

Gromov-Wasserstein Learning for Graph Matching, Partitioning, and Embedding

Hongteng Xu

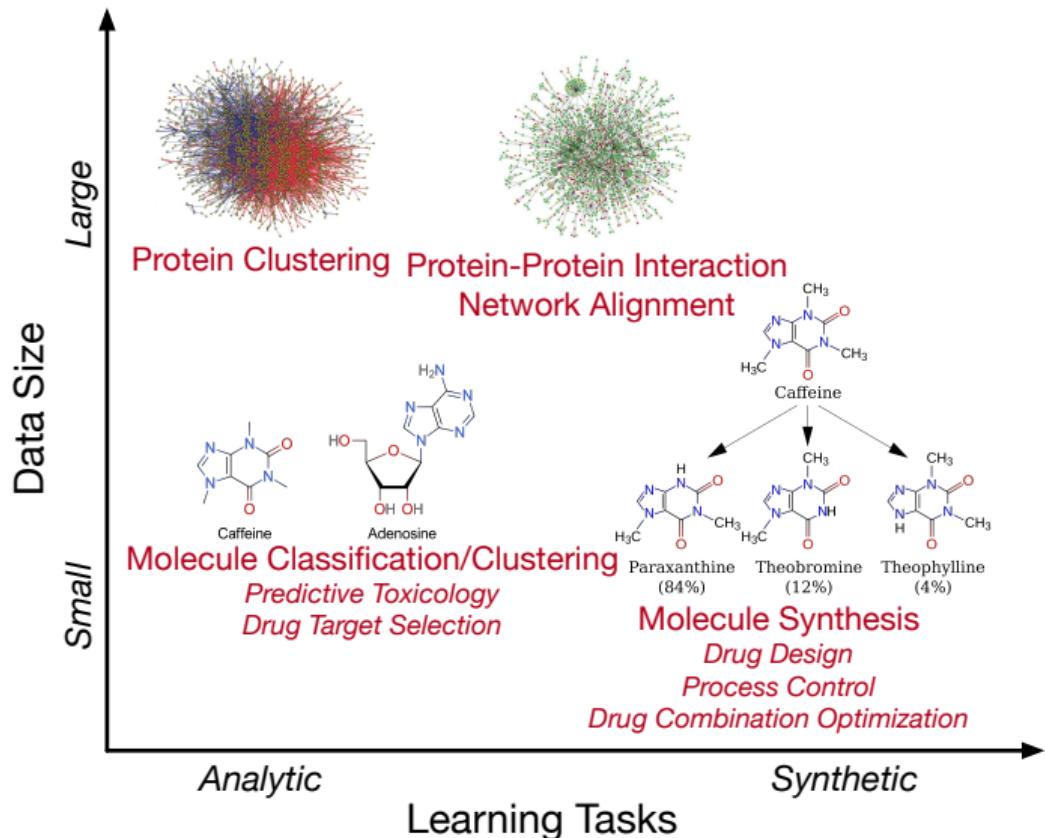
Gaoling School of Artificial Intelligence, Renmin University of China



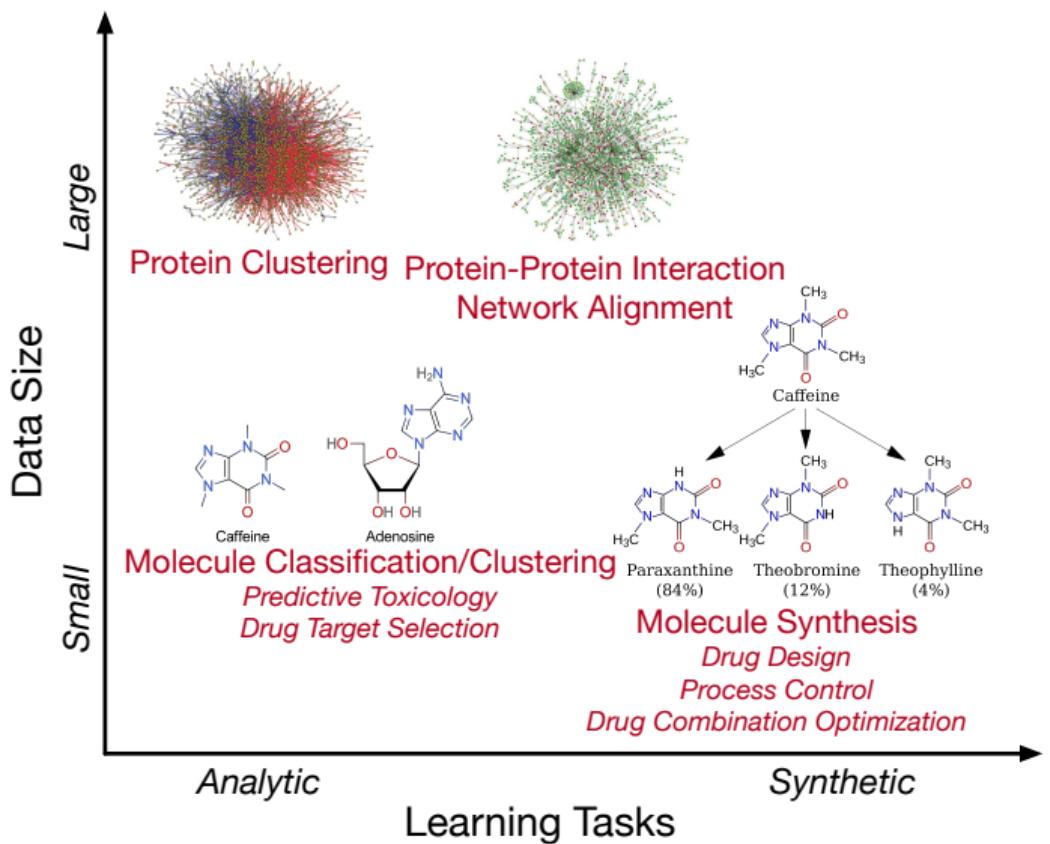
中國人民大學
RENMIN UNIVERSITY OF CHINA

高領人工智能學院
Gaoling School of Artificial Intelligence

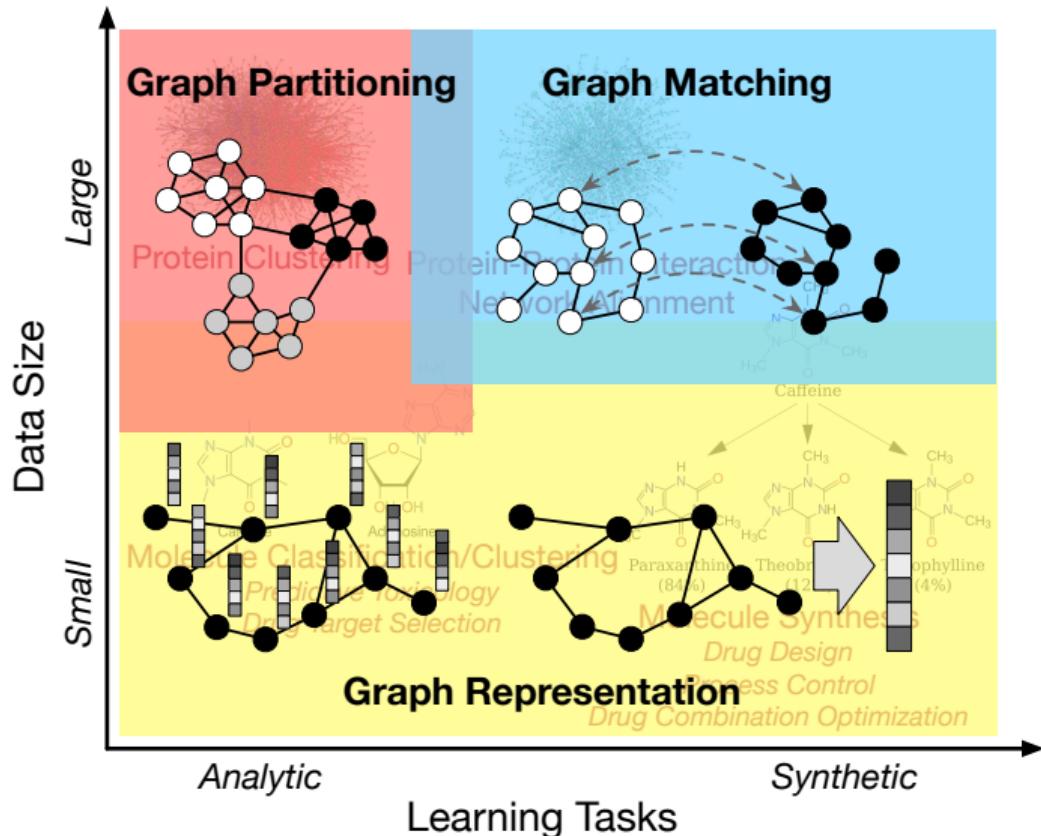
The space of graph-related problems



The space of graph-related problems



The space of graph-related problems

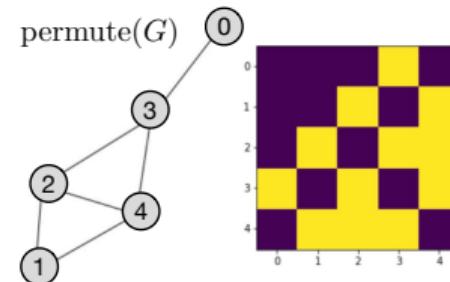
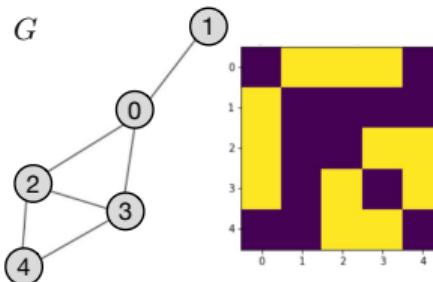


The challenges of the problems

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

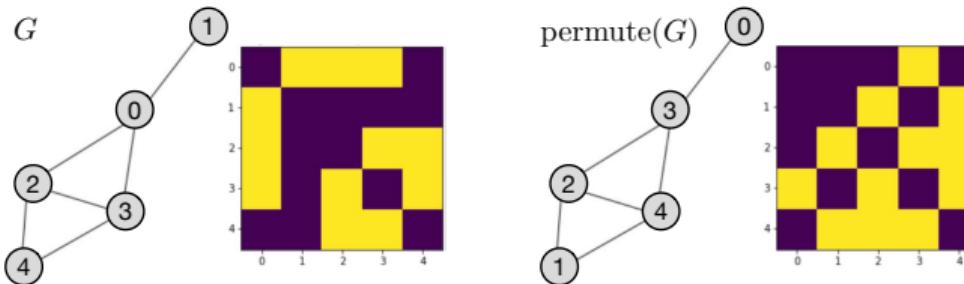
The challenges of the problems

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



The challenges of the problems

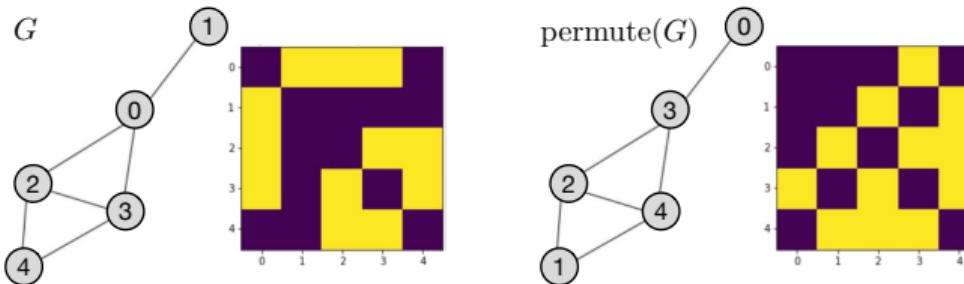
- ▶ NP-completeness
 - ▶ Approximation algorithms with high stability and scalability
- ▶ Permutation invariance



- ▶ 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

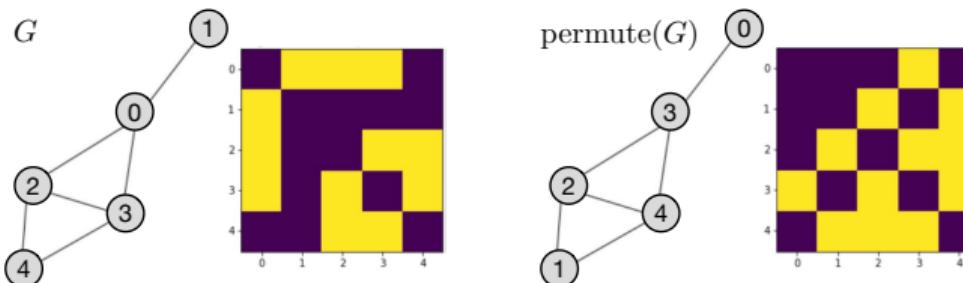
- ▶ NP-completeness
 - ▶ Approximation algorithms with high stability and scalability
- ▶ Permutation invariance



- ▶ 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

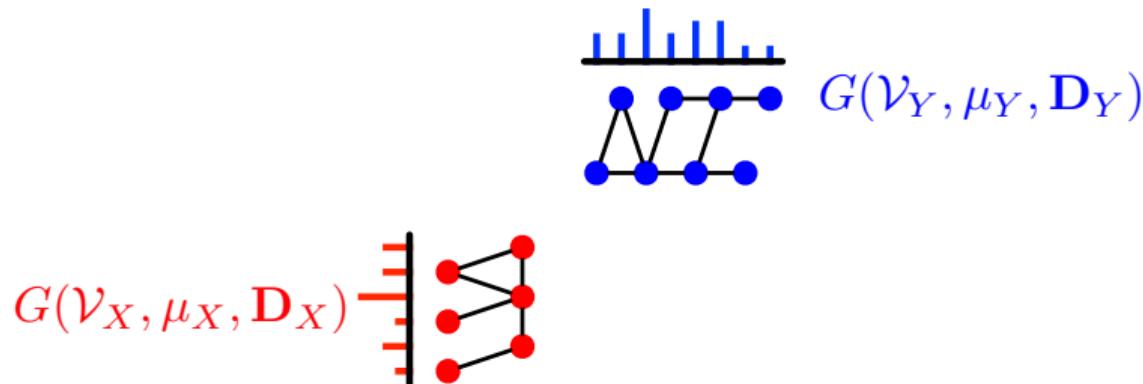
- ▶ NP-completeness
 - ▶ Approximation algorithms with high stability and scalability
- ▶ Permutation invariance



- ▶ 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

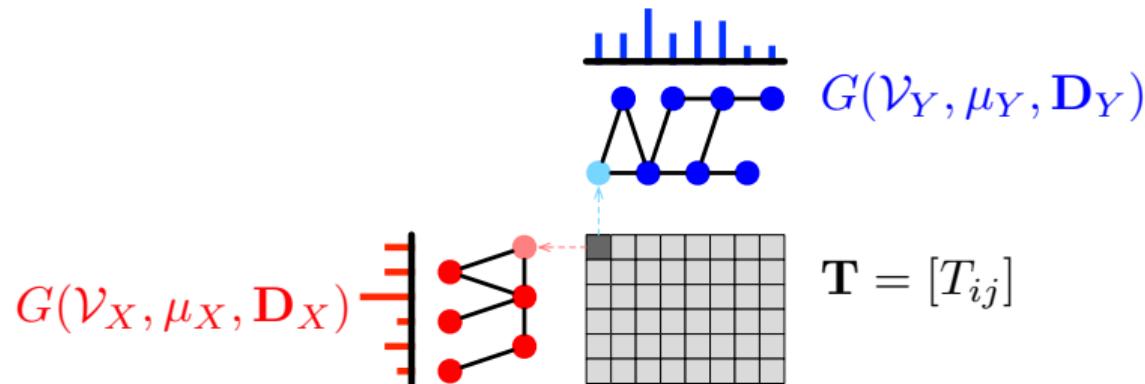
Gromov-Wasserstein Learning (GWL) provides a potential solution.
Applications: PPI network alignment, molecule clustering and classification.

Gromov-Wasserstein distance (GWD) between two graphs



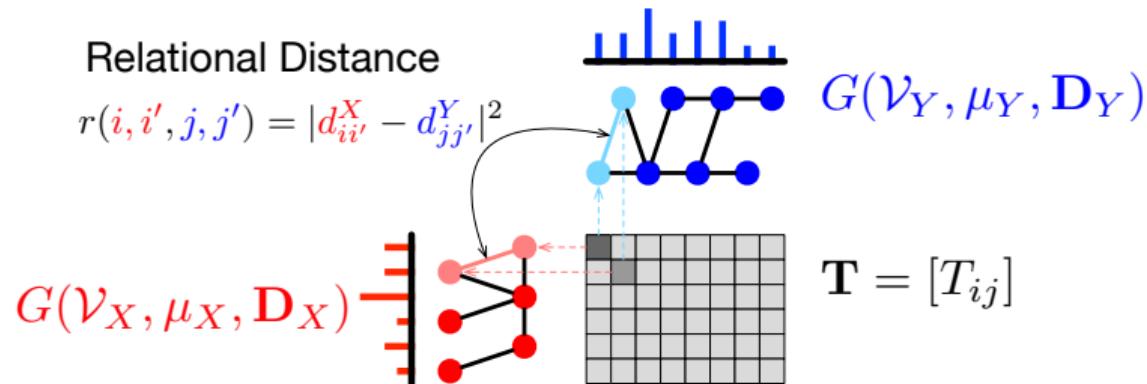
- ▶ \mathcal{V} : the node set
- ▶ μ : a predefined distribution of nodes
- ▶ $\mathbf{D} = [d_{ii'}]$: the adjacency / distance / kernel matrix

Gromov-Wasserstein distance (GWD) between two graphs



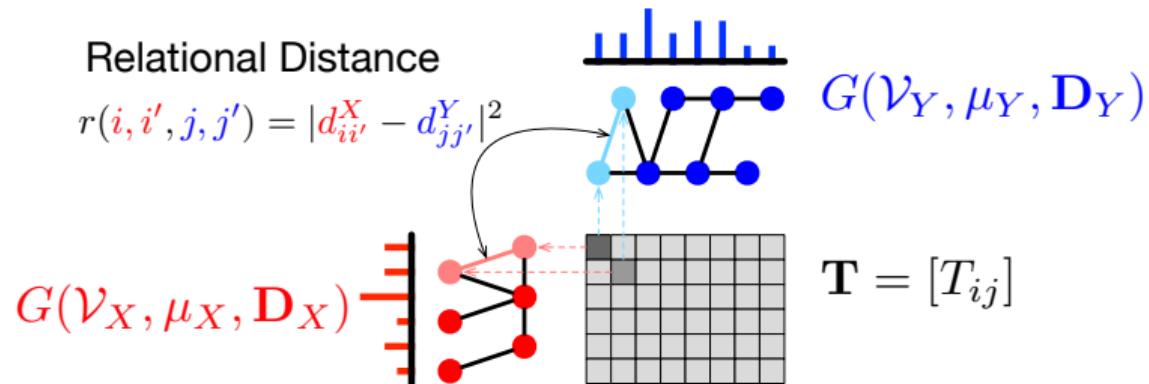
- ▶ $\mathbf{T} = [T_{ij}]$: a joint distribution of nodes
- ▶ $(i \in \mathcal{V}_X, j \in \mathcal{V}_Y) \sim \mathbf{T}$.
- ▶ $\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \boldsymbol{\mu}_Y) = \{\mathbf{T} \geq \mathbf{0} \mid \mathbf{T}\mathbf{1} = \boldsymbol{\mu}_X, \mathbf{T}^\top \mathbf{1} = \boldsymbol{\mu}_Y\}$

Gromov-Wasserstein distance (GWD) between two graphs



- ▶ $\underbrace{\mathbf{T} \otimes \mathbf{T}}$: a joint distribution of edges.
Kronecker product
- ▶ The pair of edges $(d_{ii'}^X, d_{jj'}^Y) \sim \mathbf{T} \otimes \mathbf{T}$.
- ▶ Relational distance $r(i, i', j, j')$: the difference between the edges.

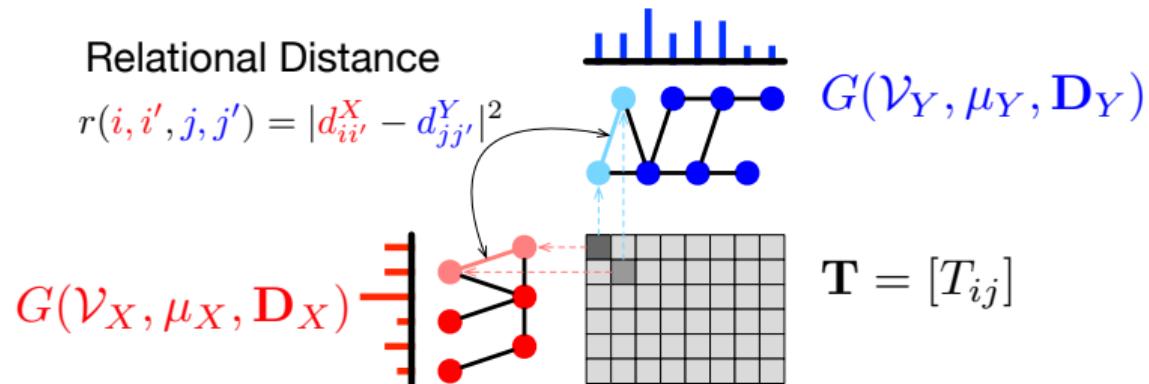
Gromov-Wasserstein distance (GWD) between two graphs



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

$$\begin{aligned} d_{gw}(G_X, G_Y) &:= \min_{\mathbf{T} \in \Pi(\mu_X, \mu_Y)} \mathbb{E}_{(i, i', j, j') \sim \mathbf{T} \otimes \mathbf{T}} [r(i, i', j, j')] \\ &= \min_{\mathbf{T} \in \Pi(\mu_X, \mu_Y)} \sum_{i, i'} \sum_{j, j'} \underbrace{|d_{ii'}^X - d_{jj'}^Y|^2}_{\text{distance } r} \underbrace{T_{ij} T_{i'j'}}_{\text{prob}(r)} \quad (1) \end{aligned}$$

Gromov-Wasserstein distance (GWD) between two graphs

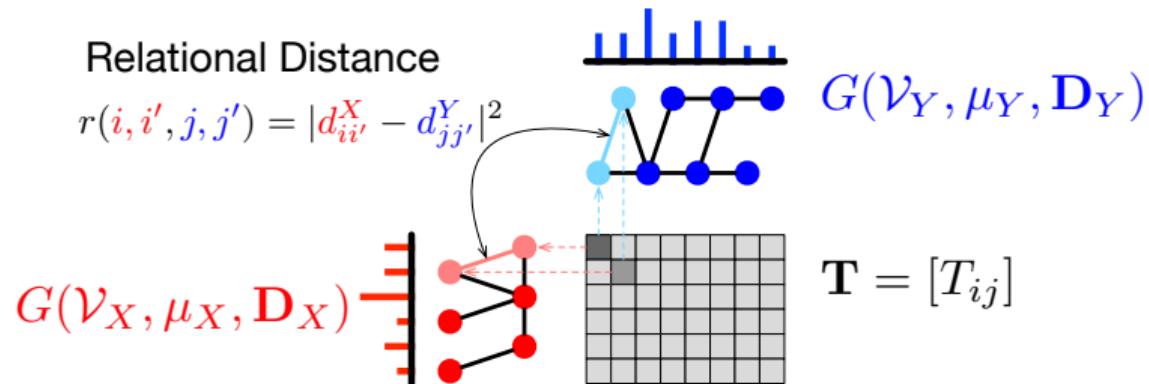


The GWD is **the minimum expectation of the relational distance**:

$$\begin{aligned} d_{gw}(G_X, G_Y) &:= \min_{\mathbf{T} \in \Pi(\mu_X, \mu_Y)} \mathbb{E}_{(i, i', j, j') \sim \mathbf{T} \otimes \mathbf{T}} [r(i, i', j, j')] \\ &= \min_{\mathbf{T} \in \Pi(\mu_X, \mu_Y)} \langle \mathbf{D}_{XY} - 2\mathbf{D}_X \mathbf{T} \mathbf{D}_Y^\top, \mathbf{T} \rangle, \end{aligned} \tag{1}$$

► $\mathbf{D}_{XY} = (\mathbf{D}_X \odot \mathbf{D}_X) \mu_X \mathbf{1}_{|\mathcal{V}_Y|}^\top + \mathbf{1}_{|\mathcal{V}_X|} \mu_Y^\top (\mathbf{D}_Y \odot \mathbf{D}_Y)^\top.$

Gromov-Wasserstein distance (GWD) between two graphs

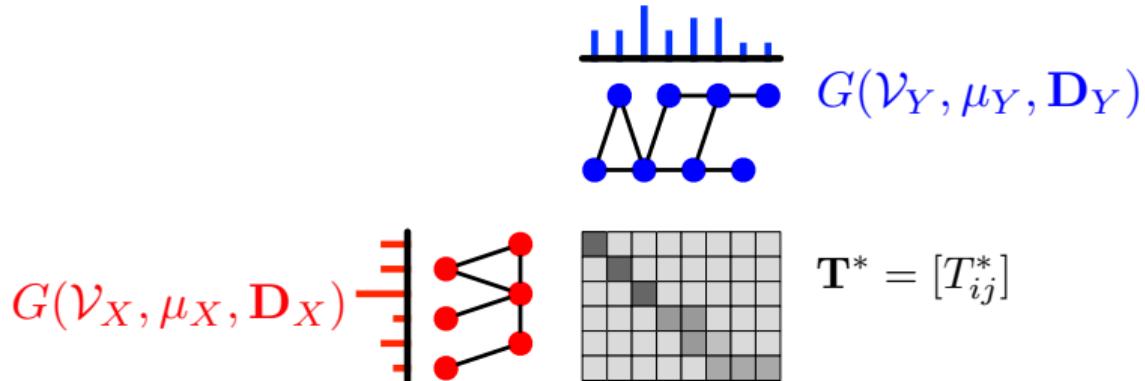


The GWD is **the minimum expectation of the relational distance**:

$$\begin{aligned} d_{gw}(G_X, G_Y) &:= \min_{\mathbf{T} \in \Pi(\mu_X, \mu_Y)} \mathbb{E}_{(i, i', j, j') \sim \mathbf{T} \otimes \mathbf{T}} [r(i, i', j, j')] \\ &= \min_{\mathbf{T} \in \Pi(\mu_X, \mu_Y)} \langle \mathbf{D}_{XY} - 2\mathbf{D}_X \mathbf{T} \mathbf{D}_Y^\top, \mathbf{T} \rangle, \end{aligned} \quad (1)$$

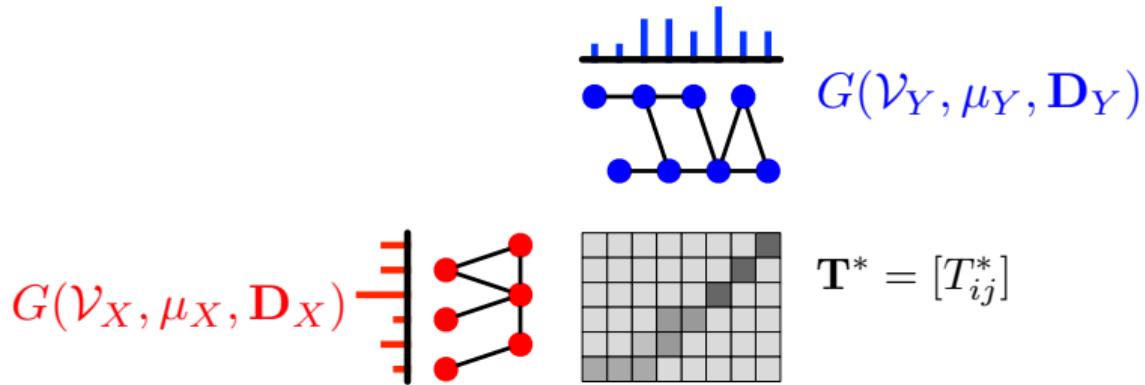
- ▶ $\mathbf{D}_{XY} = (\mathbf{D}_X \odot \mathbf{D}_X) \mu_X \mathbf{1}_{|\mathcal{V}_Y|}^\top + \mathbf{1}_{|\mathcal{V}_X|} \mu_Y^\top (\mathbf{D}_Y \odot \mathbf{D}_Y)^\top.$
- ▶ Given comparable node attributes, $\mathbf{D}_{XY} \leftarrow \mathbf{D}_{XY} + \mathbf{D}(\mathbf{F}_X, \mathbf{F}_Y)$

Advantages of GWD



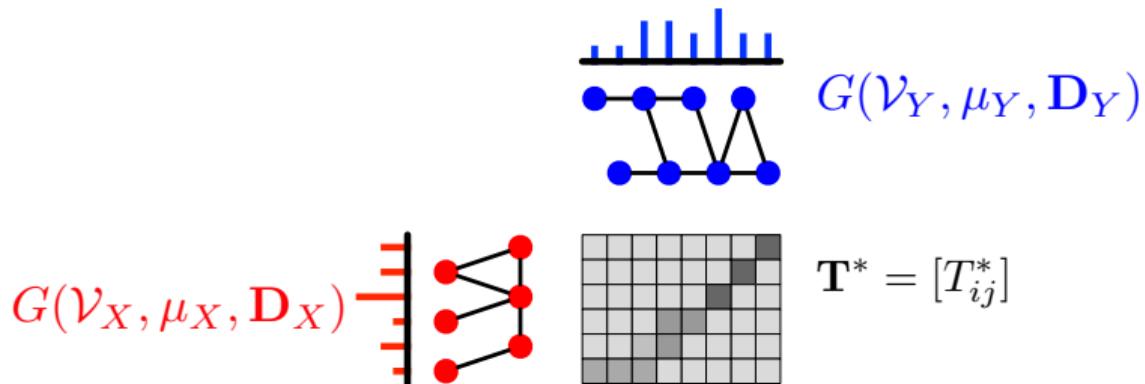
- ▶ The optimal joint distribution \mathbf{T}^* (or called “**optimal transport**” matrix) indicates the correspondence between the two graphs.

Advantages of GWD



- ▶ The optimal joint distribution \mathbf{T}^* (or called “**optimal transport**” matrix) indicates the correspondence between the two graphs.
- ▶ A **permutation-invariant** (pseudo) metric
 - ▶ $d_{gw}(G_X, G_Y) = d_{gw}(G_X, \text{permute}(G_Y))$

Advantages of GWD



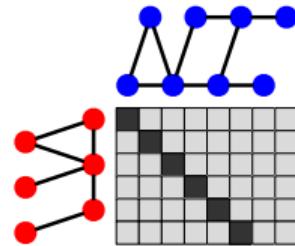
- ▶ The optimal joint distribution \mathbf{T}^* (or called “**optimal transport**” matrix) indicates the correspondence between the two graphs.
- ▶ A **permutation-invariant** (pseudo) metric
 - ▶ $d_{gw}(G_X, G_Y) = d_{gw}(G_X, \text{permute}(G_Y))$
- ▶ Applicable to the graphs with different sizes, *i.e.*, $|\mathcal{V}_X| \neq |\mathcal{V}_Y|$.
- ▶ Applicable to the graphs with/without node attributes.

Matching via learning optimal transport

Quadratic assignment problem (QAP):

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

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

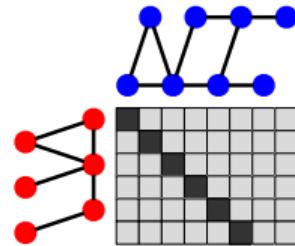


Matching via learning optimal transport

Quadratic assignment problem (QAP):

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

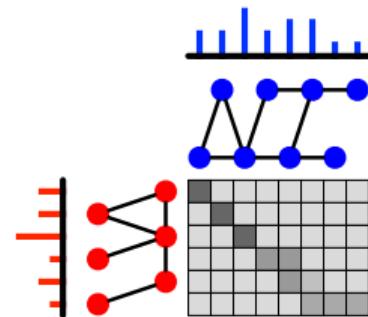
$$\mathcal{P} = \{P \in \{0, 1\}^{|\mathcal{V}_X| \times |\mathcal{V}_Y|} \mid P\mathbf{1} = \mathbf{1}, P^\top \mathbf{1} \leq \mathbf{1}\}.$$



Gromov-Wasserstein distance (GWD):

$$\min_{T \in \Pi(\mu_X, \mu_Y)} \langle \mathbf{D}_{XY} - 2\mathbf{D}_X T \mathbf{D}_Y^\top, T \rangle,$$

$$\Pi(\mu_X, \mu_Y) = \{T \geq 0 \mid T\mathbf{1} = \mu_X, T^\top \mathbf{1} = \mu_Y\}$$

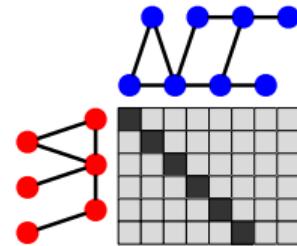


Matching via learning optimal transport

Quadratic assignment problem (QAP):

$$\max_{P \in \mathcal{P}} \langle D_X P D_Y^\top, P \rangle,$$

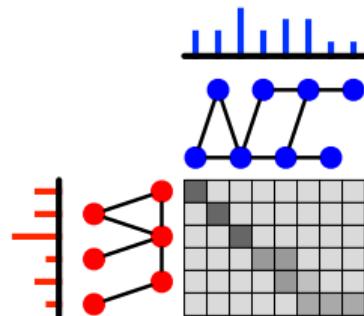
$$\mathcal{P} = \{P \in \{0, 1\}^{|\mathcal{V}_X| \times |\mathcal{V}_Y|} \mid P\mathbf{1} = \mathbf{1}, P^\top \mathbf{1} \leq \mathbf{1}\}.$$



Gromov-Wasserstein distance (GWD):

$$\min_{T \in \Pi(\mu_X, \mu_Y)} \langle D_{XY} - 2D_X T D_Y^\top, T \rangle,$$

$$\Pi(\mu_X, \mu_Y) = \{T \geq 0 \mid T\mathbf{1} = \mu_X, T^\top \mathbf{1} = \mu_Y\}$$

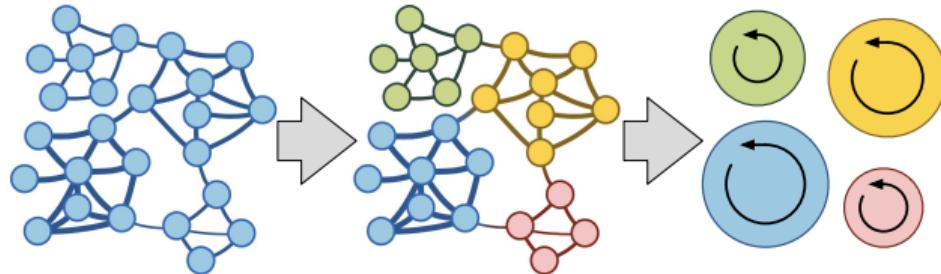


- ▶ Conditional prob. $P^* = \frac{T^*}{\mu_X \mathbf{1}^\top} = [P^*(j|i)]$.
- ▶ For each node $i \in \mathcal{V}_X$, $j^* = \arg \max_j P^*(j|i)$.

Partitioning is also matching

Modularity maximization principle

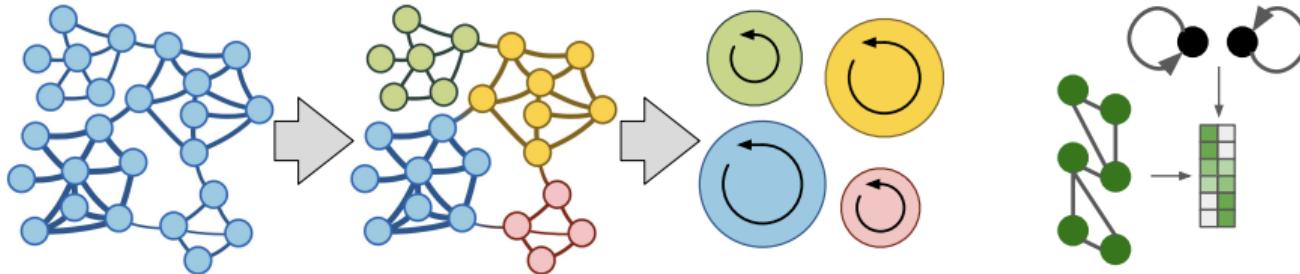
- ▶ Dense internal edges + sparse external edges.



Partitioning is also matching

Modularity maximization principle

- Dense internal edges + sparse external edges.



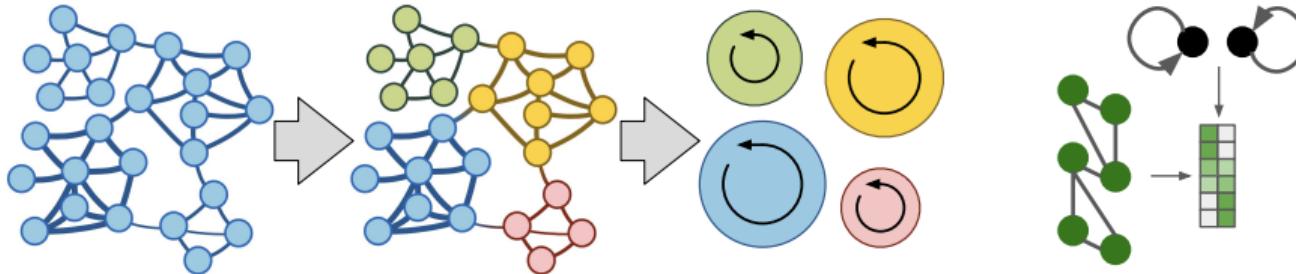
A GWD-based solution [*Xu, et al., NeurIPS 2019*]:

- $\mathbf{T}^* \in \mathbb{R}^{|\mathcal{V}| \times N} \leftarrow d_{gw}(G, G_{iso})$
- $G_{iso}(\mathcal{V}_{iso}, \frac{1}{N} \mathbf{1}_N, \mathbf{I}_{N \times N})$

Partitioning is also matching

Modularity maximization principle

- ▶ Dense internal edges + sparse external edges.



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

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

Proposed Optimization Methods

- ▶ **Entropic Regularization** [Peyré, et al., ICML 2016]
- ▶ **Proximal Point Algorithm** [Xu, et al., ICML 2019]
- ▶ **Bregman ADMM** [Xu, AAAI 2020]

Proposed Optimization Methods

- ▶ **Entropic Regularization** [Peyré, et al., ICML 2016]
- ▶ **Proximal Point Algorithm** [Xu, et al., ICML 2019]
- ▶ **Bregman ADMM** [Xu, AAAI 2020]

Convergence

- ▶ $\lim_{m \rightarrow \infty} T^{(m)}$ is a stationary point.
- ▶ Linear convergence.

Proposed Optimization Methods

- ▶ **Entropic Regularization** [Peyré, et al., ICML 2016]
- ▶ **Proximal Point Algorithm** [Xu, et al., ICML 2019]
- ▶ **Bregman ADMM** [Xu, AAAI 2020]

Convergence

- ▶ $\lim_{m \rightarrow \infty} \mathbf{T}^{(m)}$ is a stationary point.
- ▶ Linear convergence.

Computational complexity per iteration

$$\min_{\mathbf{T} \in \Pi(\mu_X, \mu_Y)} \langle \mathbf{D}_{XY} - 2\mathbf{D}_X \mathbf{T} \mathbf{D}_Y^\top, \mathbf{T} \rangle$$

- ▶ $\mathbf{D}_X, \mathbf{D}_Y$ are dense distance/kernel matrices: $\mathcal{O}(V^3)$.

Proposed Optimization Methods

- ▶ **Entropic Regularization** [Peyré, et al., ICML 2016]
- ▶ **Proximal Point Algorithm** [Xu, et al., ICML 2019]
- ▶ **Bregman ADMM** [Xu, AAAI 2020]

Convergence

- ▶ $\lim_{m \rightarrow \infty} T^{(m)}$ is a stationary point.
- ▶ Linear convergence.

Computational complexity per iteration

$$\min_{T \in \Pi(\mu_X, \mu_Y)} \langle D_{XY} - 2\mathbf{D}_X T \mathbf{D}_Y^\top, T \rangle$$

- ▶ D_X, D_Y are dense distance/kernel matrices: $\mathcal{O}(V^3)$.
- ▶ D_X, D_Y are adjacency matrices: $\mathcal{O}(VE)$.

Proposed Optimization Methods

- ▶ **Entropic Regularization** [Peyré, et al., ICML 2016]
- ▶ **Proximal Point Algorithm** [Xu, et al., ICML 2019]
- ▶ **Bregman ADMM** [Xu, AAAI 2020]

Convergence

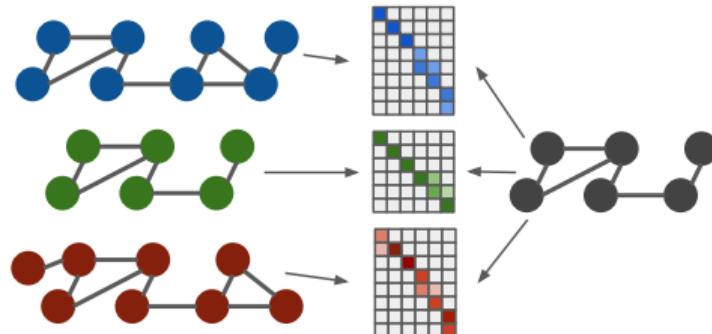
- ▶ $\lim_{m \rightarrow \infty} \mathbf{T}^{(m)}$ is a stationary point.
- ▶ Linear convergence.

Computational complexity per iteration

$$\min_{\mathbf{T} \in \Pi(\mu_X, \mu_Y)} \langle \mathbf{D}_{XY} - 2\mathbf{D}_X \mathbf{T} \mathbf{D}_Y^\top, \mathbf{T} \rangle$$

- ▶ $\mathbf{D}_X, \mathbf{D}_Y$ are dense distance/kernel matrices: $\mathcal{O}(V^3)$.
- ▶ $\mathbf{D}_X, \mathbf{D}_Y$ are adjacency matrices: $\mathcal{O}(VE)$.
- ▶ When $V = |\mathcal{V}_X| \gg |\mathcal{V}_Y| = N$ (graph partitioning): $\mathcal{O}(N(E + V))$.

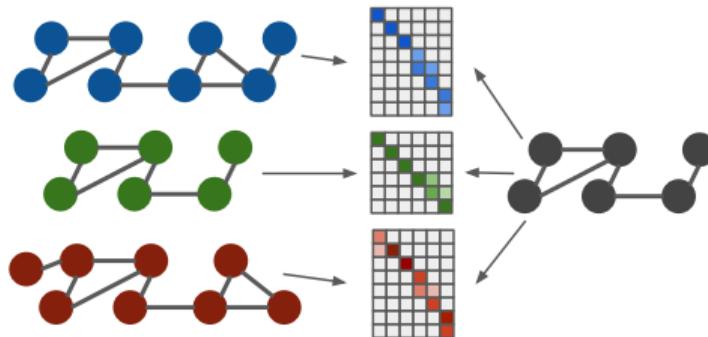
Large-scale matching based on GW barycenters



Given $\{G_k\}_{k=1}^K$, $K \geq 2$, their GW barycenter is defined as

$$\underbrace{B_{gw}(\bar{\mathcal{V}}, \bar{\boldsymbol{\mu}}, \mathbf{B}^*)}_{\text{Barycenter graph}}, \quad \underbrace{\{\mathbf{T}_k^*\}_{k=1}^K}_{\text{OT matrices}} := \arg \min_B \sum_{k=1}^K \lambda_k d_{gw}(B, G_k), \quad (2)$$

Large-scale matching based on GW barycenters



Given $\{G_k\}_{k=1}^K$, $K \geq 2$, their GW barycenter is defined as

$$\underbrace{B_{gw}(\bar{\mathcal{V}}, \bar{\boldsymbol{\mu}}, \mathbf{B}^*)}_{\text{Barycenter graph}}, \quad \underbrace{\{\mathbf{T}_k^*\}_{k=1}^K}_{\text{OT matrices}} := \arg \min_B \sum_{k=1}^K \lambda_k d_{gw}(B, G_k), \quad (2)$$

Learn $\{\mathbf{T}_k^*\}_{k=1}^K$ and \mathbf{B}^* via **alternating optimization**.

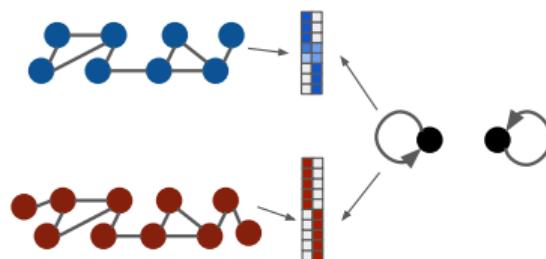
$$\mathbf{B}^* = \frac{1}{\bar{\boldsymbol{\mu}} \bar{\boldsymbol{\mu}}^\top} \sum_{k=1}^K \lambda_k (\mathbf{T}_k^*)^\top \mathbf{D}_k \mathbf{T}_k^* \quad (3)$$

Large-scale matching based on GW barycenters

Co-partition two graphs:

$$B^*, T_X^*, T_Y^* = \arg \min \frac{|\mathcal{V}_X|}{|\mathcal{V}_X| + |\mathcal{V}_Y|} d_{gw}(B, G_X) + \frac{|\mathcal{V}_Y|}{|\mathcal{V}_X| + |\mathcal{V}_Y|} d_{gw}(B, G_Y)$$

Initialize the barycenter graph by a disconnected graph.

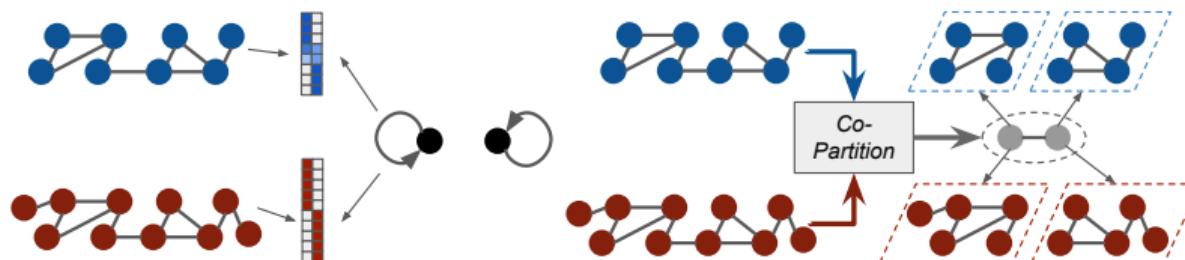


Large-scale matching based on GW barycenters

Co-partition two graphs:

$$B^*, T_X^*, T_Y^* = \arg \min \frac{|\mathcal{V}_X|}{|\mathcal{V}_X| + |\mathcal{V}_Y|} d_{gw}(B, G_X) + \frac{|\mathcal{V}_Y|}{|\mathcal{V}_X| + |\mathcal{V}_Y|} d_{gw}(B, G_Y)$$

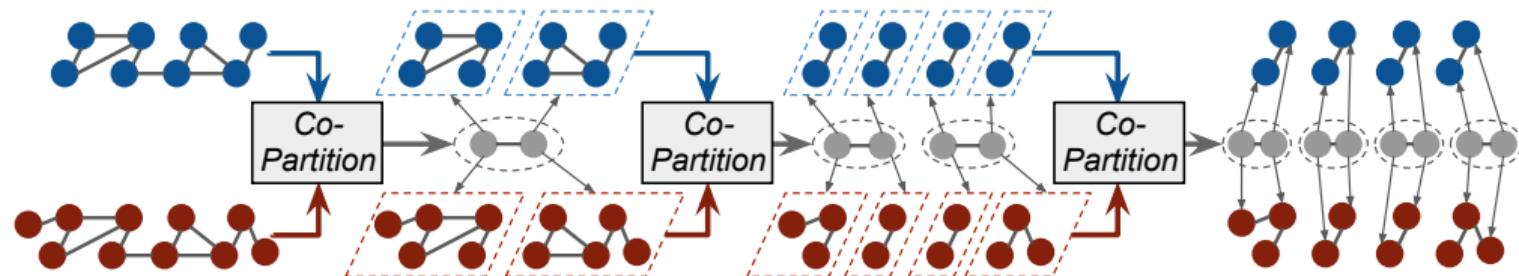
Initialize the barycenter graph by a disconnected graph.



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

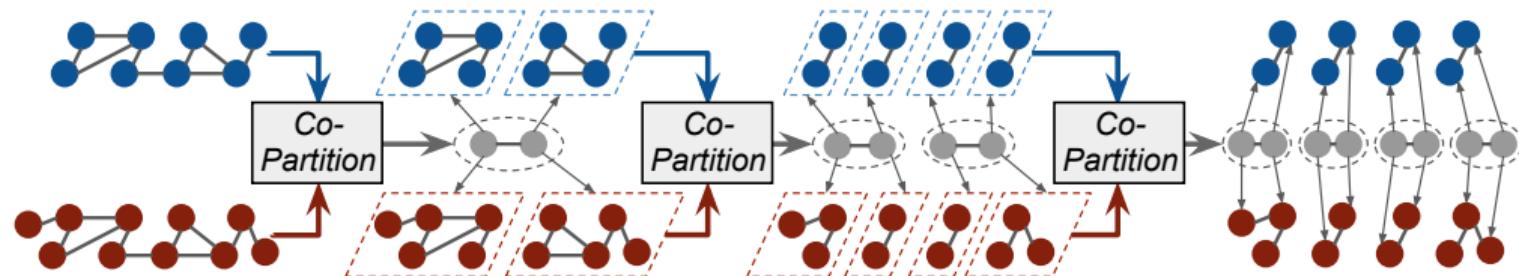
Large-scale matching based on GW barycenters

A “Divide and Conquer” strategy based on recursive co-partitioning [Xu, et al., NeurIPS 2019]



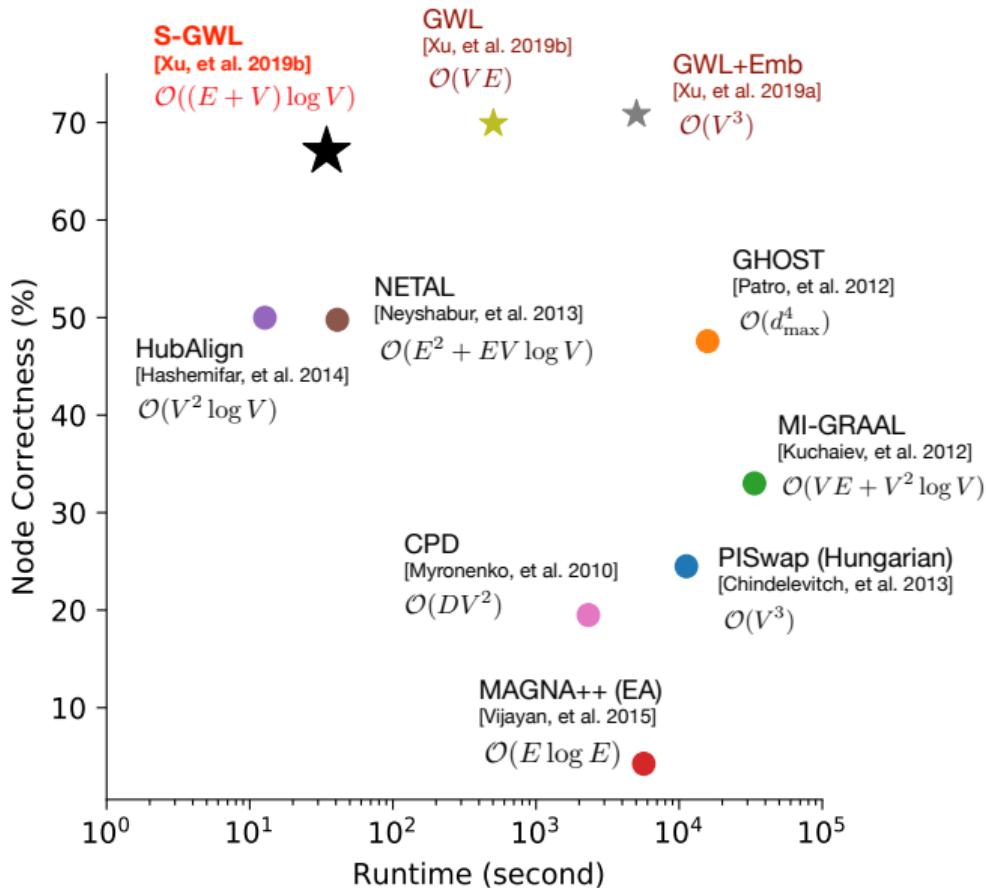
Large-scale matching based on GW barycenters

A “Divide and Conquer” strategy based on recursive co-partitioning [Xu, et al., NeurIPS 2019]



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

Matching synthetic graphs

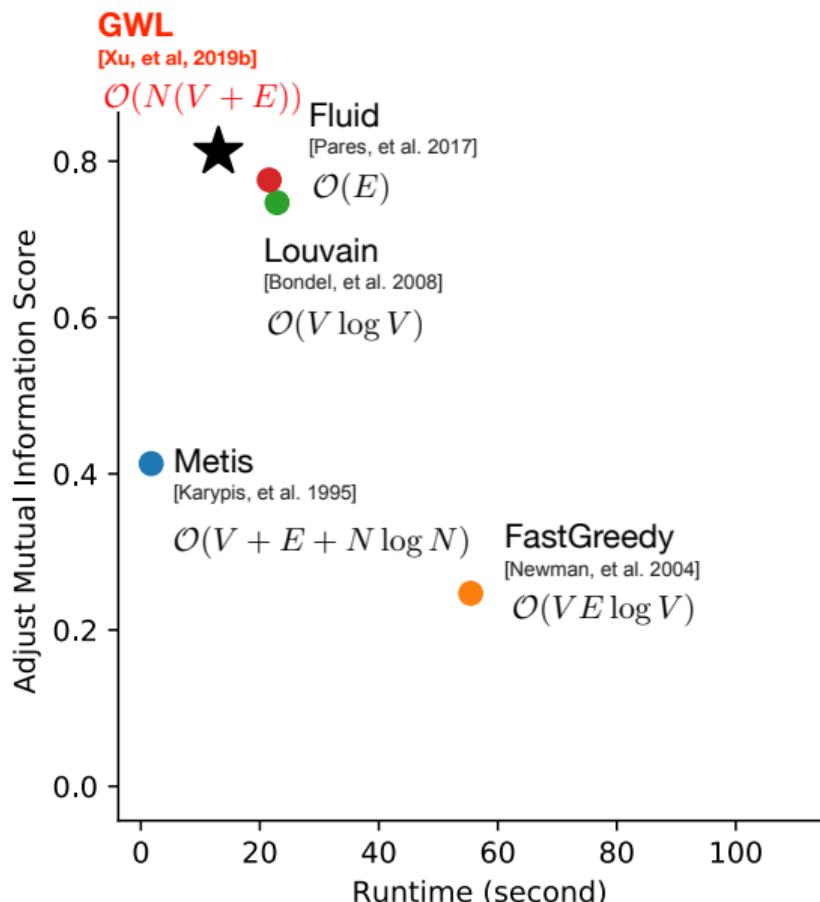
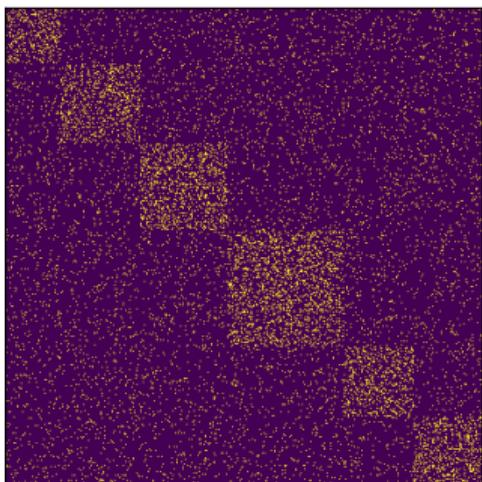


Partitioning synthetic graphs

$V = 4,000$

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

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

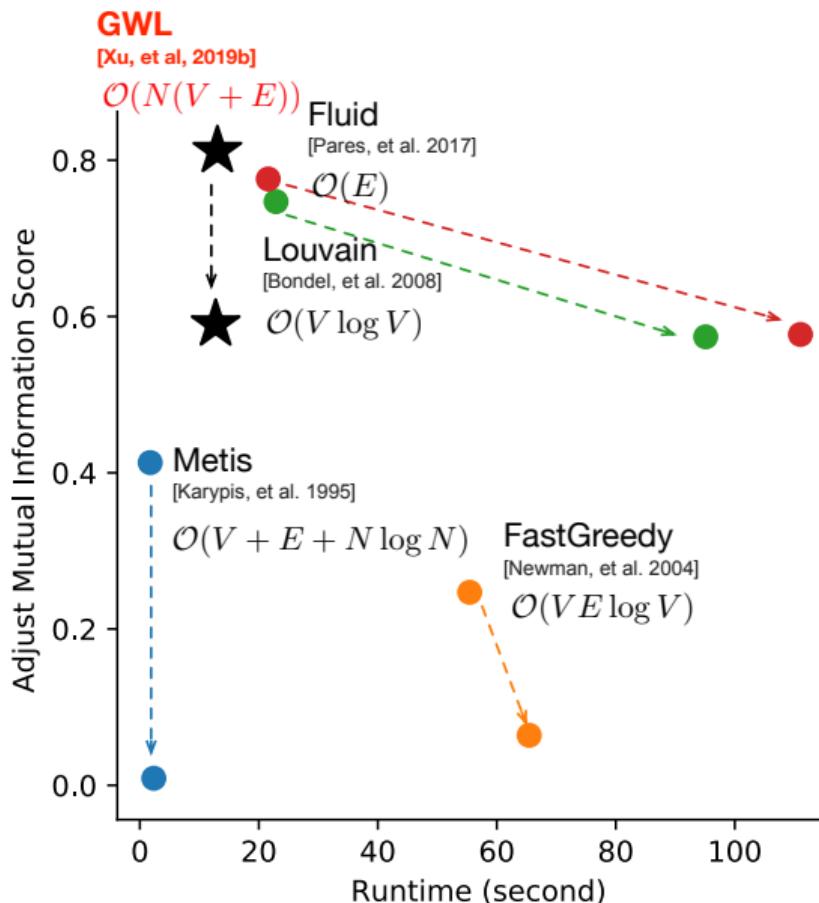
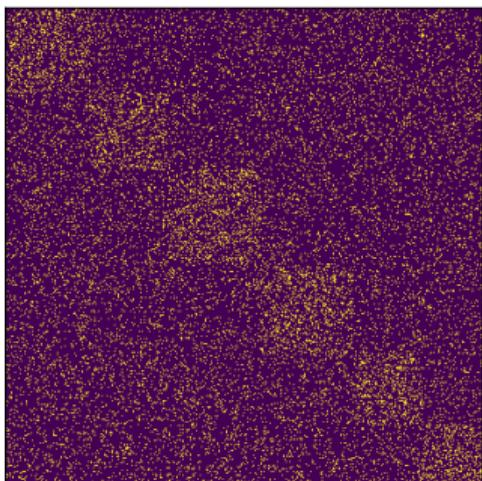


Partitioning synthetic graphs

$V = 4,000$

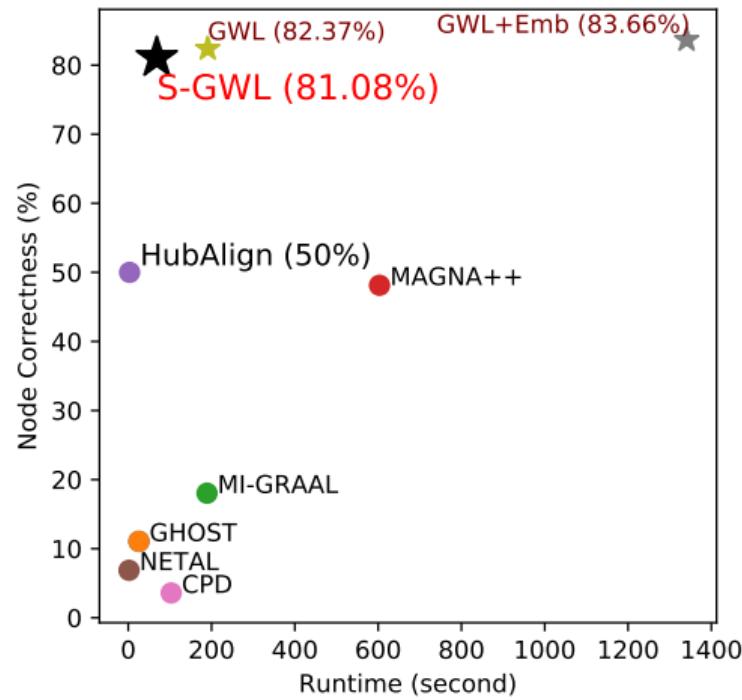
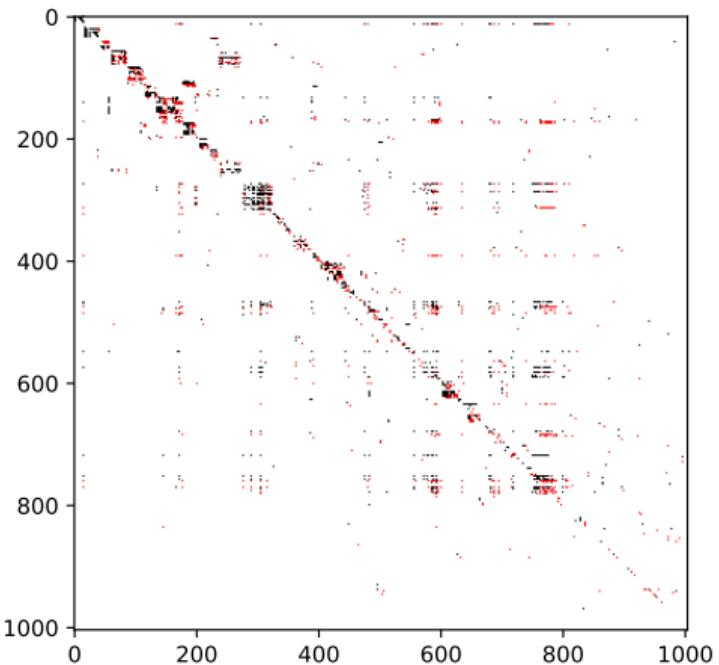
$p_{\text{within}} = 0.2$

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



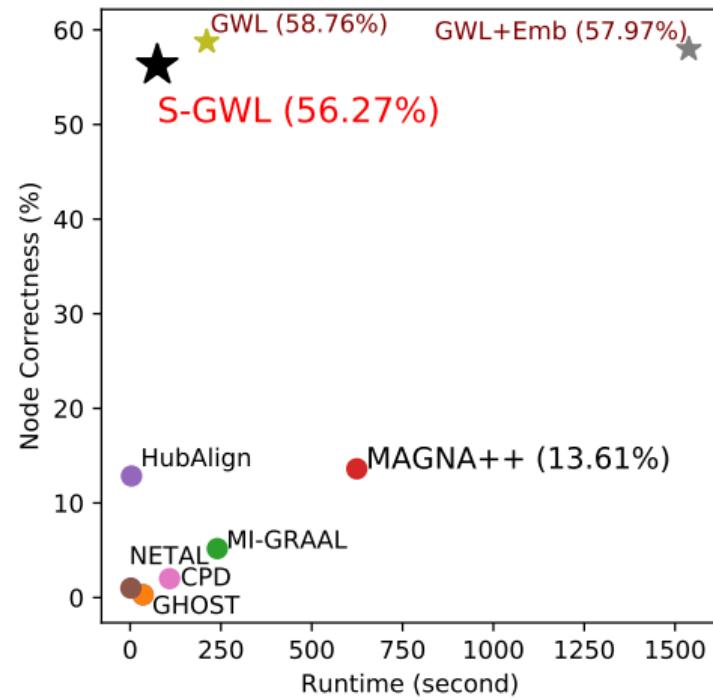
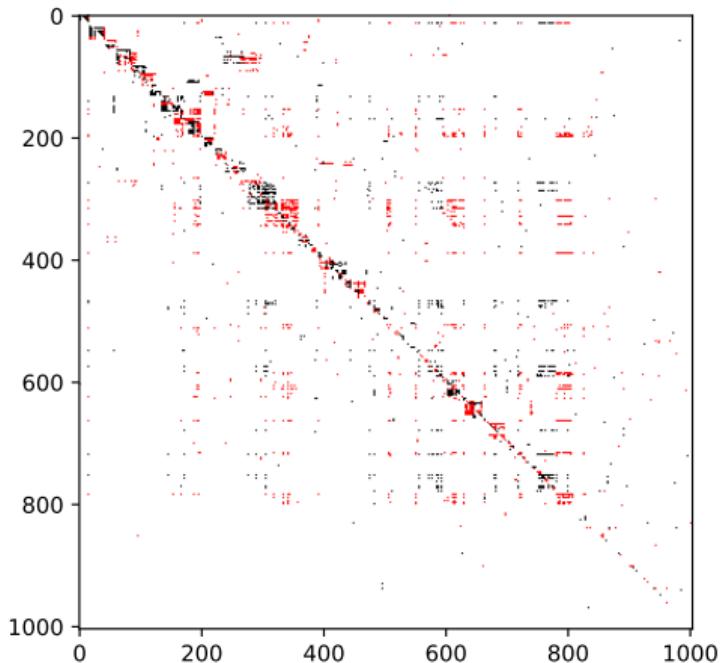
Real-world PPI network alignment

Yeast PPI \leftrightarrow Yeast PPI + 5% LC edges

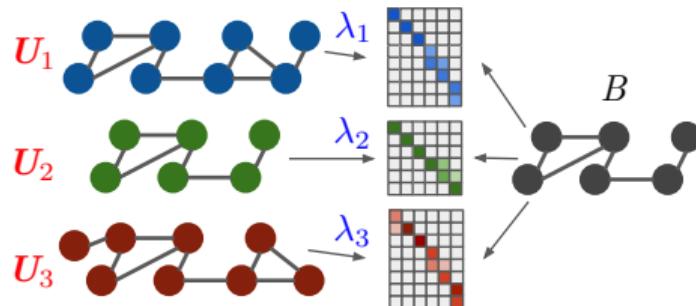


Real-world PPI network alignment

Yeast PPI \leftrightarrow Yeast PPI + 25% LC edges



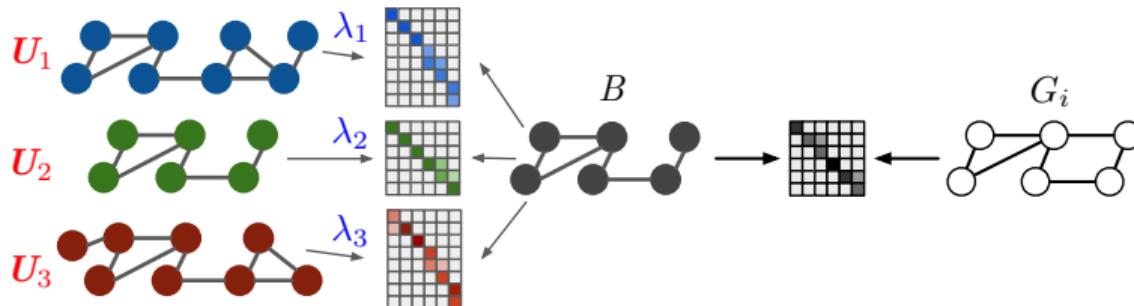
Graph representation via Gromov-Wasserstein factorization



$$B_{gw}(\mathbf{U}_{1:K}, \boldsymbol{\lambda}) := \arg \min_B \sum_{k=1}^K \lambda_k d_{gw}(B, \mathbf{G}_k(\mathbf{U}_k)). \quad (4)$$

- ▶ $\{\mathbf{G}_k(\mathbf{U}_k)\}_{k=1}^K$: a set of graph factors.
- ▶ $\boldsymbol{\lambda} = [\lambda_k] \in \Delta^{K-1}$: the coefficients of the graph factors.

Graph representation via Gromov-Wasserstein factorization



$$B_{gw}(\mathbf{U}_{1:K}, \boldsymbol{\lambda}) := \arg \min_B \sum_{k=1}^K \lambda_k d_{gw}(B, \mathbf{G}_k(\mathbf{U}_k)). \quad (4)$$

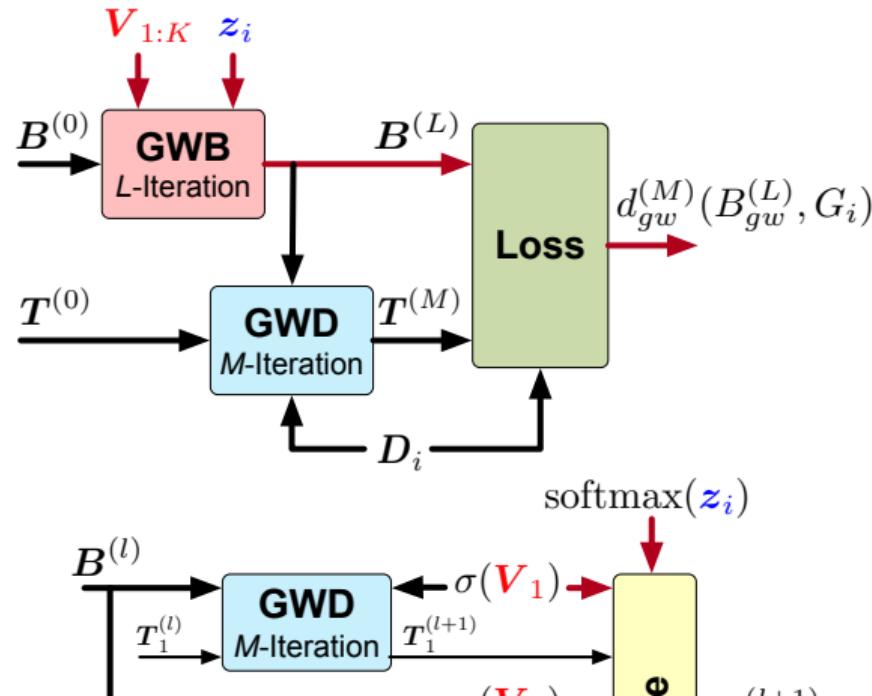
- ▶ $\{\mathbf{G}_k(\mathbf{U}_k)\}_{k=1}^K$: a set of graph bases.
- ▶ $\boldsymbol{\lambda} = [\lambda_k] \in \Delta^{K-1}$: the coefficients of the graph basis.
- ▶ Estimate each graph by a GW barycenter graph [Xu, AAAI 2020]:

$$\min_{\mathbf{1} \geq \mathbf{U}_{1:K} \geq \mathbf{0}, \ \boldsymbol{\lambda}_{1:I} \in \Delta^{K-1}} \sum_{i=1}^I d_{gw}(B_{gw}(\mathbf{U}_{1:K}, \underbrace{\boldsymbol{\lambda}_i}_{\text{Rep. of } G_i}), G_i). \quad (5)$$

Graph representation via Gromov-Wasserstein factorization

Reparameterize the problem to an unconstrained optimization problem:

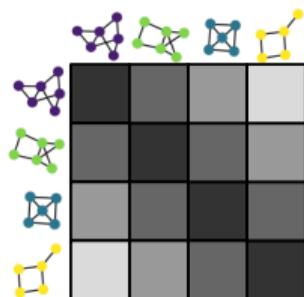
$$\min_{\mathbf{V}_{1:K}, \mathbf{z}_{1:I}} \sum_{i=1}^I d_{gw}(B_{gw}(\underbrace{\sigma(\mathbf{V}_{1:K})}_{\mathbf{U}_{1:K}}, \underbrace{\text{softmax}(\mathbf{z}_i)}_{\lambda_i}), G_i). \quad (6)$$



Experiments on molecule clustering

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

GWD Kernel



GWD+Kmeans

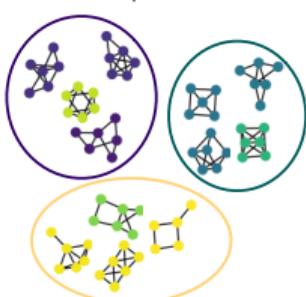


Table: Comparisons on clustering accuracy (%)

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

Summary of the GWL framework

Theoretical
Fundamentals

Gromov-Wasserstein Distance for
Structured Data

Summary of the GWL framework

Optimization
Theoretical
Fundamentals

Proximal Gradient	ADMM	Alternating Opt.	...
Constrained Non-convex Optimization			
Gromov-Wasserstein Distance for Structured Data			

Summary of the GWL framework

Models	Graph convolution networks		Factorization Model
	Unsupervised and Semi-supervised Learning		
Optimization	Proximal Gradient	ADMM	Alternating Opt.
	Constrained Non-convex Optimization		
Theoretical Fundamentals	Gromov-Wasserstein Distance for Structured Data		

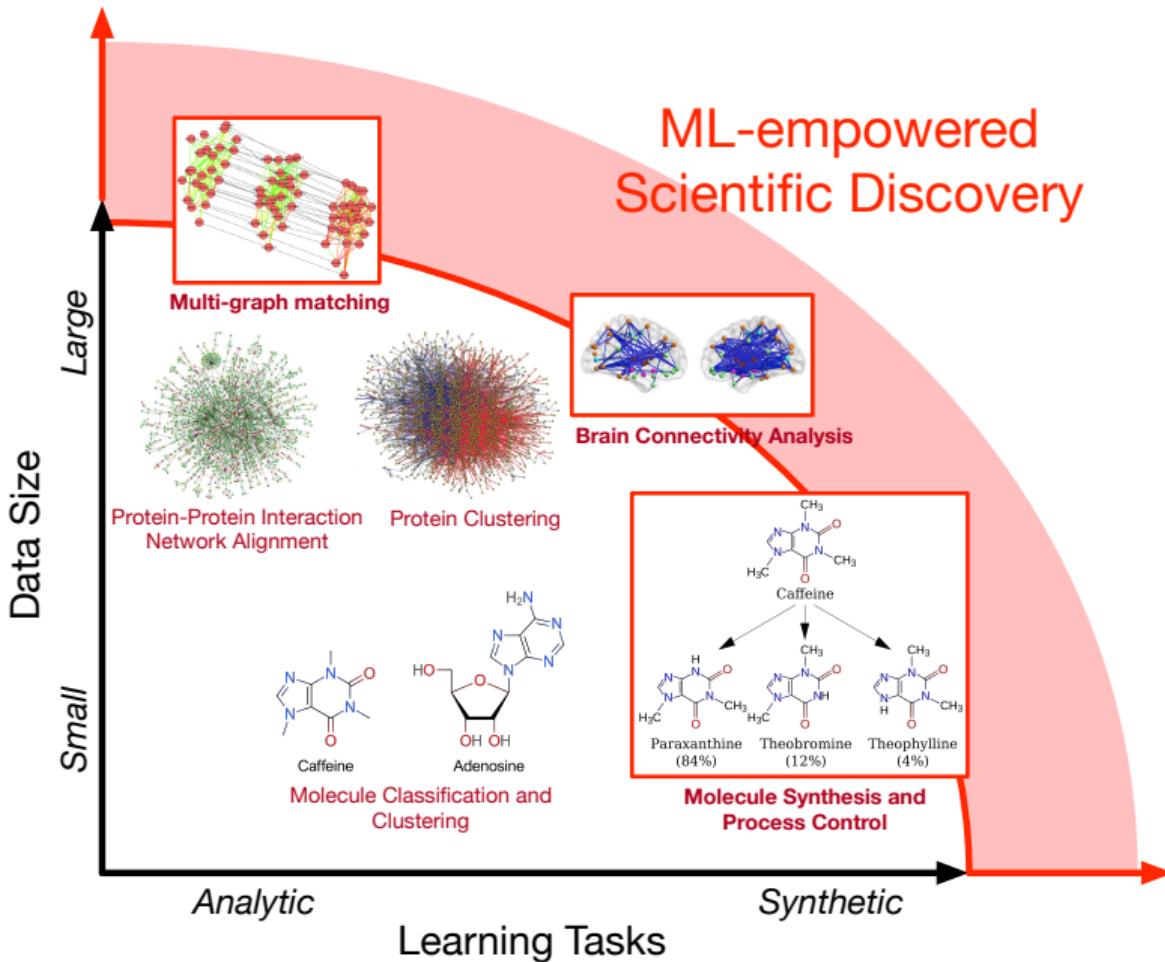
Summary of the GWL framework

Applications	Graph Matching	Graph Partitioning	Graph Representation
Models	Graph convolution networks		Factorization Model
Unsupervised and Semi-supervised Learning			
Optimization	Proximal Gradient	ADMM	Alternating Opt.
Constrained Non-convex Optimization			
Theoretical Fundamentals	Gromov-Wasserstein Distance for Structured Data		

Summary of the GWL framework

Tasks	PPI Network Alignment	Molecule Clustering and Classification	...
Applications	Graph Matching	Graph Partitioning	Graph Representation
Models	Graph convolution networks		Factorization Model
	Unsupervised and Semi-supervised Learning		
Optimization	Proximal Gradient	ADMM	Alternating Opt.
	Constrained Non-convex Optimization		
Theoretical Fundamentals	Gromov-Wasserstein Distance for Structured Data		

Future work



Thanks! Q&A

<https://hongtengxu.github.io>

<https://github.com/HongtengXu>

hongtengxu@ruc.edu.cn

AAAI 2022: OT-SDM Workshop <https://ot-sdm.github.io> + Tutorial on GWL