

# **MOREL: MULTI-OMICS RELATIONAL LEARNING**

**ICLR 2022**

Arman Hasanzadeh Ehsan Hajiramezanali, Nick Duffiel and Xiaoning Qian

Presenter: Qingmei Wang

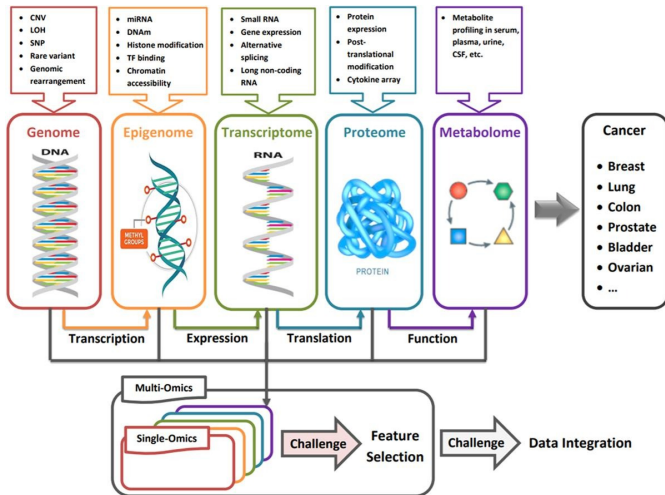
October 19, 2023

# Overview

- ▶ Preliminaries
- ▶ Problem Formulation and Notations
- ▶ Method
- ▶ Experiments
- ▶ Conclusion

# What's Multi-omics Data

“Multi-omics” refers to multiple “omics” datasets: Genomics, Epigenomics, Transcriptomics, Proteomics, and Metabolomics.



# Molecular Relational Learning

- ▶ Molecular Interaction Prediction: Crucial for discovering or designing new molecules.
- ▶ Drug-Drug Interaction Prediction: Learning whether a combination of two drugs will produce side effects.

# Wasserstein Distance & Gromov-Wasserstein Distance

## Couplings and Optimal Transport (EMD)

Input distributions  $\mu = \sum_i \mu_i \delta_{x_i}$   
 $\nu = \sum_j \nu_j \delta_{y_j}$

Points  $(x_i)_i, (y_j)_j$

Weights  $\mu_i \geq 0, \nu_j \geq 0$ .

$\sum_{i=1}^{N_1} \mu_i = \sum_{j=1}^{N_2} \nu_j = 1$   $d_{i,j} = d(x_i, y_j)$

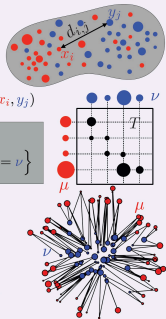
Def. Couplings

$$C_{\mu,\nu} \stackrel{\text{def.}}{=} \left\{ T \in \mathbb{R}_+^{N_1 \times N_2} ; T \mathbb{1}_{N_1} = \mu, T^\top \mathbb{1}_{N_2} = \nu \right\}$$

Def. Wasserstein Distance / EMD

$$W_p^p(\mu, \nu) \stackrel{\text{def.}}{=} \min \left\{ \sum_{i,j} T_{i,j} d_{i,j}^p ; T \in C_{\mu,\nu} \right\}$$

[Kantorovich 1942]



## Gromov-Wasserstein

Inputs:  $\{(\text{similarity/kernel matrix, histogram})\}$

$$(d, \mu) \quad \mu = \sum_i \mu_i \delta_{x_i} \quad d_{i,i'} = d(x_i, x_{i'})$$

$$(\bar{d}, \nu) \quad \nu = \sum_j \nu_j \delta_{y_j} \quad \bar{d}_{j,j'} = \bar{d}(y_j, y_{j'})$$

Def. Gromov-Wasserstein distance:

$$GW_p^p(d, \mu, \bar{d}, \nu) \stackrel{\text{def.}}{=} \min_{T \in C_{\mu,\nu}} \mathcal{E}_{d,\bar{d}}^p(T)$$

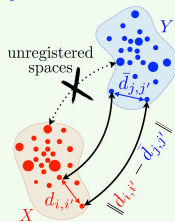
$$\mathcal{E}_{d,\bar{d}}^p(T) \stackrel{\text{def.}}{=} \sum_{i,i',j,j'} |d_{i,i'} - \bar{d}_{j,j'}|^p T_{i,j} T_{i',j'}$$

[Memoli 2011]

Computation of GW is a QAP:

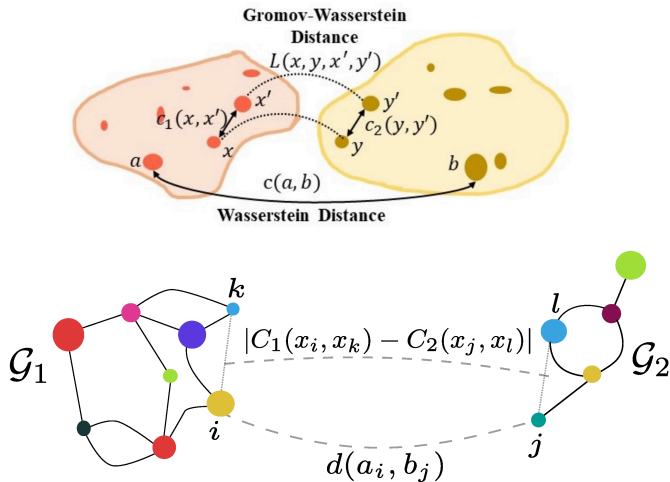
→ NP-hard in general.

→ need for a fast approximate solver.



- ▶ WD quantifies the geometric discrepancy between two probability distributions.
- ▶ GWD has been proposed as a natural extension of WD when a meaningful transportation cost between the distributions cannot be defined.

# Fused Gromov-Wasserstein



FGW distance is an optimization problem consisting of a Wasserstein term and a GW term

# Problem Formulation and Notation

Aims to infer the interrelations among entities, i.e. features, across all of the views.

- ▶ multiple views,  $\mathcal{V}$ , of data are given
- ▶ assume that the structure information is available for some of the views  $\mathcal{V}_s \subset \mathcal{V}$
- ▶  $\mathcal{V}_u = \mathcal{V} \setminus \mathcal{V}_s$  are unstructured
- ▶  $\mathfrak{G}_s = \{\mathcal{G}^{(v)}\}_{v \in \mathcal{V}_s}$  as the set of graphs for structured views
- ▶  $\mathfrak{A}_s = \{\mathbf{A}^{(v)}\}_{v \in \mathcal{V}_s}$  as adjacency matrices
- ▶  $\mathcal{X}_s = \{\mathbf{X}^{(v)}\}_{v \in \mathcal{V}_s}$  as the set of node attributes for structured views
- ▶  $\mathcal{X}_u = \{\mathbf{X}^{(v)}\}_{v \in \mathcal{V}_u}$  as the set of data for unstructured views
- ▶  $N_v$  denotes the number of nodes in structured views and number of features for unstructured views

# Problem Formulation and Notation

- ▶ MoReL infers the interactions among the nodes in  $\mathfrak{G}_s$  and features in  $\mathcal{X}_u$
- ▶ Inter-relations by a multi-partite graph with  $\sum_{v \in \mathcal{V}} N_v$  nodes and a multi-adjacency tensor  $\mathcal{A} = \left\{ \mathbf{A}^{(vv')} \right\}_{v, v' \in \mathcal{V}, v \neq v'}$ , where  $\mathbf{A}^{(vv')}$  is the  $N_v \times N_{v'}$  bi-adjacency matrix between views  $v$  and  $v'$



# MOREL Generative Model

- ▶ Define a hierarchical Bayesian model for MoReL with three sets of latent variables:
- ▶  $\mathcal{H} = \mathcal{H}_s \cup \mathcal{H}_u = \{\mathbf{H}^{(v)}\}_{v \in \mathcal{V}_s \cup \mathcal{V}_u}$ , which captures the (hidden) structural information
- ▶  $\mathcal{A}$ , which encodes the interaction among features across views
- ▶  $\mathcal{Z} = \mathcal{Z}_s \cup \mathcal{Z}_u = \{\mathbf{Z}^{(v)}\}_{v \in \mathcal{V}_s \cup \mathcal{V}_u}$ , which summarizes the feature/attribute specific information
- ▶ The joint probability of observations and latent variables factorizes as follows:

$$p_{\theta}(\mathcal{X}_u, \mathcal{X}_s, \mathcal{A}_s, \mathcal{H}, \mathcal{A}, \mathcal{Z}) = \\ p_{\theta_x}(\mathcal{X}_u | \mathcal{Z}_u) p_{\theta_x}(\mathcal{X}_s | \mathcal{Z}_s) p_{\theta_g}(\mathcal{A}_s | \mathcal{H}_s) p_{\theta_z}(\mathcal{Z} | \mathcal{H}, \mathcal{A}) p_{\theta_a}(\mathcal{A} | \mathcal{H}) p(\mathcal{H}).$$

## FGW Distance.

- Given two structured probability distributions,  $\mathbf{\Lambda} \in \mathcal{P}(\mathbb{X})$  and  $\mathbf{\Delta} \in \mathcal{P}(\mathbb{Y})$ , FGW is defined as follows:

$$\begin{aligned}\mathcal{D}_{\text{FGW}}(\mathbf{\Lambda}, \mathbf{\Delta}) &= \alpha \mathcal{D}_{\text{W}}(\mathbf{\Lambda}, \mathbf{\Delta}) + \beta \mathcal{D}_{\text{GW}}(\mathbf{\Lambda}, \mathbf{\Delta}) \\ &= \alpha \inf_{\pi_w \in \Pi(\mathbb{X} \times \mathbb{Y})} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \pi_w} [c^{(\mathbb{X}\mathbb{Y})}(\mathbf{x}, \mathbf{y})] \\ &\quad + \beta \inf_{\pi_{gw} \in \Pi(\mathbb{X} \times \mathbb{Y})} \mathbb{E}_{(\mathbf{x}, \mathbf{y}), (\mathbf{x}', \mathbf{y}') \sim \pi_{gw}} [\| c^{(\mathbb{X})}(\mathbf{x}, \mathbf{x}') - c^{(\mathbb{Y})}(\mathbf{y}, \mathbf{y}') \|],\end{aligned}\tag{3}$$

where  $\alpha, \beta \in [0, 1]$  are scalar hyper-parameters,  $\Pi(\mathbb{X} \times \mathbb{Y})$  is the set of all admissible couplings between  $\mathbf{\Lambda}$  and  $\mathbf{\Delta}$ , and  $c^{(\mathbb{X}\mathbb{Y})}$ ,  $c^{(\mathbb{X})}$ , and  $c^{(\mathbb{Y})}$  are corresponding transportation cost functions.  $\mathcal{D}_{\text{FGW}}$  can be further simplified by choosing  $\pi_w$  to be equal to  $\pi_{gw}$  (Chen et al., 2020).

# Relational learning via FGW

- ▶ FGW distance based decoder for every pair of views, in which each view independently belongs to either structured or unstructured views, i.e.  $(v, v') \in \mathcal{V}$
- ▶ define the transportation cost functions  $c^{(vv')}$  and  $c^{(v)}$ , and then approximate  $\mathcal{D}_{\text{FGW}}$
- ▶ define the inter-cost function for the first term of FGW, i.e.  $\mathcal{D}_{\text{W}}$ , as follows:

$$c^{(vv')} \left( \mathbf{H}_{i,:}^{(v)}, \mathbf{H}_{j,:}^{(v')} \right) = 1 - \sigma \left( \mathbf{H}_{i,:}^{(v)} \left( \mathbf{H}_{j,:}^{(v')} \right)^T \right); \quad v, v' \in \mathcal{V}$$

where  $\sigma$  denotes the sigmoid function, and  $\mathbf{H}_{i,:}^{(v)}$  represents the structural latent variable of node/feature  $i$  in view  $v$

## Relational learning via FGW

- For the structured views, define the cost function as a combination of the shortest path distance from the graph and the distance between structural latent variables. Given the normalized shortest path distance matrix between every pair of nodes in the input graph  $\mathbf{D}^{(v)}$  :

$$c^{(v)} \left( \mathbf{H}_{i,:}^{(v)}, \mathbf{H}_{j,:}^{(v)} \right) = \mathbf{D}^{(v)} \odot \left( 1 - \sigma \left( \mathbf{H}_{i,:}^{(v)} \left( \mathbf{H}_{j,:}^{(v)} \right)^T \right) \right); \quad \text{for } v \in \mathcal{V}_s$$

where  $\odot$  denotes the Hadamard product.

- For unstructured views, we define the cost function between two features as follows:

$$c^{(v)} \left( \mathbf{H}_{i,:}^{(v)}, \mathbf{H}_{j,:}^{(v)} \right) = 1 - \sigma \left( \mathbf{H}_{i,:}^{(v)} \left( \mathbf{H}_{j,:}^{(v)} \right)^T \right); \quad \text{for } v \in \mathcal{V}_u$$

# Relational learning via FGW

- Rewrite  $\mathcal{D}_{\text{FGW}}$  between two views of data with shared transport matrix as follows:

$$\begin{aligned} & \mathcal{D}_{\text{FGW}} \left( p \left( \mathbf{H}^{(v)} \right), p \left( \mathbf{H}^{(v')} \right) \right) \\ & \sum_{i=1}^{N_v} \sum_{j=1}^{N_{v'}} \min_{\mathbf{T}_{gw}^{(vv')} \in \Pi} \sum_{\mathbf{H}_{i,:}^{(v)}, \mathbf{H}_{j,:}^{(v')}, \mathbf{H}_{i,:}^{(v)'}, \mathbf{H}_{j,:}^{(v')'}} \left[ \alpha c^{(vv')} \left( \mathbf{H}_{i,:}^{(v)}, \mathbf{H}_{j,:}^{(v')} \right) + \right. \\ & \left. \beta c^{(v)} \left( \mathbf{H}_{i,:}^{(v)}, \mathbf{H}_{i,:}^{(v)'} \right) - c^{(v')} \left( \mathbf{H}_{j,:}^{(v')}, \mathbf{H}_{j,:}^{(v')'} \right) \right]. \end{aligned}$$

- To approximate the FGW distance, firstly deploy GW algorithm to obtain  $\mathbf{T}_{gw}^{(vv')}$  and  $\mathcal{D}_{\text{GW}}$ , and then utilize  $\mathbf{T}_{gw}^{(vv')}$  along with the defined transportation cost  $c^{(vv')}$  to calculate Wasserstein distance term in  $\mathcal{D}_{\text{FGW}}$ .

## Prior Construction

- ▶ The paper imposes independent zero-mean unit-variance Gaussian priors on elements of  $\mathcal{H}$ . The prior for  $\mathcal{Z}$  is a multivariate Gaussian distribution whose mean and diagonal covariance matrix are constructed from the inferred multi-partite graph and the structural latent variable  $\mathcal{H}$ .
- ▶ The paper uses two graph neural networks (GNNs)  $g_{pz}^{(\mu)}$  and  $g_{pz}^{(\sigma)}$  to map  $\mathcal{H}$  and  $\mathcal{A}$  to the parameters of  $p_{\theta_z}(\mathcal{Z})$
- ▶ Specifically,

$$p_{\theta_z}(\mathcal{Z} \mid \mathcal{H}, \mathcal{A}) = \prod_{v \in \mathcal{V}_s} \prod_{i=1}^{N_v} p_{\theta_z} \left( \mathbf{z}_{i,:}^{(v)} \mid \mathcal{H}, \mathcal{A} \right); \quad p_{\theta_z} \left( \mathbf{z}_{i,:}^{(v)} \mid \mathcal{H}, \mathcal{A} \right) = \mathcal{N} \left( \boldsymbol{\mu}_{pz}^{(v,i)}, \boldsymbol{\sigma}_{pz}^{(v,i)} \right),$$

with  $\left[ \boldsymbol{\mu}_{pz}^{(v,i)} \right]_{v,i} = g_{pz}^{(\mu)}(\mathcal{H}, \mathcal{A}), \quad \left[ \boldsymbol{\sigma}_{pz}^{(v,i)} \right]_{v,i} = g_{pz}^{(\sigma)}(\mathcal{H}, \mathcal{A}).$

# Likelihood of observations

**Likelihood of observations.** To reconstruct the input graphs in the structured views, we assume that the views and edges are conditionally independent. More specifically, we employ an inner-product decoder as follows:

$$p_{\theta_g}(\mathcal{A}_s | \mathcal{H}_s) = \prod_{v \in \mathcal{V}_s} \prod_{i,j=1}^{N_v} p_{\theta_g} \left( \mathbf{A}_{i,j}^{(v)} | \mathbf{H}_{i,:}^{(v)}, \mathbf{H}_{j,:}^{(v)} \right);$$
$$p_{\theta_g} \left( \mathbf{A}_{i,j}^{(v)} | \mathbf{H}_{i,:}^{(v)}, \mathbf{H}_{j,:}^{(v)} \right) = \text{Ber} \left( \sigma(\mathbf{H}_{i,:}^{(v)} (\mathbf{H}_{j,:}^{(v)})^T) \right).$$

To generate the features in unstructured views and node attributes in structured views, we assume that views are conditionally independent. Hence we can expand the feature reconstruction terms in the equation (2) as follows:

$$p_{\theta_x}(\mathcal{X}_u | \mathcal{Z}_u) = \prod_{v \in \mathcal{V}_u} p_{\theta_x}(\mathbf{X}^{(v)} | \mathbf{Z}^{(v)}), \quad p_{\theta_x}(\mathcal{X}_s | \mathcal{Z}_s) = \prod_{v \in \mathcal{V}_s} p_{\theta_x}(\mathbf{X}^{(v)} | \mathbf{Z}^{(v)}).$$

We note that  $p_{\theta_x}$  could also be view specific depending on whether the node attributes/features in a view are discrete or continuous. In our experiments, we have deployed the Gaussian likelihood with the unit variance. The mapping from  $\mathcal{Z}$  to the parameters of  $p_{\theta_x}(\mathcal{X})$ , in our case, the mean of the Gaussian distribution, can be any highly expressive function such as neural networks. We denote these functions by  $f_{px}^{(v,s)}$  and  $f_{px}^{(v,u)}$ .

# Posterior

**Posterior.** We model the posterior of the structural latent variables as a Gaussian distribution and infer its parameters independently for each view. More specifically,

$$q_{\phi_h}(\mathcal{H}_u | \mathcal{X}_u) = \prod_{v \in \mathcal{V}_u} q_{\phi_h}(\mathbf{H}^{(v)} | \mathbf{X}^{(v)}), \quad q_{\phi_h}(\mathcal{H}_s | \mathcal{X}_u, \mathcal{A}_s) = \prod_{v \in \mathcal{V}_s} q_{\phi_h}(\mathbf{H}^{(v)} | \mathbf{X}^{(v)}, \mathbf{A}^{(v)}).$$

We use two GNNs for each structured view,  $\{g_{qh}^{(\mu,v)}(\mathbf{X}^{(v)}, \mathbf{A}^{(v)}), g_{qh}^{(\sigma,v)}(\mathbf{X}^{(v)}, \mathbf{A}^{(v)})\}_{v \in \mathcal{V}_s}$ , and two fully connected neural networks per unstructured view,  $\{f_{qh}^{(\mu,v)}(\mathbf{X}^{(v)}), f_{qh}^{(\sigma,v)}(\mathbf{X}^{(v)})\}_{v \in \mathcal{V}_u}$ , to map inputs to the mean and variance of the posteriors. We consider the variational distribution of  $\mathcal{Z}$  to be a multivariate Gaussian distribution, and it is factorized as follows:

$$q_{\phi_z}(\mathcal{Z}_u | \mathcal{X}_u) = \prod_{v \in \mathcal{V}_u} q_{\phi_z}(\mathbf{Z}^{(v)} | \mathbf{X}^{(v)}), \quad q_{\phi_z}(\mathcal{Z}_s | \mathcal{X}_u, \mathcal{A}_s) = \prod_{v \in \mathcal{V}_s} q_{\phi_z}(\mathbf{Z}^{(v)} | \mathbf{X}^{(v)}, \mathbf{A}^{(v)}).$$

We use two GNNs per structured view,  $\{g_{qz}^{(\mu,v)}(\mathbf{X}^{(v)}, \mathbf{A}^{(v)}), g_{qz}^{(\sigma,v)}(\mathbf{X}^{(v)}, \mathbf{A}^{(v)})\}_{v \in \mathcal{V}_s}$ , and two fully connected neural networks for each unstructured view,  $\{f_{qz}^{(\mu,v)}(\mathbf{X}^{(v)}), f_{qz}^{(\sigma,v)}(\mathbf{X}^{(v)})\}_{v \in \mathcal{V}_u}$ , in the same fashion as  $q_{\phi_h}$  to infer parameters of  $q_{\phi_z}$ .



# Objective Function

$$\begin{aligned}\mathcal{L} &= -\text{ELBO} + \mathcal{L}_{\text{FGW}} \\ &= \mathbb{E}_{q_{\phi_z}(\mathcal{Z}_u, \mathcal{H}_u | \mathcal{X}_u)} \log p_{\theta}(\mathcal{Z}_u | \mathcal{A}, \mathcal{H}) + \mathbb{E}_{q_{\phi_z}(\mathcal{Z}_s, \mathcal{H}_s | \mathcal{X}_s, \mathfrak{A}_s)} \log p_{\theta}(\mathcal{Z}_s | \mathcal{A}, \mathcal{H}) \\ &\quad - \mathbb{E}_{q_{\phi_z}(\mathcal{Z}_u | \mathcal{X}_u)} \log q_{\phi_z}(\mathcal{Z}_u | \mathcal{X}_u) - \mathbb{E}_{q_{\phi_z}(\mathcal{Z}_s | \mathcal{X}_s, \mathfrak{A}_s)} \log q_{\phi_z}(\mathcal{Z}_s | \mathcal{X}_s, \mathfrak{A}_s) \\ &\quad + \mathbb{E}_{q_{\phi_h}(\mathcal{H}_u | \mathcal{X}_u)} \log p(\mathcal{H}_u) + \mathbb{E}_{q_{\phi_h}(\mathcal{H}_s | \mathcal{X}_s, \mathfrak{A}_s)} \log p(\mathcal{H}_s) \\ &\quad - \mathbb{E}_{q_{\phi_h}(\mathcal{H}_u | \mathcal{X}_u)} \log q_{\phi_h}(\mathcal{H}_u | \mathcal{X}_u) - \mathbb{E}_{q_{\phi_h}(\mathcal{H}_s | \mathcal{X}_s, \mathfrak{A}_s)} \log q_{\phi_h}(\mathcal{H}_s | \mathcal{X}_u, \mathfrak{A}_s) \\ &\quad + \mathbb{E}_{q_{\phi_z}(\mathcal{Z}_u | \mathcal{X}_u)} \log p_{\theta_x}(\mathcal{X}_u | \mathcal{Z}_u) + \mathbb{E}_{q_{\phi_z}(\mathcal{Z}_s | \mathcal{X}_s, \mathfrak{A}_s)} \log p_{\theta_x}(\mathcal{X}_s | \mathcal{Z}_s) \\ &\quad + \mathbb{E}_{q_{\phi_h}(\mathcal{H}_s | \mathcal{X}_s, \mathfrak{A}_s)} \log p_{\theta_g}(\mathfrak{A}_s | \mathcal{H}_s) + \sum_{v \in \mathcal{V}} \sum_{\substack{v' \in \mathcal{V} \\ v' \neq v}} \mathcal{D}_{\text{FGW}} \left( p(\mathbf{H}^{(v)}), p(\mathbf{H}^{(v')}) \right).\end{aligned}$$

**Figure 1:** Overall loss function as the sum of the negative variational ELBO and FGW regularization terms

# Datasets and Evaluation Metrics

## Datasets:

- ▶ microbiome-metabolite interactions in cystic fibrosis (CF)
  - ▶ Positive accuracy refers to the accuracy of identifying validated interactions with *P. aeruginosa*.
  - ▶ Negative accuracy exploits the fact that there should not be any common metabolite targets between known anaerobic microbes (*Veillonella*, *Fusobacterium*, *Prevotella*, and *Streptococcus*) and notable pathogen *P. aeruginosa*.
  - ▶ Having both higher positive and negative accuracy is desired
- ▶ gene-drug interactions in precision medicine
  - ▶ prediction sensitivity of identifying known interactions in the test sets
  - ▶ the average density of the overall constructed graphs

# Experiments

Table 1: Comparison of positive accuracy (in %) on CF dataset at negative accuracy of  $> 97\%$ .

	SRCA	BCCA	MoReL <sub>uu</sub>	MoReL <sub>us</sub>
Positive accuracy	26.41	$28.30 \pm 3.21$	$56.16 \pm 1.85$	$63.77 \pm 1.11$

Table 2: Comparison of prediction sensitivity (in %) in the precision medicine experiment.

Avg. degree	0.10	0.15	0.20	0.25	0.30	0.40	0.50
SRCA	8.03	12.00	17.15	20.70	26.85	34.93	45.79
BCCA	$9.65 \pm 0.75$	$14.34 \pm 0.06$	$18.96 \pm 0.42$	$23.29 \pm 0.52$	$28.22 \pm 0.66$	$38.02 \pm 2.15$	$46.88 \pm 1.88$
MoReL <sub>uu</sub>	$11.29 \pm 0.16$	$15.74 \pm 0.62$	$21.21 \pm 0.81$	$26.20 \pm 1.10$	$30.47 \pm 1.07$	$39.05 \pm 0.75$	$50.19 \pm 0.19$
MoReL <sub>us</sub>	$12.79 \pm 0.39$	$17.51 \pm 2.21$	$22.82 \pm 1.01$	$29.58 \pm 1.08$	$35.05 \pm 1.27$	$45.74 \pm 1.75$	$53.16 \pm 0.96$

# Comparison with Bayrel

Table 3: Positive accuracy (%) on CF dataset.

	BayReL	MoReL <sub>ss</sub>
Positive Acc.	$82.70 \pm 4.70$	$89.50 \pm 3.29$

Table 4: Prediction sensitivity (%) in the precision medicine experiment.

Avg. degree	BayReL	MoReL <sub>ss</sub>
0.4	$47.90 \pm 0.43$	$49.24 \pm 1.64$
0.5	$56.76 \pm 0.50$	$58.92 \pm 0.40$

Table 5: Positive accuracy (%) on CF dataset with unpaired samples.

	BayReL	MoReL <sub>us</sub>
Positive Acc.	31.56	$63.24 \pm 2.13$
Negative Acc.	72	97

# Conclusion

- ▶ A novel Bayesian deep generative model that efficiently infers hidden molecular relations across heterogeneous views of data.
- ▶ To handle unpaired samples across the views of data;
- ▶ To combine multiple views from different data sources with any number of missing samples.