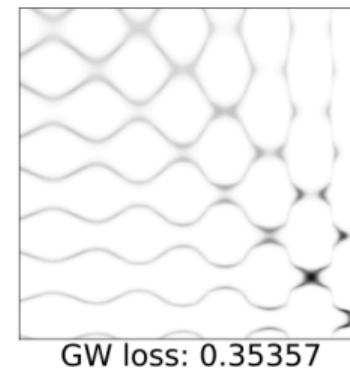
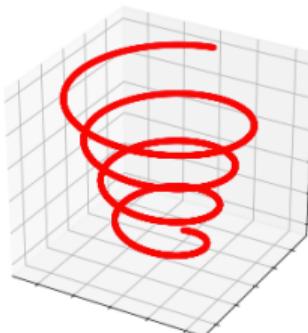
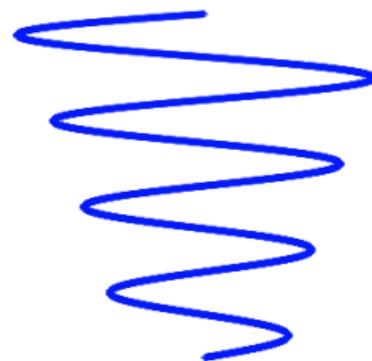
# Acceleration Strategy 1: Divide and Conquer



## Acceleration Strategy 2: Low-rank Structured GW

### Motivation:

- ▶ The optimal transport matrix often has a low-rank approximation.



An illustration of low-rank structured GW shown in [Scetbon, et al., 2021].

[Scetbon, et al., 2021] Scetbon, M., Peyré, G., & Cuturi, M. Linear-time Gromov-Wasserstein distances using low rank couplings and costs. arXiv 2021.

## Acceleration Strategy 2: Low-rank Structured GW

### Principle:

- ▶ A low-rank factorization of proposed optimal transport matrix

$$\mathbf{T} = \mathbf{Q} \text{diag}(1/\mathbf{g}) \mathbf{R}^\top, \text{ where}$$

$$\underbrace{\mathbf{Q} \in \Pi(\boldsymbol{\mu}_X, \mathbf{g}), \quad \mathbf{R} \in \Pi(\boldsymbol{\mu}_Y, \mathbf{g}), \quad \mathbf{g} \in \Delta^{r-1}}_{\mathcal{C}(\boldsymbol{\mu}_X, \boldsymbol{\mu}_Y, r)}, \quad r \ll M, N. \quad (24)$$

## Acceleration Strategy 2: Low-rank Structured GW

### Principle:

- ▶ A low-rank factorization of proposed optimal transport matrix

$$\mathbf{T} = \mathbf{Q} \text{diag}(1/\mathbf{g}) \mathbf{R}^\top, \text{ where}$$

$$\underbrace{\mathbf{Q} \in \Pi(\mu_X, \mathbf{g}), \quad \mathbf{R} \in \Pi(\mu_Y, \mathbf{g}), \quad \mathbf{g} \in \Delta^{r-1}}_{\mathcal{C}(\mu_X, \mu_Y, r)}, \quad r \ll M, N. \quad (24)$$

- ▶ In each iteration, solving the following entropic low-rank GW problem via the Dykstra's Algorithm shown in [Scetbon, et al., ICML 2021]

$$\min_{\mathbf{T}} -\langle \mathbf{D}_X \underbrace{\mathbf{T}^{(m)}}_{\text{Constant}} \mathbf{D}_Y^\top, \mathbf{T} \rangle - \gamma \underbrace{H(\mathbf{Q}, \mathbf{R}, \mathbf{g})}_{\text{Entropy}} \quad (25)$$

$$s.t. \quad \mathbf{T} = \mathbf{Q} \text{diag}(1/\mathbf{g}) \mathbf{R}^\top, \text{ and } \mathbf{Q}, \mathbf{R}, \mathbf{g} \in \mathcal{C}(\mu_X, \mu_Y, r).$$

- ▶ The complexity per iteration becomes  $\mathcal{O}(M^2r + N^2r)$ .

[Scetbon, et al., ICML 2021] Scetbon, M., Cuturi, M., & Peyré, G. Low-rank Sinkhorn factorization. ICML 2021.

## Acceleration Strategy 2: Low-rank Structured GW

- ▶ In each step, the convergence to a unique solution is guaranteed (because of the convergence of Dykstra's algorithm).

## Acceleration Strategy 2: Low-rank Structured GW

- ▶ In each step, the convergence to a unique solution is guaranteed (because of the convergence of Dykstra's algorithm).
- ▶ Furthermore, when the matrices  $\mathbf{D}_X$  and  $\mathbf{D}_Y$  own low-rank structures, i.e.,

$$\mathbf{D}_X = \mathbf{XX}^T, \quad \mathbf{D}_Y = \mathbf{YY}^T, \text{ where } \mathbf{X} \in \mathbb{R}^{M \times d}, \quad \mathbf{Y} \in \mathbb{R}^{N \times d'}, \quad (26)$$

the complexity can be further reduced to  $\mathcal{O}(Mr(r + d) + Nr(r + d'))$  (generally,  $d \ll M$  and  $d' \ll N$ ).

## Acceleration Strategy 3: Tree Structured GW

### Motivation:

- ▶ In practice, graphs or point clouds often own tree structures.
- ▶ The tree structures are also beneficial for accelerating GW distance.

# Acceleration Strategy 3: Tree Structured GW

## Motivation:

- ▶ In practice, graphs or point clouds often own tree structures.
- ▶ The tree structures are also beneficial for accelerating GW distance.

## Principle:

- ▶ Given  $\mathbf{X} = \{x_i\}_{i=1}^M$  and  $\mathbf{Y} = \{y_j\}_{j=1}^N$ , define a **flow-based alignment discrepancy** [Le, et al., AISTATS 2021]:

$$d_A(\mathbf{X}, \mathbf{Y}) = \underbrace{\min_{r_x, r_y, \mathbf{T} \in \Pi(\mu_X, \mu_Y)} \sum_{i,j} |d_{T_X}(r_x, x_i) - d_{T_Y}(r_y, y_j)|^2 T_{ij}}_{d_w(T_X(r_x, \mu_X), T_Y(r_y, \mu_Y))}, \quad (27)$$

where explore optimal transport under optimal root pair  $(r_x, r_y)$ .

[Le, et al., AISTATS 2021] Le, T., Ho, N., & Yamada, M. Flow-based alignment approaches for probability measures in different spaces. AISTATS 2021.

## Acceleration Strategy 3: Tree Structured GW

$$d_A(\mathbf{X}, \mathbf{Y}) = \min_{r_x, r_z, \mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \boldsymbol{\mu}_Y)} \sum_{i,j} |d_{T_X}(r_x, x_i) - d_{T_Y}(r_y, y_j)|^2 T_{ij},$$

- ▶ The roots and tree structures can be generated by exploring the hierarchical clustering structures within  $\mathbf{X}$  and  $\mathbf{Y}$ , respectively.
- ▶ Based on the tree structure, the distance from a root to a node (i.e.,  $d_{T_X}(r_x, x_i)$ ) can be calculated in a recursive manner.
- ▶ Complexity becomes  $\mathcal{O}(MN)$ .

## Acceleration Strategy 3: Tree Structured GW

$$d_A(\mathbf{X}, \mathbf{Y}) = \min_{r_x, r_z, \mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \boldsymbol{\mu}_Y)} \sum_{i,j} |d_{T_X}(r_x, x_i) - d_{T_Y}(r_y, y_j)|^2 T_{ij},$$

- ▶ The roots and tree structures can be generated by exploring the hierarchical clustering structures within  $\mathbf{X}$  and  $\mathbf{Y}$ , respectively.
- ▶ Based on the tree structure, the distance from a root to a node (i.e.,  $d_{T_X}(r_x, x_i)$ ) can be calculated in a recursive manner.
- ▶ Complexity becomes  $\mathcal{O}(MN)$ .
- ▶ **Theorem 1** in [Le, et al., AISTATS 2021]:  $d_A$  is a pseudo-distance satisfying symmetry and the triangle inequality, and  $d_{gw} \leq 2d_A$  for the trees with  $depth \leq 2$ .

## Acceleration Strategy 4: Sliced Gromov-Wasserstein

### Motivation:

- ▶ The optimal transport between one-dimensional distributions is relatively easy to solve.

## Acceleration Strategy 4: Sliced Gromov-Wasserstein

### Motivation:

- ▶ The optimal transport between one-dimensional distributions is relatively easy to solve.

### Sliced-Wasserstein distance [Kolouri, et al., NeurIPS 2019]:

- ▶ Given two distributions  $\mu$  and  $\eta$  defined on a metric space  $(\mathcal{X} \subset \mathbb{R}^d, d_X)$ , projecting  $\forall x \in \mathcal{X}$  through a linear projection  $\theta \in \mathcal{S}^{d-1}$ , i.e.,  $R_\theta(x) = \langle x, \theta \rangle$ , leads to
  - ▶ A one-dimensional metric space  $(R_\theta(\mathcal{X}), d_{R_\theta(\mathcal{X})})$ .
  - ▶ The one-dimensional distributions after projection  $R_\theta \#\mu$  and  $R_\theta \#\eta$ .

## Acceleration Strategy 4: Sliced Gromov-Wasserstein

### Motivation:

- ▶ The optimal transport between one-dimensional distributions is relatively easy to solve.

### Sliced-Wasserstein distance [Kolouri, et al., NeurIPS 2019]:

- ▶ Given two distributions  $\mu$  and  $\eta$  defined on a metric space  $(\mathcal{X} \subset \mathbb{R}^d, d_X)$ , projecting  $\forall x \in \mathcal{X}$  through a linear projection  $\theta \in \mathcal{S}^{d-1}$ , i.e.,  $R_\theta(x) = \langle x, \theta \rangle$ , leads to
  - ▶ A one-dimensional metric space  $(R_\theta(\mathcal{X}), d_{R_\theta(\mathcal{X})})$ .
  - ▶ The one-dimensional distributions after projection  $R_\theta \#\mu$  and  $R_\theta \#\eta$ .

$$d_{sw}(\mu, \eta) = \mathbb{E}_{\theta \sim p_{\mathcal{S}^{d-1}}} [d_w(R_\theta \#\mu, R_\theta \#\eta)] = \int_{\theta \in \mathcal{S}^{d-1}} p(\theta) d_w(R_\theta \#\mu, R_\theta \#\eta) d\theta \quad (28)$$

[Kolouri, et al., NeurIPS 2019] Kolouri, S., Nadjahi, K., Simsekli, U., Badeau, R., & Rohde, G. Generalized sliced wasserstein distances. NeurIPS, 2019.

## Acceleration Strategy 4: Sliced Gromov-Wasserstein

- ▶ Practical implementation:
  - ▶ Samples  $\mathbf{X} = \{x_i\}_{i=1}^N \sim \mu$  and  $\mathbf{Y} = \{y_i\}_{i=1}^N \sim \eta$  are provided.
  - ▶ Finite number of projections are sampled based on a distribution  $\{\theta_l\}_{l=1}^L \sim p_{\mathcal{S}^{d-1}}$ .

## Acceleration Strategy 4: Sliced Gromov-Wasserstein

- ▶ Practical implementation:
  - ▶ Samples  $\mathbf{X} = \{x_i\}_{i=1}^N \sim \mu$  and  $\mathbf{Y} = \{y_i\}_{i=1}^N \sim \eta$  are provided.
  - ▶ Finite number of projections are sampled based on a distribution  $\{\theta_l\}_{l=1}^L \sim p_{\mathcal{S}^{d-1}}$ .
- ▶ Sample-based sliced Wasserstein distance:

$$\hat{d}_{sw}(\mu, \eta) = \frac{1}{L} \sum_{l=1}^L \left( \min_{T \in \Pi(\frac{1}{N}\mathbf{1}_N, \frac{1}{N}\mathbf{1}_N)} \sum_{i,j=1}^N |\theta_l^\top x_i - \theta_l^\top y_j|^2 T_{ij} \right)$$

## Acceleration Strategy 4: Sliced Gromov-Wasserstein

- ▶ Practical implementation:
  - ▶ Samples  $\mathbf{X} = \{x_i\}_{i=1}^N \sim \mu$  and  $\mathbf{Y} = \{y_i\}_{i=1}^N \sim \eta$  are provided.
  - ▶ Finite number of projections are sampled based on a distribution  $\{\theta_l\}_{l=1}^L \sim p_{\mathcal{S}^{d-1}}$ .
- ▶ Sample-based sliced Wasserstein distance:

$$\begin{aligned}\hat{d}_{sw}(\mu, \eta) &= \frac{1}{L} \sum_{l=1}^L \left( \min_{\mathbf{T} \in \Pi(\frac{1}{N} \mathbf{1}_N, \frac{1}{N} \mathbf{1}_N)} \sum_{i,j=1}^N |\theta_l^\top x_i - \theta_l^\top y_j|^2 T_{ij} \right) \\ &= \frac{1}{L} \sum_{l=1}^L \left( \frac{1}{N} \min_{\sigma \in \mathcal{P}_N} \sum_{i=1}^N |\theta_l^\top x_i - \theta_l^\top y_{\sigma(i)}|^2 \right)\end{aligned}\tag{29}$$

**Theorem:** For one dimensional  $x_1 \leq \dots \leq x_N$  and  $y_1 \leq \dots \leq y_N$ , identity permutation ( $\sigma(i) = i$  for  $i = 1, \dots, N$ ) leads to the optimal transport between them.

## Acceleration Strategy 4: Sliced Gromov-Wasserstein

- ▶ Similarly, given  $(\mathcal{X} \in \mathbb{R}^d, d_X, \mu_X)$  and  $(\mathcal{Y} \in \mathbb{R}^d, d_Y, \mu_Y)$ , we can project the mm-spaces to one-dimensional mm-spaces by a linear projection  $\theta \in \mathcal{S}^{d-1}$ :

$$\begin{aligned}\mathcal{X}_\theta &:= R_\theta(\mathcal{X})_{d_{R_\theta(X)}, R_\theta \# \mu_X}, \\ \mathcal{Y}_\theta &:= R_\theta(\mathcal{Y})_{d_{R_\theta(Y)}, R_\theta \# \mu_Y}.\end{aligned}\tag{30}$$

- ▶ When  $\dim(\mathcal{X}) \neq \dim(\mathcal{Y})$ , padding zeros to the samples in the lower-dimensional space.

## Acceleration Strategy 4: Sliced Gromov-Wasserstein

- ▶ Similarly, given  $(\mathcal{X} \in \mathbb{R}^d, d_X, \mu_X)$  and  $(\mathcal{Y} \in \mathbb{R}^d, d_Y, \mu_Y)$ , we can project the mm-spaces to one-dimensional mm-spaces by a linear projection  $\theta \in \mathcal{S}^{d-1}$ :

$$\begin{aligned}\mathcal{X}_\theta &:= R_\theta(\mathcal{X})_{d_{R_\theta(X)}, R_\theta \# \mu_X}, \\ \mathcal{Y}_\theta &:= R_\theta(\mathcal{Y})_{d_{R_\theta(Y)}, R_\theta \# \mu_Y}.\end{aligned}\tag{30}$$

- ▶ When  $\dim(\mathcal{X}) \neq \dim(\mathcal{Y})$ , padding zeros to the samples in the lower-dimensional space.
- ▶ **Sliced Gromov-Wasserstein distance** [Titouan, et al., NeurIPS 2019]:

$$d_{sgw}(\mathcal{X}_{d_X, \mu_X}, \mathcal{Y}_{d_Y, \mu_Y}) = \mathbb{E}_{\theta \sim p_{\mathcal{S}^{d-1}}} [d_{gw}(\mathcal{X}_\theta, \mathcal{Y}_\theta)].\tag{31}$$

[Titouan, et al., NeurIPS 2019] Titouan, V., Flamary, R., Courty, N., Tavenard, R., & Chapel, L. Sliced gromov-wasserstein. NeurIPS, 2019.

## Acceleration Strategy 4: Sliced Gromov-Wasserstein

- ▶ Practical implementation:
  - ▶ Samples  $\mathbf{X} = \{x_i\}_{i=1}^N \sim \mu_X$  and  $\mathbf{Y} = \{y_i\}_{i=1}^N \sim \mu_Y$  are provided.
  - ▶ Finite number of projections are sampled based on a distribution  $\{\theta_l\}_{l=1}^L \sim p_{\mathcal{S}^{d-1}}$ .
- ▶ **Sample-based sliced Gromov-Wasserstein distance:**

$$\hat{d}_{sgw}(\mathbf{X}, \mathbf{Y})$$

$$= \frac{1}{L} \sum_{l=1}^L \left( \min_{T \in \Pi(\frac{1}{N}\mathbf{1}_N, \frac{1}{N}\mathbf{1}_N)} \sum_{i,i,j,j'=1}^N (|\theta_l^\top x_i - \theta_l^\top x_{i'}|^2 - |\theta_l^\top y_j - \theta_l^\top y_{j'}|^2)^2 T_{ij} T_{i'j'} \right)$$

## Acceleration Strategy 4: Sliced Gromov-Wasserstein

- ▶ Practical implementation:
  - ▶ Samples  $\mathbf{X} = \{x_i\}_{i=1}^N \sim \mu_X$  and  $\mathbf{Y} = \{y_i\}_{i=1}^N \sim \mu_Y$  are provided.
  - ▶ Finite number of projections are sampled based on a distribution  $\{\theta_l\}_{l=1}^L \sim p_{\mathcal{S}^{d-1}}$ .
- ▶ **Sample-based sliced Gromov-Wasserstein distance:**

$$\begin{aligned}\hat{d}_{sgw}(\mathbf{X}, \mathbf{Y}) &= \frac{1}{L} \sum_{l=1}^L \left( \min_{T \in \Pi(\frac{1}{N} \mathbf{1}_N, \frac{1}{N} \mathbf{1}_N)} \sum_{i,i',j,j'=1}^N (|\theta_l^\top x_i - \theta_l^\top x_{i'}|^2 - |\theta_l^\top y_j - \theta_l^\top y_{j'}|^2)^2 T_{ij} T_{i'j'} \right) \quad (32) \\ &= \frac{1}{L} \sum_{l=1}^L \left( \frac{1}{N^2} \min_{\sigma \in \mathcal{P}_N} \sum_{i,i'=1}^N (|\theta_l^\top x_i - \theta_l^\top x_{i'}|^2 - |\theta_l^\top y_{\sigma(i)} - \theta_l^\top y_{\sigma(i')}|^2)^2 \right)\end{aligned}$$

**Theorem 3.2** in [Titouan, et al., NeurIPS 2019]: If  $x_1 < \dots < x_n$  and  $y_1 < \dots < y_n$ , this result is achieved either by the identity or the anti-identity permutation.

# Facing to More Complicated Scenarios

## **Challenges:**

- ▶ The nodes of graphs might be attributed.
- ▶ The graphs might have clustering structures.
- ▶ The node distributions might be unavailable or unreliable.

# Facing to More Complicated Scenarios

## Challenges:

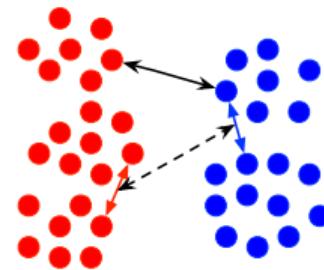
- ▶ The nodes of graphs might be attributed.
- ▶ The graphs might have clustering structures.
- ▶ The node distributions might be unavailable or unreliable.

## More variants of GW:

- ▶ Consider node attributes  $\Rightarrow$  Fused GW Distance
- ▶ Consider clustering structures  $\Rightarrow$  Hierarchical (Fused) GW Distance
- ▶ Relax doubly-stochastic constraint  $\Rightarrow$  Unbalance/Partial GW Distance

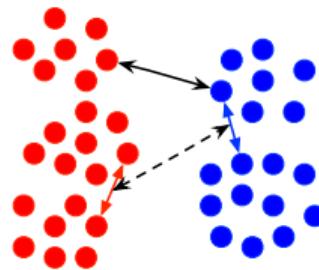
## Fused Gromov-Wasserstein Distance

For graphs/point clouds, consider the GW distance between their topological structure/nodes' relations and the Wasserstein distance between their node attributions jointly [Titouan, et al., ICML 2019].



## Fused Gromov-Wasserstein Distance

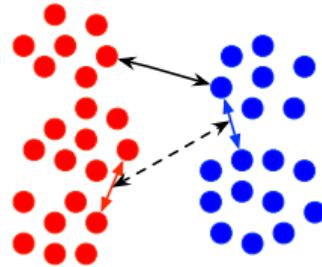
For graphs/point clouds, consider the GW distance between their topological structure/nodes' relations and the Wasserstein distance between their node attributions jointly [Titouan, et al., ICML 2019].



$$d_{fgw}(\mathcal{X}_{d_X, \mu_X}, \mathcal{Y}_{d_Y, \mu_Y}; \beta) = \inf_{\pi \in \Pi(\mu_X, \mu_Y)} \left( (1 - \beta) \underbrace{\int_{\mathcal{X} \times \mathcal{Y}} d_{x,y} \pi(x, y) dx dy}_{\text{Wasserstein term}} + \beta \underbrace{\int_{\mathcal{X} \times \mathcal{Y}} \int_{\mathcal{X} \times \mathcal{Y}} r_{x,y,x',y'} \pi(x, y) \pi(x', y') dx dx' dy dy' }_{\text{Gromov-Wasserstein term}} \right). \quad (33)$$

[Titouan, et al., ICML 2019] Titouan, V., Courty, N., Tavenard, R., & Flamary, R. Optimal transport for structured data with application on graphs. ICML 2019.

## Fused Gromov-Wasserstein Distance

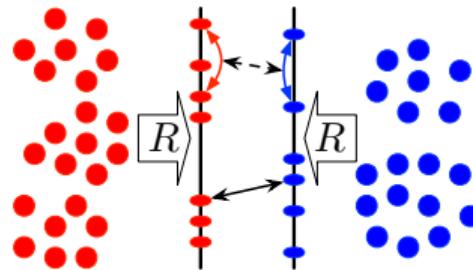


In practice, given  $\mathbf{X} = \{x_i\}_{i=1}^M$  and  $\mathbf{Y} = \{y_j\}_{j=1}^N$ , we have

$$\begin{aligned}\hat{d}_{fgw}(\mathbf{X}, \mathbf{Y}; \beta) &= \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \boldsymbol{\mu}_Y)} \left( (1 - \beta) \sum_{i=1}^M \sum_{j=1}^N d(x_i, y_j) T_{ij} \right. \\ &\quad \left. + \beta \sum_{i,i'=1}^M \sum_{j,j'=1}^N r(x_i, x_{i'}, y_j, y_{j'}) T_{ij} T_{i'j'} \right) \quad (34) \\ &= \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \boldsymbol{\mu}_Y)} \langle \mathbf{D}_{XY} - 2\mathbf{D}_X \mathbf{T} \mathbf{D}_Y^\top, \mathbf{T} \rangle\end{aligned}$$

## Sliced Fused Gromov-Wasserstein Distance

The Sliced FGW distance in [Xu, et al., ICML 2020] considers the expectation of the FGW distance between one-dimensional mm-spaces derived by random projections:



$$d_{sfgw}(\mathcal{X}_{d_X, \mu_X}, \mathcal{Y}_{d_Y, \mu_Y}; \beta) = \mathbb{E}_{\theta \sim p_{S^{d-1}}} [d_{fgw}(\mathcal{X}_\theta, \mathcal{Y}_\theta; \beta)]. \quad (35)$$

[Xu, et al., ICML 2020] Xu, H., Luo, D., Henao, R., Shah, S., & Carin, L. Learning autoencoders with relational regularization. ICML 2020.

## Sliced Fused Gromov-Wasserstein Distance

Given  $\mathbf{X} = \{x_i\}_{i=1}^N$  and  $\mathbf{Y} = \{y_j\}_{j=1}^N$ , and  $L$  linear projections  $\{\theta_l \in \mathcal{S}^{d-1}\}$ , we have

$$\begin{aligned} \hat{d}_{sf gw}(\mathbf{X}, \mathbf{Y}; \beta) &= \frac{1}{L} \sum_{l=1}^L \left( \min_{\sigma \in \mathcal{P}_N} (1 - \beta) \sum_{i=1}^N |\theta_l^\top x_i - \theta_l^\top y_{\sigma(i)}|^2 \right. \\ &\quad \left. + \beta \sum_{i,i'=1}^N (|\theta_l^\top x_i - \theta_l^\top x_{i'}|^2 - |\theta_l^\top y_{\sigma(i)} - \theta_l^\top y_{\sigma(i')}|^2)^2 \right). \end{aligned} \tag{36}$$

## Sliced Fused Gromov-Wasserstein Distance

Given  $\mathbf{X} = \{x_i\}_{i=1}^N$  and  $\mathbf{Y} = \{y_j\}_{j=1}^N$ , and  $L$  linear projections  $\{\theta_l \in \mathcal{S}^{d-1}\}$ , we have

$$\begin{aligned} \hat{d}_{sf gw}(\mathbf{X}, \mathbf{Y}; \beta) = & \frac{1}{L} \sum_{l=1}^L \left( \min_{\sigma \in \mathcal{P}_N} (1 - \beta) \sum_{i=1}^N |\theta_l^\top x_i - \theta_l^\top y_{\sigma(i)}|^2 \right. \\ & \left. + \beta \sum_{i,i'=1}^N (|\theta_l^\top x_i - \theta_l^\top x_{i'}|^2 - |\theta_l^\top y_{\sigma(i)} - \theta_l^\top y_{\sigma(i')}|^2)^2 \right). \end{aligned} \quad (36)$$

- ▶ **Theorem 3.4** in [Xu, et al., ICML 2020]: For one-dimensional samples  $x_1 < \dots < x_n$  and  $y_1 < \dots < y_n$ , we denote their zero-mean translations as  $x' = \{x'_i\}$  and  $y' = \{y'_i\}$ , respectively.

1. When  $(\sum_i x'_i y'_i + \frac{1-\beta}{8\beta})^2 \geq (\sum_i x'_i y'_{n+1-i} + \frac{1-\beta}{8\beta})^2$ , the solution is the identity permutation  $\sigma(i) = i$ .
2. Otherwise, the solution is the anti-identity permutation  $\sigma(i) = n + 1 - i$ .

## Sliced Fused Gromov-Wasserstein Distance

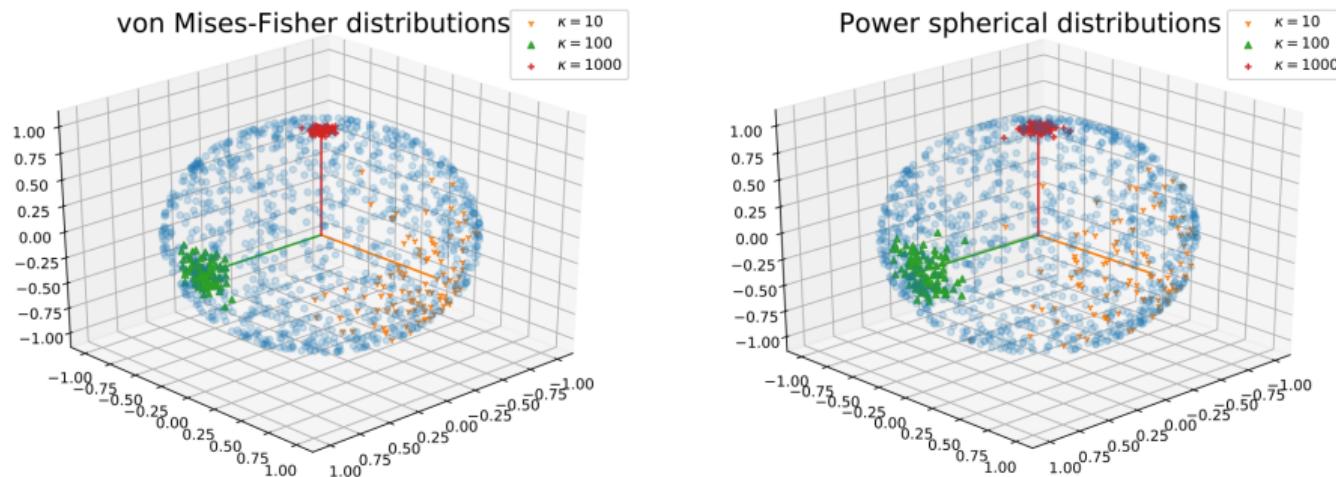
Given  $\mathbf{X} = \{x_i\}_{i=1}^N$  and  $\mathbf{Y} = \{y_j\}_{j=1}^N$ , and  $L$  linear projections  $\{\theta_l \in \mathcal{S}^{d-1}\}$ , we have

$$\begin{aligned} \hat{d}_{sf gw}(\mathbf{X}, \mathbf{Y}; \beta) &= \frac{1}{L} \sum_{l=1}^L \left( \min_{\sigma \in \mathcal{P}_N} (1 - \beta) \sum_{i=1}^N |\theta_l^\top x_i - \theta_l^\top y_{\sigma(i)}|^2 \right. \\ &\quad \left. + \beta \sum_{i,i'=1}^N (|\theta_l^\top x_i - \theta_l^\top x_{i'}|^2 - |\theta_l^\top y_{\sigma(i)} - \theta_l^\top y_{\sigma(i')}|^2)^2 \right). \end{aligned} \quad (36)$$

- ▶ **Theorem 3.4** in [Xu, et al., ICML 2020]: For one-dimensional samples  $x_1 < \dots < x_n$  and  $y_1 < \dots < y_n$ , we denote their zero-mean translations as  $x' = \{x'_i\}$  and  $y' = \{y'_i\}$ , respectively.
  1. When  $(\sum_i x'_i y'_i + \frac{1-\beta}{8\beta})^2 \geq (\sum_i x'_i y'_{n+1-i} + \frac{1-\beta}{8\beta})^2$ , the solution is the identity permutation  $\sigma(i) = i$ .
  2. Otherwise, the solution is the anti-identity permutation  $\sigma(i) = n + 1 - i$ .
- ▶ Obviously, when  $\beta = 1$ , we obtain sliced GW distance. Theorem 3.4 extends the Theorem 3.2 in [Titouan, et al, NeurIPS 2019], which provides an explicit condition to compute the sliced GW distance.

# Sampling Strategies for Sliced (F)GW Distance

- ▶ Uniform sampling:  $p_{\mathcal{S}^{d-1}} = \text{Uniform}(\mathcal{S}^{d-1})$ .
- ▶ Spherical distributed sampling [Nguyen, et al., ICLR 2020]



The von Mises-Fisher distribution and the power spherical distribution illustrated in [Nguyen, et al., ICLR 2020].

[Nguyen, et al., ICLR 2020] Nguyen, K., Nguyen, S., Ho, N., Pham, T., & Bui, H. Improving Relational Regularized Autoencoders with Spherical Sliced Fused Gromov Wasserstein. ICLR 2020.

# Sampling Strategies for Sliced (F)GW Distance

- ▶ Projection pursuit driven by diversity [Meng, et al., NeurIPS 2019]
  - ▶ Given current projections, finding the new projection confusing a discriminator.
  - ▶ Connect with the idea of GAN — the projections aim at cheating the discriminator classifying point sources.

[Meng, et al., NeurIPS 2019] Meng, C., Ke, Y., Zhang, J., Zhang, M., Zhong, W., & Ma, P. Large-scale optimal transport map estimation using projection pursuit. NeurIPS 2019.

## Hierarchical Fused Gromov-Wasserstein Distance

For  $(\mathcal{X}, d_X, \mu_X)$  and  $(\mathcal{Y}, d_Y, \mu_Y)$ , suppose that

$$\begin{aligned}\mu_X &= \sum_{i=1}^I a_i p_i, \quad \mathbf{a} \in \Delta^{I-1}, \quad p_i = \mathcal{N}(\mathbf{u}_i, \Sigma_i), \\ \mu_Y &= \sum_{j=1}^J b_j q_j, \quad \mathbf{b} \in \Delta^{J-1}, \quad q_j = \mathcal{N}(\mathbf{v}_j, \Phi_j).\end{aligned}\tag{37}$$

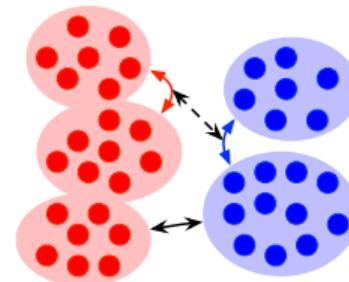
# Hierarchical Fused Gromov-Wasserstein Distance

For  $(\mathcal{X}, d_X, \mu_X)$  and  $(\mathcal{Y}, d_Y, \mu_Y)$ , suppose that

$$\begin{aligned}\mu_X &= \sum_{i=1}^I a_i p_i, \quad \mathbf{a} \in \Delta^{I-1}, \quad p_i = \mathcal{N}(\mathbf{u}_i, \Sigma_i), \\ \mu_Y &= \sum_{j=1}^J b_j q_j, \quad \mathbf{b} \in \Delta^{J-1}, \quad q_j = \mathcal{N}(\mathbf{v}_j, \Phi_j).\end{aligned}\tag{37}$$

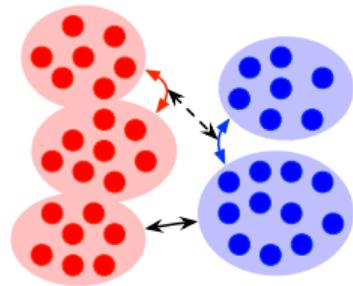
## Principle:

- ▶ Consider the fused GW distance between the GMM components in different mm-spaces.



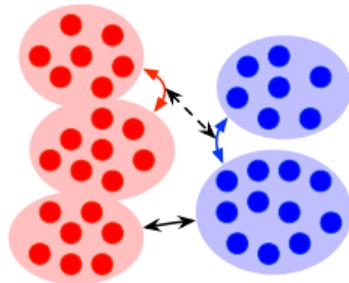
[Xu, et al., ICML 2020] Xu, H., Luo, D., Henao, R., Shah, S., & Carin, L. Learning autoencoders with relational regularization. ICML 2020.

## Hierarchical Fused Gromov-Wasserstein Distance



$$\begin{aligned} & d_{hfgw}(\mathcal{X}_{d_X, \mu_X}, \mathcal{Y}_{d_Y, \mu_Y}; \beta) \\ &= d_{fgw}(\mathbf{a}, \mathbf{b}; \beta) \text{ with Wasserstein ground distance} \\ &= \min_{\mathbf{T} \in \Pi(\mathbf{a}, \mathbf{b})} (1 - \beta) \sum_{i,j} d_w(p_i, q_j) T_{ij} + \\ & \quad \beta \sum_{i,i',j,j'} |d_w(p_i, p_{i'}) - d_w(q_j, q_{j'})|^2 T_{ij} T_{i'j'}. \end{aligned} \tag{38}$$

## Hierarchical Fused Gromov-Wasserstein Distance



$$\begin{aligned} d_{hfgw}(\mathcal{X}_{d_X, \mu_X}, \mathcal{Y}_{d_Y, \mu_Y}; \beta) \\ = & d_{fgw}(\mathbf{a}, \mathbf{b}; \beta) \text{ with Wasserstein ground distance} \\ = & \min_{\mathbf{T} \in \Pi(\mathbf{a}, \mathbf{b})} (1 - \beta) \sum_{i,j} d_w(p_i, q_j) T_{ij} + \\ & \beta \sum_{i,i',j,j'} |d_w(p_i, p_{i'}) - d_w(q_j, q_{j'})|^2 T_{ij} T_{i'j'}. \end{aligned} \tag{38}$$

- ▶ Closed form:  $d_w(p, q) = \|\mathbf{u}_p - \mathbf{u}_q\|_2^2 + \text{trace}(\Sigma_p + \Sigma_q - 2(\Sigma_p^{\frac{1}{2}} \Sigma_q \Sigma_p^{\frac{1}{2}})^{\frac{1}{2}})$  [Takatsu, 2011].
- ▶  $\mathbf{T} = [T_{ij}]$  is the optimal transport between the components of the two GMMs.
- ▶ Complexity:  $\mathcal{O}(I^2J + J^2I)$ .

[Takatsu, 2011] Takatsu, A. Wasserstein geometry of Gaussian measures. Osaka Journal of Mathematics, 48(4), 1005-1026, 2011.

# Partial and Unbalanced Gromov-Wasserstein Distance

## Motivation:

- ▶ For  $(\mathcal{X}, d_X, \mu_X)$  and  $(\mathcal{Y}, d_Y, \mu_Y)$ , the  $\mu_X$  and  $\mu_Y$  are prior distributions.

# Partial and Unbalanced Gromov-Wasserstein Distance

## Motivation:

- ▶ For  $(\mathcal{X}, d_X, \mu_X)$  and  $(\mathcal{Y}, d_Y, \mu_Y)$ , the  $\mu_X$  and  $\mu_Y$  are prior distributions.

Given  $\mathbf{X} \subset \mathcal{X}$  and  $\mathbf{Y} \subset \mathcal{Y}$ , partial Gromov-Wasserstein Distance [Chapel, et al., NeurIPS 2020]:

$$d_{pgw}(\mathbf{X}, \mathbf{Y}) = \min_{\mathbf{T} \in \Pi_s(\boldsymbol{\mu}_X, \boldsymbol{\mu}_Y)} \sum_{i,i'=1}^M \sum_{j,j'=1}^N r(x_i, x_{i'}, y_j, y_{j'}) T_{ij} T_{i'j'} \quad (39)$$

where  $\Pi_s(\boldsymbol{\mu}_X, \boldsymbol{\mu}_Y) = \{\mathbf{T} > \mathbf{0} | \mathbf{T}\mathbf{1}_N \leq \boldsymbol{\mu}_X, \mathbf{T}^\top \mathbf{1}_M \leq \boldsymbol{\mu}_Y, \mathbf{1}_M^\top \mathbf{T}\mathbf{1}_N = s\}$ .

# Partial and Unbalanced Gromov-Wasserstein Distance

## Motivation:

- ▶ For  $(\mathcal{X}, d_X, \mu_X)$  and  $(\mathcal{Y}, d_Y, \mu_Y)$ , the  $\mu_X$  and  $\mu_Y$  are prior distributions.

Given  $\mathbf{X} \subset \mathcal{X}$  and  $\mathbf{Y} \subset \mathcal{Y}$ , partial Gromov-Wasserstein Distance [Chapel, et al., NeurIPS 2020]:

$$d_{pgw}(\mathbf{X}, \mathbf{Y}) = \min_{\mathbf{T} \in \Pi_s(\mu_X, \mu_Y)} \sum_{i,i'=1}^M \sum_{j,j'=1}^N r(x_i, x_{i'}, y_j, y_{j'}) T_{ij} T_{i'j'} \quad (39)$$

where  $\Pi_s(\mu_X, \mu_Y) = \{\mathbf{T} > \mathbf{0} | \mathbf{T}\mathbf{1}_N \leq \mu_X, \mathbf{T}^\top \mathbf{1}_M \leq \mu_Y, \mathbf{1}_M^\top \mathbf{T}\mathbf{1}_N = s\}$ .

- ▶ Introduce virtual points  $\tilde{\mathbf{X}} = \mathbf{X} \cup x_v$  and  $\tilde{\mathbf{Y}} = \mathbf{Y} \cup y_v$ , and reformulate the problem to a classic GW problem:

$$\tilde{\mathbf{T}} \in \Pi(\tilde{\mu}_X, \tilde{\mu}_Y), \quad \text{where } \tilde{\mu}_X = [\mu_X; 1 - \|\mu_X\|_1], \quad \tilde{\mu}_Y = [\mu_Y; 1 - \|\mu_Y\|_1]. \quad (40)$$

[Chapel, et al., NeurIPS 2020] Chapel, L., Alaya, M. Z., & Gasso, G. Partial optimal transport with applications on positive-unlabeled learning. NeurIPS 2020.

## Partial and Unbalanced Gromov-Wasserstein Distance

Unbalanced Gromov-Wasserstein Distance — Regularize the marginal distributions of optimal transport directly [Séjourné, et al., NeurIPS 2021]:

$$d_{ugw}(\mathcal{X}_{d_X, \mu_X}, \mathcal{Y}_{d_Y, \mu_Y}) = \inf_{\pi} \int_{(x, x', y, y') \in \mathcal{X}^2 \times \mathcal{Y}^2} r(x, x', y, y') \pi(x, y) \pi(x', y') dx dx' dy dy' + D_\psi(\pi_X \otimes \pi_X \| \mu_X \otimes \mu_X) + D_\psi(\pi_Y \otimes \pi_Y \| \mu_Y \otimes \mu_Y). \quad (41)$$

- ▶  $\pi_X$  and  $\pi_Y$  are marginals of  $\pi$ .
- ▶ The quadratic  $\psi$ -divergence (which can be specialized as KL-divergence, Total-variance) makes  $d_{ugw}$  2-homogeneous.
- ▶ When  $\psi$ -divergence = TV,  $d_{ugw} \Rightarrow d_{pgw}$ .
- ▶ An alternating Sinkhorn algorithm can be used to solve this problem.

[Séjourné, et al., NeurIPS 2021] Séjourné, T., Vialard, F. X., & Peyré, G. The unbalanced Gromov Wasserstein distance: Conic formulation and relaxation. NeurIPS 2021.

## Application Scenarios of GWD and Its Variances

- ▶ **Input:** Relation matrices and marginal distributions;
- ▶ **Output:** OT matrix
  - ▶ (Fused) GWD (CG, ProxGrad, BADMM, ...)
  - ▶ Acceleration: S-GWL (Divide and conquer), Low-rank structured GW.
  - ▶ Unreliable distribution: PGW or UGW

# Application Scenarios of GWD and Its Variances

- ▶ **Input:** Relation matrices and marginal distributions;
- ▶ **Output:** OT matrix
  - ▶ (Fused) GWD (CG, ProxGrad, BADMM, ...)
  - ▶ Acceleration: S-GWL (Divide and conquer), Low-rank structured GW.
  - ▶ Unreliable distribution: PGW or UGW
- ▶ **Input:** Point clouds and marginal distributions
- ▶ **Output:** (Fused) GW distance or its surrogates.
  - ▶ (Fused) GWD (CG, ProxGrad, BADMM, ...), Sliced (Fused) GWD.
  - ▶ Acceleration: S-GWL (Divide and conquer), Low-rank structured GW.
  - ▶ With clustering structure: Tree-structured GWD, Hierarchical (Fused) GWD.
  - ▶ Unreliable distribution: PGW or UGW

## Summary

- ▶ GW distance can be calculated through various iterative optimization methods.
- ▶ Many strategies are applicable to accelerate and/or approximate the GW distance.
- ▶ With simple extensions, GW distance becomes applicable to complicated scenarios, e.g., attributed graph, point clouds with clustering structures, etc.

Next,

- ▶ **Beyond graph matching and partitioning, apply GW distance and its variants to more machine learning tasks.**
- ▶ **Discuss ongoing research and open problems.**

5-min break for Q & A

# Outline

## Part 1 Introduction of Gromov-Wasserstein Distance

- ▶ Preliminary and basic concepts
- ▶ Connections to structured data modeling and analysis

## Part 2 The Computation of Gromov-Wasserstein Distance

- ▶ Typical optimization algorithms
- ▶ The variants of Gromov-Wasserstein distance

## Part 3 GW-based Models and Their ML Applications

- ▶ Generative modeling
- ▶ Graph representation
- ▶ Graph generation

# Overview of GWL and Its Applications

## Methodology

- ▶ Take GW distance and/or its variants as objective functions or regularizers.
- ▶ Leverage the optimization of GW distance to construct computational modules of learning systems.

# Overview of GWL and Its Applications

## Methodology

- ▶ Take GW distance and/or its variants as objective functions or regularizers.
- ▶ Leverage the optimization of GW distance to construct computational modules of learning systems.

## Applications

- ▶ Graph matching and partitioning (See Part 1)
- ▶ Generative modeling: Wasserstein Autoencoder with Relational Regularization
- ▶ Graph representation: GW factorization model
- ▶ Graph generation: GW barycenter-based Graphon estimation

# Generative Modeling: Relationally-Regularized Wasserstein Autoencoder

**Wasserstein autoencoder (WAE)** [Tolstikhin, et al., ICLR 2018] fits the model distribution  $p_G$  by minimizing its Wasserstein distance to the data distribution  $p_x$ :

$$\min D_w(p_x, p_G) \implies \min_{G, Q} \underbrace{\mathbb{E}_{p_x} \mathbb{E}_{q_{z|x;Q}} [d(x, G(z))]}_{\text{reconstruction loss}} + \gamma \underbrace{D(\mathbb{E}_{p_x} [q_{z|x;Q}], p_z)}_{\text{distance(posterior, prior)}}, \quad (42)$$

- ▶  $G : \mathcal{Z} \mapsto \mathcal{X}$  is the target generative model (**decoder**).
- ▶  $q_{z|x;Q}$  is the posterior of  $z$  given  $x$ , parameterized by an **encoder**  $Q : \mathcal{X} \mapsto \mathcal{Z}$ .
- ▶  $q_{z;Q} = \mathbb{E}_{p_x} [q_{z|x;Q}]$  is the expectation of the posterior distributions.
- ▶  $p_z$  is the prior of  $z$ .

[Tolstikhin, et al., ICLR 2018] Tolstikhin, I., Bousquet, O., Gelly, S., & Schoelkopf, B. (2018, February). Wasserstein Auto-Encoders. ICLR 2018.

## Autoencoders with Predefined or Learnable Prior

$$\min_{G,Q} \underbrace{\mathbb{E}_{p_x} \mathbb{E}_{q_{z|x;Q}} [d(x, G(z))]}_{\text{reconstruction loss}} + \underbrace{\gamma D(\overbrace{\mathbb{E}_{p_x} [q_{z|x;Q}]}^{q_{z;Q}}, \textcolor{red}{p_z})}_{\text{distance(posterior, prior)}}, \quad (43)$$

Predefined prior, e.g.,  $p_z = \mathcal{N}(z; 0, I)$ .

- ▶ Over-regularization
- ▶ Difficult for conditional generation

## Autoencoders with Predefined or Learnable Prior

$$\min_{G,Q} \underbrace{\mathbb{E}_{p_x} \mathbb{E}_{q_{z|x;Q}} [d(x, G(z))]}_{\text{reconstruction loss}} + \underbrace{\gamma D(\overbrace{\mathbb{E}_{p_x} [q_{z|x;Q}]}^{q_{z;Q}}, \textcolor{red}{p_z})}_{\text{distance(posterior, prior)}}, \quad (43)$$

Predefined prior, e.g.,  $p_z = \mathcal{N}(z; 0, I)$ .

- ▶ Over-regularization
- ▶ Difficult for conditional generation

Learnable prior, e.g.,  $p_z = \frac{1}{K} \sum_{k=1}^K \mathcal{N}(z; u_k, \Sigma_k)$ .

- ▶ Under-regularization
- ▶ Sensitive to hyperparameters and initialization

## Proposed Method: Learning with A Relational Regularization

Introduce a relational regularizer based on the Gromov-Wasserstein distance between distributions:

$$\min_{G,Q,p_z \in \mathcal{P}} \mathbb{E}_{p_x} \mathbb{E}_{q_{z|x};Q} [d(x, G(z))] + \gamma \underbrace{( (1 - \beta) D(q_{z|Q}, p_z) )}_{\text{direct comparison}} + \underbrace{\beta D_{\text{gw}}(q_{z|Q}, p_z)}_{\text{relational comparison}}. \quad (44)$$

## Proposed Method: Learning with A Relational Regularization

Introduce a relational regularizer based on the Gromov-Wasserstein distance between distributions:

$$\min_{G,Q,p_z \in \mathcal{P}} \mathbb{E}_{p_x} \mathbb{E}_{q_{z|x;Q}} [d(x, G(z))] + \gamma \underbrace{((1 - \beta) D(q_{z;Q}, p_z)}_{\text{direct comparison}} + \underbrace{\beta D_{\text{gw}}(q_{z;Q}, p_z)}_{\text{relational comparison}}. \quad (44)$$

We consider the direct comparison  $D(q_{z;Q}, p_z)$  as the Wasserstein distance, and reformulate the regularizers as **fused Gromov-Wasserstein (FGW)** distance [Titouan, et al., ICML 2019].

$$\min_{G,Q,p_z \in \mathcal{P}} \mathbb{E}_{p_x} \mathbb{E}_{q_{z|x;Q}} [d(x, G(z))] + \gamma d_{fgw}(q_{z;Q}, p_z; \beta). \quad (45)$$

## Proposed Method: Learning with A Relational Regularization

Introduce a relational regularizer based on the Gromov-Wasserstein distance between distributions:

$$\min_{G,Q,p_z \in \mathcal{P}} \mathbb{E}_{p_x} \mathbb{E}_{q_{z|x;Q}} [d(x, G(z))] + \gamma \underbrace{((1 - \beta) D(q_{z;Q}, p_z)}_{\text{direct comparison}} + \underbrace{\beta D_{\text{gw}}(q_{z;Q}, p_z)}_{\text{relational comparison}}. \quad (44)$$

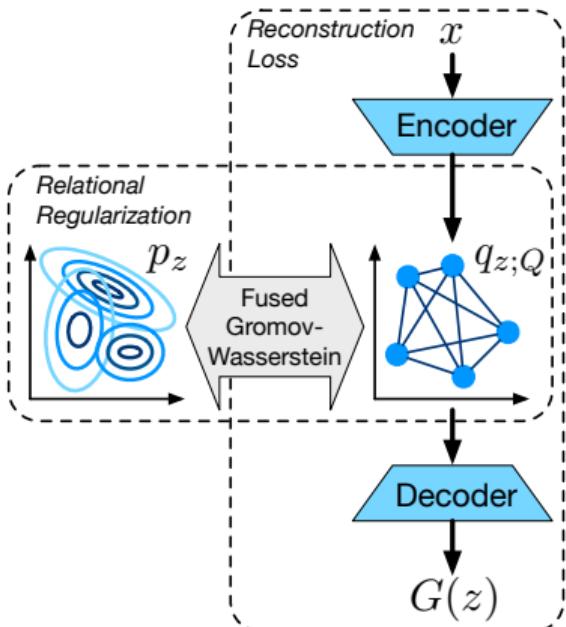
We consider the direct comparison  $D(q_{z;Q}, p_z)$  as the Wasserstein distance, and reformulate the regularizers as **fused Gromov-Wasserstein (FGW)** distance [Titouan, et al., ICML 2019].

$$\min_{G,Q,p_z \in \mathcal{P}} \mathbb{E}_{p_x} \mathbb{E}_{q_{z|x;Q}} [d(x, G(z))] + \gamma d_{fgw}(q_{z;Q}, p_z; \beta). \quad (45)$$

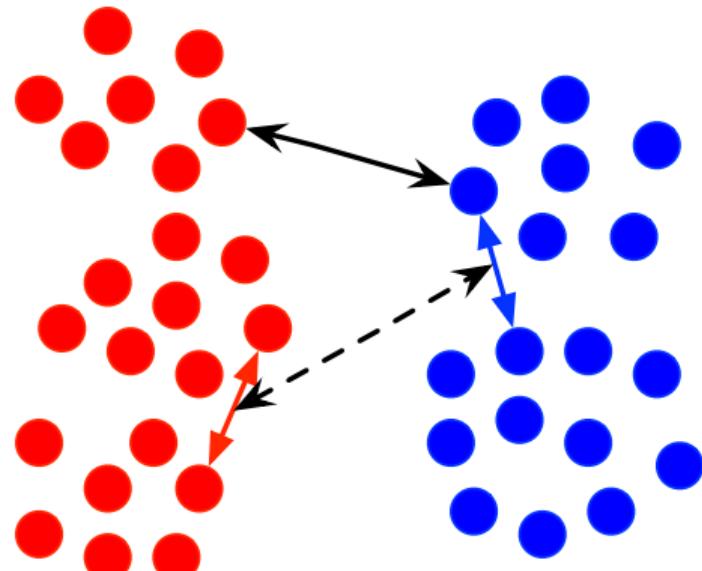
- ▶ A shared optimal transport  $\pi$  encourages the consistency of the relations between  $q_{z;Q}$  and  $p_z$  under different criteria.
- ▶ Minimization of an upper bound.

# Relationally-Regularized Wasserstein Autoencoder

$$\min_{G, Q, p_z} \mathbb{E}_{p_x} \mathbb{E}_{q_{z|x;Q}} [d(x, G(z))] + \gamma d_{\text{fgw}}(q_{z;Q}, p_z; \beta), \text{ where } p_z = \frac{1}{K} \sum_{k=1}^K \mathcal{N}(z; u_k, \Sigma_k). \quad (46)$$



The scheme of the method



The principle of fused GW distance.

## Two Computationally-friendly Implementations

### Probabilistic Autoencoder + Hierarchical FGW Distance:

- ▶ Recall that the prior  $p_z$  is a GMM.
- ▶ The encoder  $Q$  outputs the mean and the log(variance) of  $q_{z|x;Q}$  given a sample  $x$ . Given a batch of  $x$ 's, the expected posterior  $q_{z;Q}$  becomes a GMM as well.

## Two Computationally-friendly Implementations

### Probabilistic Autoencoder + Hierarchical FGW Distance:

- ▶ Recall that the prior  $p_z$  is a GMM.
- ▶ The encoder  $Q$  outputs the mean and the log(variance) of  $q_{z|x;Q}$  given a sample  $x$ . Given a batch of  $x$ 's, the expected posterior  $q_{z;Q}$  becomes a GMM as well.

$$d_{hfgw}(p, q; \beta) = \min_{T \in \Pi(a, b)} \langle D_{pq} - 2\beta D_p T D_q^\top, T \rangle, \quad (47)$$

- ▶ We solve (47) efficiently by the proximal gradient method in [Xu, et al., ICML 2019].
- ▶ The computational complexity is  $\mathcal{O}(N^2K + K^2N)$ .
- ▶ For each batch, we first learn the  $T$ , and then update autoencoders with fixed  $T$  by backpropagation.

## Two Computationally-friendly Implementations

### Deterministic Autoencoder + Sliced FGW Distance:

- ▶ The encoder  $Q$  outputs  $N$  samples of  $q_{z;Q}$  directly.
- ▶ For the GMM prior  $p_z$ , we can generate  $N$  (differentiable) samples by reparameterization trick.

## Two Computationally-friendly Implementations

### Deterministic Autoencoder + Sliced FGW Distance:

- ▶ The encoder  $Q$  outputs  $N$  samples of  $q_{z;Q}$  directly.
- ▶ For the GMM prior  $p_z$ , we can generate  $N$  (differentiable) samples by reparameterization trick.

$$\hat{d}_{\text{sfgw}}(p_x, p_y; \beta) = \frac{1}{L} \sum_{l=1}^L \min_{\sigma \in \mathcal{P}_N} \left( \frac{1-\beta}{N} \sum_{i=1}^N (x_{i,\theta_l} - y_{\sigma(i),\theta_l})^2 + \frac{\beta}{N} \sum_{i,j=1}^N ((x_{i,\theta_l} - x_{j,\theta_l})^2 - (y_{\sigma(i),\theta_l} - y_{\sigma(j),\theta_l})^2)^2 \right). \quad (48)$$

- ▶ For projected latent codes, sort them and compute the 1D FGW distance by their identity permutation or anti-identity permutation.
- ▶ The average of all the 1D FGW distances leads to the SFGW distance.
- ▶ The computational complexity is  $\mathcal{O}(L(\underbrace{Nd}_{\text{proj.}} + \underbrace{N \log N}_{\text{sort}} + \underbrace{N^2}_{\text{1D fgw}}))$ .

## Comparisons with Other Autoencoders

Table: Comparisons for different autoencoders

Method	$Q : \mathcal{X} \mapsto \mathcal{Z}$	$p_z$	$d(q_z; Q, p_z)$
VAE	Probabilistic	$\mathcal{N}(z; 0, I)$	KL
WAE	Deterministic	$\mathcal{N}(z; 0, I)$	MMD/GAN
SWAE	Deterministic	$\mathcal{N}(z; 0, I)$	$d_w$
GMVAE	Probabilistic	$\frac{1}{K} \sum_k \mathcal{N}(z; u_k, \Sigma_k)$	KL
VampPrior	Probabilistic	$\frac{1}{K} \sum_k \mathcal{N}(z; Q(x_k))$	KL
<b>Our RAE</b>	Probabilistic	$\frac{1}{K} \sum_k \mathcal{N}(z; u_k, \Sigma_k)$	$d_{\text{hfgw}}$
	Deterministic		$\hat{d}_{\text{sfgw}}$

## Experiments: Random Image Generation



Figure: WAE v.s. RAE

## Experiments: Conditional Image Generation



(a) VampPrior



(b) GMVAE



(c) Probabilistic RAE



(d) Deterministic RAE

Figure: Comparisons on conditional digit generation.

## Experiments: Conditional Image Generation



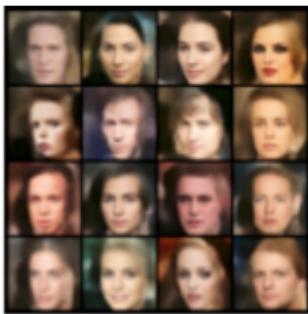
Front Face  
Short Hair



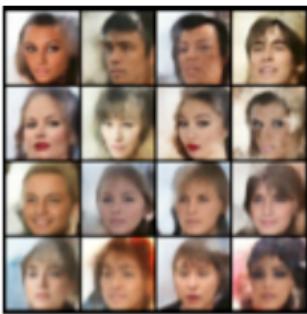
Front Face  
Long Hair



Front Face  
Dark Background



Left Cheek  
Dark Background



Left Cheek  
Light Background



Right Cheek  
Dark Background



Right Cheek  
Light Background



White Hair  
Blonde Hair



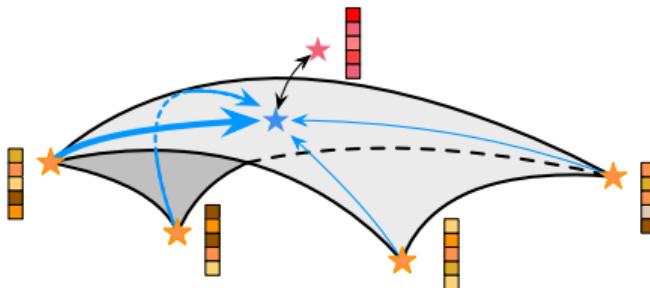
Dark Skin



Smile With Teeth

Figure: Deterministic RAE for conditional face generation.

## GW-based Graph Representation: GW Factorization

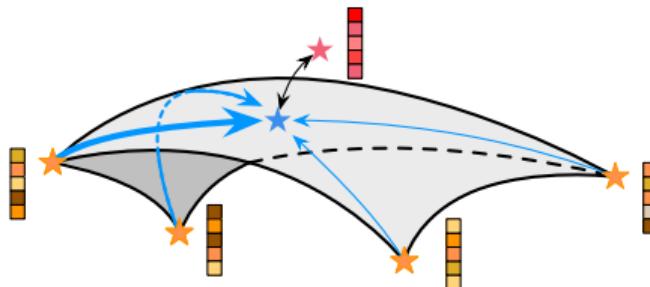


### Traditional Factorization Models:

- Given  $\{y_1, \dots, y_I\}$ , we would like to learn the basis  $A = [a_1, \dots, a_K]$  and the coefficient vector (representation)  $\lambda_i$  for each  $y_i$ .

$$\min_{\{A, \lambda_{1:I}\} \in \Omega} \sum_{i=1}^I d_{\text{loss}}(A\lambda_i, y_i). \quad (49)$$

# GW-based Graph Representation: GW Factorization



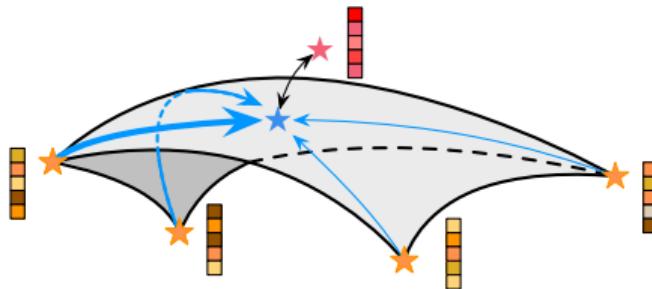
## Traditional Factorization Models:

- Given  $\{\mathbf{y}_1, \dots, \mathbf{y}_I\}$ , we would like to learn the basis  $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_K]$  and the coefficient vector (representation)  $\boldsymbol{\lambda}_i$  for each  $\mathbf{y}_i$ .

$$\min_{\{\mathbf{A}, \boldsymbol{\lambda}_{1:I}\} \in \Omega} \sum_{i=1}^I d_{\text{loss}}(\mathbf{A}\boldsymbol{\lambda}_i, \mathbf{y}_i). \quad (49)$$

- PCA:  $d_{\text{loss}} = \ell_2$ . (MSE)
- Robust PCA:  $d_{\text{loss}} = \ell_1$ . (MAE)
- NMF:  $d_{\text{loss}} = \ell_2$ ,  $\Omega = \text{Nonnegativeness}$ .
- LDA:  $d_{\text{loss}} = \text{KL}$ ,  $\Omega = \text{Simplex}$ .
- Wasserstein dictionary learning:  $d_{\text{loss}} = \text{Wasserstein}$ ,  $\Omega = \text{Simplex}$ .

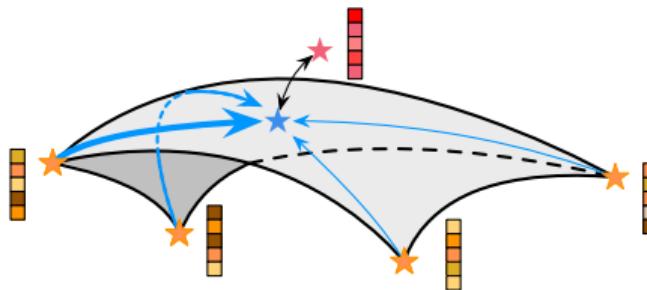
## GW-based Graph Representation: GW Factorization



- When  $\Omega = \text{Simplex}$ , we have

$$\mathbf{A}\boldsymbol{\lambda} = \underbrace{\arg \min_{\mathbf{y}} \sum_{k=1}^K \lambda_k \|\mathbf{y} - \mathbf{a}_k\|_2^2}_{\text{Euclidean barycenter}} = \mathbf{b}(\mathbf{A}, \boldsymbol{\lambda}; \ell_2). \quad (50)$$

## GW-based Graph Representation: GW Factorization



- When  $\Omega = \text{Simplex}$ , we have

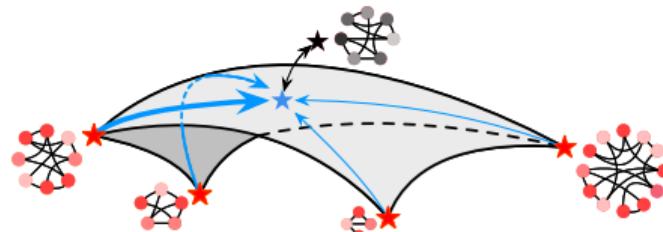
$$\mathbf{A}\boldsymbol{\lambda} = \underbrace{\arg \min_{\mathbf{y}} \sum_{k=1}^K \lambda_k \|\mathbf{y} - \mathbf{a}_k\|_2^2}_{\text{Euclidean barycenter}} = \mathbf{b}(\mathbf{A}, \boldsymbol{\lambda}; \ell_2). \quad (50)$$

- Accordingly, a generalized factorization model can be written as

$$\min_{\{\mathbf{A}, \boldsymbol{\lambda}_{1:I}\} \in \Omega} \sum_{i=1}^I d_{\text{loss}}(\mathbf{b}(\mathbf{A}, \boldsymbol{\lambda}_i; \underbrace{d_b}_{\text{b's metric}}), \mathbf{y}_i). \quad (51)$$

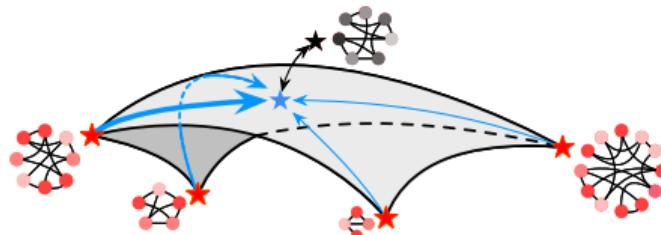
## GW-based Graph Representation: GW Factorization

GW factorization [Xu, AAAI 2020]: Learning interpretable factorization model to represent unaligned graphs.



## GW-based Graph Representation: GW Factorization

GW factorization [Xu, AAAI 2020]: Learning interpretable factorization model to represent unaligned graphs.



- ▶ Estimate each graph by a GW barycenter graph:

$$B_{gw}(\mathbf{U}_{1:K}, \boldsymbol{\lambda}) := \arg \min_B \sum_{k=1}^K \boldsymbol{\lambda}_k d_{gw}(B, G_k(\mathbf{U}_k)). \quad (52)$$

- ▶  $\{G_k(\mathbf{U}_k)\}_{k=1}^K$ : a set of graph factors.
- ▶  $\boldsymbol{\lambda}_{1:I} = \{\boldsymbol{\lambda}_i \in \Delta^{K-1}\}_{i=1}^I$ : the coefficients of the graph factors corresponding to  $\{G_i\}_{i=1}^I$  (The representations of the observed graphs).

[Xu, AAAI 2020] Xu, H. Gromov-Wasserstein factorization models for graph clustering. AAAI 2020.

## GW-based Graph Representation: GW Factorization

Learning task:

$$\min_{\mathbf{1} \geq \mathbf{U}_{1:K} \geq \mathbf{0}, \boldsymbol{\lambda}_{1:I} \in \Delta^{K-1}} \sum_{i=1}^I d_{\text{loss}}(B_{gw}(\mathbf{U}_{1:K}, \boldsymbol{\lambda}_i), G_i). \quad (53)$$

## GW-based Graph Representation: GW Factorization

Learning task:

$$\min_{\mathbf{1} \geq \mathbf{U}_{1:K} \geq \mathbf{0}, \lambda_{1:I} \in \Delta^{K-1}} \sum_{i=1}^I d_{\text{loss}}(B_{gw}(\mathbf{U}_{1:K}, \lambda_i), G_i). \quad (53)$$

Reparameterize the problem to an unconstrained optimization problem:

$$\min_{\mathbf{V}_{1:K}, \mathbf{z}_{1:I}} \sum_{i=1}^I d_{\text{loss}}(B_{gw}(\sigma(\mathbf{V}_{1:K}), \text{softmax}(\mathbf{z}_i)), G_i). \quad (54)$$

# GW-based Graph Representation: GW Factorization

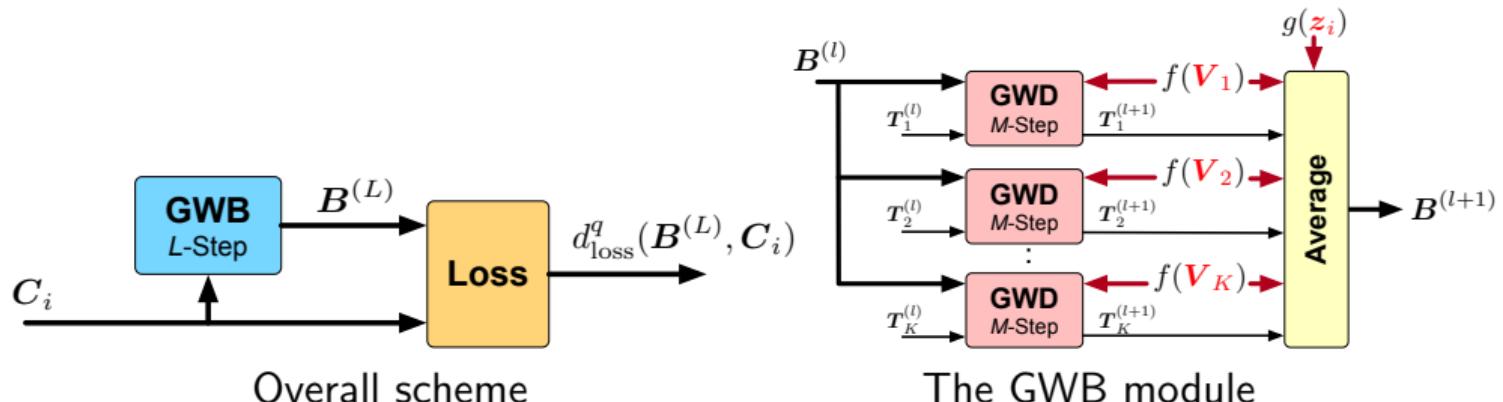
Learning task:

$$\min_{\mathbf{1} \geq \mathbf{U}_{1:K} \geq \mathbf{0}, \lambda_{1:I} \in \Delta^{K-1}} \sum_{i=1}^I d_{\text{loss}}(B_{gw}(\mathbf{U}_{1:K}, \lambda_i), G_i). \quad (53)$$

Reparameterize the problem to an unconstrained optimization problem:

$$\min_{\mathbf{V}_{1:K}, \mathbf{z}_{1:I}} \sum_{i=1}^I d_{\text{loss}}(B_{gw}(\sigma(\mathbf{V}_{1:K}), \text{softmax}(\mathbf{z}_i)), G_i). \quad (54)$$

(1) Compute optimal transport matrices by GWD modules; (2) Fix the OT matrices and learn the model parameters via SGD:



# GWF v.s. Graph Dictionary Learning (GDL)

## GW Factorization

- ▶ Learn arbitrarily-sized graph factors, and reconstruct graphs via the weighted GW barycenter of the factors (Nonlinear)
- ▶ Unsupervised learning under arbitrary loss functions (MSE, Cross-Entropy, GWD, ...) because of the permutation invariance of GW barycenter.
- ▶ High computational complexity

# GWF v.s. Graph Dictionary Learning (GDL)

## GW Factorization

- ▶ Learn arbitrarily-sized graph factors, and reconstruct graphs via the weighted GW barycenter of the factors (Nonlinear)
- ▶ Unsupervised learning under arbitrary loss functions (MSE, Cross-Entropy, GWD, ...) because of the permutation invariance of GW barycenter.
- ▶ High computational complexity

## Graph Dictionary Learning [Vincent, et al., ICML 2021]

- ▶ Learn graph factors with fixed sizes, and reconstruct graphs via a linear model:

$$\hat{G} = \sum_{k=1}^K \lambda_k \mathbf{U}_k \quad (55)$$

- ▶ Unsupervised learning, but the loss function has to be GW distance:

$$\min_{\mathbf{U}_{1:\mathcal{K}}, \boldsymbol{\lambda}_{1:I}} \sum_{i=1}^I d_{gw}(\hat{G}_i, G_i). \quad (56)$$

- ▶ Low computational complexity

[Vincent, et al., ICML 2021] Vincent-Cuaz, C., Vayer, T., Flamary, R., Corneli, M., & Courty, N. Online graph dictionary learning. ICML 2021.

# GWF v.s. Random-Walk Neural Network (RWNN)

## GW Factorization

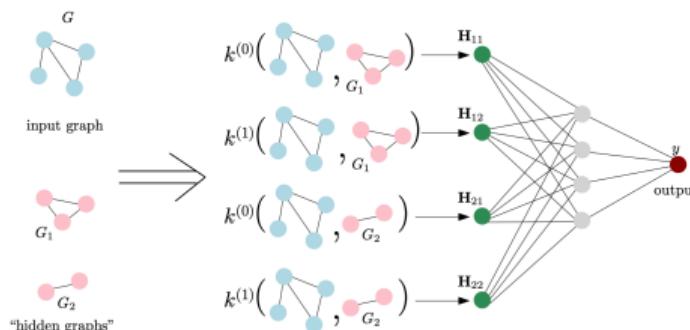
- ▶ Consider the GWD between graph factors and observed graphs, and learn the factors unsupervisedly via minimizing reconstruction errors.

# GWF v.s. Random-Walk Neural Network (RWNN)

## GW Factorization

- Consider the GWD between graph factors and observed graphs, and learn the factors unsupervisedly via minimizing reconstruction errors.

## Random-Walk Neural Network in [Nikolentzos, et al., NeurIPS 2020]



The scheme of RWNN shown in [Nikolentzos, et al., NeurIPS 2020]

- Consider the Random-Walk kernel between graph factors and observed graphs, and learn the factors supervisedly driven by graph classification labels.

[Nikolentzos, et al., NeurIPS 2020] Nikolentzos, G., & Vazirgiannis, M. Random walk graph neural networks. NeurIPS 2020.

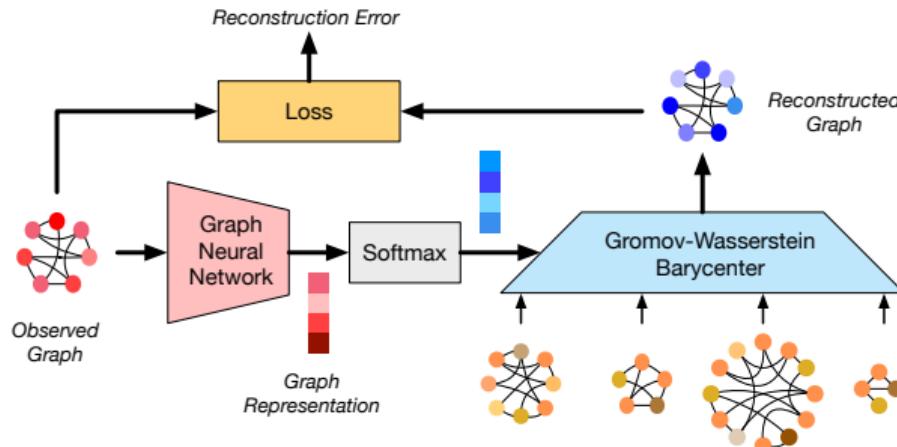
# An Inductive Graph Representation Combining GWF with GNN

Achieve an **inductive** GW factorization model by introducing graph neural networks (i.e., GWF-GNN) [Xu, et al, 2022]:

$$\min_{\mathbf{V}_{1:K}, \theta} \sum_{i=1}^I d_{gw}(B_{gw}(\sigma(\mathbf{V}_{1:K}), \text{softmax}(\text{GNN}_\theta(G_i))), G_i). \quad (57)$$

[Xu, et al, 2022] Representing graphs via Gromov-Wasserstein factorization. IEEE TPAMI, 2022  
(Accepted)

## GWF-GNN v.s. VGAE



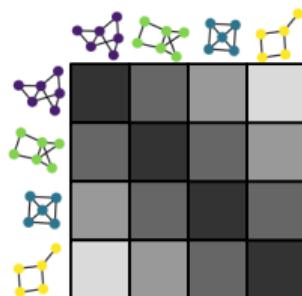
- ▶ **GWF-GNN:** A new auto-encoding architecture for learning graph-level representations.
- ▶ **VGAE in [Kipf, et al., 2016]:** An auto-encoding architecture for node-level representations.

[Kipf, et al., 2016] Kipf, T. N., & Welling, M. Variational graph auto-encoders. arXiv preprint arXiv:1611.07308, 2016.

# Experiments on Graph Clustering

- ▶ AIDS: 2,000 compounds active/inactive to anti-HIV
- ▶ PROTEIN: 1,113 enzymatic/non-enzymatic proteins

GWD Kernel



GWD+Kmeans

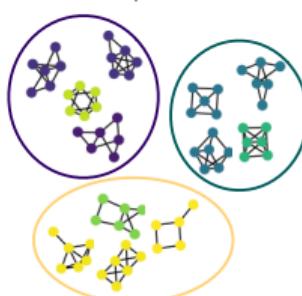
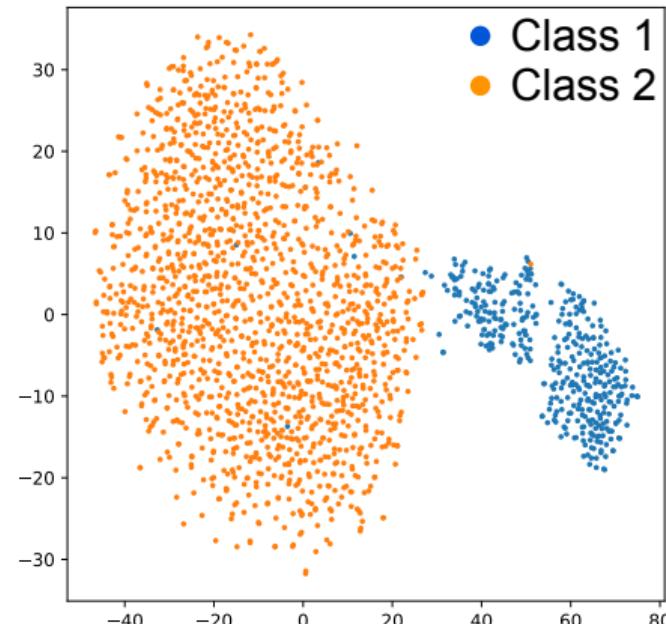


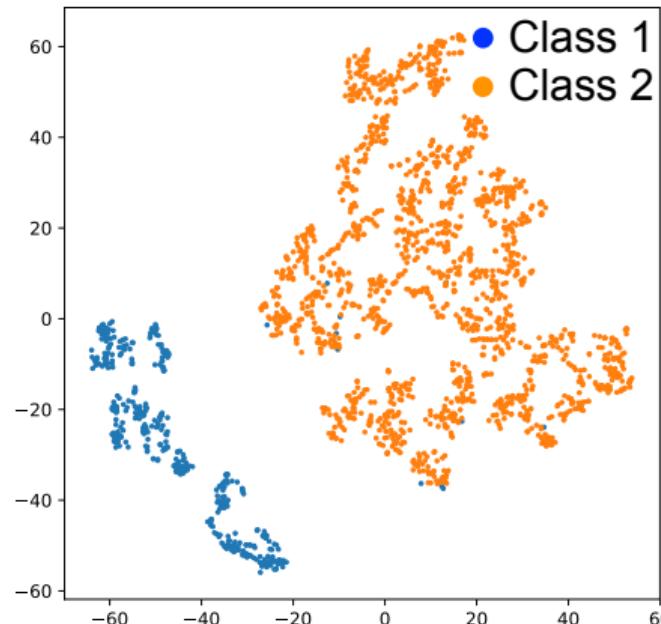
Table: Comparisons on clustering accuracy (%)

Method	AIDS	PROTEIN
GWD Kernel + SC	$91.0 \pm 0.7$	$66.4 \pm 0.8$
GWD + Kmeans	$95.2 \pm 0.9$	$64.7 \pm 1.1$
GWF + Kmeans	$99.5 \pm 0.4$	$70.7 \pm 0.7$

# Experiments on Graph Clustering



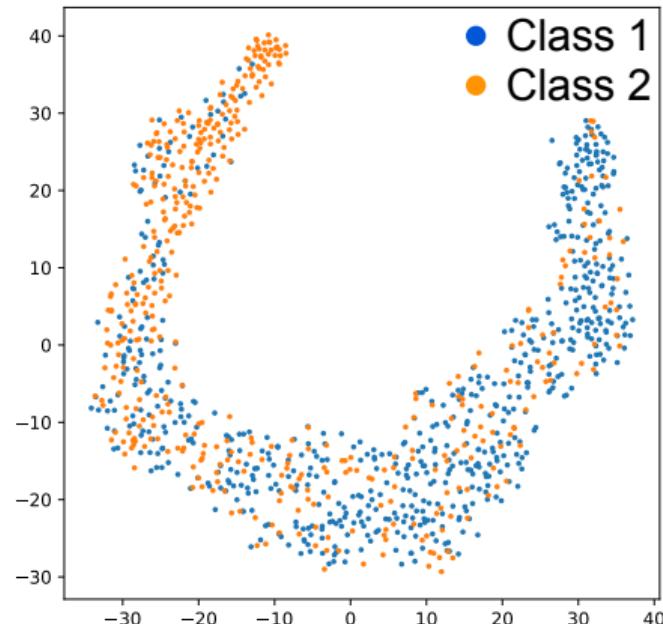
(a) AIDS (GWF)



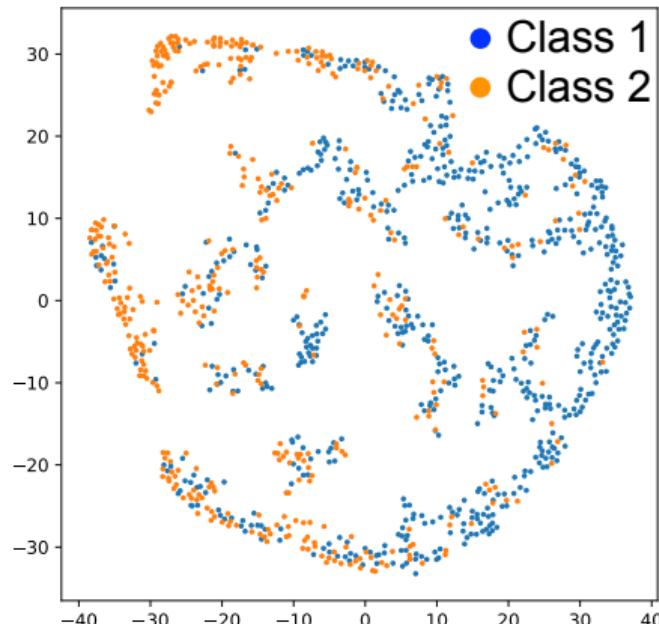
(b) AIDS (GWF-GNN)

Figure: The t-SNE visualizations of  $z_{1:I}$  of AIDS dataset derived by the GWF and GWF-GNN.

# Experiments on Graph Clustering



(a) PROTEIN (GWF)



(b) PROTEIN (GWF-GNN)

Figure: The t-SNE visualizations of  $z_{1:I}$  of Protein dataset derived by the GWF and GWF-GNN.

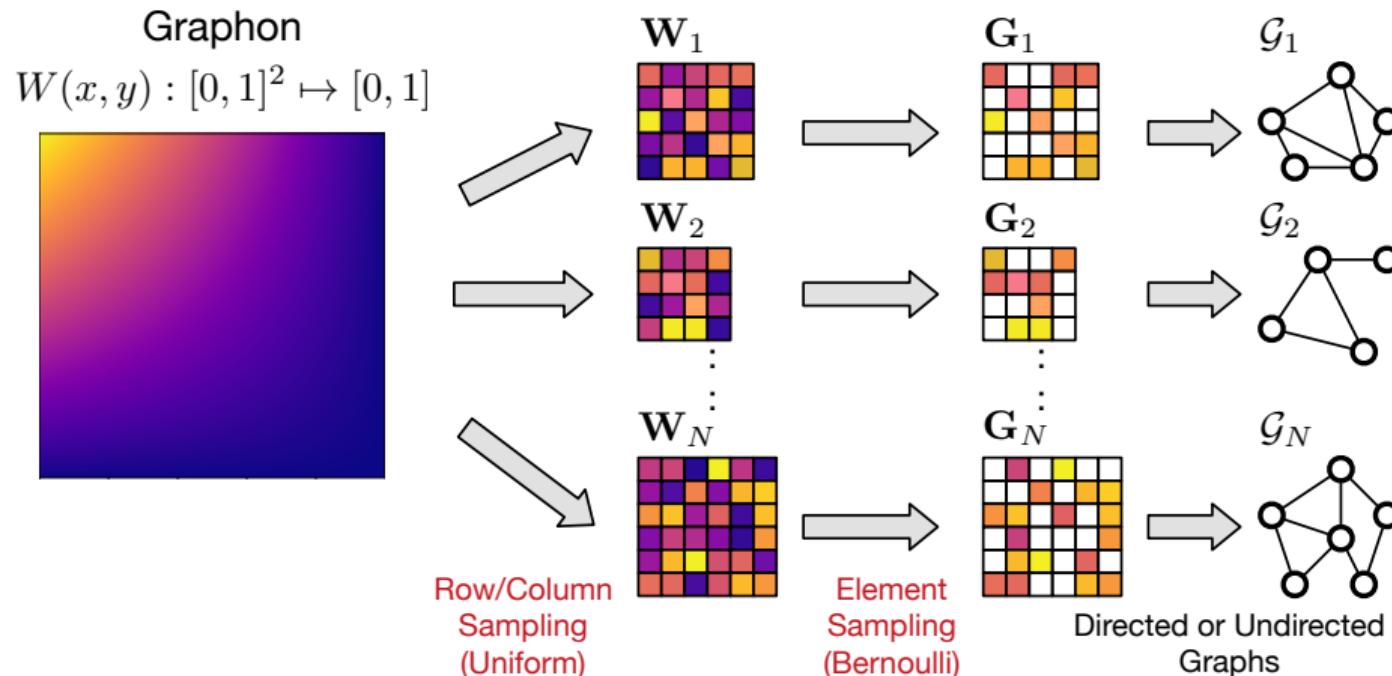
# Experiments on Graph Classification

Table: Comparison on classification accuracy (%).

	Dataset	MUTAG	PTC-MR	IMDB-B	IMDB-M
Method Category	# Graphs	188	344	1000	1500
	# Classes	2	2	2	3
	Avg. # Nodes	17.93	14.29	19.77	13
	Avg. # Edges	19.79	14.69	96.53	65.94
Graph kernel dim=512	RW	83.71	57.85	50.68	34.65
	SP	85.22	58.24	55.60	37.99
	GK	81.66	57.26	65.87	43.89
	WL	80.72	57.97	72.30	46.95
	DGK	87.44	60.08	66.96	44.55
	MLG	87.94	63.26	66.55	41.17
Recent Embedding dim=512	node2vec	72.63	58.58	-	-
	sub2vec	61.05	59.99	55.26	36.67
	graph2vec	83.15	60.17	71.1	<b>50.44</b>
GWL dim=15	GWF	82.93	62.57	63.1	40.67
	GWF-GNN	<b>90.11</b>	<b>63.58</b>	<b>73.50</b>	49.40

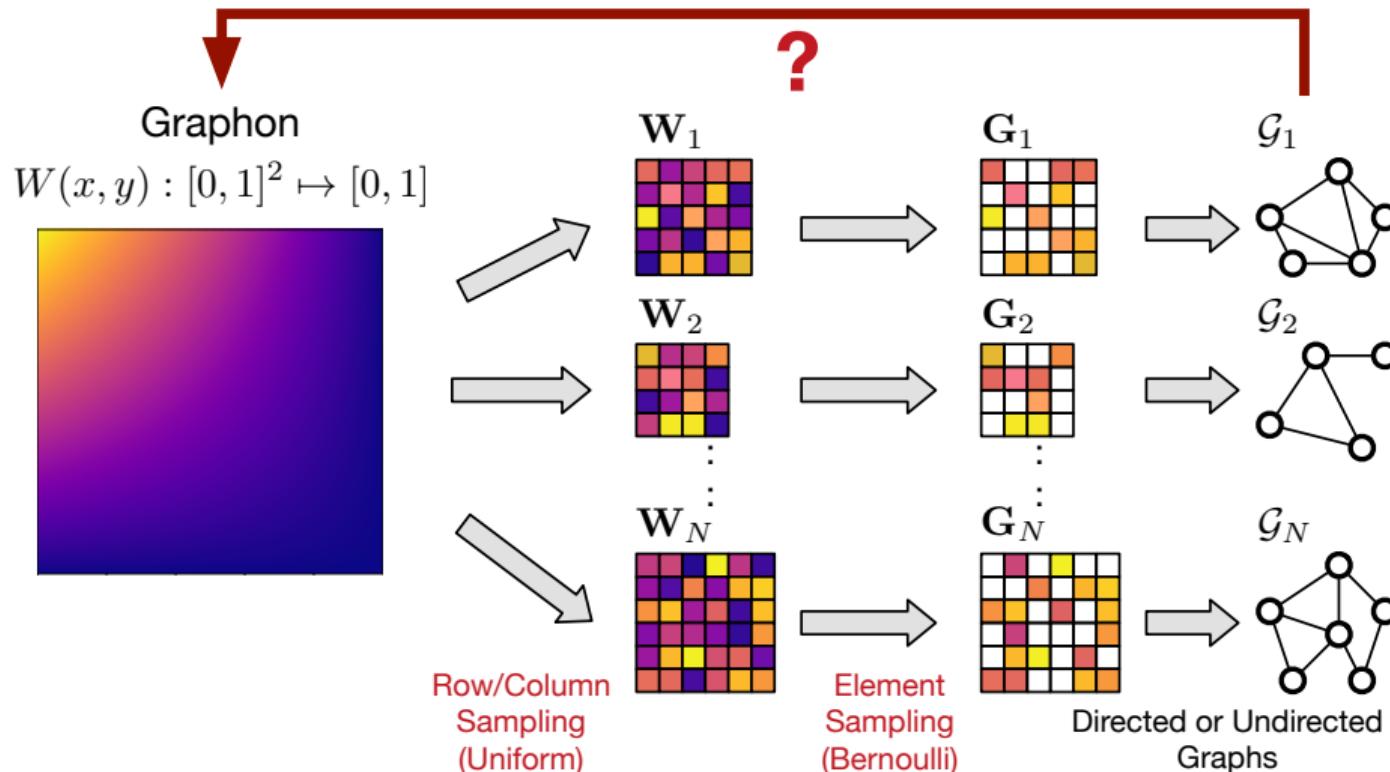
# GW-based Graph Generation: GWB-based Graphon Estimation

## Graphon: A Nonparametric Graph Generative Model



# GW-based Graph Generation: GWB-based Graphon Estimation

## Graphon: A Nonparametric Graph Generative Model



# GW-based Graph Generation: GWB-based Graphon Estimation

Two learning paradigms:

- ▶ **Traditional paradigm:** Given a single large-scale graph, estimate a graphon by a step-function and make it as smooth as possible.
  - ▶ Collecting and processing a large-scale graph are challenging.
  - ▶ The estimation and its smoothness depend on the sorting of nodes (according to their degrees).
- ▶ **The proposed paradigm:** Given a set of unaligned but small graphs, estimate a graphon by solving a GW barycenter problem [Xu, et al., AAAI 2021].
  - ▶ Reduce the difficulty on data collection and processing.
  - ▶ Robust to the challenging cases where the graph nodes are hard to sort.

[Xu, et al., AAAI 2021] Xu, H., Luo, D., Carin, L., & Zha, H. Learning Graphons via Structured Gromov-Wasserstein Barycenters. AAAI 2021.

# Oracle Estimator

**Cut norm:**  $\|W\|_{\square} := \sup_{\mathcal{X}, \mathcal{Y} \subset \Omega} \left| \int_{\mathcal{X} \times \mathcal{Y}} W(x, y) dx dy \right|,$  (58)

**Cut distance:**  $\delta_{\square}(W_1, W_2) := \inf_{\substack{\phi \in \mathcal{S}_{\Omega} \\ \text{Measure-preserved map}}} \left\| \underbrace{W_1 - W_2^{\phi}}_{\text{Residual}} \right\|_{\square}.$

Graphon      Step function

$W(x, y)$



$W_{\mathcal{P}}$

$\approx$



Weak Regularity Lemma

$$\|W - W_{\mathcal{P}}\|_{\square} \leq \frac{2}{\sqrt{\log K}} \|W\|_{L_2}$$

# Oracle Estimator

**Cut norm:**  $\|W\|_{\square} := \sup_{\mathcal{X}, \mathcal{Y} \subset \Omega} \left| \int_{\mathcal{X} \times \mathcal{Y}} W(x, y) dx dy \right|,$

**Cut distance:**  $\delta_{\square}(W_1, W_2) := \inf_{\phi \in \mathcal{S}_{\Omega}} \underbrace{\|W_1 - W_2^{\phi}\|_{\square}}_{\text{Residual}}$ . (58)

$$\delta_{\square}(W, W_O) \leq \frac{C}{|\mathcal{P}|}$$

Graphon  
 $W(x, y)$



Step function  
 $W_{\mathcal{P}}$

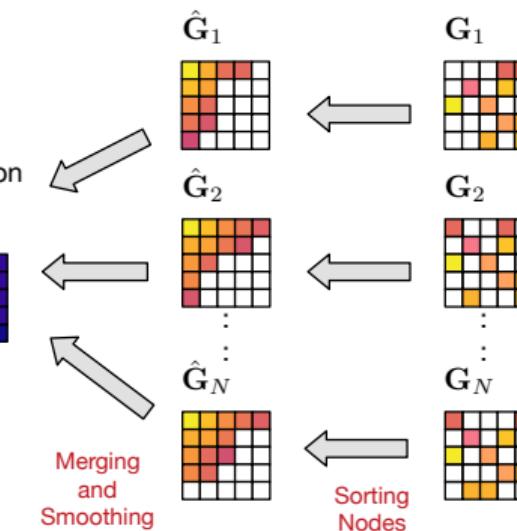


Oracle  
Estimation  
 $W_O$



Weak Regularity Lemma

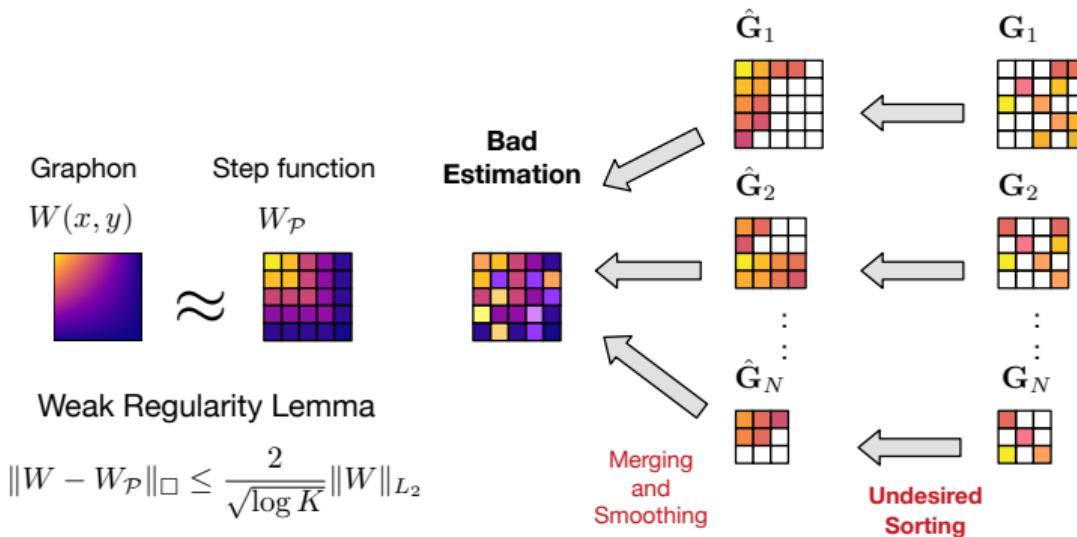
$$\|W - W_{\mathcal{P}}\|_{\square} \leq \frac{2}{\sqrt{\log K}} \|W\|_{L_2}$$



# Oracle Estimator

**Cut norm:**  $\|W\|_{\square} := \sup_{\mathcal{X}, \mathcal{Y} \subset \Omega} \left| \int_{\mathcal{X} \times \mathcal{Y}} W(x, y) dx dy \right|,$

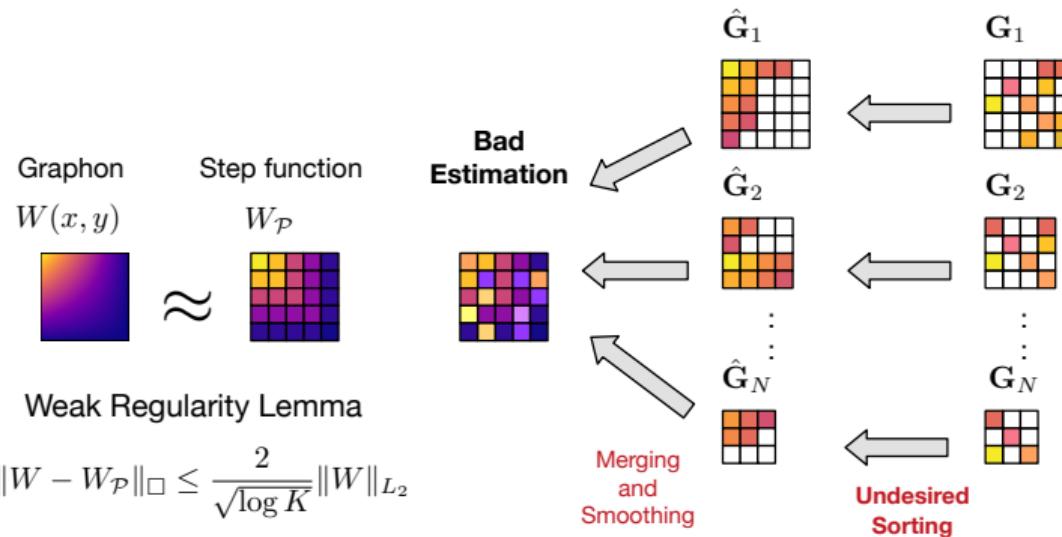
**Cut distance:**  $\delta_{\square}(W_1, W_2) := \inf_{\substack{\phi \in \mathcal{S}_{\Omega} \\ \text{Measure-preserved map}}} \|W_1 - W_2^{\phi}\|_{\square}.$  (58)



# Oracle Estimator

**Cut norm:**  $\|W\|_{\square} := \sup_{\mathcal{X}, \mathcal{Y} \subset \Omega} \left| \int_{\mathcal{X} \times \mathcal{Y}} W(x, y) dx dy \right|,$

**Cut distance:**  $\delta_{\square}(W_1, W_2) := \inf_{\substack{\phi \in \mathcal{S}_{\Omega} \\ \text{Measure-preserved map}}} \|W_1 - W_2^{\phi}\|_{\square}.$  (58)



Can we achieve graph alignment and graphon learning jointly?

## Optimizing An Upper Bound of Estimation Error

$$\delta_{\square}(W, W_{\mathcal{P}}) \leq \delta_{\square}(W, W_O) + \delta_{\square}(W_O, W_{\mathcal{P}}) \quad (\text{Triangle Inequality})$$

$$= \delta_{\square}(W, W_O) + \delta_{\square}\left(\frac{1}{M} \sum_{m=1}^M \hat{G}_m, W_{\mathcal{P}}\right)$$

# Optimizing An Upper Bound of Estimation Error

$$\delta_{\square}(W, W_{\mathcal{P}}) \leq \delta_{\square}(W, W_O) + \delta_{\square}(W_O, W_{\mathcal{P}}) \quad (\text{Triangle Inequality})$$

$$= \delta_{\square}(W, W_O) + \delta_{\square}\left(\frac{1}{M} \sum_{m=1}^M \hat{G}_m, W_{\mathcal{P}}\right)$$

$$\leq \delta_{\square}(W, W_O) + \frac{1}{M} \sum_{m=1}^M \delta_{\square}(\hat{G}_m, W_{\mathcal{P}}) \quad (\text{Triangle Inequality})$$

# Optimizing An Upper Bound of Estimation Error

$$\delta_{\square}(W, W_{\mathcal{P}}) \leq \delta_{\square}(W, W_O) + \delta_{\square}(W_O, W_{\mathcal{P}}) \quad (\text{Triangle Inequality})$$

$$= \delta_{\square}(W, W_O) + \delta_{\square}\left(\frac{1}{M} \sum_{m=1}^M \hat{G}_m, W_{\mathcal{P}}\right)$$

$$\leq \delta_{\square}(W, W_O) + \frac{1}{M} \sum_{m=1}^M \delta_{\square}(\hat{G}_m, W_{\mathcal{P}}) \quad (\text{Triangle Inequality})$$

$$= \delta_{\square}(W, W_O) + \frac{1}{M} \sum_{m=1}^M \delta_{\square}(G_m, W_{\mathcal{P}}) \quad (\delta_{\square}(G_m, \hat{G}_m) = 0)$$

# Optimizing An Upper Bound of Estimation Error

$$\begin{aligned}\delta_{\square}(W, W_{\mathcal{P}}) &\leq \delta_{\square}(W, W_O) + \delta_{\square}(W_O, W_{\mathcal{P}}) && \text{(Triangle Inequality)} \\ &= \delta_{\square}(W, W_O) + \delta_{\square}\left(\frac{1}{M} \sum_{m=1}^M \hat{G}_m, W_{\mathcal{P}}\right) \\ &\leq \delta_{\square}(W, W_O) + \frac{1}{M} \sum_{m=1}^M \delta_{\square}(\hat{G}_m, W_{\mathcal{P}}) && \text{(Triangle Inequality)} \quad (59) \\ &= \delta_{\square}(W, W_O) + \frac{1}{M} \sum_{m=1}^M \delta_{\square}(G_m, W_{\mathcal{P}}) && (\delta_{\square}(G_m, \hat{G}_m) = 0) \\ &\leq \delta_{\square}(W, W_O) + \frac{1}{M} \sum_{m=1}^M \delta_1(G_m, W_{\mathcal{P}}). && (\delta_{\square}(W, W') \leq \delta_1(W, W'))\end{aligned}$$

- **$\delta_1$  distance:**  $\delta_1(W_1, W_2) := \inf_{\phi \in \mathcal{S}_{\Omega}} \|W_1 - W_2^{\phi}\|_{L_1}$ .
- **Task:**  $\min_{W_{\mathcal{P}} \in [0,1]^{K \times K}} \sum_{m=1}^M \delta_1(G_m, W_{\mathcal{P}})$ .

# Learning Graphons as Gromov-Wasserstein Barycenters

$$\underbrace{\delta_1(W_1, W_2) = d_{\text{gw}, \ell_1}(W_1, W_2)}_{[\text{Janson, 2013}]}$$

## Learning Graphons as Gromov-Wasserstein Barycenters

$$\underbrace{\delta_1(W_1, W_2) = d_{\text{gw}, \ell_1}(W_1, W_2)}_{[\text{Janson, 2013}]} \cong \underbrace{d_{\text{gw}, \ell_2}(W_1, W_2)}_{[\text{Mémoli, 2011}]} . \quad (60)$$

# Learning Graphons as Gromov-Wasserstein Barycenters

$$\underbrace{\delta_1(W_1, W_2) = d_{\text{gw}, \ell_1}(W_1, W_2)}_{[\text{Janson, 2013}]} \cong \underbrace{d_{\text{gw}, \ell_2}(W_1, W_2)}_{[\text{Mémoli, 2011}]} . \quad (60)$$

The task becomes a **Gromov-Wasserstein barycenter (GWB)** problem:

$$\min_{W_{\mathcal{P}} \in [0,1]^{K \times K}} \sum_{m=1}^M d_{\text{gw}, \ell_2}(G_m, W_{\mathcal{P}}) . \quad (61)$$

where

$$d_{\text{gw}, \ell_2}(G_m, W_{\mathcal{P}}) = \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_1, \boldsymbol{\mu}_2)} \underbrace{\sum_{i,i',j,j'} r_{iji'j'} \overbrace{T_{ii'} T_{jj'}}^{p(r)}}_{\mathbb{E}[r]} = \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_1, \boldsymbol{\mu}_2)} -\langle \mathbf{A}_1 \mathbf{T} \mathbf{A}_2^\top, \mathbf{T} \rangle . \quad (62)$$

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

[Mémoli, 2011] Mémoli, F. GromovWasserstein distances and the metric approach to object matching. Foundations of computational mathematics, 2011.

## Two Variants for Structured GWBs

Can we achieve a smoothed graphon?

## Two Variants for Structured GWBs

Can we achieve a smoothed graphon?

- ▶ **Smoothed GW Barycenter (SGWB)** [Xu, et al., AAAI 2021]:

$$\min_{W_{\mathcal{P}} \in [0,1]^{K \times K}} \sum_{m=1}^M d_{\text{gw},2}^2(G_m, W_{\mathcal{P}}) + \alpha \|\Delta W_{\mathcal{P}}\|_F^2. \quad (63)$$

## Two Variants for Structured GWBs

Can we achieve a smoothed graphon?

- ▶ **Smoothed GW Barycenter (SGWB)** [Xu, et al., AAAI 2021]:

$$\min_{W_{\mathcal{P}} \in [0,1]^{K \times K}} \sum_{m=1}^M d_{\text{gw},2}^2(G_m, W_{\mathcal{P}}) + \alpha \|\Delta W_{\mathcal{P}}\|_F^2. \quad (63)$$

- ▶ Similar to classic GWB problem, an alternating optimization strategy works well.

## Two Variants for Structured GWBs

Can we achieve a smoothed graphon?

- ▶ **Smoothed GW Barycenter (SGWB)** [Xu, et al., AAAI 2021]:

$$\min_{W_{\mathcal{P}} \in [0,1]^{K \times K}} \sum_{m=1}^M d_{\text{gw},2}^2(G_m, W_{\mathcal{P}}) + \alpha \|\Delta W_{\mathcal{P}}\|_F^2. \quad (63)$$

- ▶ Similar to classic GWB problem, an alternating optimization strategy works well.

When the graphs are from multiple graphons, how to learn the model?

## Two Variants for Structured GWBs

Can we achieve a smoothed graphon?

- ▶ **Smoothed GW Barycenter (SGWB)** [Xu, et al., AAAI 2021]:

$$\min_{W_{\mathcal{P}} \in [0,1]^{K \times K}} \sum_{m=1}^M d_{\text{gw},2}^2(G_m, W_{\mathcal{P}}) + \alpha \|\Delta W_{\mathcal{P}}\|_F^2. \quad (63)$$

- ▶ Similar to classic GWB problem, an alternating optimization strategy works well.

When the graphs are from multiple graphons, how to learn the model?

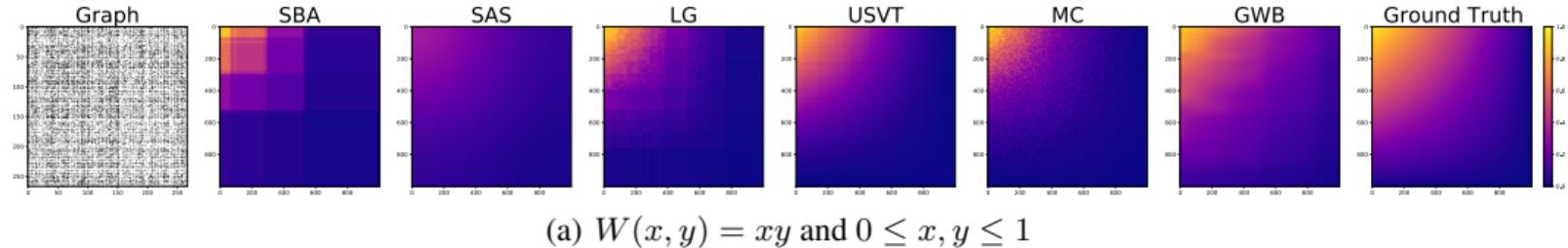
- ▶ **A Mixture of GWBs** [Xu, et al., AAAI 2021]:

$$\min_{\{W_c\}_{c=1}^C, \mathbf{P} \in \Pi(\frac{1}{C}\mathbf{1}_C, \frac{1}{M}\mathbf{1}_M)} \underbrace{\sum_{c=1}^C \sum_{m=1}^M p_{cm} d_{\text{gw},2}^2(G_m, W_c)}_{\langle \mathbf{P}, \mathbf{D}_{\text{gw}} \rangle} \quad (64)$$

- ▶  $p_{cm}$ : the probability of generating the  $m$ -th graph from the  $c$ -th graphon.
- ▶ Learn a graphon set to minimize its hierarchical GW distance to the observed graphs.

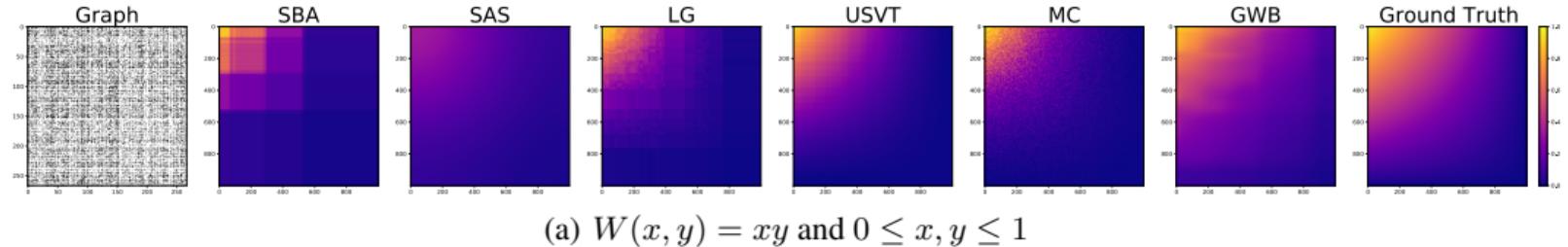
# Experiments

**Easy Case:** The node degrees provide strong evidence for sorting nodes.

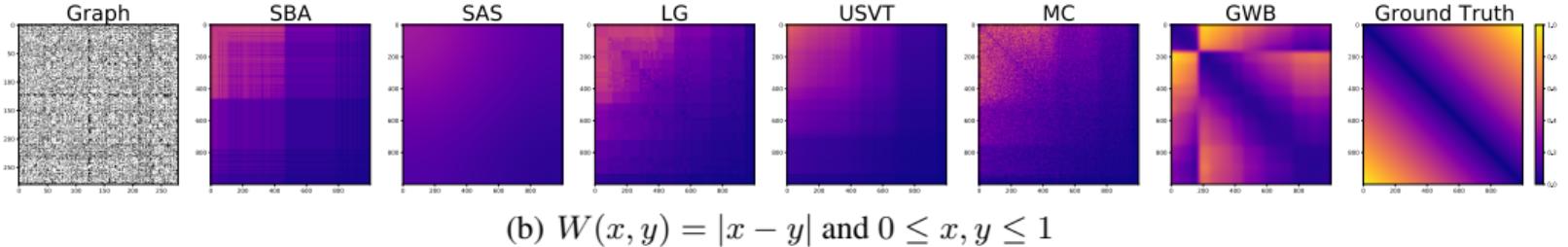


# Experiments

**Easy Case:** The node degrees provide strong evidence for sorting nodes.



**Hard Case:** The nodes of each graph have comparable degrees.



# Summary of the GWL framework

Tasks	PPI Network Alignment	Molecule and Protein Clustering and Classification	Network Simulation	...
Applications	Graph Matching	Graph Partitioning	Graph Representation	Graph Generation
Models	Autoencoder	Factorization Model		Graphon
Optimization	Unsupervised and Semi-supervised Learning			
Optimization	Conditional Gradient	ADMM	Alternating Opt.	...
Theoretical Fundamentals	Constrained Non-convex Optimization			
Theoretical Fundamentals	Gromov-Wasserstein Distance and Its Variants for Structured Data Analysis			

# Open Problems and Ongoing Research

## ► Gromov-Monge Formulation and Its Applications

- ▶ Inverse Problems for MM-Spaces [Mémoli, et al., 2021].
- ▶ Reversible Gromov-Monge Sampler [Hur, et al., 2021].

## ► Connections to Stochastic Control and Dynamic Systems

- ▶ OT problems and algorithms are relevant to the Schrödinger bridge problem [Léonard, 2014] (A diffusion model predicting and tracking particles over time.)
- ▶ How to connect GW distance and its algorithms to Stochastic Control?

## ► Large-scale Heterogeneous Matching and Learning

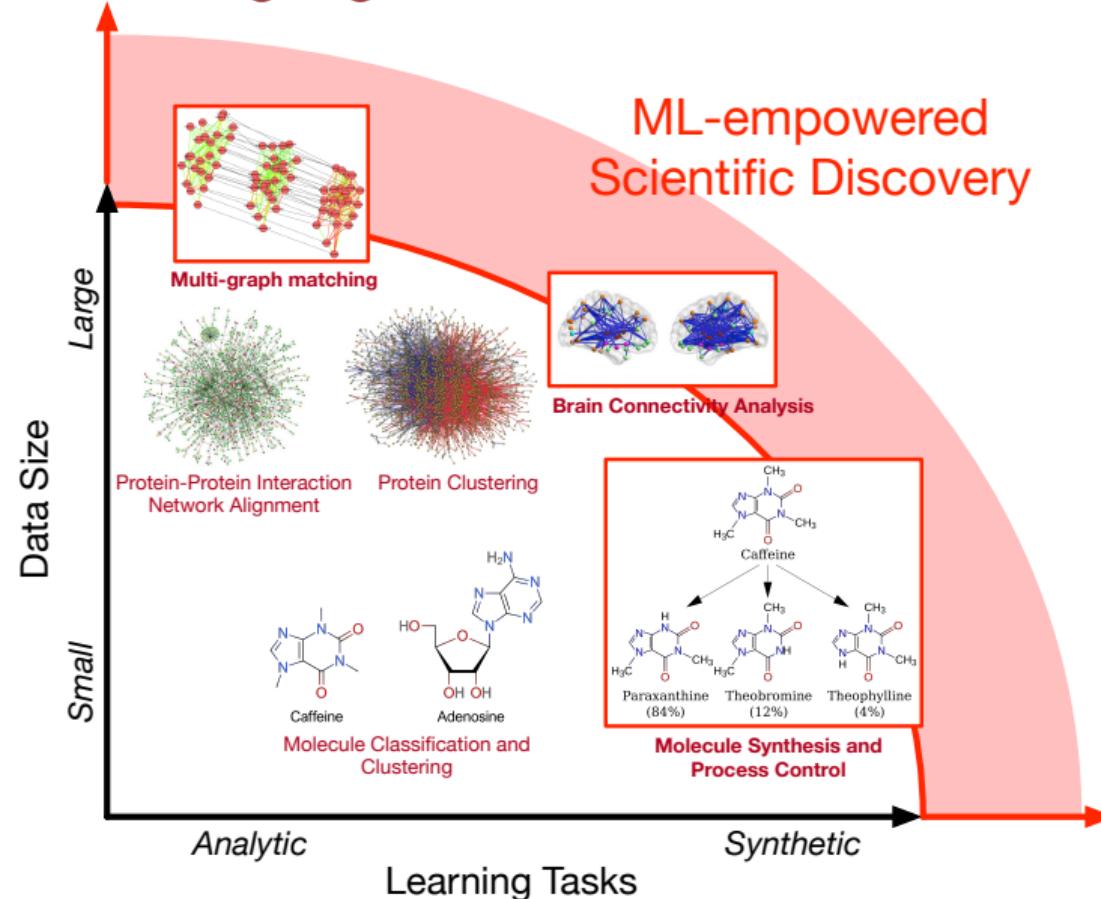
- ▶ Algorithms with lower complexity, higher stability, better approximation.
- ▶ Parallel and distributed learning.

[Mémoli, et al., 2021] Mémoli, F., & Needham, T. Distance distributions and inverse problems for metric measure spaces. arXiv 2021.

[Hur, et al., 2021] Hur, Y., Guo, W., & Liang, T. Reversible Gromov-Monge Sampler for Simulation-Based Inference. arXiv 2021.

[Léonard, 2014] Léonard, C. A survey of the Schrödinger problem and some of its connections with optimal transport. Discrete and Continuous Dynamical Systems-Series A, 2014.

# Open Problems and Ongoing Research



Thank you!

<https://hongtengxu.github.io>

<https://github.com/HongtengXu>

[hongtengxu@ruc.edu.cn](mailto:hongtengxu@ruc.edu.cn)

The 1st Workshop on Optimal Transport and Structured Data Modeling (OT-SDM)

<https://ot-sdm.github.io> (Feb. 28, co-organized with AAAI)