

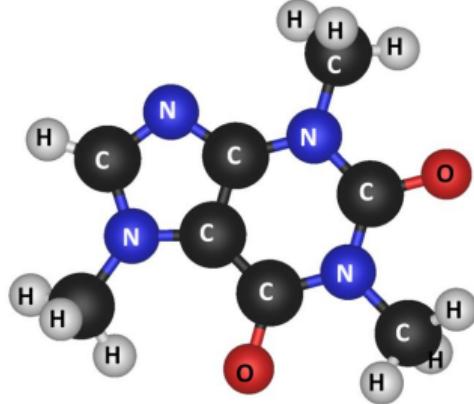
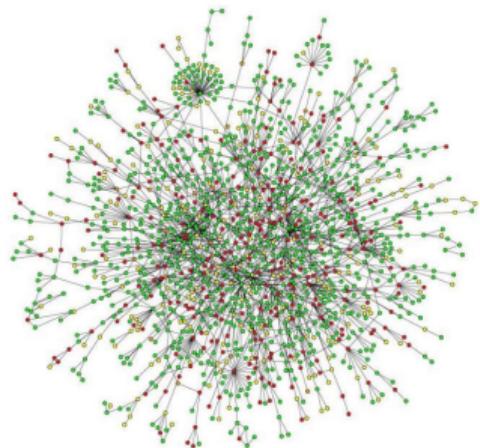
Gromov-Wasserstein Learning: A New Machine Learning Framework for Structured Data Analysis

Hongteng Xu

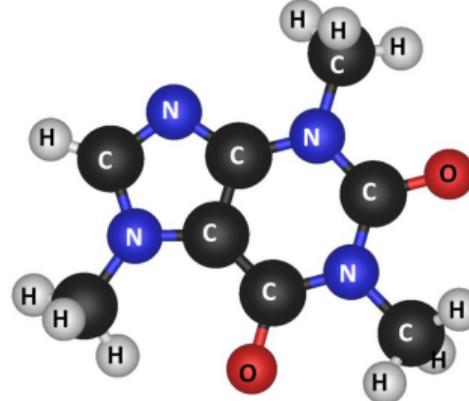
Infinia ML Inc. and Duke University



Structured data analysis



Structured data analysis



≡ Google Scholar

"structured data" OR "graph" OR "network" OR "molecule" AND "learning"

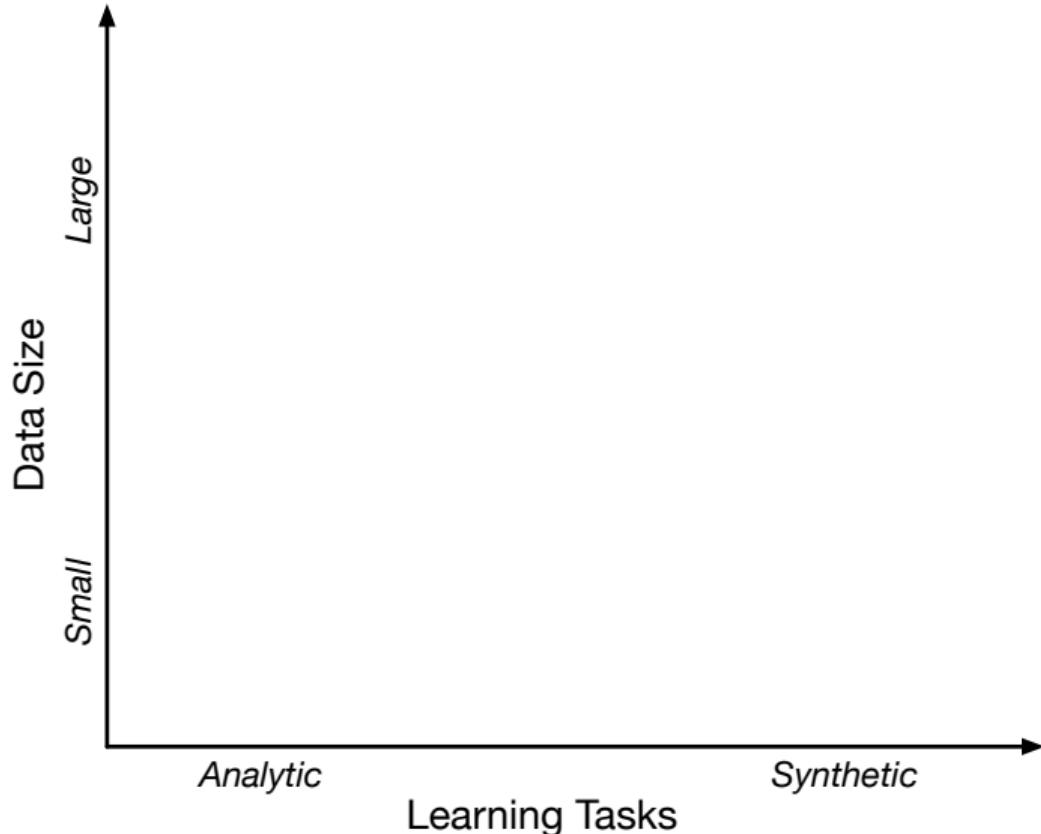


Articles

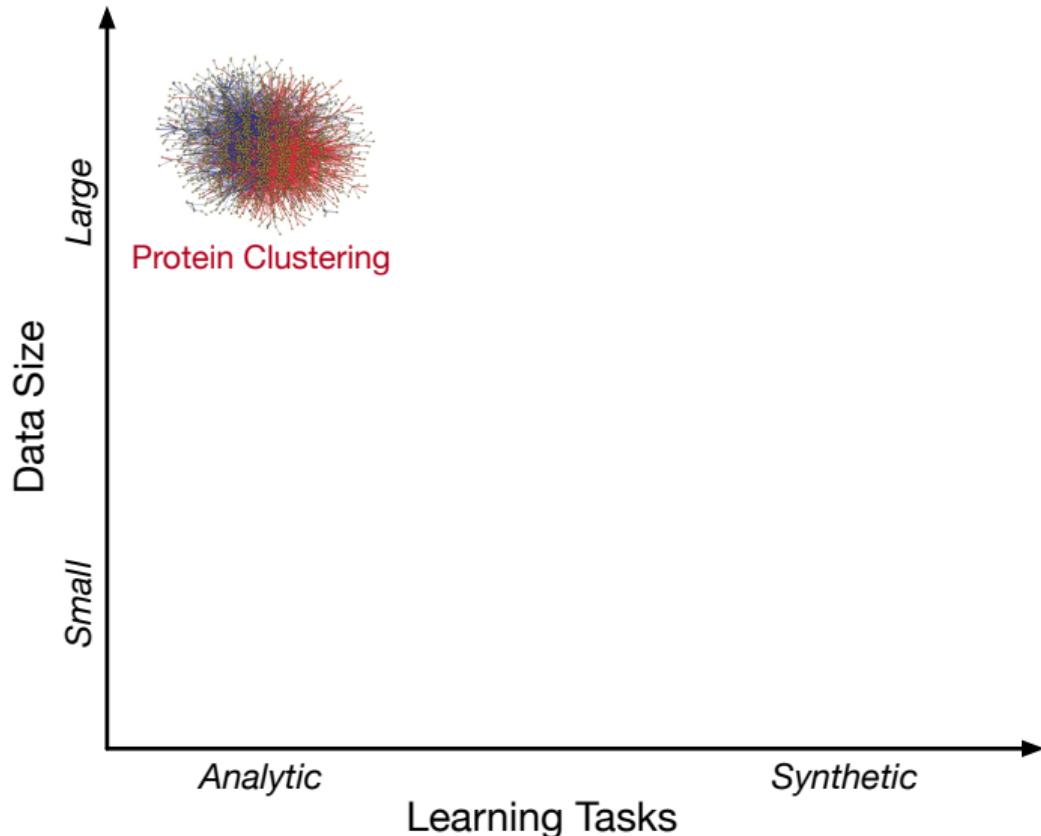
About 4,700,000 results (0.13 sec)

Extract knowledge from relations effectively and efficiently.

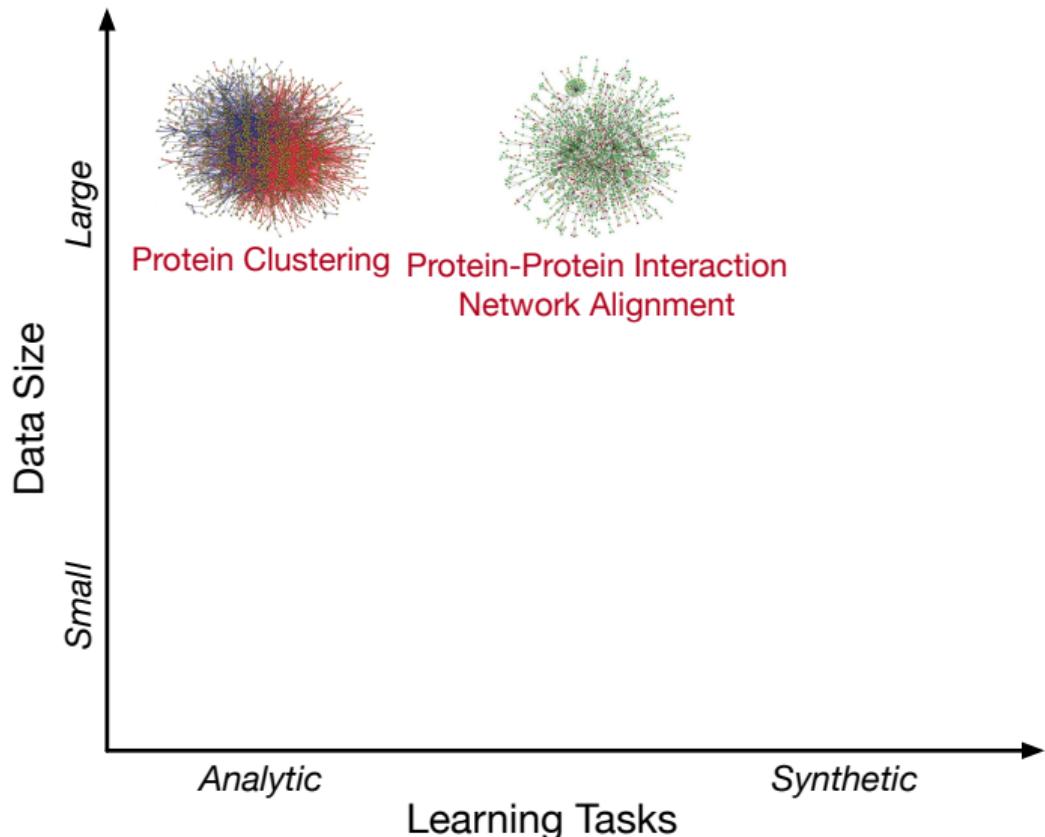
The space of problems



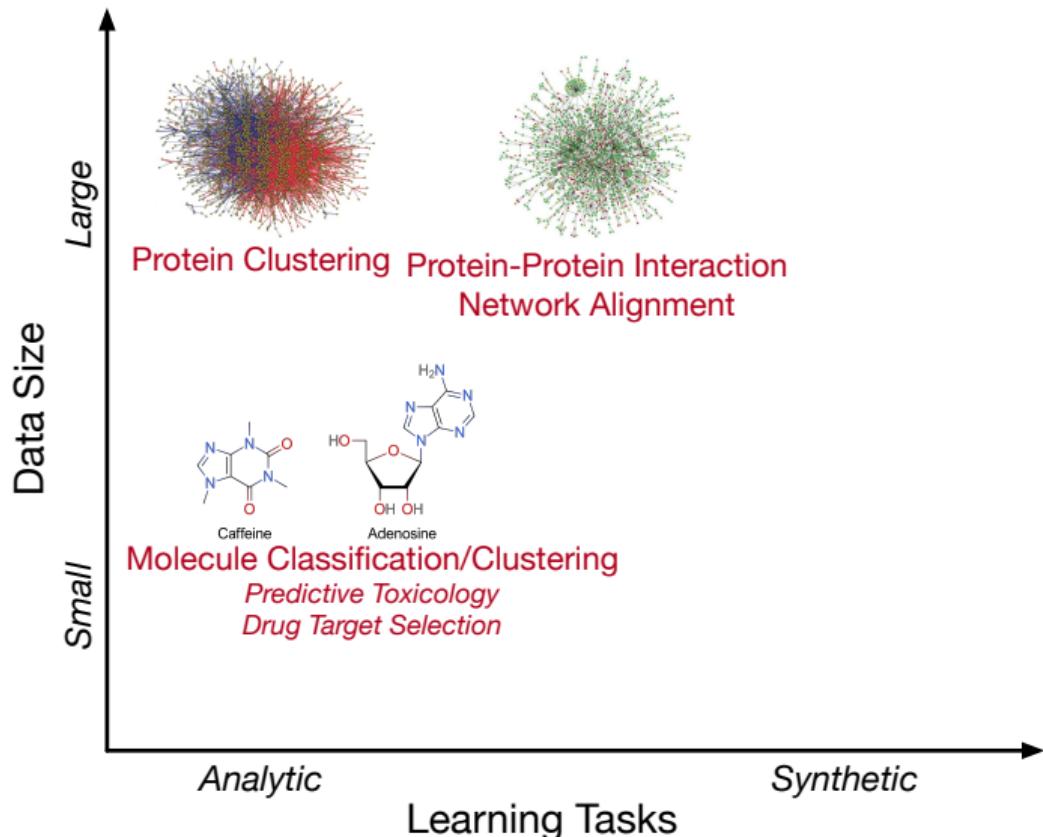
The space of problems



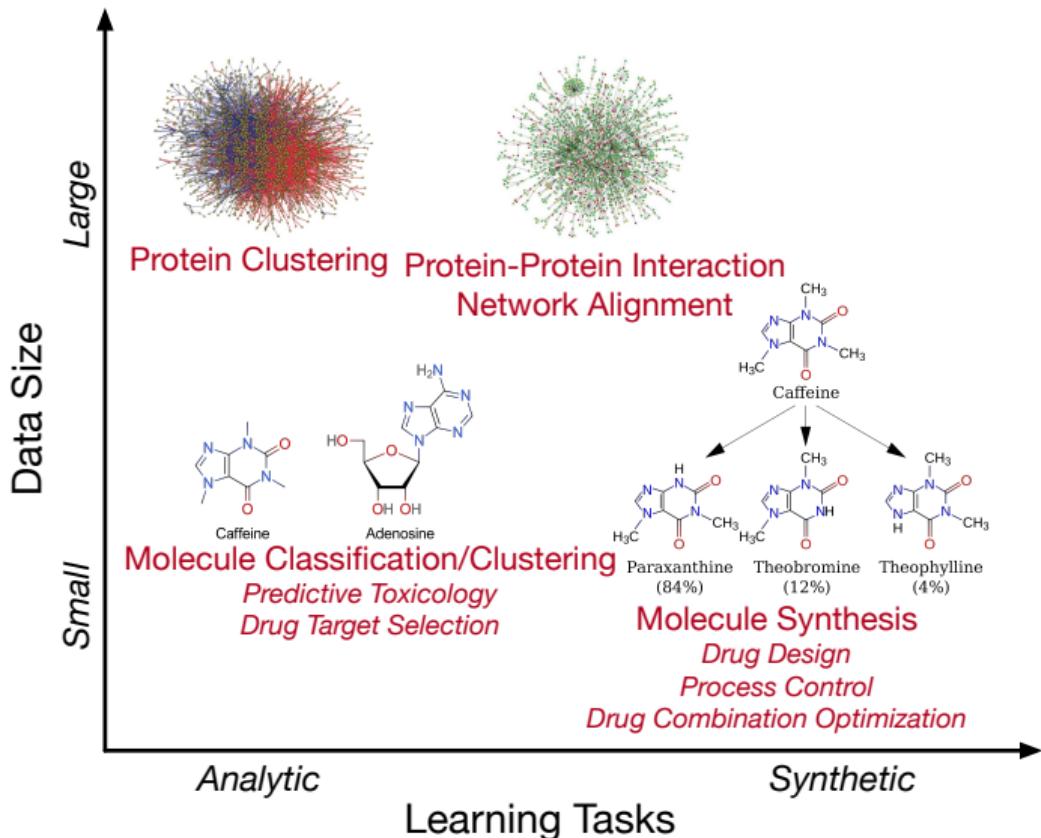
The space of problems



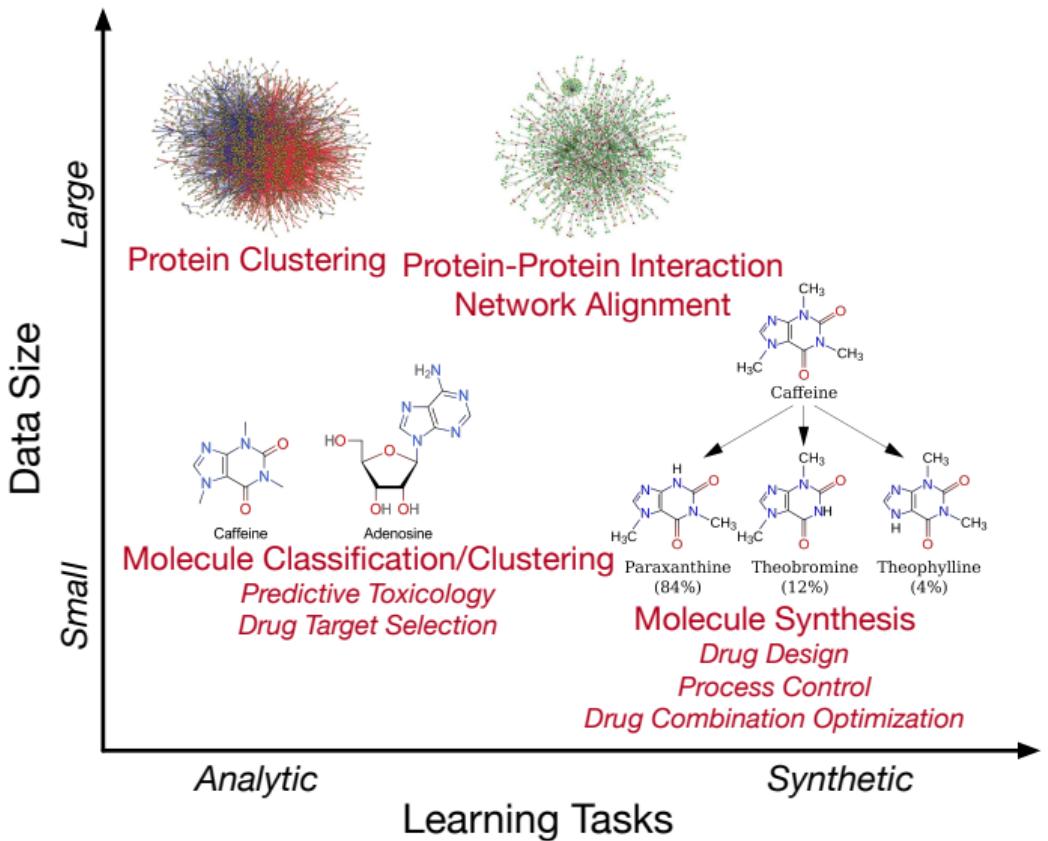
The space of problems



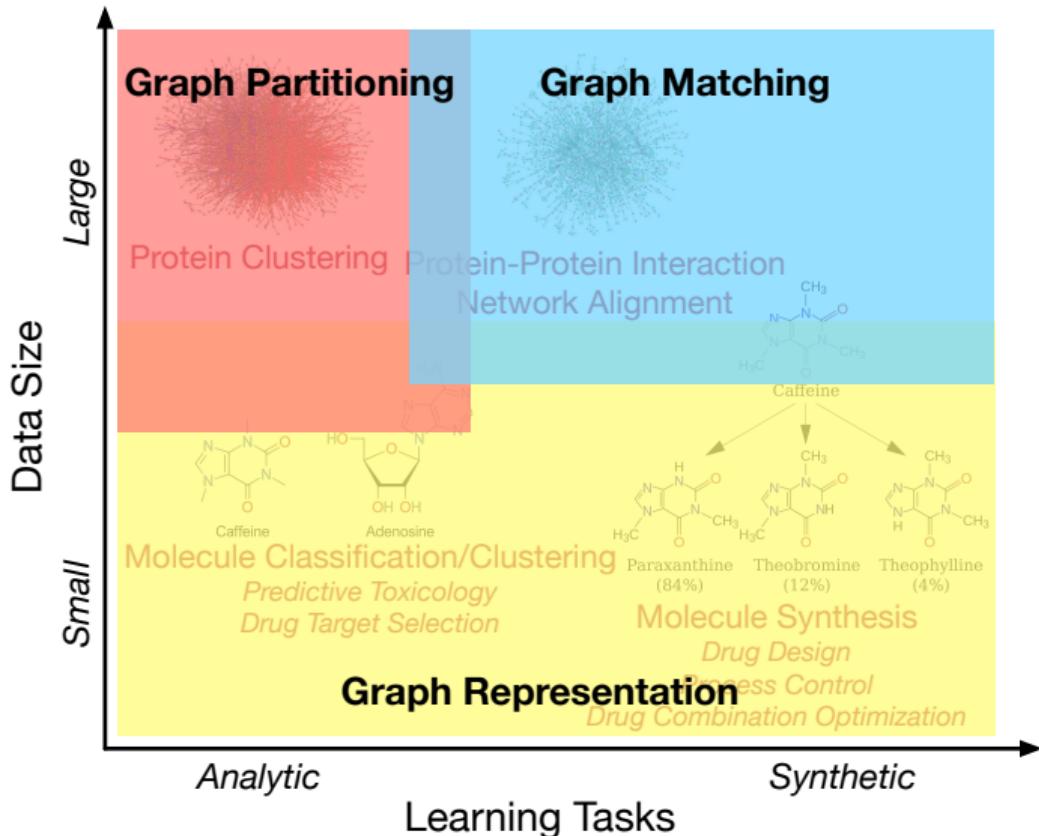
The space of problems



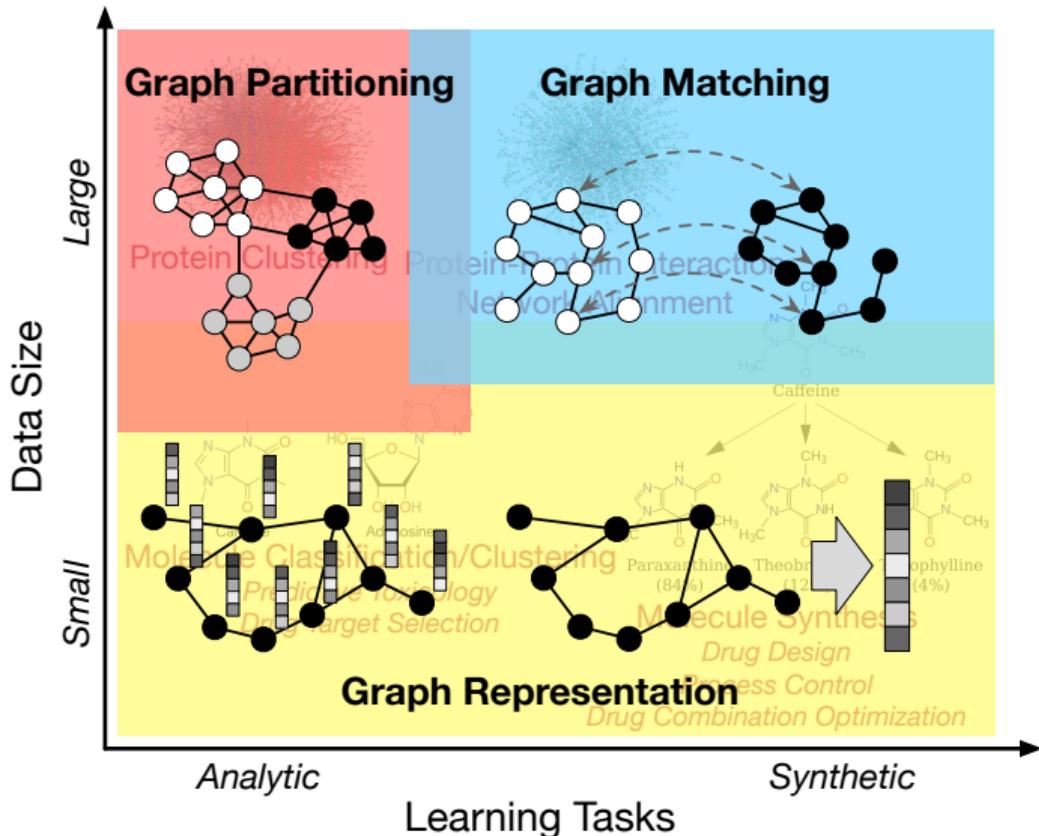
The space of problems



The space of problems



The space of problems



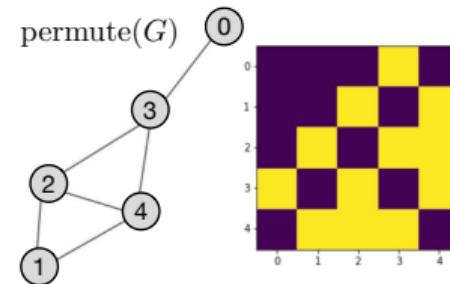
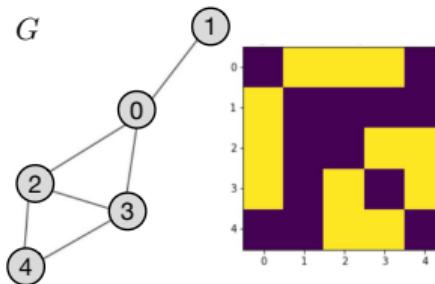
The challenges of the problems

The challenges of the problems

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

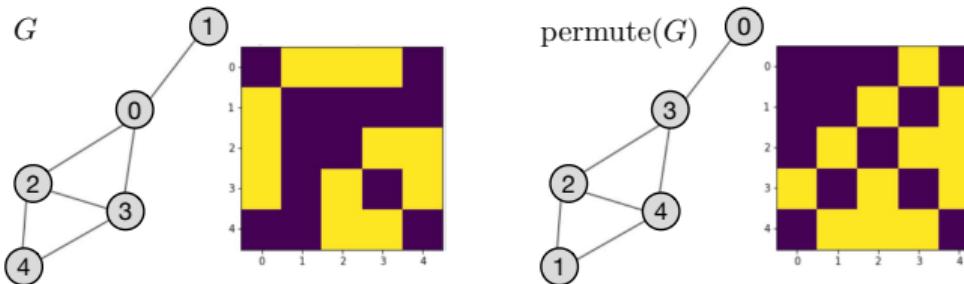
The challenges of the problems

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



The challenges of the problems

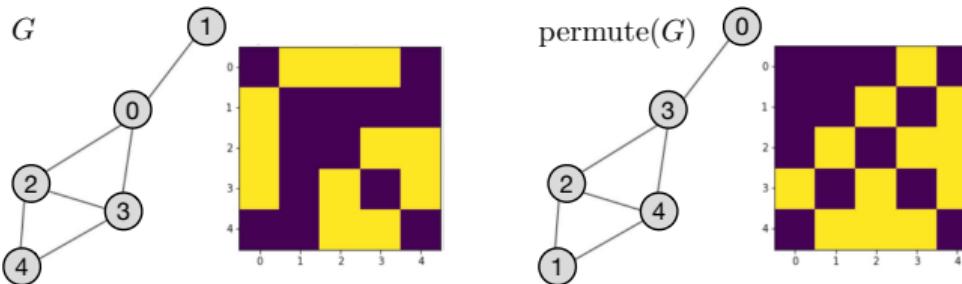
- ▶ NP-completeness
 - ▶ Approximation algorithms with high stability and scalability
- ▶ Permutation invariance



- ▶ 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

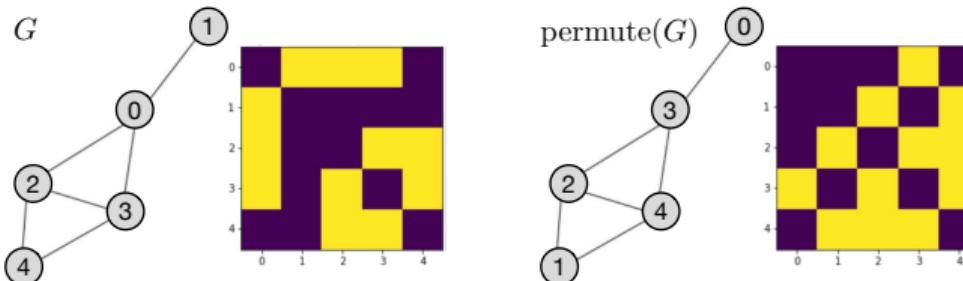
- ▶ NP-completeness
 - ▶ Approximation algorithms with high stability and scalability
- ▶ Permutation invariance



- ▶ 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

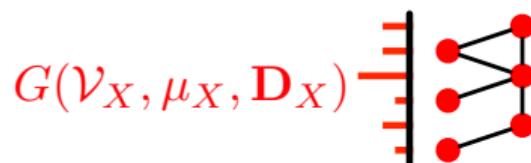
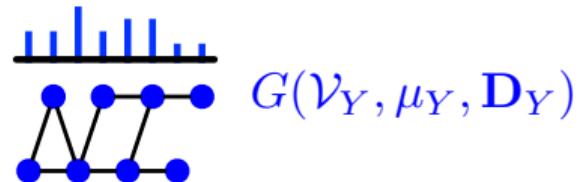


- ▶ A permutation-invariant metric d : $d(G_X, G_Y) = d(G_X, \text{permute}(G_Y))$
- ▶ A permutation-invariant representation model f : $f(G) = f(\text{permute}(G))$
- ▶ (Often) No labels
 - ▶ Unsupervised or semi-supervised learning

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

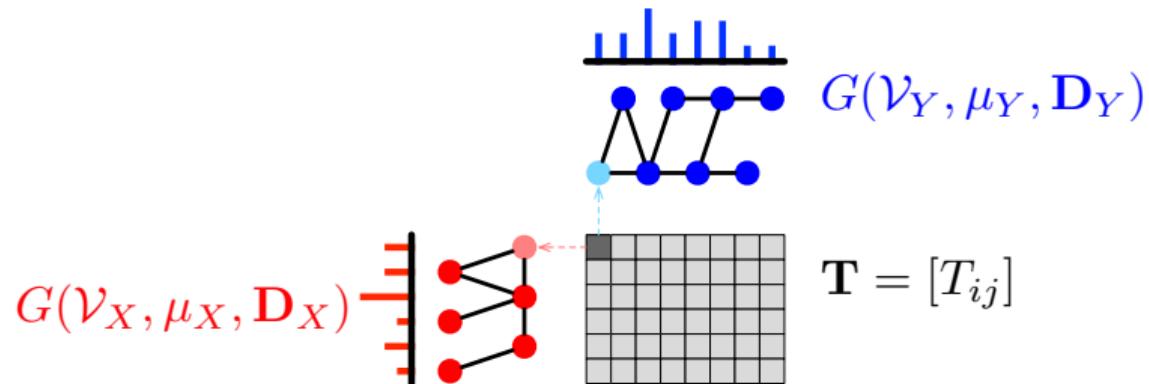
Gromov-Wasserstein distance (GWD) for structured data

Gromov-Wasserstein distance (GWD) for structured data



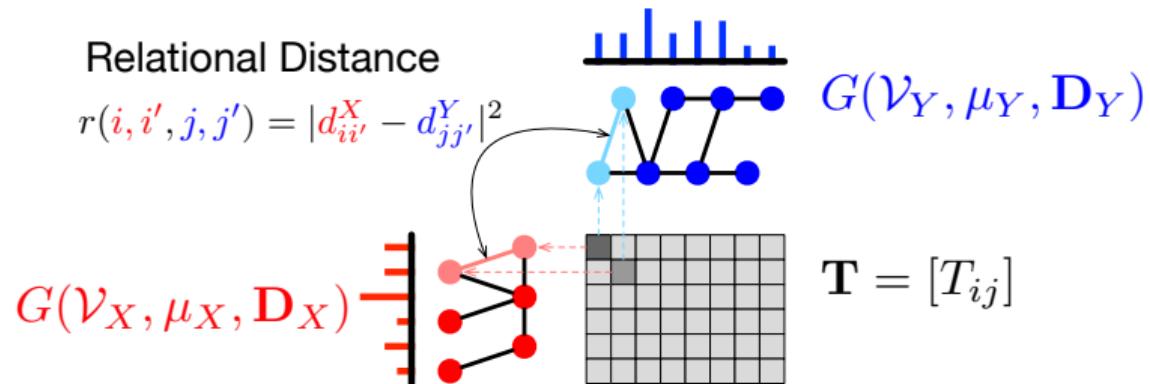
- ▶ \mathcal{V} : the node set
- ▶ μ : a predefined distribution of nodes
- ▶ $\mathbf{D} = [d_{ii'}]$: the adjacency / distance / kernel matrix

Gromov-Wasserstein distance (GWD) for structured data



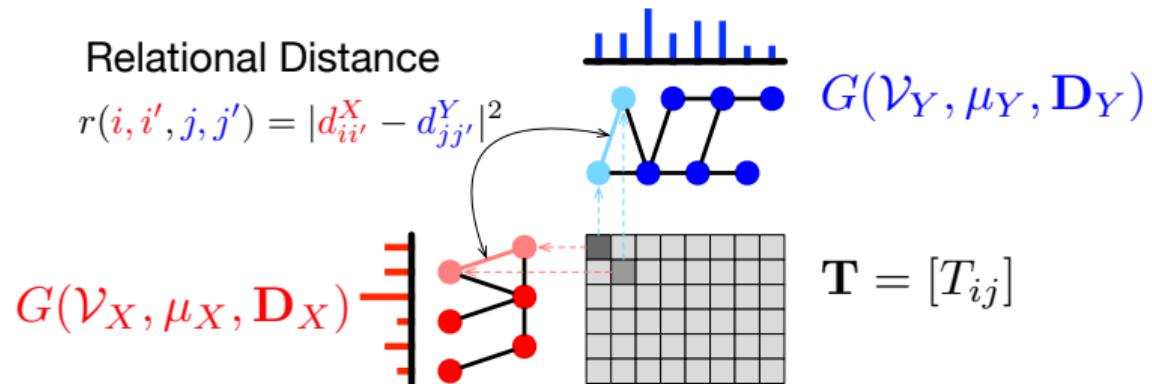
- ▶ $\mathbf{T} = [T_{ij}]$: a joint distribution of nodes
- ▶ $(i \in \mathcal{V}_X, j \in \mathcal{V}_Y) \sim \mathbf{T}$.
- ▶ $\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \boldsymbol{\mu}_Y) = \{\mathbf{T} \geq \mathbf{0} \mid \mathbf{T}\mathbf{1} = \boldsymbol{\mu}_X, \mathbf{T}^\top \mathbf{1} = \boldsymbol{\mu}_Y\}$

Gromov-Wasserstein distance (GWD) for structured data



- ▶ $\underbrace{\mathbf{T} \otimes \mathbf{T}}$: 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.

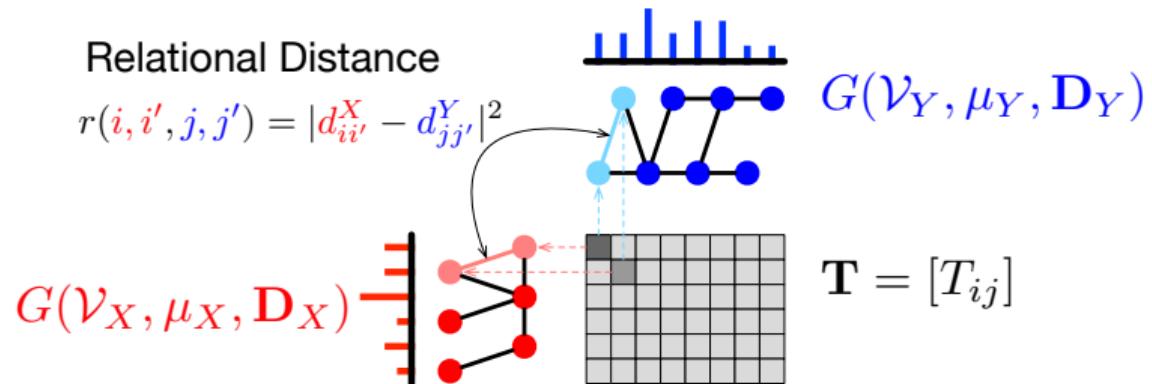
Gromov-Wasserstein distance (GWD) for structured data



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) for structured data

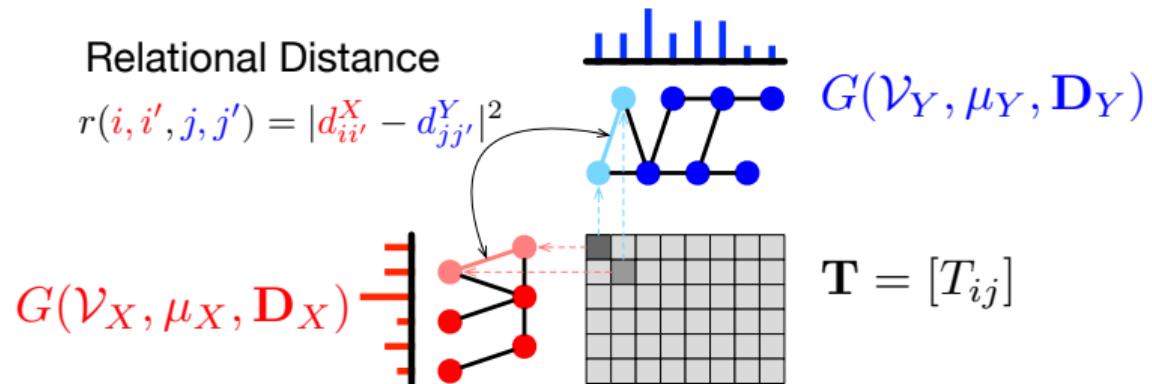


The GWD is **the minimum expectation of the relational distance**:

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

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

Gromov-Wasserstein distance (GWD) for structured data

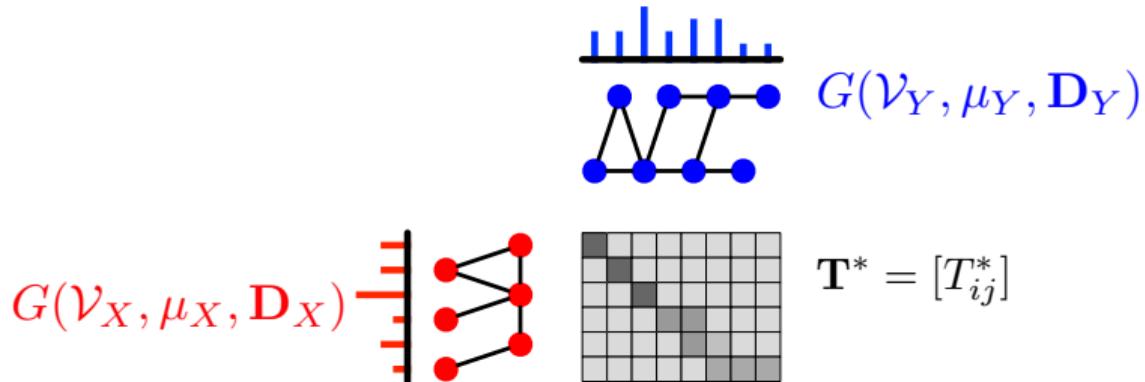


The GWD is **the minimum expectation of the relational distance**:

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

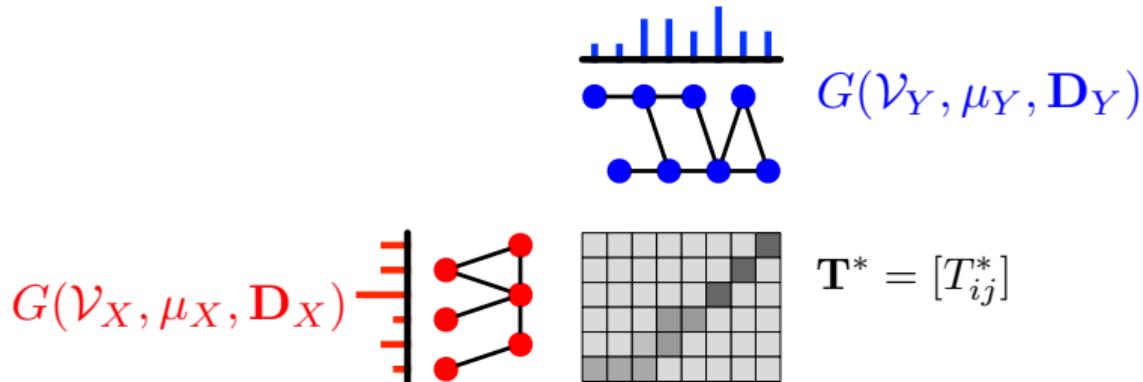
- ▶ $\mathbf{D}_{XY} = (\mathbf{D}_X \odot \mathbf{D}_Y) \mu_X \mathbf{1}_{|\mathcal{V}_Y|}^\top + \mathbf{1}_{|\mathcal{V}_X|} \mu_Y^\top (\mathbf{D}_Y \odot \mathbf{D}_Y)^\top.$
- ▶ Given comparable node attributes, $\mathbf{D}_{XY} \leftarrow \mathbf{D}_{XY} + \mathbf{D}(\mathbf{F}_X, \mathbf{F}_Y)$

Advantages of GWD



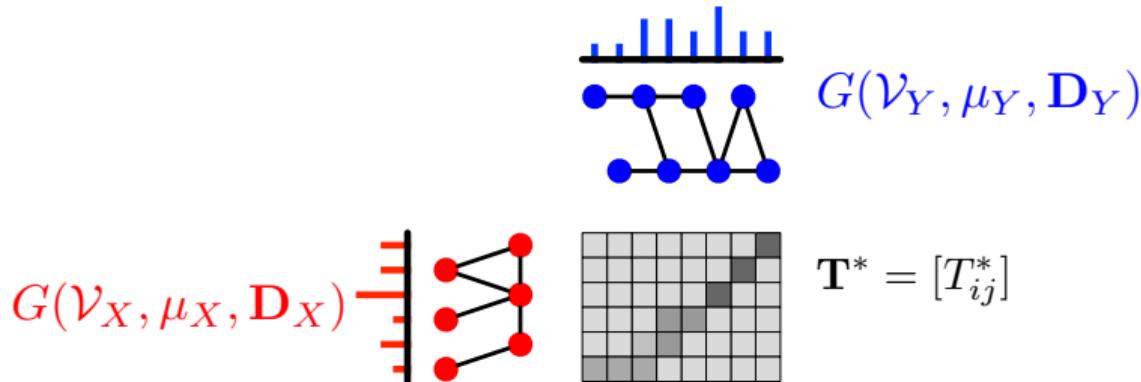
- ▶ The optimal joint distribution \mathbf{T}^* (or called “**optimal transport**” matrix) indicates the correspondence between the two graphs.

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.

Straightforward applications

Graph matching and partitioning

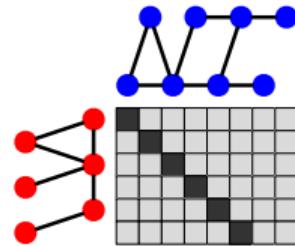
Apply the Gromov-Wasserstein distance as objective functions

Matching via learning optimal transport

Quadratic assignment problem (QAP):

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

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

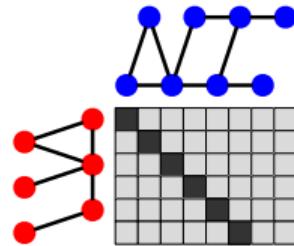


Matching via learning optimal transport

Quadratic assignment problem (QAP):

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

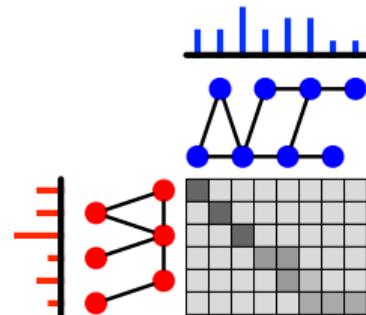
$$\mathcal{P} = \{P \in \{0, 1\}^{|\mathcal{V}_X| \times |\mathcal{V}_Y|} \mid P\mathbf{1} = \mathbf{1}, P^\top \mathbf{1} \leq \mathbf{1}\}.$$



Gromov-Wasserstein distance (GWD):

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

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

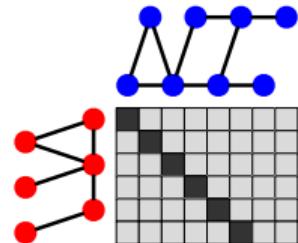


Matching via learning optimal transport

Quadratic assignment problem (QAP):

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

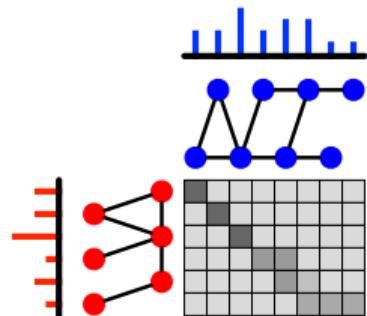
$$\mathcal{P} = \{P \in \{0, 1\}^{|\mathcal{V}_X| \times |\mathcal{V}_Y|} \mid P\mathbf{1} = \mathbf{1}, P^\top \mathbf{1} \leq \mathbf{1}\}.$$



Gromov-Wasserstein distance (GWD):

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

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



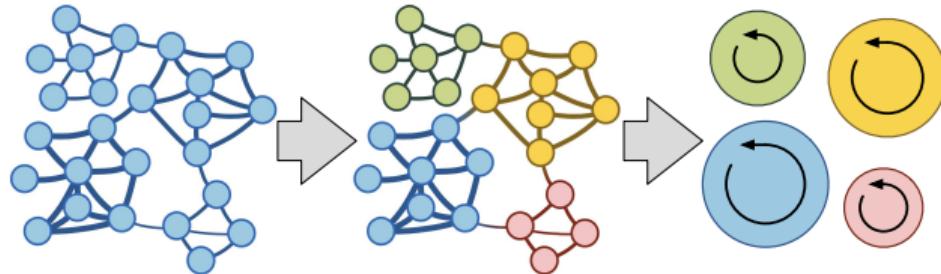
- ▶ Conditional prob. $P^* = \frac{T^*}{\mu_X \mathbf{1}^\top} = [P^*(j|i)]$.
- ▶ For each node $i \in \mathcal{V}_X$, $j^* = \arg \max_j P^*(j|i)$.

Partitioning is also matching

Partitioning is also matching

Modularity maximization principle

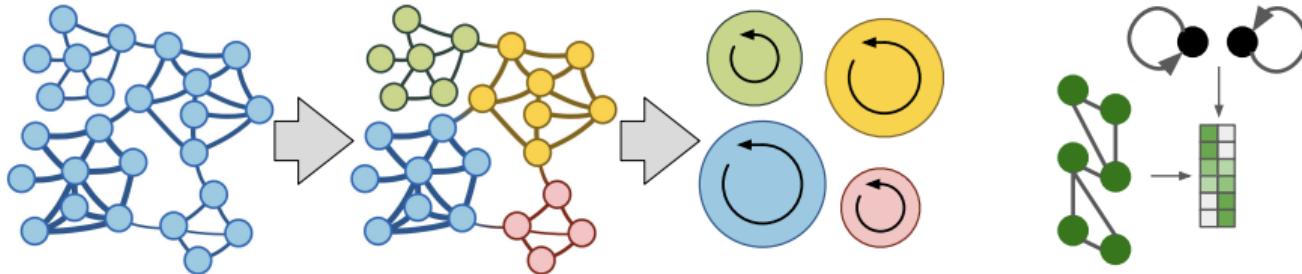
- ▶ Dense internal edges + sparse external edges.



Partitioning is also matching

Modularity maximization principle

- ▶ Dense internal edges + sparse external edges.



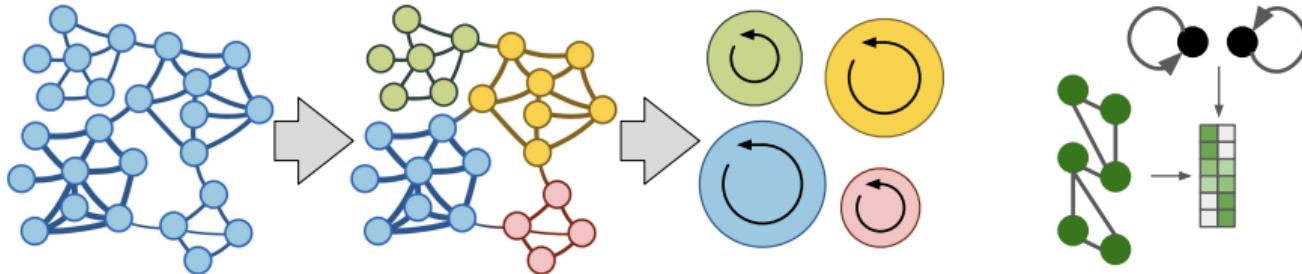
A GWD-based solution [*Xu, et al., NeurIPS 2019*]:

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

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 algorithms

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

Proposed algorithms

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

Proximal Gradient Algorithm [Xu, et al., ICML 2019]

$$\min_{T \in \Pi(\mu_X, \mu_Y)} \underbrace{\langle D_{XY} - 2D_X \mathbf{T}^{(m)} D_Y^\top, T \rangle}_{\text{GW term}} + \gamma \underbrace{\text{KL}(T \| \mathbf{T}^{(m)})}_{\text{Proximal term}}$$

Proposed algorithms

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

Proximal Gradient Algorithm [Xu, et al., ICML 2019]

$$\begin{aligned} & \min_{T \in \Pi(\mu_X, \mu_Y)} \underbrace{\langle D_{XY} - 2D_X \mathbf{T}^{(m)} D_Y^\top, T \rangle}_{\text{GW term}} + \gamma \underbrace{\text{KL}(T \| \mathbf{T}^{(m)})}_{\text{Proximal term}} \\ \rightarrow & \underbrace{\min_{T \in \Pi(\mu_X, \mu_Y)} \langle D_{XY} - 2D_X \mathbf{T}^{(m)} D_Y^\top - \gamma \log \mathbf{T}^{(m)}, T \rangle}_{\text{Sinkhorn distance}} - \gamma \underbrace{\mathsf{H}(T)}_{\text{Entropy}} \end{aligned} \quad (2)$$

Proposed algorithms

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

Proximal Gradient Algorithm [Xu, et al., ICML 2019]

$$\begin{aligned} & \min_{T \in \Pi(\mu_X, \mu_Y)} \underbrace{\langle D_{XY} - 2D_X \mathbf{T}^{(m)} D_Y^\top, T \rangle}_{\text{GW term}} + \gamma \underbrace{\text{KL}(T \| \mathbf{T}^{(m)})}_{\text{Proximal term}} \\ \rightarrow & \underbrace{\min_{T \in \Pi(\mu_X, \mu_Y)} \langle D_{XY} - 2D_X \mathbf{T}^{(m)} D_Y^\top - \gamma \log \mathbf{T}^{(m)}, T \rangle}_{\text{Constant}} - \gamma \underbrace{\mathsf{H}(T)}_{\text{Entropy}} \end{aligned} \quad (2)$$

Sinkhorn distance

Bregman ADMM [Xu, AAAI 2020]

$$\underbrace{\min_{\substack{T \in \Pi(\mu_X, \cdot), S \in \Pi(\cdot, \mu_Y), \mathbf{T} = \mathbf{S}}} \langle D_{XY} - 2D_X S D_Y^\top, T \rangle}_{\text{Decoupling constraints by Aux.}}$$

Proposed algorithms

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

Proximal Gradient Algorithm [Xu, et al., ICML 2019]

$$\begin{aligned} & \min_{T \in \Pi(\mu_X, \mu_Y)} \underbrace{\langle D_{XY} - 2D_X \mathbf{T}^{(m)} D_Y^\top, T \rangle}_{\text{GW term}} + \gamma \underbrace{\text{KL}(T \| \mathbf{T}^{(m)})}_{\text{Proximal term}} \\ \rightarrow & \min_{T \in \Pi(\mu_X, \mu_Y)} \underbrace{\langle D_{XY} - 2D_X \mathbf{T}^{(m)} D_Y^\top - \gamma \log \mathbf{T}^{(m)}, T \rangle}_{\text{Constant}} - \gamma \underbrace{\mathsf{H}(T)}_{\text{Entropy}} \\ & \qquad \qquad \qquad \underbrace{\qquad \qquad \qquad}_{\text{Sinkhorn distance}} \end{aligned} \quad (2)$$

Bregman ADMM [Xu, AAAI 2020]

$$\begin{aligned} & \underbrace{\min_{\substack{T \in \Pi(\mu_X, \cdot), S \in \Pi(\cdot, \mu_Y), \mathbf{T} = S}}}_{\text{Decoupling constraints by Aux.}} \langle D_{XY} - 2D_X S D_Y^\top, \mathbf{T} \rangle \\ \rightarrow & \min_{\substack{T \in \Pi(\mu_X, \cdot), S \in \Pi(\cdot, \mu_Y), Z}} \langle D_{XY} - 2D_X S D_Y^\top, \mathbf{T} \rangle + \underbrace{\langle Z, \mathbf{T} - S \rangle}_{\text{Augmented Lagrangian}} + \underbrace{\gamma \mathcal{B}(\mathbf{T}, S)}_{\text{Bregman Div.}} \end{aligned} \quad (3)$$

Further analysis

Convergence

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

Further analysis

Convergence

- ▶ $\lim_{m \rightarrow \infty} T^{(m)}$ is a stationary point.
- ▶ Linear convergence.
- ▶ PGA requires fewer iterations.
- ▶ B-ADMM works better on directed graphs.

Further analysis

Convergence

- ▶ $\lim_{m \rightarrow \infty} \mathbf{T}^{(m)}$ is a stationary point.
- ▶ Linear convergence.
- ▶ PGA requires fewer iterations.
- ▶ B-ADMM works better on directed graphs.

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: $\mathcal{O}(V^3)$.

Further analysis

Convergence

- ▶ $\lim_{m \rightarrow \infty} \mathbf{T}^{(m)}$ is a stationary point.
- ▶ Linear convergence.
- ▶ PGA requires fewer iterations.
- ▶ B-ADMM works better on directed graphs.

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: $\mathcal{O}(VE)$.

Further analysis

Convergence

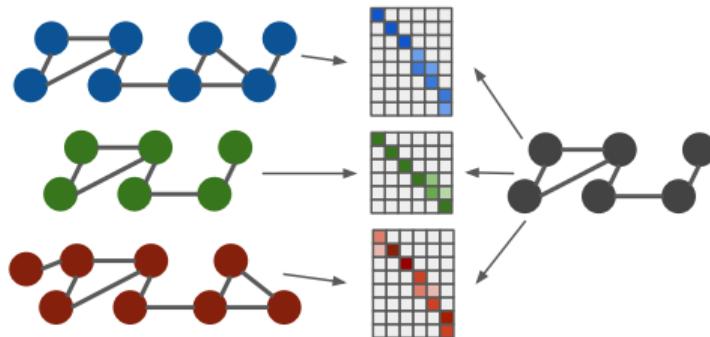
- ▶ $\lim_{m \rightarrow \infty} \mathbf{T}^{(m)}$ is a stationary point.
- ▶ Linear convergence.
- ▶ PGA requires fewer iterations.
- ▶ B-ADMM works better on directed graphs.

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): $\mathcal{O}(N(E + V))$.

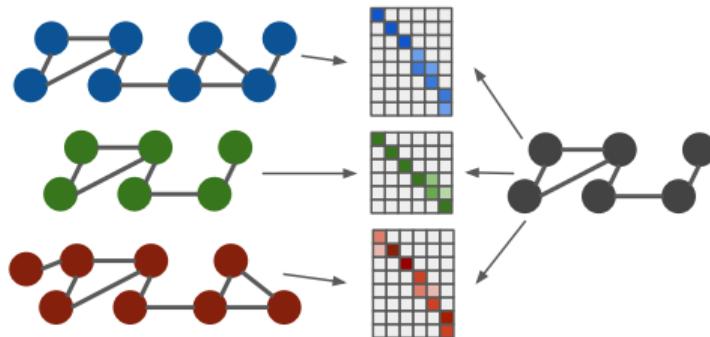
Large-scale matching based on GW barycenters



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

$$\underbrace{B_{gw}(\bar{\mathcal{V}}, \bar{\boldsymbol{\mu}}, \bar{\mathbf{B}}^*)}_{\text{Barycenter graph}}, \quad \underbrace{\{\mathbf{T}_k^*\}_{k=1}^K}_{\text{OT matrices}} := \arg \min_B \sum_{k=1}^K \lambda_k d_{gw}(B, G_k), \quad (4)$$

Large-scale matching based on GW barycenters



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

$$\underbrace{B_{gw}(\bar{\mathcal{V}}, \bar{\boldsymbol{\mu}}, \mathbf{B}^*)}_{\text{Barycenter graph}}, \quad \underbrace{\{\mathbf{T}_k^*\}_{k=1}^K}_{\text{OT matrices}} := \arg \min_B \sum_{k=1}^K \lambda_k d_{gw}(B, G_k), \quad (4)$$

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

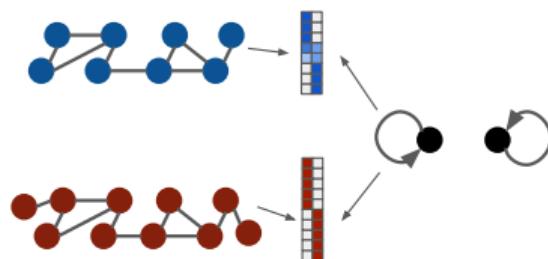
$$\mathbf{B}^* = \frac{1}{\bar{\boldsymbol{\mu}} \bar{\boldsymbol{\mu}}^\top} \sum_{k=1}^K \lambda_k (\mathbf{T}_k^*)^\top \mathbf{D}_k \mathbf{T}_k^* \quad (5)$$

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.

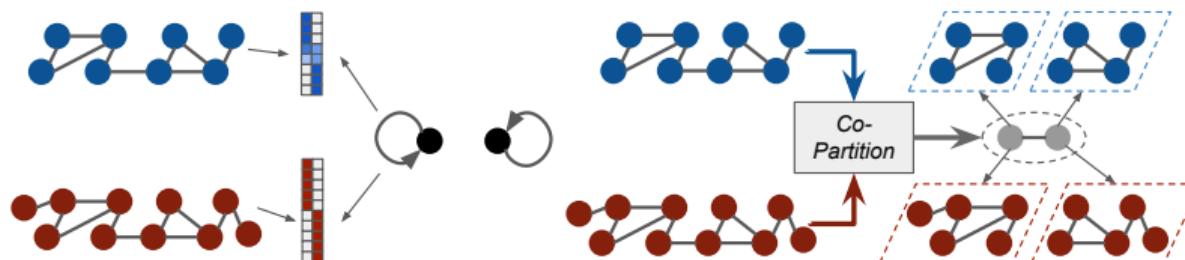


Large-scale matching based on GW barycenters

Co-partition two graphs:

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

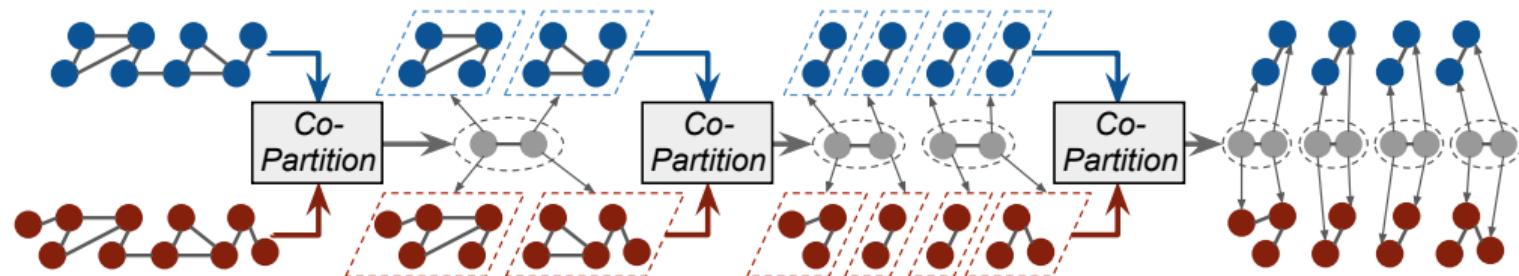
Initialize the barycenter graph by a disconnected graph.



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

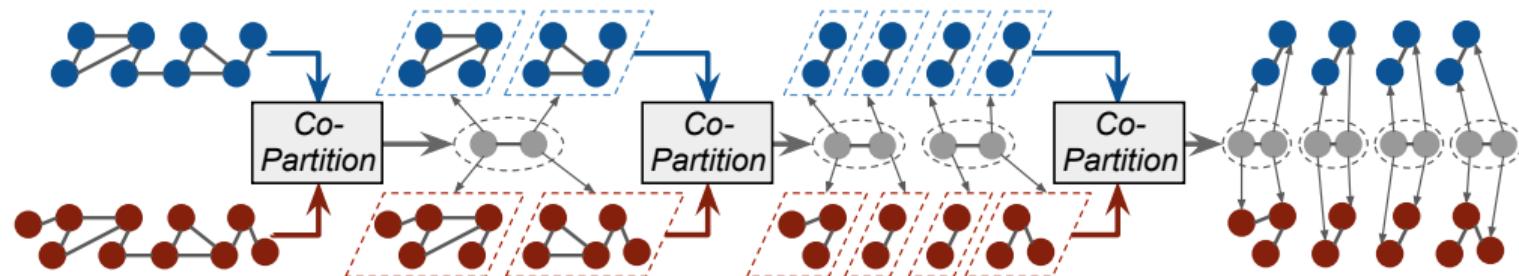
Large-scale matching based on GW barycenters

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



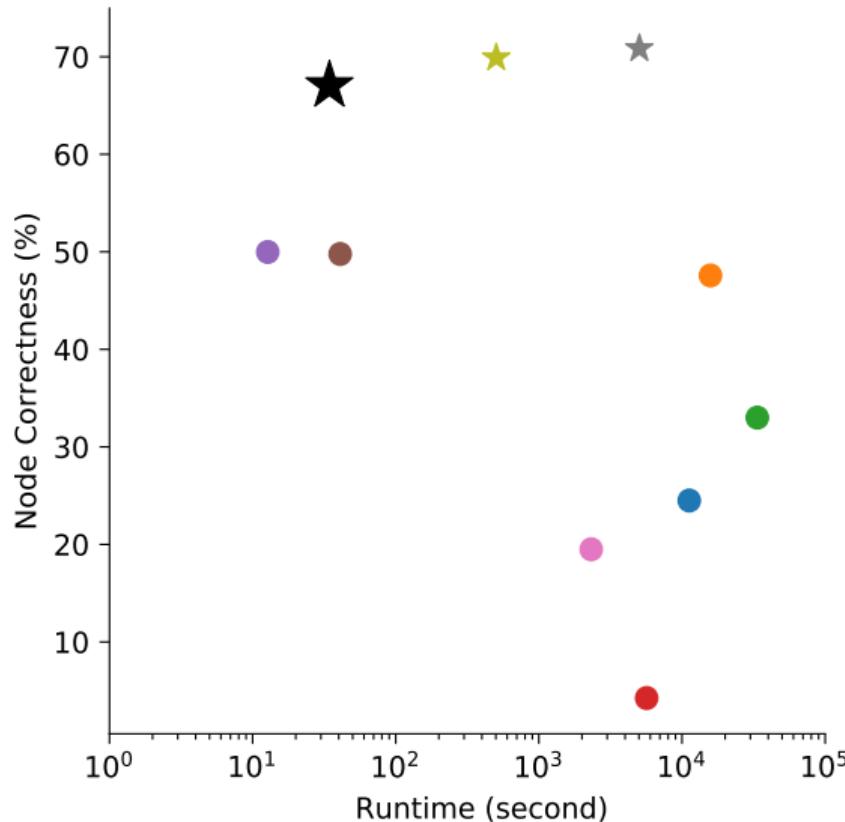
Large-scale matching based on GW barycenters

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

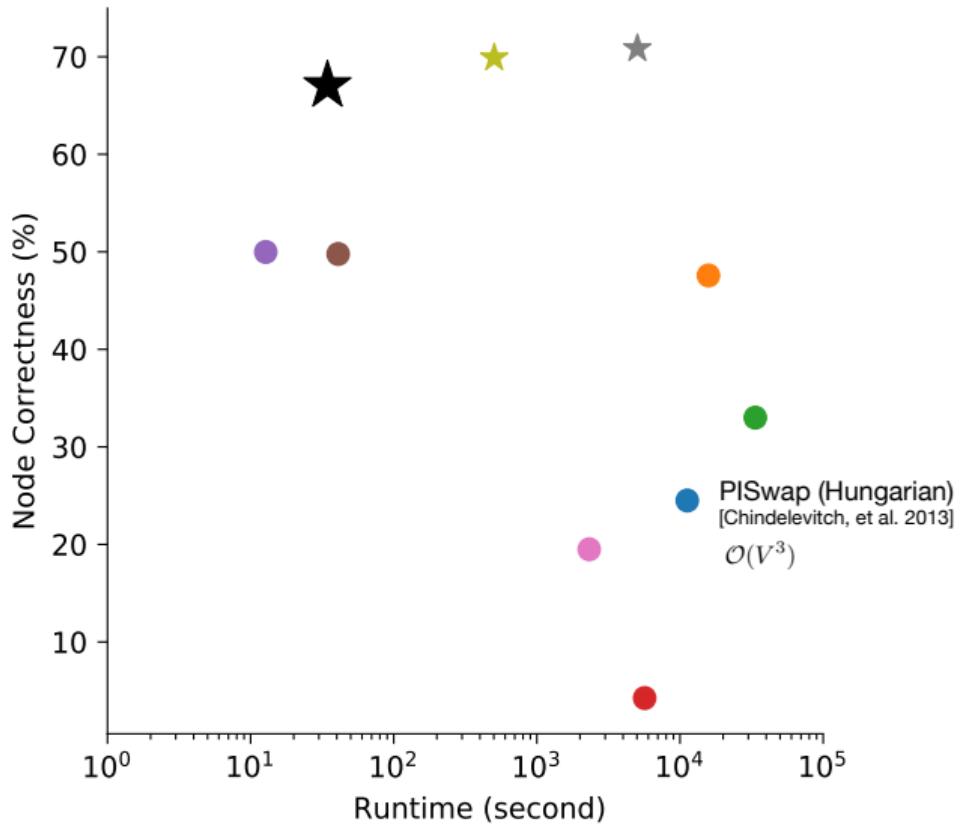


Computational complexity: $\mathcal{O}((E + V) \log V)$

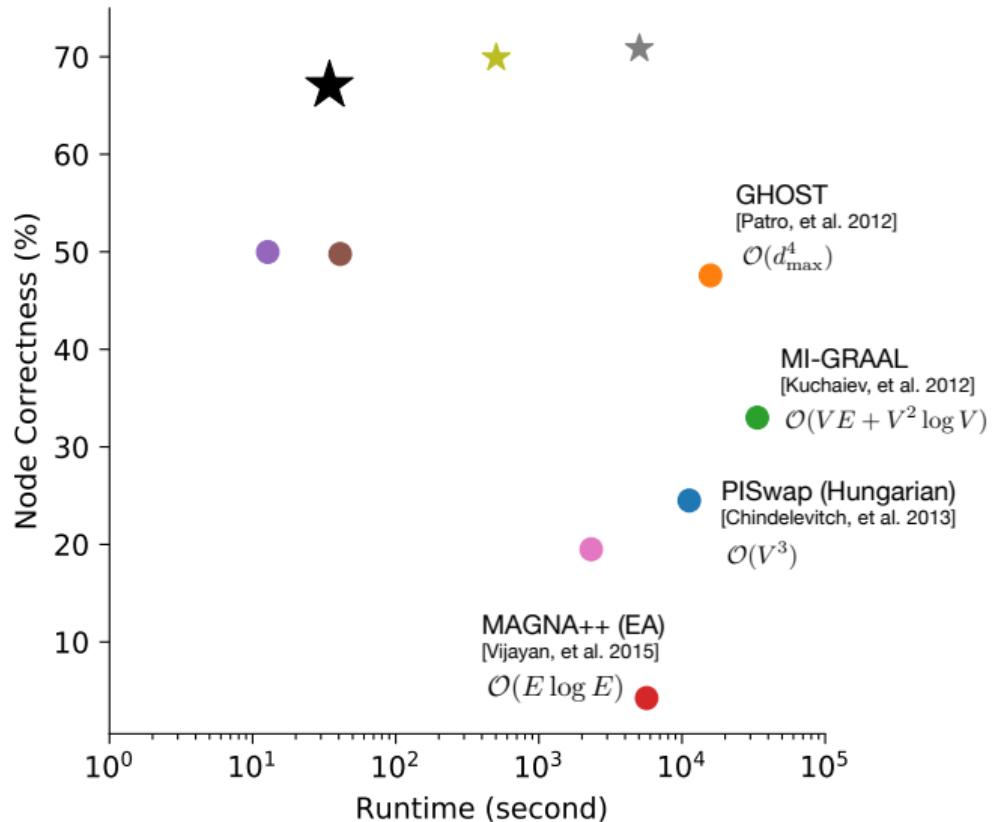
Matching synthetic graphs



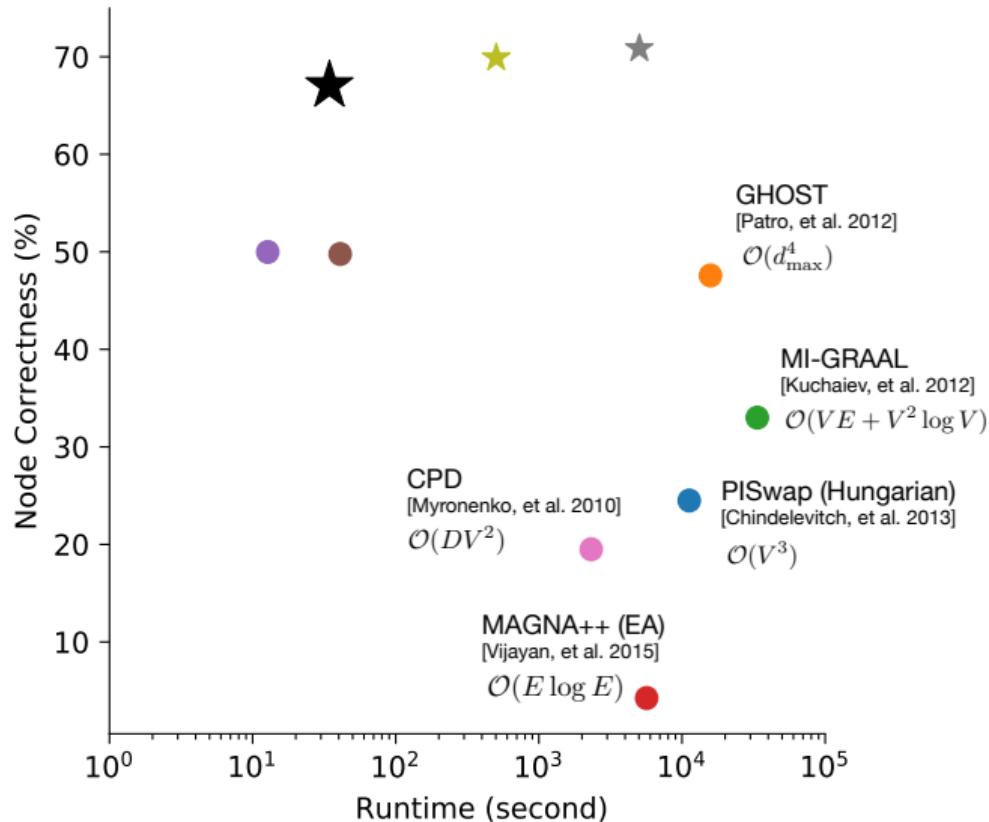
Matching synthetic graphs



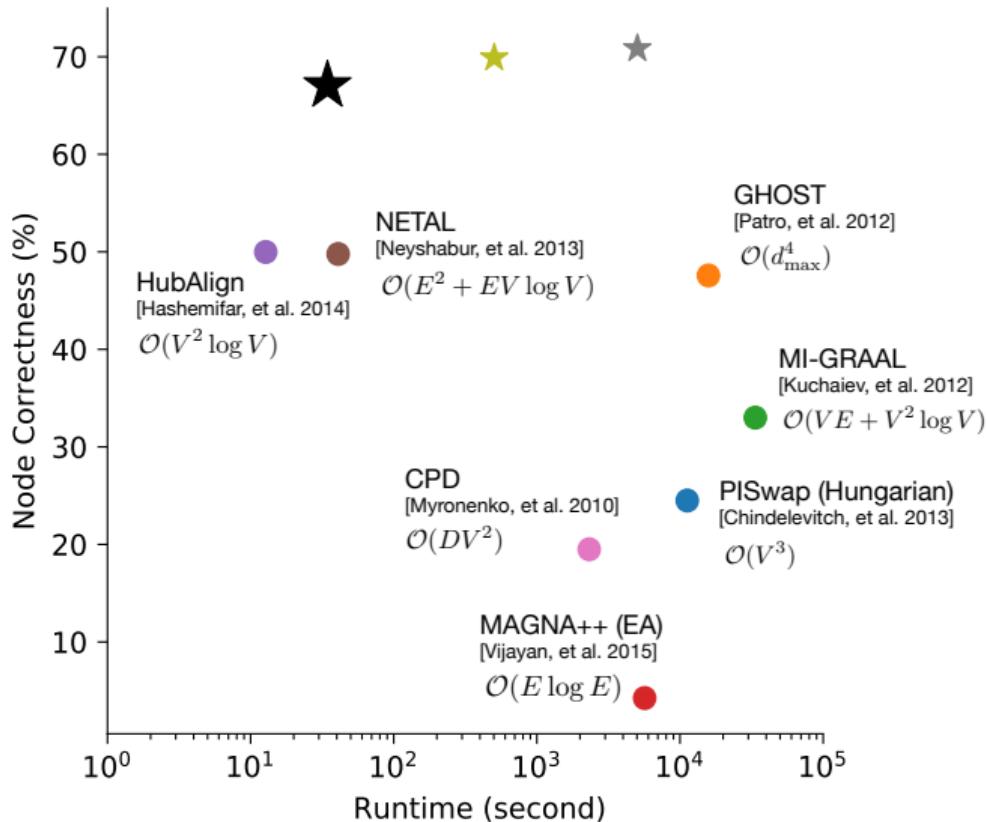
Matching synthetic graphs



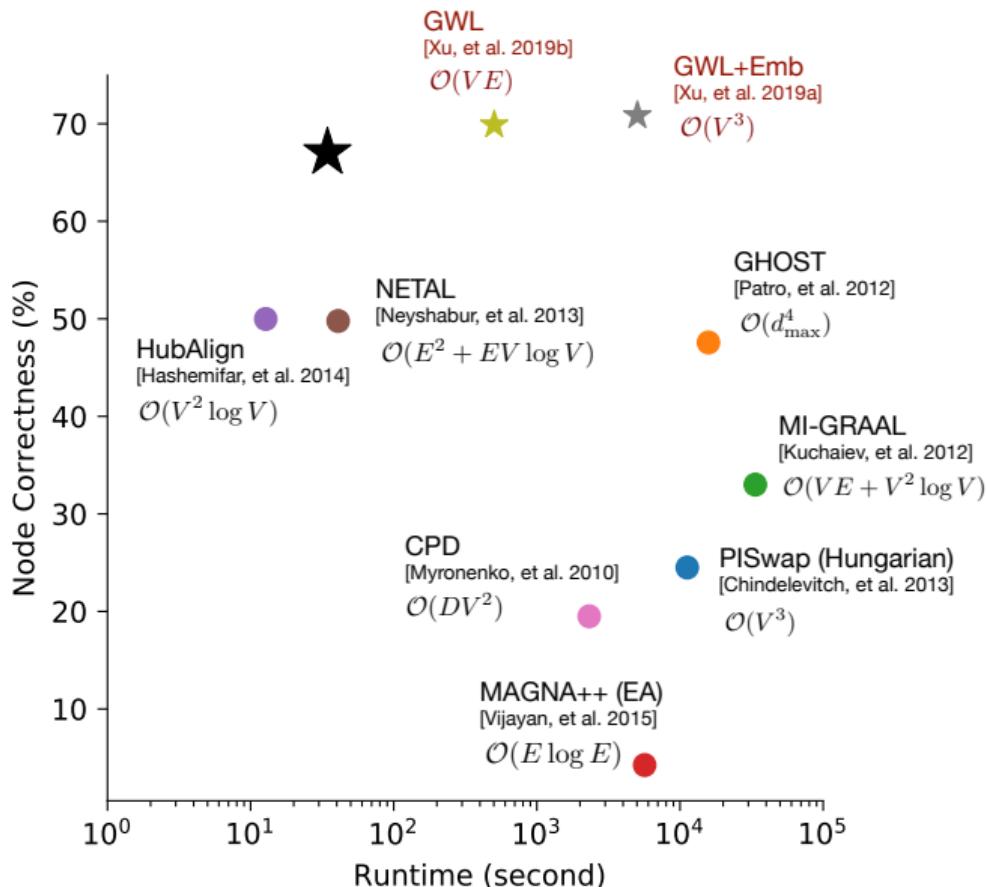
Matching synthetic graphs



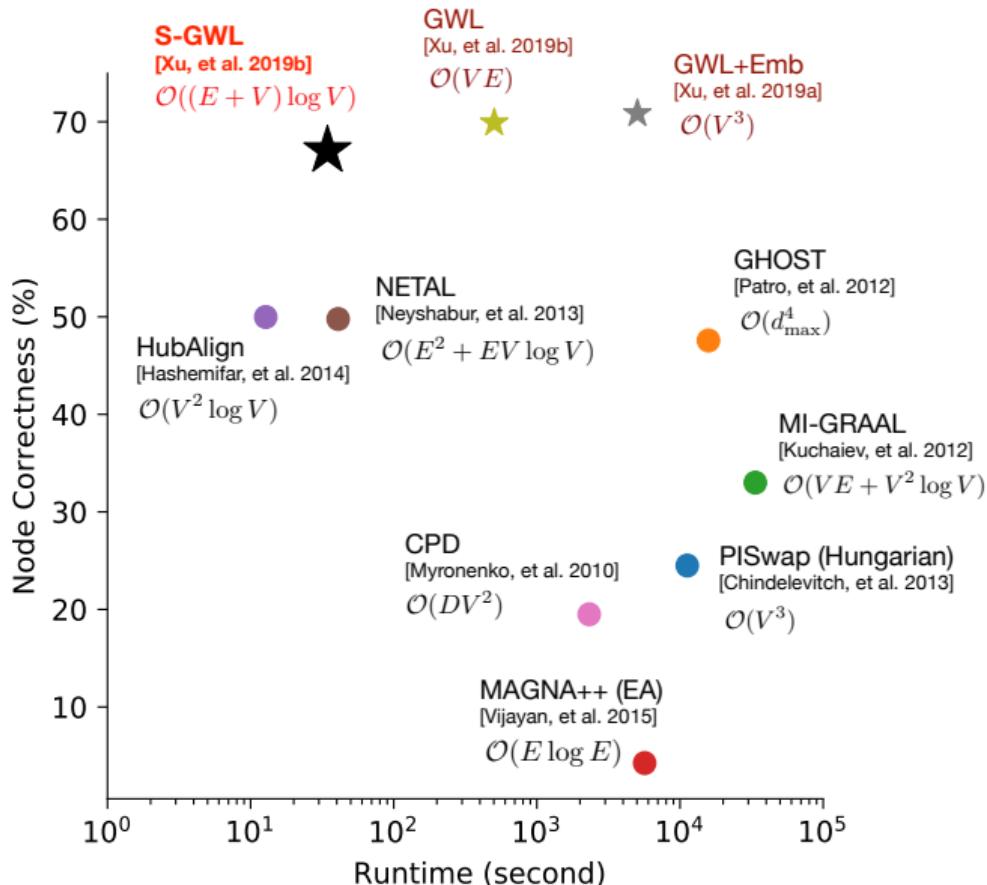
Matching synthetic graphs



Matching synthetic graphs



Matching synthetic graphs

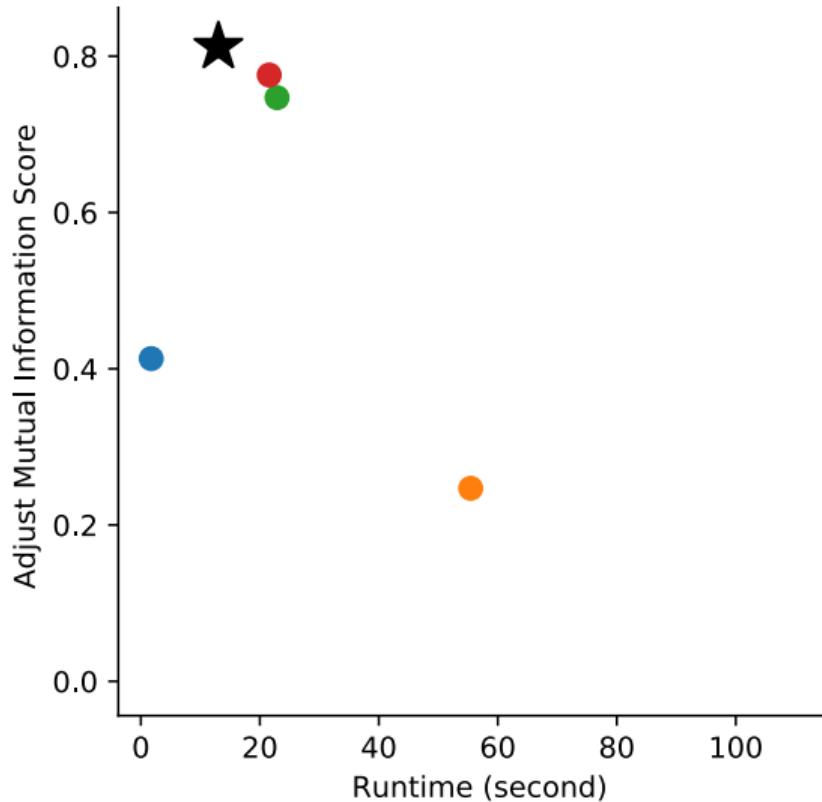
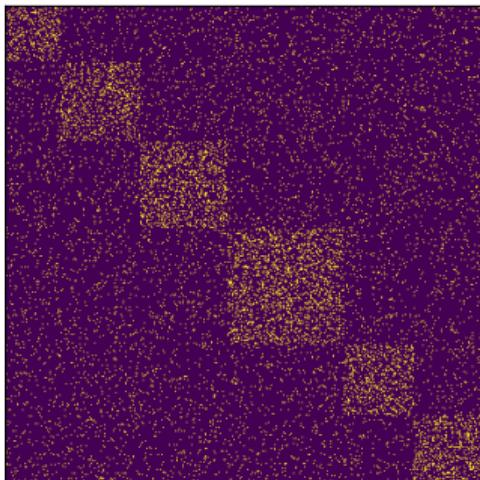


Partitioning synthetic graphs

$V = 4,000$

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

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

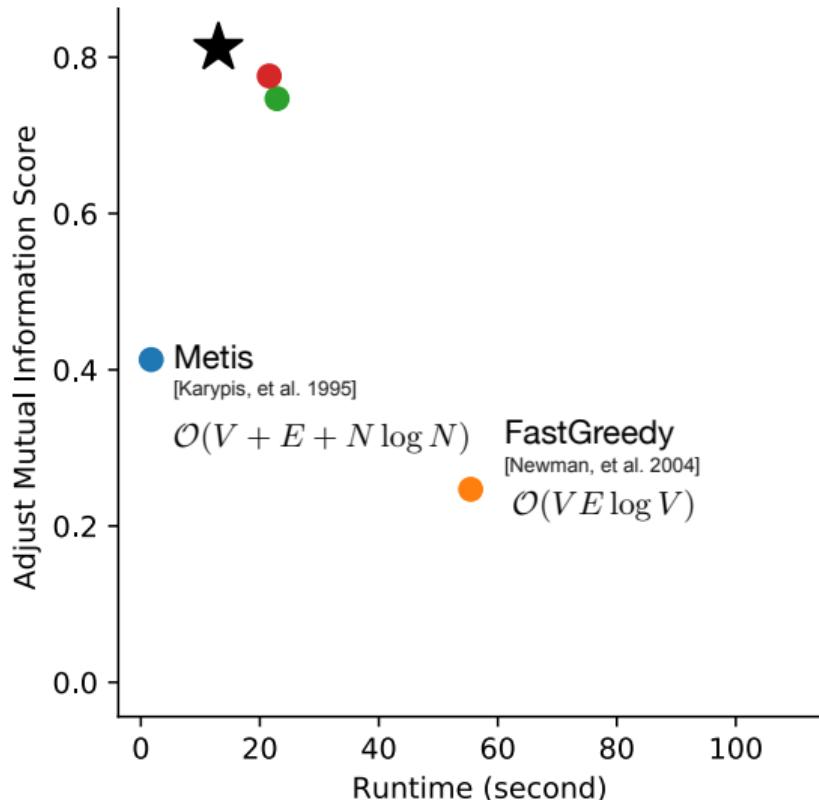
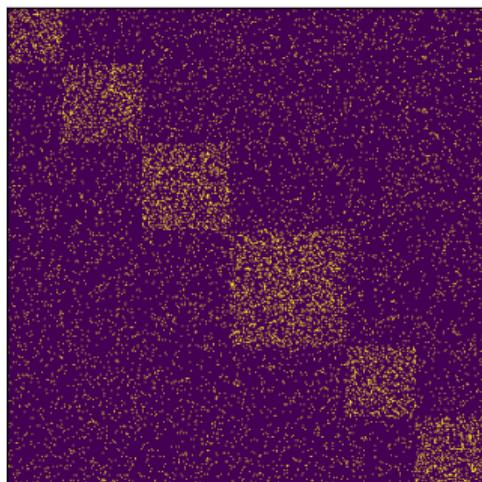


Partitioning synthetic graphs

$V = 4,000$

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

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

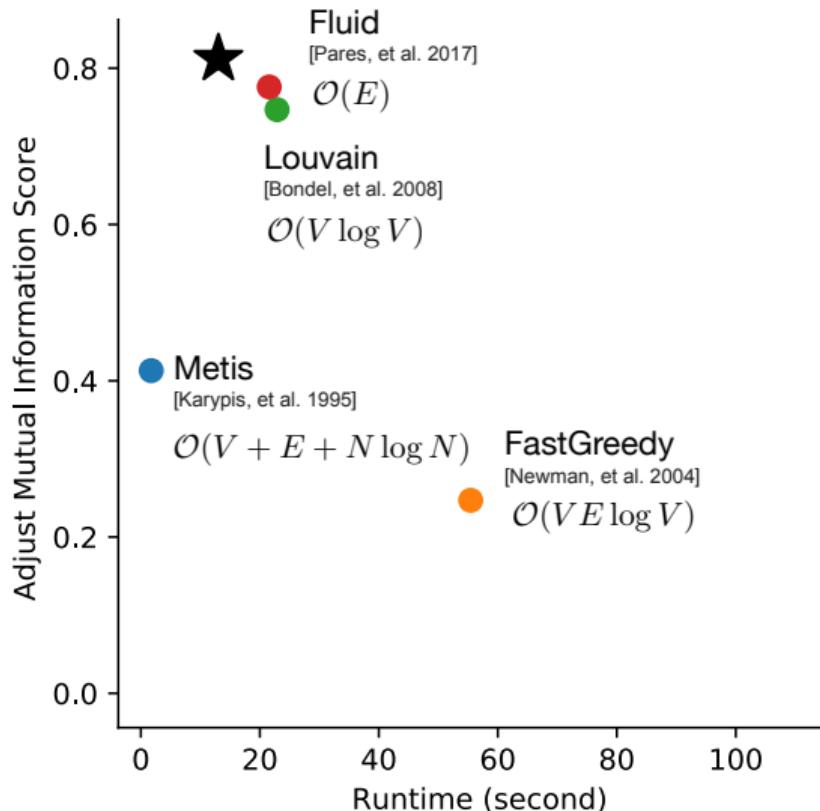
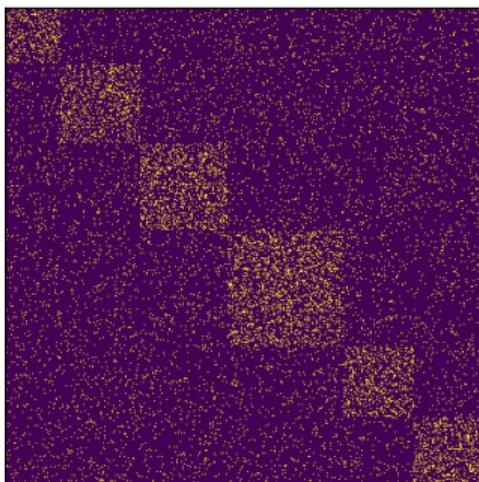


Partitioning synthetic graphs

$V = 4,000$

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

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

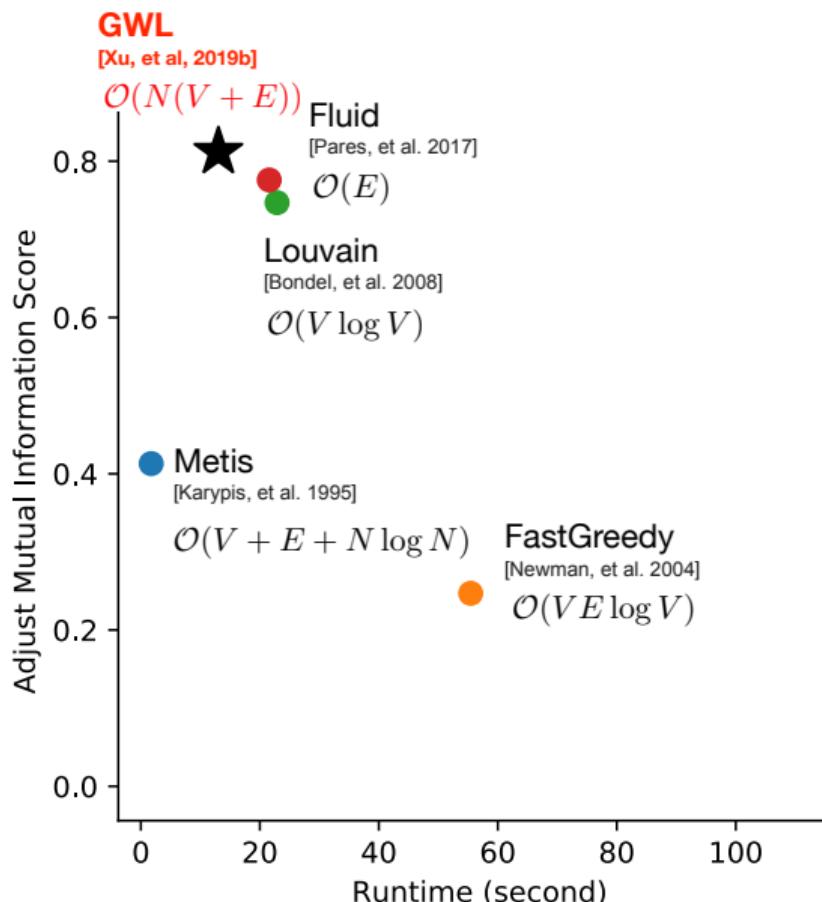
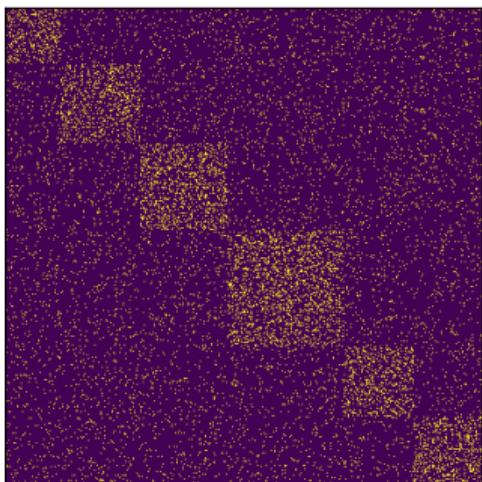


Partitioning synthetic graphs

$V = 4,000$

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

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

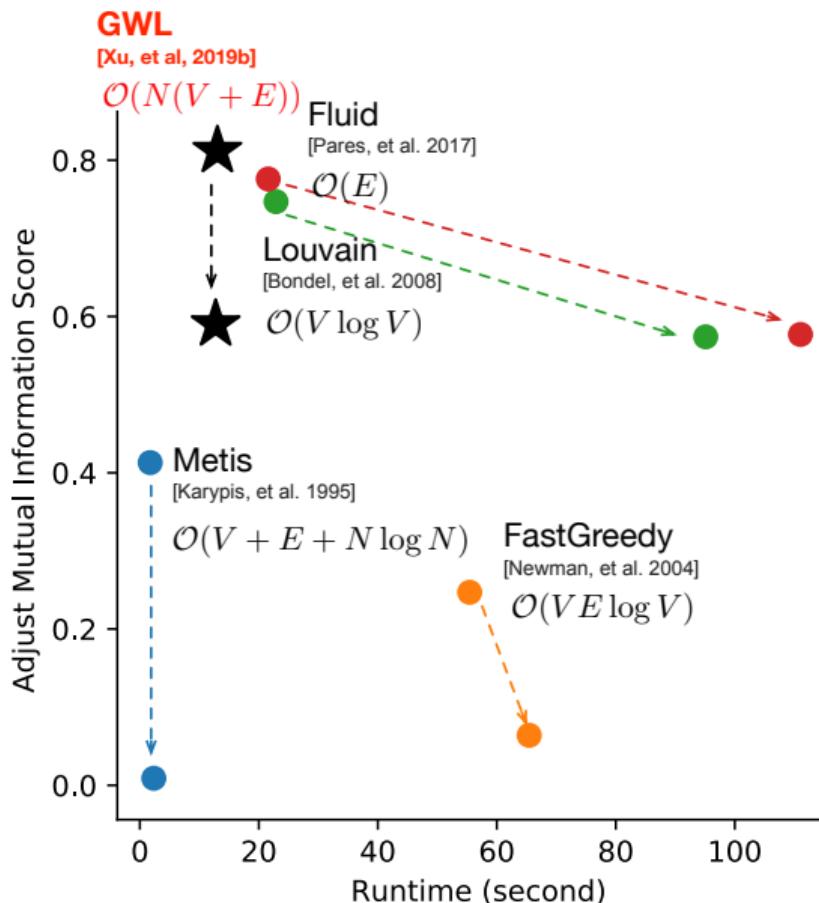
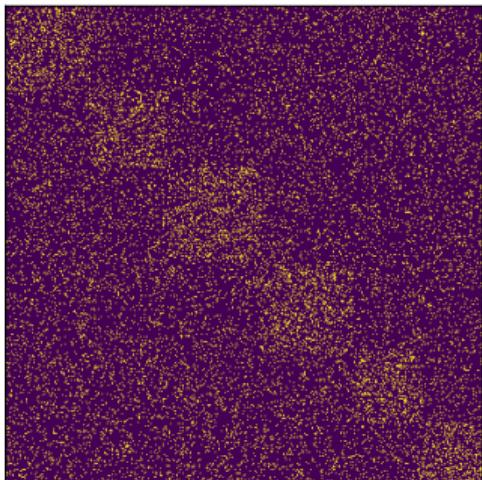


Partitioning synthetic graphs

$V = 4,000$

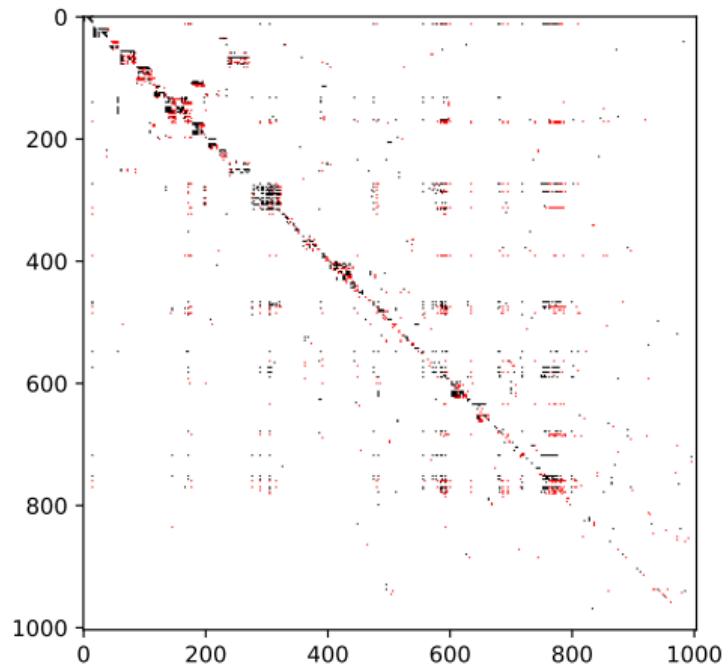
$p_{\text{within}} = 0.2$

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



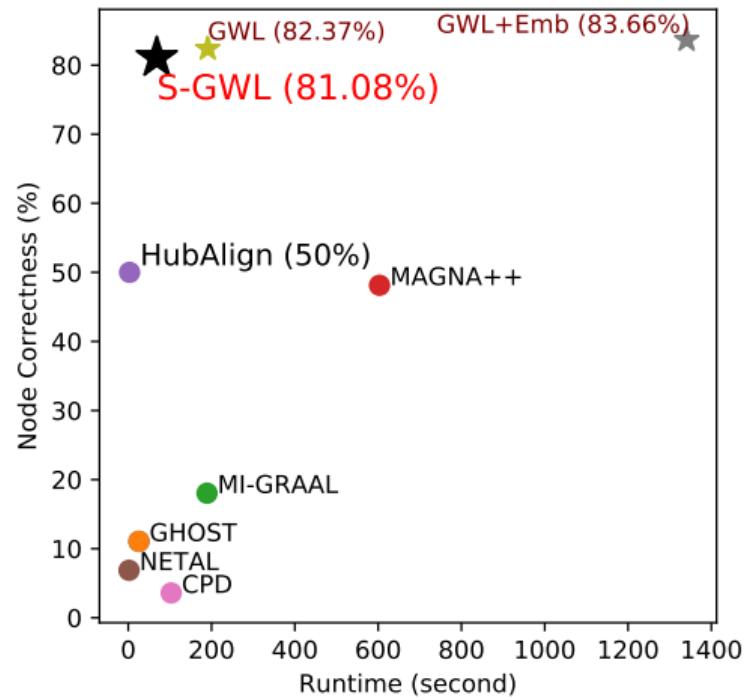
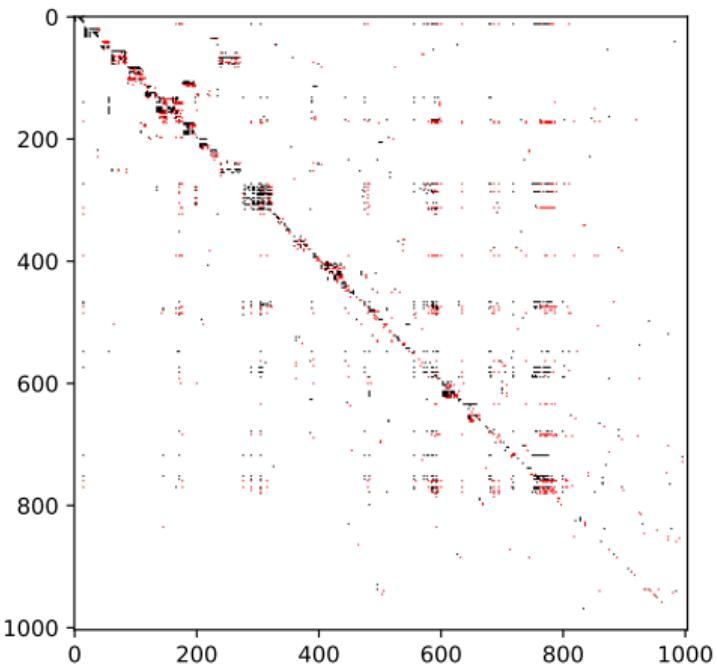
Real-world PPI network alignment

Yeast PPI \leftrightarrow Yeast PPI + 5% LC edges



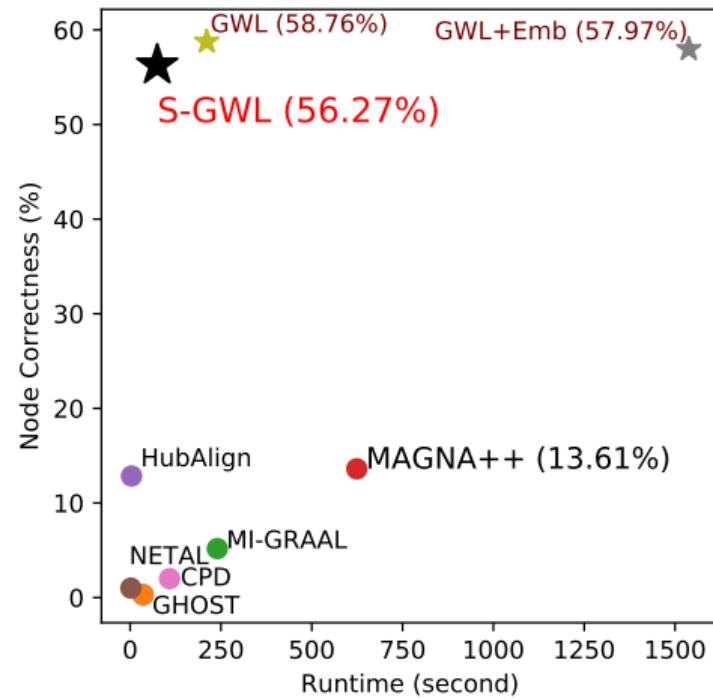
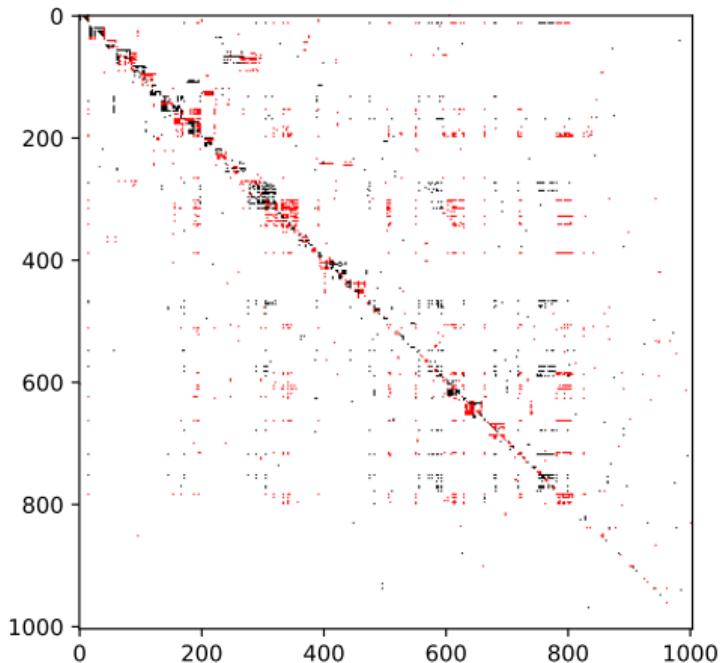
Real-world PPI network alignment

Yeast PPI \leftrightarrow Yeast PPI + 5% LC edges



Real-world PPI network alignment

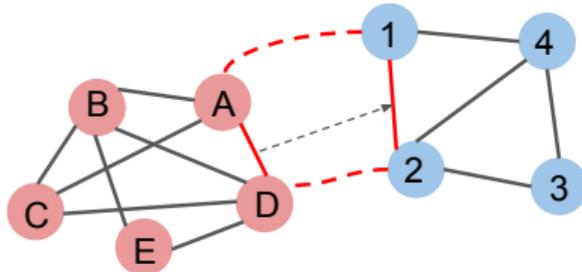
Yeast PPI \leftrightarrow Yeast PPI + 25% LC edges



Matching real-world PPI networks

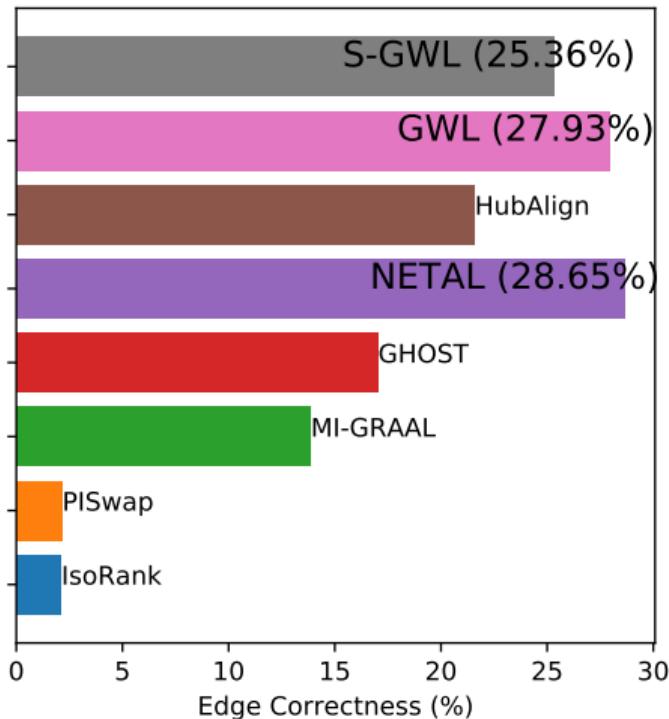
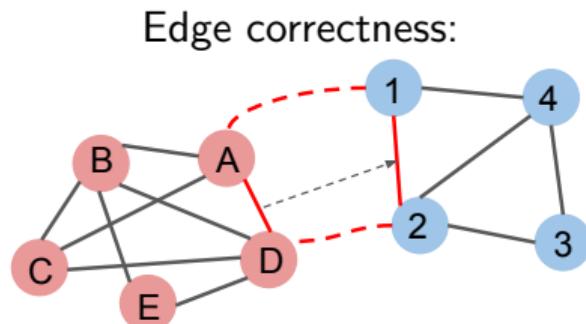
Matching the PPI network of yeast (2,340 proteins) to that of human (9,141 proteins)

Edge correctness:



Matching real-world PPI networks

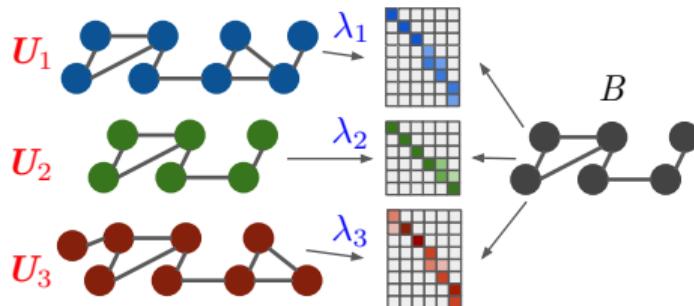
Matching the PPI network of yeast (2,340 proteins) to that of human (9,141 proteins)



Graph representation

Extend existing representation models based on Gromov-Wasserstein distance

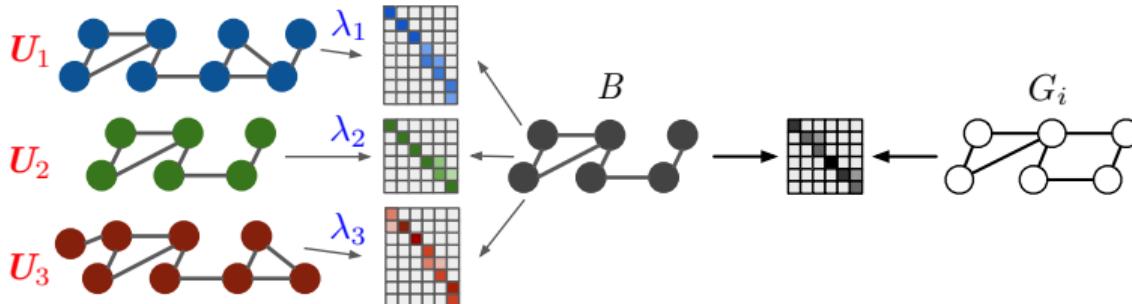
Gromov-Wasserstein factorization model



$$B_{gw}(\mathbf{U}_{1:K}, \boldsymbol{\lambda}) := \arg \min_B \sum_{k=1}^K \lambda_k d_{gw}(B, G_k(\mathbf{U}_k)). \quad (6)$$

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

Gromov-Wasserstein factorization model



$$B_{gw}(\mathbf{U}_{1:K}, \boldsymbol{\lambda}) := \arg \min_B \sum_{k=1}^K \lambda_k d_{gw}(B, \mathbf{G}_k(\mathbf{U}_k)). \quad (6)$$

- ▶ $\{\mathbf{G}_k(\mathbf{U}_k)\}_{k=1}^K$: 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]:

$$\min_{\mathbf{1} \geq \mathbf{U}_{1:K} \geq \mathbf{0}, \ \boldsymbol{\lambda}_{1:I} \in \Delta^{K-1}} \sum_{i=1}^I d_{gw}(B_{gw}(\mathbf{U}_{1:K}, \underbrace{\boldsymbol{\lambda}_i}_{\text{Rep. of } G_i}), G_i). \quad (7)$$

Gromov-Wasserstein factorization model

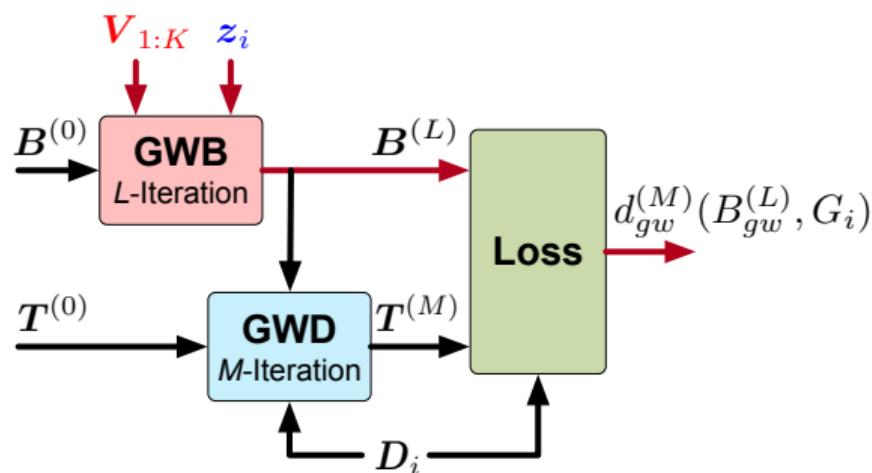
Reparameterize the problem to an unconstrained optimization problem:

$$\min_{\mathbf{V}_{1:K}, \mathbf{z}_{1:I}} \sum_{i=1}^I d_{gw}(B_{gw}(\underbrace{\sigma(\mathbf{V}_{1:K})}_{\mathbf{U}_{1:K}}, \underbrace{\text{softmax}(\mathbf{z}_i)}_{\lambda_i}), G_i). \quad (8)$$

Gromov-Wasserstein factorization model

Reparameterize the problem to an unconstrained optimization problem:

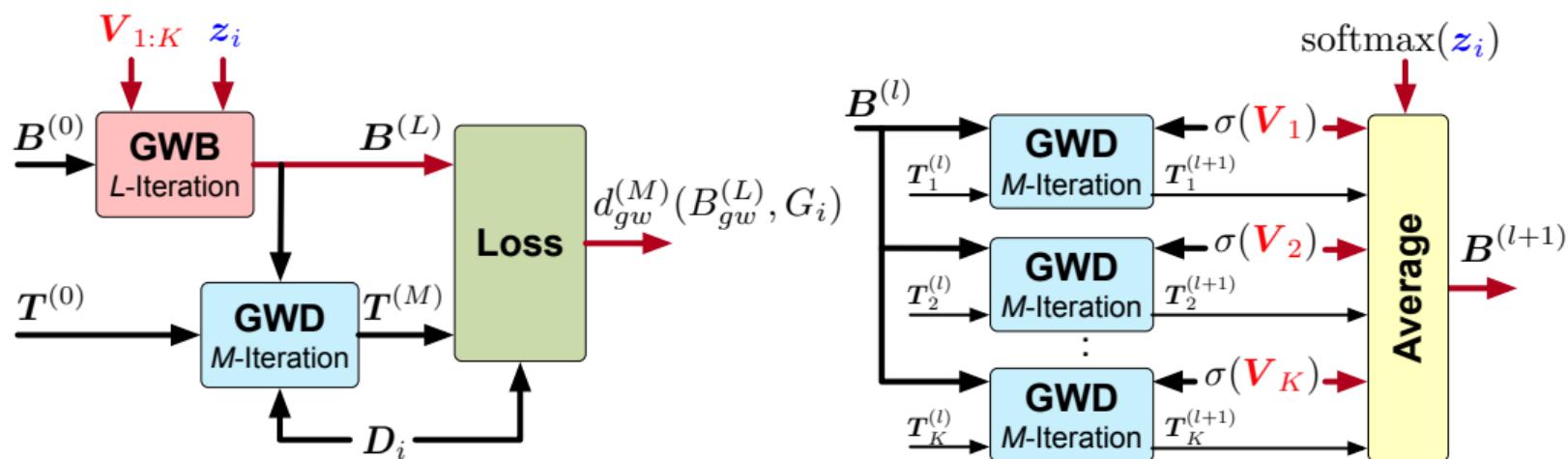
$$\min_{\mathbf{V}_{1:K}, \mathbf{z}_{1:I}} \sum_{i=1}^I d_{gw}(B_{gw}(\underbrace{\sigma(\mathbf{V}_{1:K})}_{\mathbf{U}_{1:K}}, \underbrace{\text{softmax}(\mathbf{z}_i)}_{\lambda_i}), G_i). \quad (8)$$



Gromov-Wasserstein factorization model

Reparameterize the problem to an unconstrained optimization problem:

$$\min_{\mathbf{V}_{1:K}, \mathbf{z}_{1:I}} \sum_{i=1}^I d_{gw}(B_{gw}(\underbrace{\sigma(\mathbf{V}_{1:K})}_{\mathbf{U}_{1:K}}, \underbrace{\text{softmax}(\mathbf{z}_i)}_{\lambda_i}), G_i). \quad (8)$$



From transductive to inductive

The GWF is a transductive model:

$$z_{\text{new}} = \arg \min_z d_{gw}(B_{gw}(\sigma(V_{1:K}), \text{softmax}(z)), G_{\text{new}}). \quad (9)$$

From transductive to inductive

The GWF is a transductive model:

$$\mathbf{z}_{\text{new}} = \arg \min_{\mathbf{z}} d_{gw}(B_{gw}(\sigma(\mathbf{V}_{1:K}), \text{softmax}(\mathbf{z})), G_{\text{new}}). \quad (9)$$

Make an **inductive GWF** model [*Xu, et al. ICML 2020 submission*]:

$$\min_{\mathbf{V}_{1:K}, \theta} \sum_{i=1}^I d_{gw}(B_{gw}(\sigma(\mathbf{V}_{1:K}), \text{softmax}(\underbrace{\text{GCN}_{\theta}(G_i)}_{\mathbf{z}_i})), G_i). \quad (10)$$

From transductive to inductive

The GWF is a transductive model:

$$\mathbf{z}_{\text{new}} = \arg \min_{\mathbf{z}} d_{gw}(B_{gw}(\sigma(\mathbf{V}_{1:K}), \text{softmax}(\mathbf{z})), G_{\text{new}}). \quad (9)$$

Make an **inductive GWF** model [Xu, et al. ICML 2020 submission]:

$$\min_{\mathbf{V}_{1:K}, \theta} \sum_{i=1}^I d_{gw}(B_{gw}(\sigma(\mathbf{V}_{1:K}), \text{softmax}(\underbrace{\text{GCN}_\theta(G_i)}_{\mathbf{z}_i})), G_i). \quad (10)$$

$$\mathbf{z}_{\text{new}} = \text{GCN}_\theta(G_{\text{new}}). \quad (11)$$

Experiments on molecule clustering

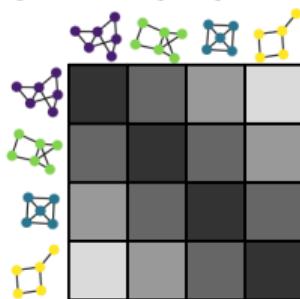
Experiments on molecule clustering

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

Experiments on molecule clustering

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

GWD Kernel



GWD+Kmeans

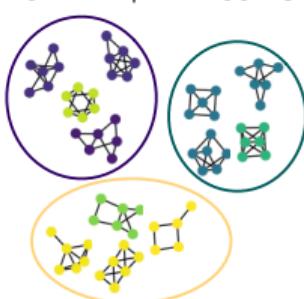


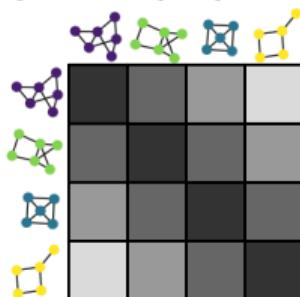
Table: Comparisons on clustering accuracy (%)

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

Experiments on molecule clustering

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

GWD Kernel



GWD+Kmeans

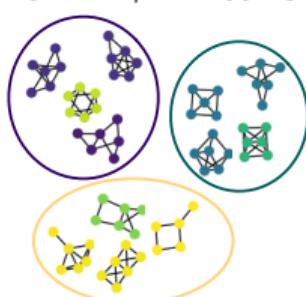
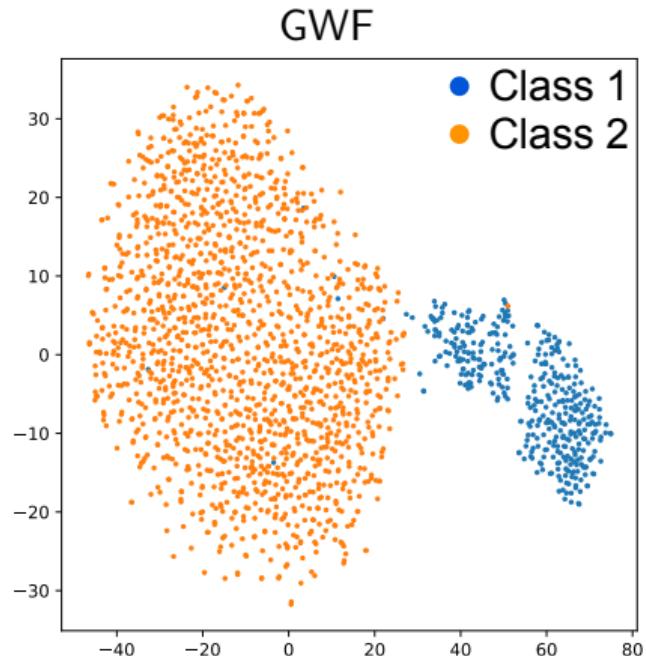


Table: Comparisons on clustering accuracy (%)

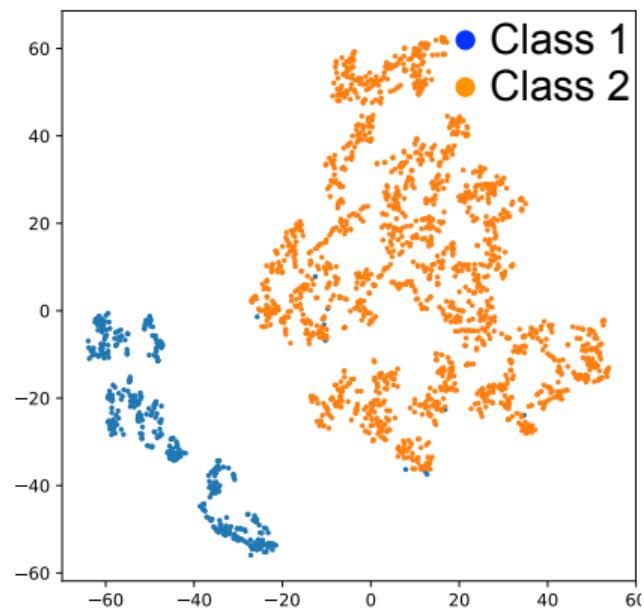
Method	AIDS	PROTEIN
GWD Kernel + SC	91.0 ± 0.7	66.4 ± 0.8
GWD + Kmeans	95.2 ± 0.9	64.7 ± 1.1
GWF + Kmeans	99.5 ± 0.4	70.7 ± 0.7
I-GWF + Kmeans	99.2 ± 0.5	73.3 ± 0.6

Visualization of coefficient vectors

AIDS



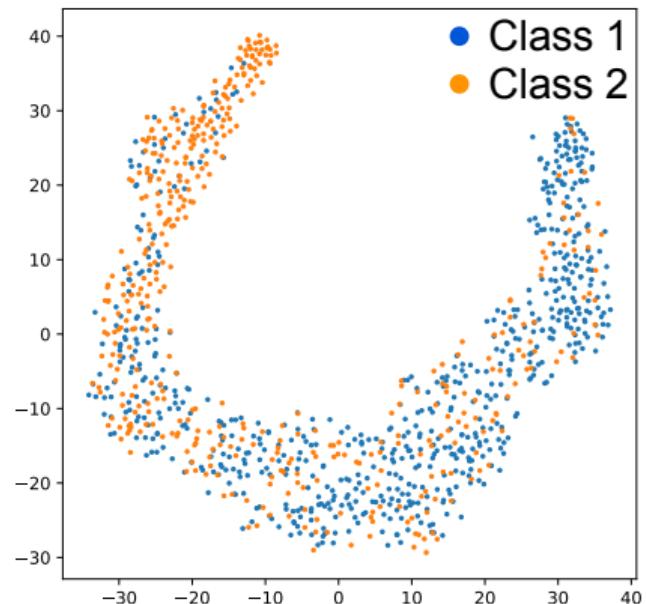
I-GWF



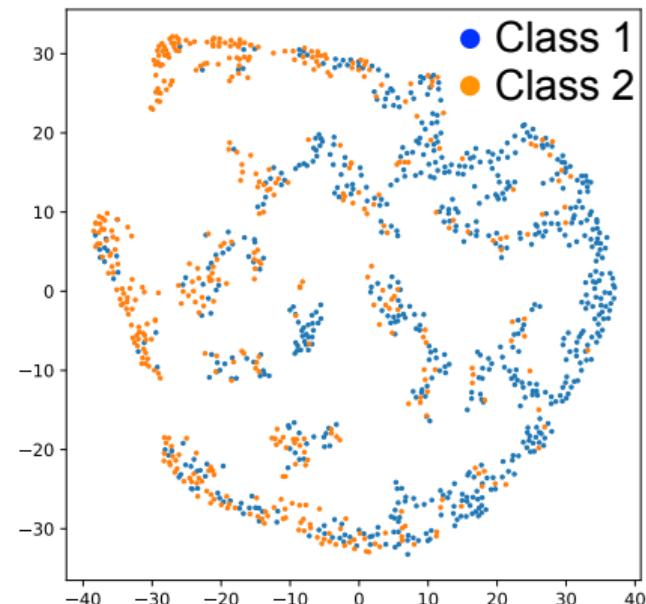
Visualization of coefficient vectors

PROTEIN

GWF



I-GWF



Semi-supervised learning of the I-GWF model

$$\min_{\mathbf{V}_{1:K}, \theta, \phi} \underbrace{\sum_{i=1}^I d_{gw}(B_{gw}(\sigma(\mathbf{V}_{1:K}), \text{softmax}(\text{GCN}_\theta(G_i))), G_i) +}_{\text{Reconstruction loss for both unlabeled and labeled graphs}} \\ \underbrace{\sum_{i \in \mathcal{G}_{\text{labeled}}} \text{loss}(\overbrace{\phi(\text{GCN}_\theta(G_i))}^{\text{classifier}}, l_i)}_{\text{MLE for labeled graphs}}. \quad (12)$$

Experiments on molecule classification

Experiments on molecule classification

- ▶ MUTAG: 188 compounds with high/low mutagenicity
- ▶ PTC-MR: 344 drugs with/without rodent carcinogenicity

Experiments on molecule classification

- ▶ MUTAG: 188 compounds with high/low mutagenicity
- ▶ PTC-MR: 344 drugs with/without rodent carcinogenicity

Table: Comparison on classification accuracy (%).

Category	Method	MUTAG	PTC-MR
Graph kernel	Random Walk	83.71	57.85
	Shortest Path	85.22	58.24
	Graphlet	81.66	57.26
	Weisfeiler-Lehman	80.72	57.97
	Deep Graph	87.44	60.08
	Multi-Scale Laplacian	<u>87.94</u>	<u>63.26</u>

Experiments on molecule classification

- ▶ MUTAG: 188 compounds with high/low mutagenicity
- ▶ PTC-MR: 344 drugs with/without rodent carcinogenicity

Table: Comparison on classification accuracy (%).

Category	Method	MUTAG	PTC-MR
Graph kernel	Random Walk	83.71	57.85
	Shortest Path	85.22	58.24
	Graphlet	81.66	57.26
	Weisfeiler-Lehman	80.72	57.97
	Deep Graph	87.44	60.08
	Multi-Scale Laplacian	<u>87.94</u>	<u>63.26</u>
Graph embedding dimension=512	node2vec	72.63	58.58
	sub2vec	61.05	59.99
	graph2vec	83.15	60.17
	InfoGraph [ICLR 2020]	<u>89.01</u>	<u>61.65</u>

Experiments on molecule classification

- ▶ MUTAG: 188 compounds with high/low mutagenicity
- ▶ PTC-MR: 344 drugs with/without rodent carcinogenicity

Table: Comparison on classification accuracy (%).

Category	Method	MUTAG	PTC-MR
Graph kernel	Random Walk	83.71	57.85
	Shortest Path	85.22	58.24
	Graphlet	81.66	57.26
	Weisfeiler-Lehman	80.72	57.97
	Deep Graph	87.44	60.08
	Multi-Scale Laplacian	87.94	63.26
Graph embedding dimension=512	node2vec	72.63	58.58
	sub2vec	61.05	59.99
	graph2vec	83.15	60.17
	InfoGraph [ICLR 2020]	89.01	61.65
GWL	I-GWF	90.11	63.58

Experiments on molecule classification

- ▶ MUTAG: 188 compounds with high/low mutagenicity
- ▶ PTC-MR: 344 drugs with/without rodent carcinogenicity

Table: Comparison on classification accuracy (%).

Category	Method	MUTAG	PTC-MR
Graph kernel	Random Walk	83.71	57.85
	Shortest Path	85.22	58.24
	Graphlet	81.66	57.26
	Weisfeiler-Lehman	80.72	57.97
	Deep Graph	87.44	60.08
	Multi-Scale Laplacian	87.94	63.26
Graph embedding dimension=512	node2vec	72.63	58.58
	sub2vec	61.05	59.99
	graph2vec	83.15	60.17
	InfoGraph [ICLR 2020]	89.01	61.65
GWL dimension=15	I-GWF	90.11	63.58
	I-GWF + SSL	91.36	65.21

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		

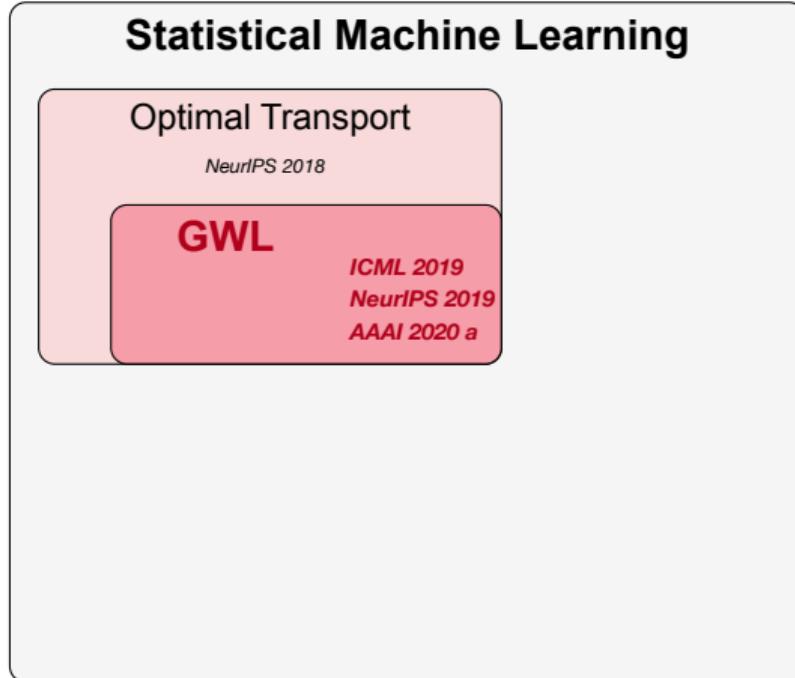
Summary of the GWL framework

Applications	Graph Matching	Graph Partitioning	Graph Representation
	Graph convolution networks		Factorization Model
Models	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

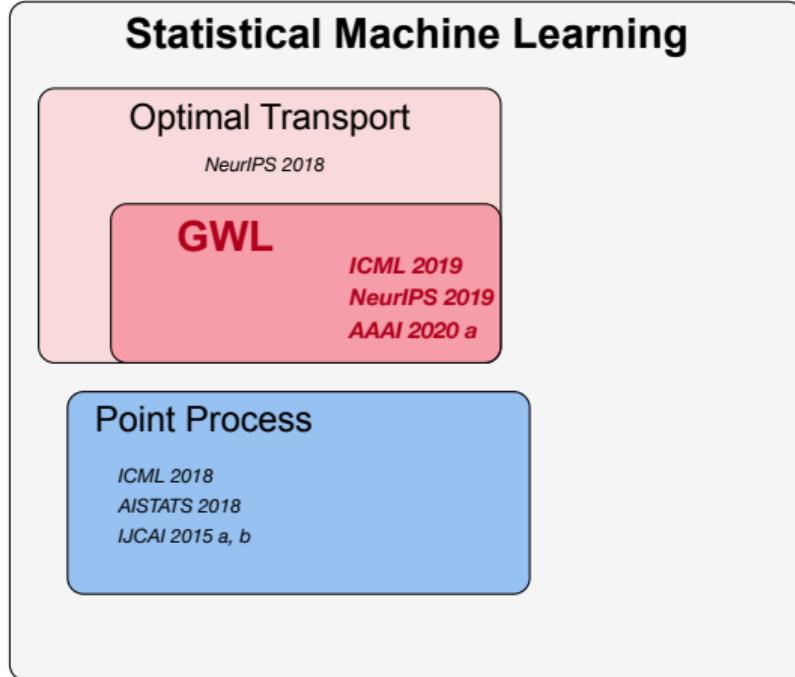
Tasks	PPI Network Alignment	Molecule Clustering and Classification	...
Applications	Graph Matching	Graph Partitioning	Graph Representation
Models	Graph convolution networks		Factorization Model
	Unsupervised and Semi-supervised Learning		
Optimization	Proximal Gradient	ADMM	Alternating Opt.
	Constrained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

Existing and ongoing work



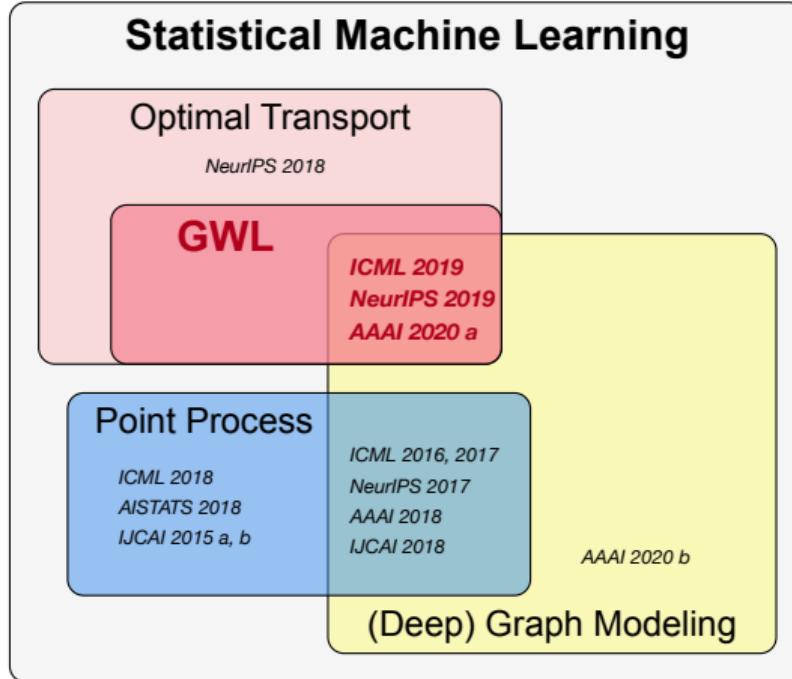
1. [ICML 2019] Hongteng Xu, D. Luo, H. Zha, and L. Carin, *Gromov-Wasserstein Learning for Graph Matching and Node Embedding*.
2. [NeurIPS 2019] Hongteng Xu, D. Luo, L. Carin, *Scalable Gromov-Wasserstein Learning for Graph Partitioning and Matching*.
3. [AAAI 2020 a] Hongteng Xu, *Gromov-Wasserstein Factorization Models for Graph Clustering*.
4. [NeurIPS 2018] Hongteng Xu, W. Wang, W. Liu, and L. Carin, *Distilled wasserstein learning for word embedding and topic modeling*.

Existing and ongoing work



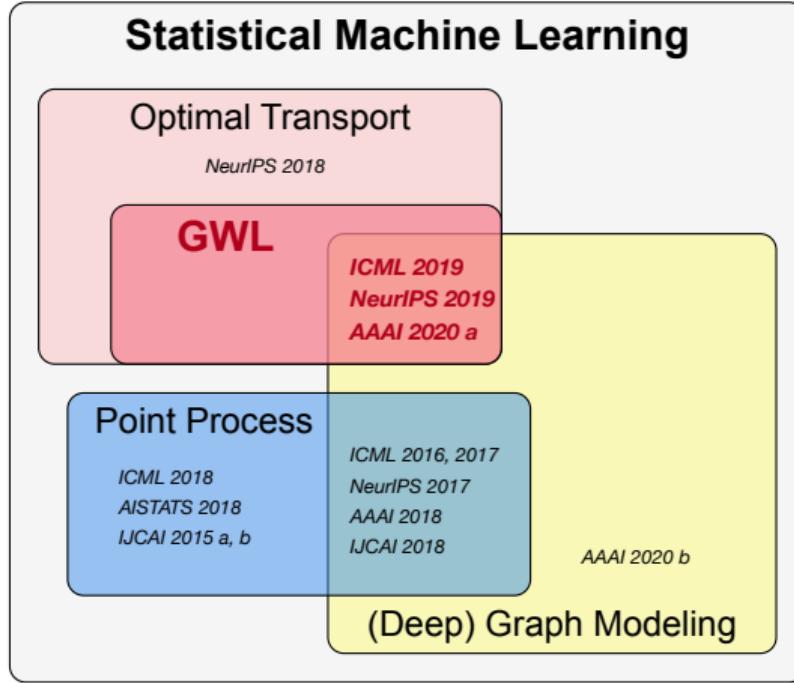
5. [ICML 2018] Hongteng Xu, H. Zha, and L. Carin, *Learning registered point processes from idiosyncratic observations*.
6. [AISTATS 2018] Hongteng Xu, D. Luo, X. Chen, and L. Carin, *Benefits from superposed hawkes processes*.
7. [IJCAI 2015 a] Hongteng Xu, D. Luo, et al., *Multi-Task Multi-Dimensional Hawkes Processes for Modeling Event Sequences*.
8. [IJCAI 2015 b] Hongteng Xu, Y. Zhen, and H. Zha, *Trailer Generation via a Point Process-Based Visual Attractiveness Model*.

Existing and ongoing work



9. [ICML 2016] Hongteng Xu, M. Farajtabar, and H. Zha, *Learning Granger Causality for Hawkes Processes*.
10. [ICML 2017] Hongteng Xu, D. Luo, and H. Zha, *Learning Hawkes processes from short doubly-censored event sequences*.
11. [NeurIPS 2017] Hongteng Xu and H. Zha, *A Dirichlet mixture model of Hawkes processes for event sequence clustering*.
12. [AAAI 2020 b] W. Wang, Hongteng Xu, et al., *Graph-Driven Generative Models for Heterogeneous Multi-Task Learning*.

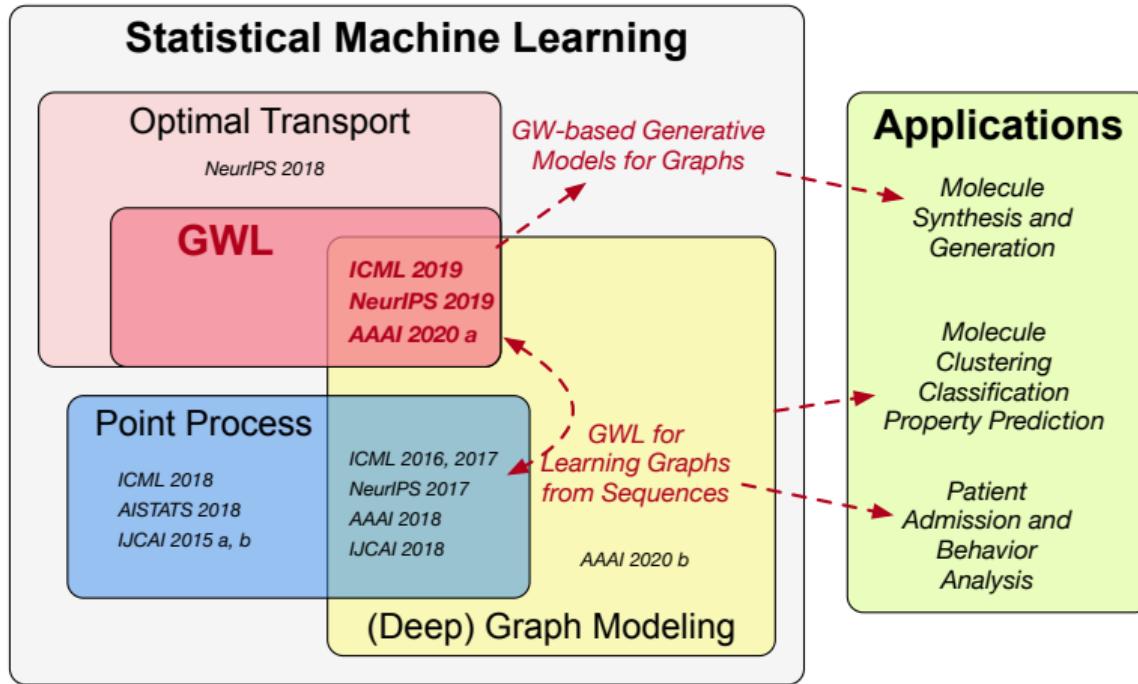
Existing and ongoing work



Miscellaneous (Manifold learning, vision, recommendation, ...)

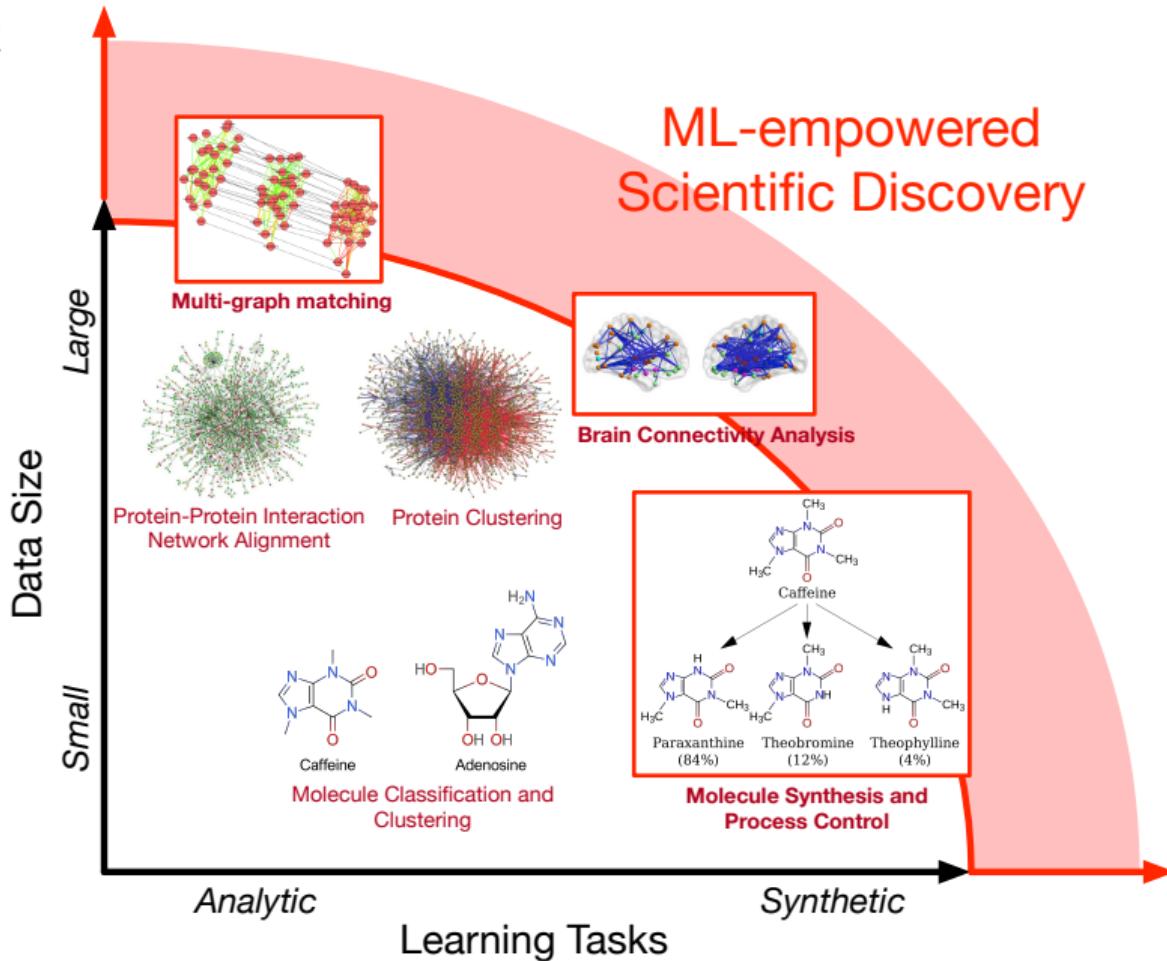
13. [AAAI 2015] Hongteng Xu, M. Davenport, et al., *Active Manifold Learning via Gershgorin Circle Guided Sample Selection*.
14. [ICCV 2015] Hongteng Xu, Y. Zhou, et al., *Unsupervised trajectory clustering via adaptive multi-kernel-based shrinkage*.
15. [CVPR 2017] Hongteng Xu, J. Yan, et al., *Fractal dimension invariant filtering and its CNN-based implementation*.
16. [WSDM 2018] X. Chen, Hongteng Xu, et al., *Sequential recommendation with user memory networks*.

Future work



17. [NeurIPS workshop 2019] D. Luo, **Hongteng Xu**, and L. Carin, *Fused Gromov-Wasserstein Alignment for Hawkes Processes*.
18. [TKDE 2017] **Hongteng Xu**, W. Wu, S. Nemati, et al., *Patient flow prediction via discriminative learning of mutually-correcting processes*.
19. [MLHC 2018] M. Engelhard, **Hongteng Xu**, et al., *Predicting smoking events with a time-varying semi-parametric hawkes process model*.
20. [TechReport 2019] D. Luo, **Hongteng Xu**, and L. Carin, *Interpretable ICD Code Embeddings with Self-and Mutual-Attention Mechanisms*.

Future work



Acknowledgements

- ▶ Lawrence Carin (Duke)
- ▶ Hongyuan Zha (Georgia Tech)
- ▶ Jun Liu (Infinia ML)
- ▶ Dixin Luo (Duke)
- ▶ Mark Davenport (Georgia Tech)
- ▶ Ricardo Henao (Duke)
- ▶ Le Song (Georgia Tech)
- ▶ Shamim Nemati (UCSD)
- ▶ Xia Ning (Ohio State)
- ▶ Yongfeng Zhang (Rutgers)
- ▶ Jiachang Liu (Duke)
- ▶ Matthew Engelhard (Duke)



National Institutes
of Health



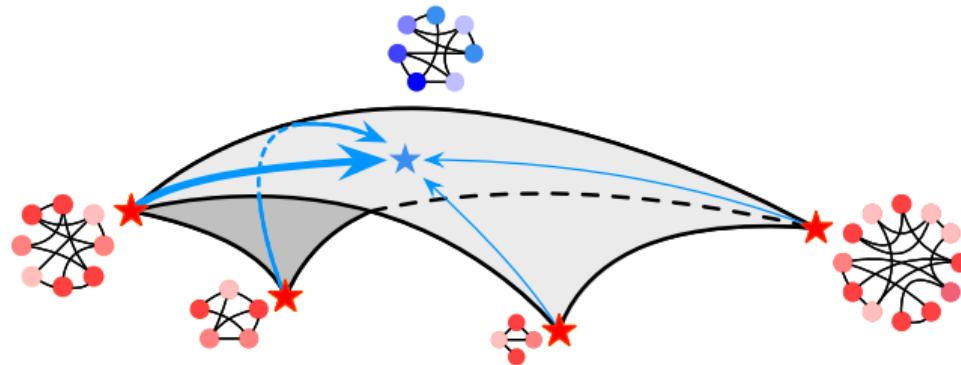
More Information

<https://sites.google.com/view/hongtengxu>

<https://github.com/HongtengXu>

hongteng.xu@duke.edu

Synthesize new graphs based on the GWF model



$$\lambda_{\text{new}} \sim P_\lambda,$$

$$G_{\text{new}} = \arg \min_B \sum_{k=1}^K \lambda_k^{\text{new}} d_{gw}(B, G_k).$$

Proposed algorithm: Proximal gradient algorithm (PGA)

Obj. $\min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \boldsymbol{\mu}_Y)} \langle \mathbf{D}_{XY} - 2\mathbf{D}_X \mathbf{T} \mathbf{D}_Y^\top, \mathbf{T} \rangle.$

Init. $m = 0, \mathbf{T}^{(m)} = \boldsymbol{\mu}_X \boldsymbol{\mu}_Y^\top.$

For $m = 0 : M$

$$\begin{aligned} & \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \boldsymbol{\mu}_Y)} \langle \mathbf{D}_{XY} - 2\mathbf{D}_X \mathbf{T}^{(m)} \mathbf{D}_Y^\top, \mathbf{T} \rangle + \gamma \underbrace{\text{KL}(\mathbf{T} \| \mathbf{T}^{(m)})}_{\text{Proximal}} \\ &= \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \boldsymbol{\mu}_Y)} \underbrace{\langle \mathbf{D}_{XY} - 2\mathbf{D}_X \mathbf{T}^{(m)} \mathbf{D}_Y^\top - \gamma \log \mathbf{T}^{(m)}, \mathbf{T} \rangle}_{\text{constant}} - \gamma \mathsf{H}(\mathbf{T}). \end{aligned} \tag{13}$$

1. Compute $\Phi = \exp\left(\frac{\mathbf{D}_{XY} - 2\mathbf{D}_X \mathbf{T}^{(m)} \mathbf{D}_Y^\top}{\gamma}\right) \odot \mathbf{T}^{(m)}, \mathbf{a} = \boldsymbol{\mu}_X.$
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. $\mathbf{T}^{(m+1)} = \text{diag}(\mathbf{a}) \Phi \text{diag}(\mathbf{b}).$

Proposed algorithm: Bregman ADMM (B-ADMM)

$$\min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \cdot), \mathbf{S} \in \Pi(\cdot, \boldsymbol{\mu}_Y), \mathbf{T} = \mathbf{S}} \langle \mathbf{D}_{XY} - 2\mathbf{D}_X \mathbf{S} \mathbf{D}_Y^\top, \mathbf{T} \rangle, \quad (14)$$

Introduce a **dual variable** \mathbf{Z} and initialize $\mathbf{Z}^{(0)} = \mathbf{0}$, $\mathbf{T}^{(0)} = \mathbf{S}^{(0)} = \boldsymbol{\mu}_X \boldsymbol{\mu}_Y^\top$:

Proposed algorithm: Bregman ADMM (B-ADMM)

$$\min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \cdot), \mathbf{S} \in \Pi(\cdot, \boldsymbol{\mu}_Y), \mathbf{T} = \mathbf{S}} \langle \mathbf{D}_{XY} - 2\mathbf{D}_X \mathbf{S} \mathbf{D}_Y^\top, \mathbf{T} \rangle, \quad (14)$$

Introduce a **dual variable** \mathbf{Z} and initialize $\mathbf{Z}^{(0)} = \mathbf{0}$, $\mathbf{T}^{(0)} = \mathbf{S}^{(0)} = \boldsymbol{\mu}_X \boldsymbol{\mu}_Y^\top$:

$$\begin{aligned} \mathbf{T}^{(m+1)} &= \arg \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \cdot)} \langle \mathbf{D}_{XY} - 2\mathbf{D}_X \mathbf{S}^{(m)} \mathbf{D}_Y^\top, \mathbf{T} \rangle + \\ &\quad \langle \mathbf{Z}^{(m)}, \mathbf{T} - \mathbf{S}^{(m)} \rangle + \gamma \mathsf{KL}(\mathbf{T} \| \mathbf{S}^{(m)}) = \left(\frac{\boldsymbol{\mu}_X}{\Phi_T \mathbf{1}_Y} \mathbf{1}_Y^\top \right) \odot \Phi_T, \end{aligned}$$

$$\text{where } \Phi_T = \exp\left(-\frac{1}{\gamma}(\mathbf{D}_{XY} - 2\mathbf{D}_X \mathbf{S}^{(m)} \mathbf{D}_Y^\top + \mathbf{Z}^{(m)})\right) \odot \mathbf{S}^{(m)}$$

Proposed algorithm: Bregman ADMM (B-ADMM)

$$\min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \cdot), \mathbf{S} \in \Pi(\cdot, \boldsymbol{\mu}_Y), \mathbf{T} = \mathbf{S}} \langle \mathbf{D}_{XY} - 2\mathbf{D}_X \mathbf{S} \mathbf{D}_Y^\top, \mathbf{T} \rangle, \quad (14)$$

Introduce a **dual variable \mathbf{Z}** and initialize $\mathbf{Z}^{(0)} = \mathbf{0}$, $\mathbf{T}^{(0)} = \mathbf{S}^{(0)} = \boldsymbol{\mu}_X \boldsymbol{\mu}_Y^\top$:

$$\begin{aligned} \mathbf{T}^{(m+1)} &= \arg \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \cdot)} \langle \mathbf{D}_{XY} - 2\mathbf{D}_X \mathbf{S}^{(m)} \mathbf{D}_Y^\top, \mathbf{T} \rangle + \\ &\quad \langle \mathbf{Z}^{(m)}, \mathbf{T} - \mathbf{S}^{(m)} \rangle + \gamma \text{KL}(\mathbf{T} \| \mathbf{S}^{(m)}) = \left(\frac{\boldsymbol{\mu}_X}{\Phi_T \mathbf{1}_Y} \mathbf{1}_Y^\top \right) \odot \Phi_T, \end{aligned}$$

$$\text{where } \Phi_T = \exp\left(-\frac{1}{\gamma}(\mathbf{D}_{XY} - 2\mathbf{D}_X \mathbf{S}^{(m)} \mathbf{D}_Y^\top + \mathbf{Z}^{(m)})\right) \odot \mathbf{S}^{(m)}$$

$$\begin{aligned} \mathbf{S}^{(m+1)} &= \arg \min_{\mathbf{S} \in \Pi(\cdot, \boldsymbol{\mu}_t)} \langle -2\mathbf{D}_X^\top \mathbf{T}^{(m+1)} \mathbf{D}_Y, \mathbf{S} \rangle + \langle \mathbf{Z}^{(m)}, \mathbf{T}^{(m+1)} - \mathbf{S} \rangle + \\ &\quad \gamma \text{KL}(\mathbf{S} \| \mathbf{T}^{(m+1)}) = (\mathbf{1}_X \left(\frac{\boldsymbol{\mu}_Y}{\Phi_S^\top \mathbf{1}_X} \right)^\top) \odot \Phi_S, \end{aligned}$$

$$\text{where } \Phi_S = \exp\left(\frac{1}{\gamma}(2\mathbf{D}_X^\top \mathbf{T}^{(m+1)} \mathbf{D}_Y + \mathbf{Z}^{(m)})\right) \odot \mathbf{T}^{(m)}$$

Proposed algorithm: Bregman ADMM (B-ADMM)

$$\min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \cdot), \mathbf{S} \in \Pi(\cdot, \boldsymbol{\mu}_Y), \mathbf{T} = \mathbf{S}} \langle \mathbf{D}_{XY} - 2\mathbf{D}_X \mathbf{S} \mathbf{D}_Y^\top, \mathbf{T} \rangle, \quad (14)$$

Introduce a **dual variable \mathbf{Z}** and initialize $\mathbf{Z}^{(0)} = \mathbf{0}$, $\mathbf{T}^{(0)} = \mathbf{S}^{(0)} = \boldsymbol{\mu}_X \boldsymbol{\mu}_Y^\top$:

$$\begin{aligned} \mathbf{T}^{(m+1)} &= \arg \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_X, \cdot)} \langle \mathbf{D}_{XY} - 2\mathbf{D}_X \mathbf{S}^{(m)} \mathbf{D}_Y^\top, \mathbf{T} \rangle + \\ &\quad \langle \mathbf{Z}^{(m)}, \mathbf{T} - \mathbf{S}^{(m)} \rangle + \gamma \text{KL}(\mathbf{T} \| \mathbf{S}^{(m)}) = \left(\frac{\boldsymbol{\mu}_X}{\Phi_T \mathbf{1}_Y} \mathbf{1}_Y^\top \right) \odot \Phi_T, \end{aligned}$$

$$\text{where } \Phi_T = \exp\left(-\frac{1}{\gamma}(\mathbf{D}_{XY} - 2\mathbf{D}_X \mathbf{S}^{(m)} \mathbf{D}_Y^\top + \mathbf{Z}^{(m)})\right) \odot \mathbf{S}^{(m)}$$

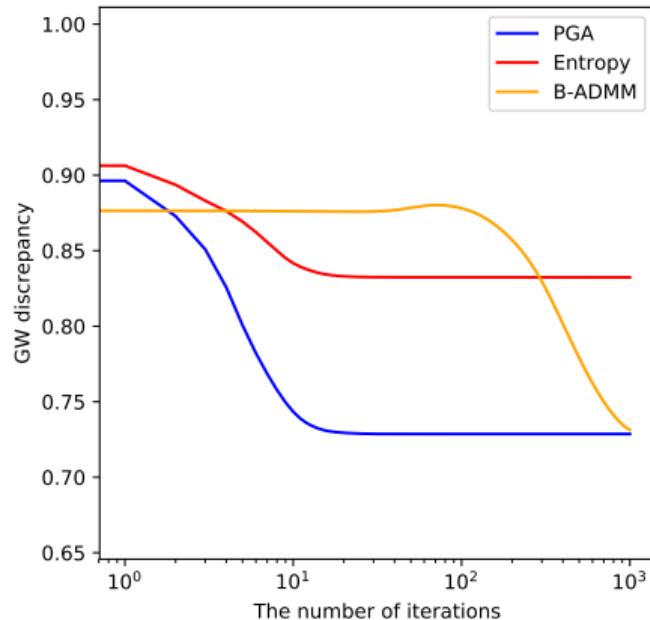
$$\begin{aligned} \mathbf{S}^{(m+1)} &= \arg \min_{\mathbf{S} \in \Pi(\cdot, \boldsymbol{\mu}_t)} \langle -2\mathbf{D}_X^\top \mathbf{T}^{(m+1)} \mathbf{D}_Y, \mathbf{S} \rangle + \langle \mathbf{Z}^{(m)}, \mathbf{T}^{(m+1)} - \mathbf{S} \rangle + \\ &\quad \gamma \text{KL}(\mathbf{S} \| \mathbf{T}^{(m+1)}) = (\mathbf{1}_X \left(\frac{\boldsymbol{\mu}_Y}{\Phi_S^\top \mathbf{1}_X} \right)^\top) \odot \Phi_S, \end{aligned}$$

$$\text{where } \Phi_S = \exp\left(\frac{1}{\gamma}(2\mathbf{D}_X^\top \mathbf{T}^{(m+1)} \mathbf{D}_Y + \mathbf{Z}^{(m)})\right) \odot \mathbf{T}^{(m)}$$

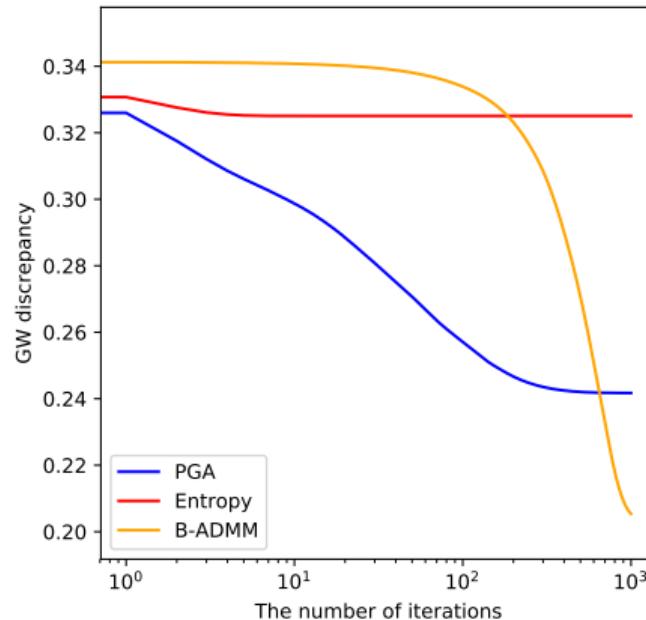
$$\mathbf{Z}^{(m+1)} = \mathbf{Z}^{(m)} + \gamma(\mathbf{T}^{(m+1)} - \mathbf{S}^{(m+1)}).$$

Empirical convergence

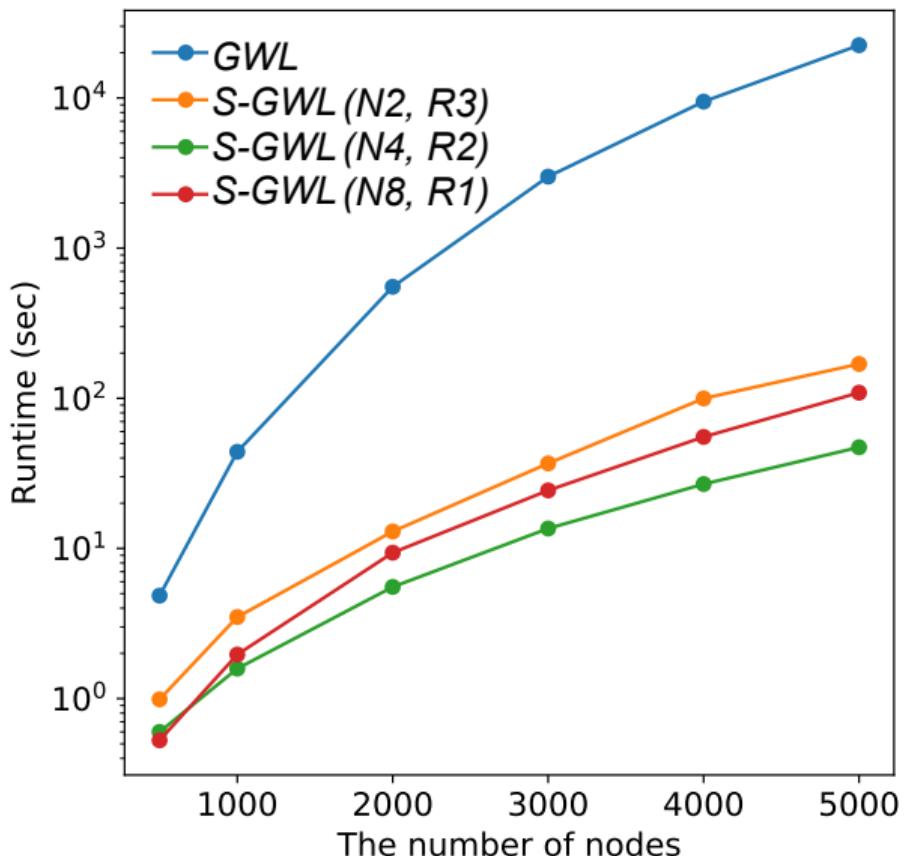
Undirected graph



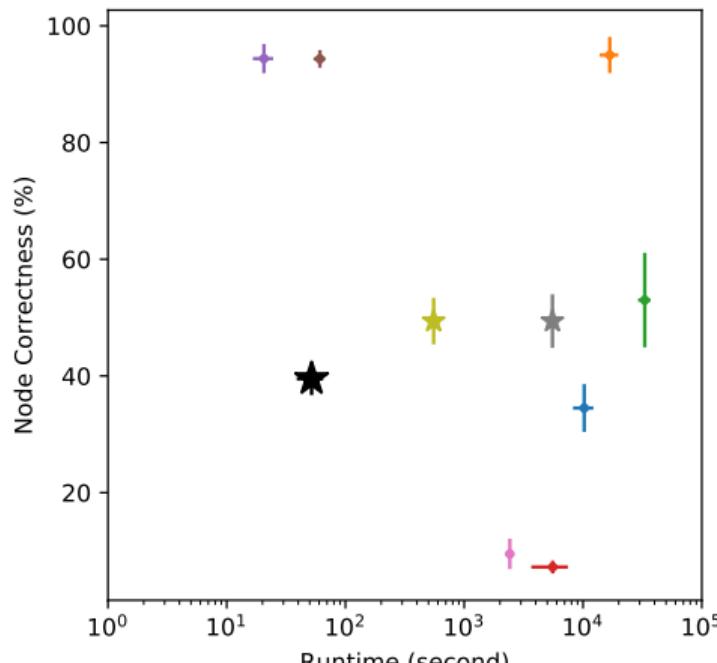
Directed graph



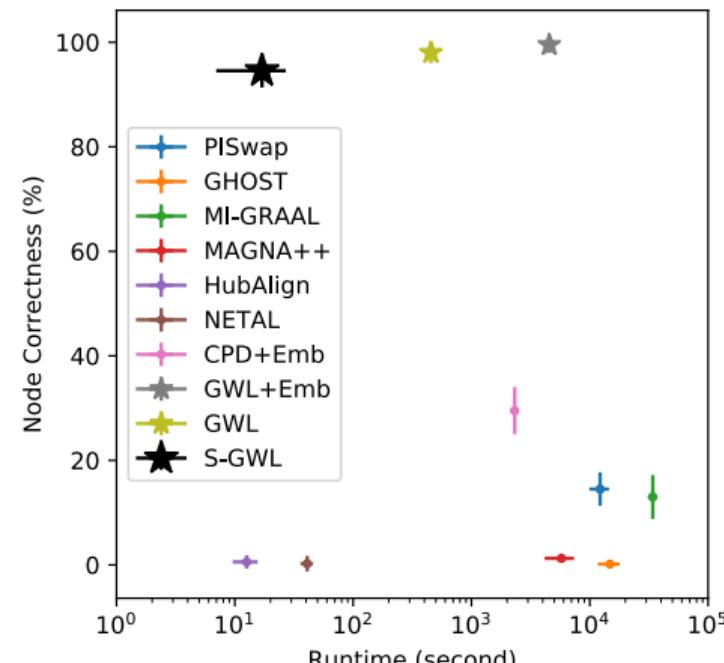
Improvements on scalability



More experimental results on graph matching



(a) BA graphs



(b) GRP graphs