

Conditional Graph Information Bottleneck for Molecular Relational Learning

Namkyeong Lee, Dongmin Hyun, Gyoung S. Na,
Sungwon Kim, Junseok Lee, Chanyoung Park



TABLE OF CONTENTS

▪ **Background**

- Molecular Relational Learning
- Functional Group
- Information Bottleneck

▪ **Motivation**

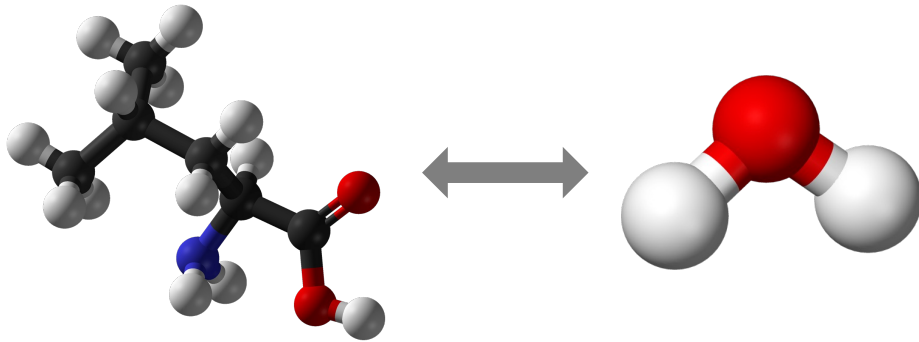
▪ **Conditional Graph Information Bottleneck**

▪ **Experiments**

▪ **Conclusion**



BACKGROUND MOLECULAR RELATIONAL LEARNING



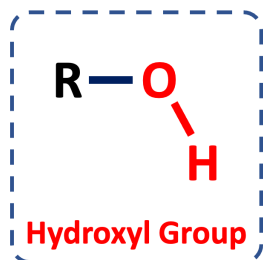
Molecular Relational Learning

Learning the interaction behavior between a pair of molecules

Examples

- Predicting optical properties when a **Chromophore** and **Solvent** react
- Predicting solubility when a **solute** and **solvent** react
- Predicting side effects when taking **two types of drugs** simultaneously

BACKGROUND FUNCTIONAL GROUP



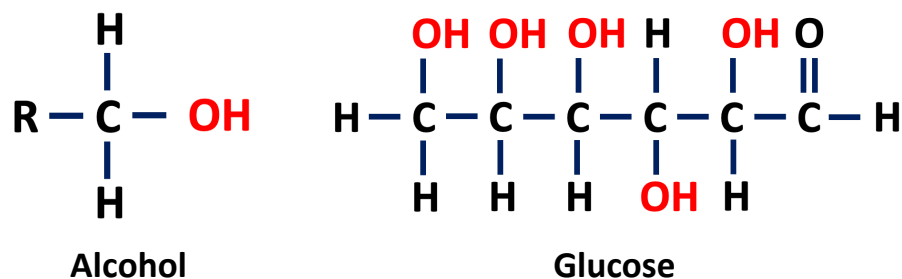
Functional Group

Specific atomic groups that play an important role in determining the chemical reactivity of organic compounds

Compounds with the same functional group generally have similar properties and undergo similar chemical reactions

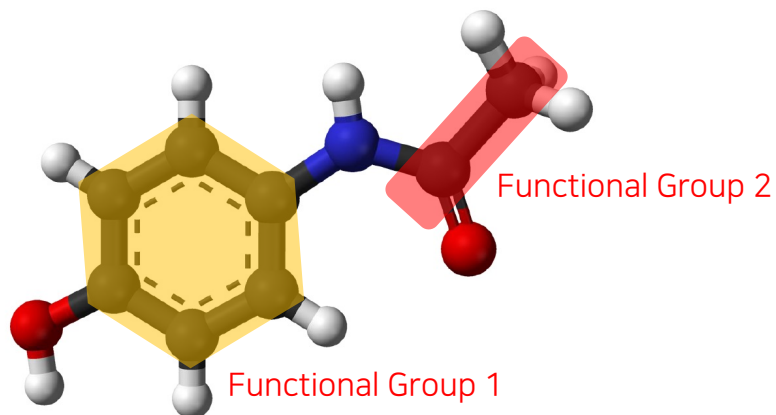
Examples

The hydroxyl group structure has the characteristic of increasing the polarity of the molecule



Hence, it is important to consider **functional group for molecular relational learning**

BACKGROUND FUNCTIONAL GROUP



Functional Group

Specific atomic groups that play an important role in determining the chemical reactivity of organic compounds

Compounds with the same functional group generally have similar properties and undergo similar chemical reactions

Molecule can be represented as a **graph**
Functional group can be represented as a **subgraph**

Recently, **information theory**-based approaches have been proposed to detect important subgraph

BACKGROUND INFORMATION BOTTLENECK

Definition 2.1. (Information Bottleneck) Given random variables X and Y , the Information Bottleneck principle aims to compress X to a bottleneck random variable T , while keeping the information relevant for predicting Y :

$$\min_T -I(Y; T) + \beta I(X; T) \quad (2)$$

where β is a Lagrangian multiplier for balancing the two mutual information terms.

Information Bottleneck Theory

A theoretical approach to trade-off between information **compression** and **preservation**

Minimize MI between X and T

→ T should contain minimal information about X

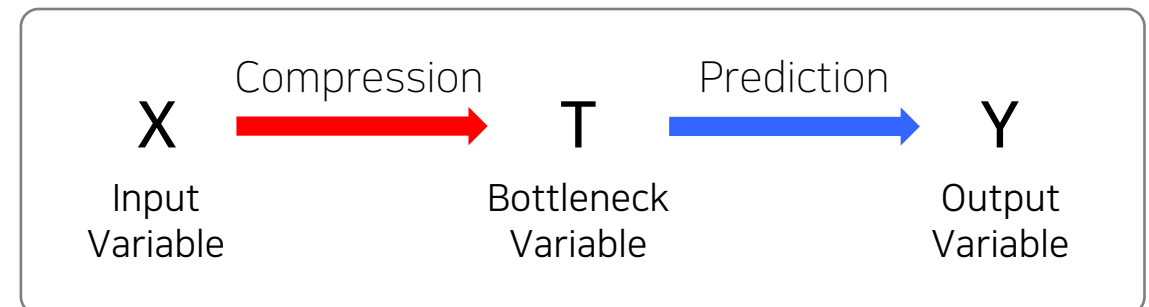
→ **Compression**

$$\min_T -I(Y; T) + \beta I(X; T)$$

Maximize MI between T and Y

→ T should contain as much information about Y as possible

→ **Prediction**



BACKGROUND GRAPH INFORMATION BOTTLENECK

Information Bottleneck Graph

Subgraph that maximally preserves the property of the original graph

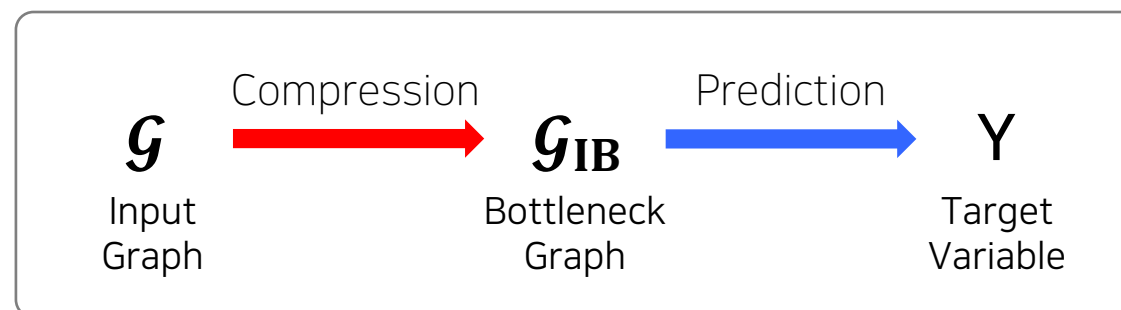
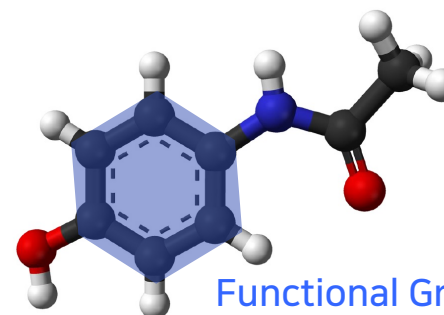
→ Motif in ordinary graphs

→ Functional group in molecules

Definition 2.2. (IB-Graph) For a graph $\mathcal{G} = (\mathbf{X}, \mathbf{A})$ and its label information \mathbf{Y} , the optimal graph $\mathcal{G}_{\text{IB}} = (\mathbf{X}_{\text{IB}}, \mathbf{A}_{\text{IB}})$ discovered under the IB principle is denoted as IB-Graph:

$$\mathcal{G}_{\text{IB}} = \arg \min_{\mathcal{G}_{\text{IB}}} -I(\mathbf{Y}; \mathcal{G}_{\text{IB}}) + \beta I(\mathcal{G}; \mathcal{G}_{\text{IB}}) \quad (3)$$

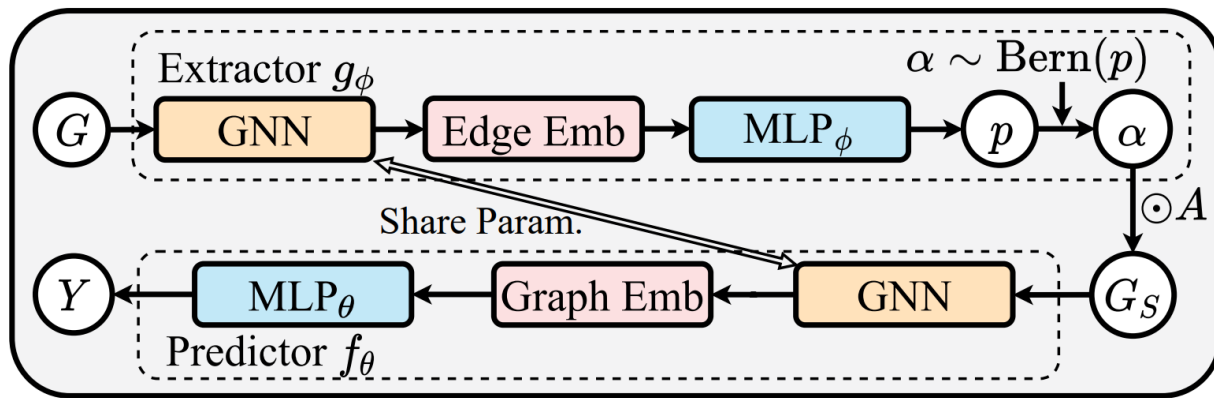
where \mathbf{X}_{IB} and \mathbf{A}_{IB} denote the task-relevant feature set and the adjacency matrix of \mathcal{G} , respectively.



BACKGROUND GRAPH INFORMATION BOTTLENECK

Extract a subgraph in terms of edges

Model an edge based Bernoulli distribution to perform graph compression



Model objective

$$\min_{\phi} -I(G_S; Y) + \beta I(G_S; G), \text{ s.t. } G_S \sim g_{\phi}(G).$$

Proposed variational upperbound

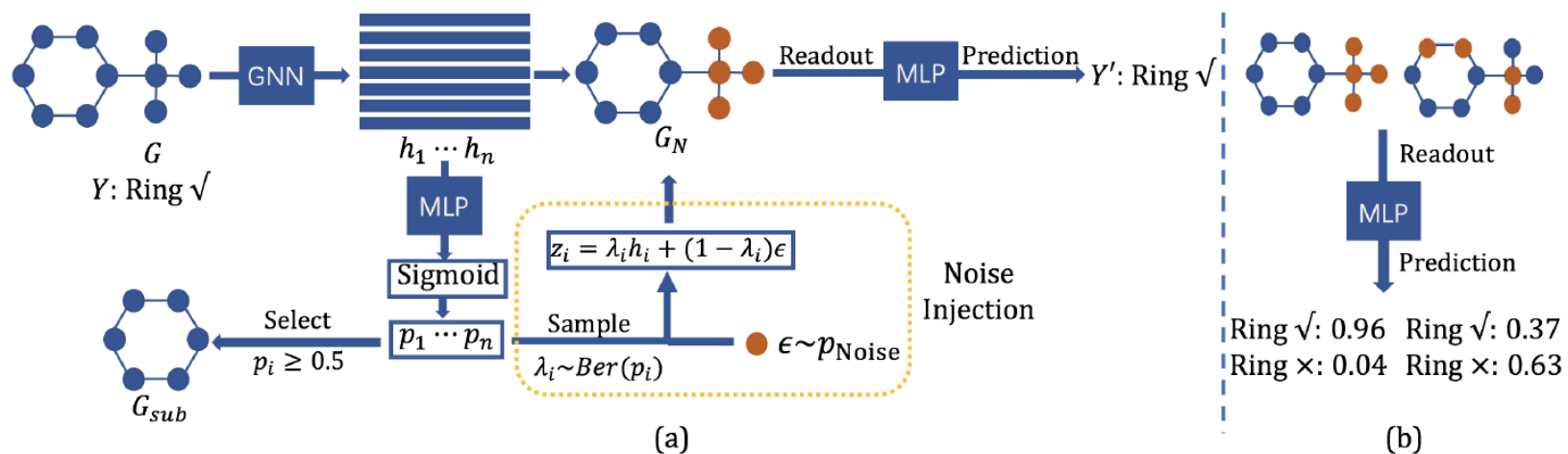
$$\min_{\theta, \phi} -\mathbb{E} [\log \mathbb{P}_{\theta}(Y|G_S)] + \beta \mathbb{E} [\text{KL}(\mathbb{P}_{\phi}(G_S|G) || \mathbb{Q}(G_S))],$$

$$\text{ s.t. } G_S \sim \mathbb{P}_{\phi}(G_S|G). \quad (8)$$

BACKGROUND GRAPH INFORMATION BOTTLENECK

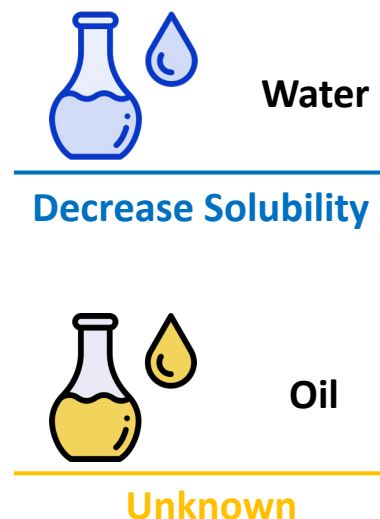
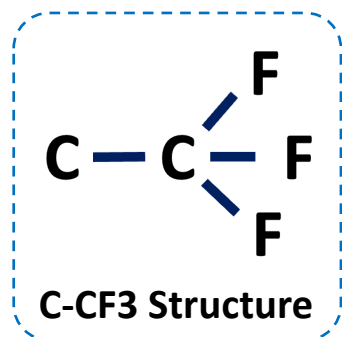
Extract a subgraph in terms of nodes

Inject noise into node embeddings to perform graph compression



Can Information bottleneck theory also benefits molecular relational learning?

MOTIVATION



Functional Group

Specific atomic groups that play an important role in determining the chemical reactivity of organic compounds

Compounds with the same functional group generally have similar properties and undergo similar chemical reactions

On the other hand, the role of functional group varies depending on which solvent the solute reacts with!

Examples: C-CF3 structure in molecules

It is important to **consider the paired solvent** when detecting the important substructure from solute
→ Existing approaches for information bottleneck cannot capture such a prior knowledge

METHODOLOGY CONDITIONAL GRAPH INFORMATION BOTTLENECK

Conditional Graph Information Bottleneck

Consider Graph 2 (Solvent) when detecting the important subgraph from Graph 1 (Solute)

$$\mathcal{G}_{\text{IB}} = \arg \min_{\mathcal{G}_{\text{IB}}} -I(\mathbf{Y}; \mathcal{G}_{\text{IB}}) + \beta I(\mathcal{G}; \mathcal{G}_{\text{IB}})$$

Graph Information Bottleneck



$$\mathcal{G}_{\text{CIB}}^1 = \arg \min_{\mathcal{G}_{\text{CIB}}^1} -I(\mathbf{Y}; \mathcal{G}_{\text{CIB}}^1 | \mathcal{G}^2) + \beta I(\mathcal{G}^1; \mathcal{G}_{\text{CIB}}^1 | \mathcal{G}^2)$$

Conditional Graph Information Bottleneck

METHODOLOGY CONDITIONAL GRAPH INFORMATION BOTTLENECK

Proof of Lemma 3.3

Assuming that \mathcal{G}^1 , $\mathcal{G}_{\text{CIB}}^1$, \mathcal{G}_n^1 , \mathcal{G}^2 , and Y satisfy the Markov condition $(Y, \mathcal{G}_n^1, \mathcal{G}^2) \rightarrow \mathcal{G}^1 \rightarrow \mathcal{G}_{\text{CIB}}^1$, we have the following inequality due to data processing inequality:

$$\begin{aligned}
 I(\mathcal{G}^1; \mathcal{G}_{\text{CIB}}^1 | \mathcal{G}^2) &= I(\mathcal{G}_{\text{CIB}}^1; \mathcal{G}^1, \mathcal{G}^2) - I(\mathcal{G}_{\text{CIB}}^1; \mathcal{G}^2) \\
 &\geq I(\mathcal{G}_{\text{CIB}}^1; Y, \mathcal{G}_n^1, \mathcal{G}^2) - I(\mathcal{G}_{\text{CIB}}^1; \mathcal{G}^2) \\
 &= I(\mathcal{G}_{\text{CIB}}^1; \mathcal{G}_n^1, \mathcal{G}^2) + I(\mathcal{G}_{\text{CIB}}^1; Y | \mathcal{G}_n^1, \mathcal{G}^2) - I(\mathcal{G}_{\text{CIB}}^1; \mathcal{G}^2) \\
 &= I(\mathcal{G}_{\text{CIB}}^1; \mathcal{G}_n^1, \mathcal{G}^2) + I(\mathcal{G}_{\text{CIB}}^1; Y | \mathcal{G}_n^1, \mathcal{G}^2) \quad (1)
 \end{aligned}$$

LEMMA 3.3. (Noise Invariance) Given a pair of graphs $(\mathcal{G}^1, \mathcal{G}^2)$ and its label information Y , let \mathcal{G}_n^1 be a task irrelevant noise in the input graph \mathcal{G}^1 . Then, the following inequality holds:

$$I(\mathcal{G}_{\text{CIB}}^1; \mathcal{G}_n^1 | \mathcal{G}^2) \leq -I(Y; \mathcal{G}_{\text{CIB}}^1 | \mathcal{G}^2) + I(\mathcal{G}^1; \mathcal{G}_{\text{CIB}}^1 | \mathcal{G}^2) \quad (6)$$

Suppose that \mathcal{G}_n^1 and Y , \mathcal{G}_n^1 and \mathcal{G}^2 , and joint random variable $(\mathcal{G}_n^1, \mathcal{G}^2)$ and Y are independent respectively. Then, for $I(\mathcal{G}_{\text{CIB}}^1; Y | \mathcal{G}_n^1, \mathcal{G}^2)$ we have:

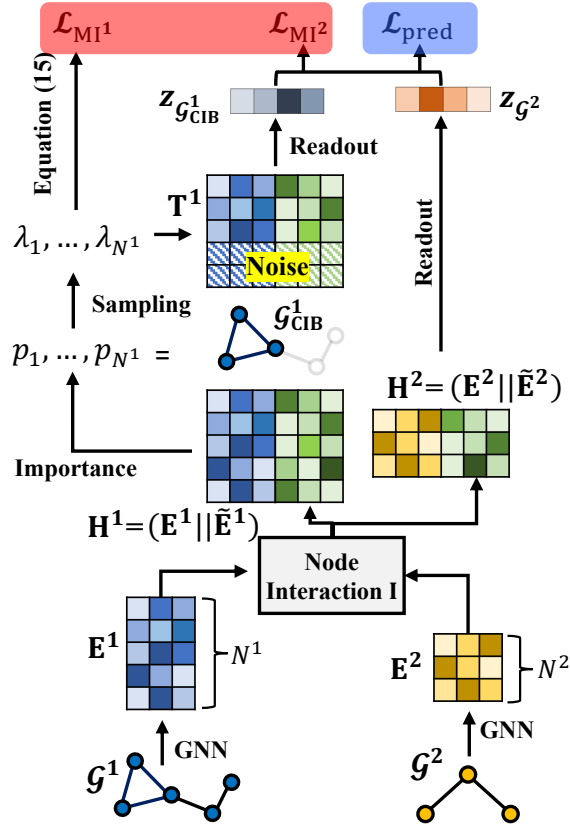
$$\begin{aligned}
 I(\mathcal{G}_{\text{CIB}}^1; Y | \mathcal{G}_n^1, \mathcal{G}^2) &= H(Y | \mathcal{G}_n^1, \mathcal{G}^2) - H(Y | \mathcal{G}_n^1, \mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2) \\
 &\geq H(Y | \mathcal{G}^2) - H(Y | \mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2) \\
 &= I(Y; \mathcal{G}_{\text{CIB}}^1 | \mathcal{G}^2) \quad (2)
 \end{aligned}$$

By plugging Equation (2) into Equation (1), we have:

$$I(\mathcal{G}^1; \mathcal{G}_{\text{CIB}}^1 | \mathcal{G}^2) \geq I(\mathcal{G}_{\text{CIB}}^1; \mathcal{G}_n^1, \mathcal{G}^2) + I(Y; \mathcal{G}_{\text{CIB}}^1 | \mathcal{G}^2)$$

By minimizing CGIB objective function,
the model learns a CIB-Graph with the smallest mutual information with task-irrelevant noise

METHODOLOGY CONDITIONAL GRAPH INFORMATION BOTTLENECK



$$\min -I(Y; \mathcal{G}_{CIB}^1 | \mathcal{G}^2) + \beta I(\mathcal{G}^1; \mathcal{G}_{CIB}^1 | \mathcal{G}^2)$$

Overall procedure

Decompose the conditional MI based on the chain rule of MI, and then derive the upper bound of the decomposed terms

$$-I(Y; \mathcal{G}_{CIB}^1 | \mathcal{G}^2) = -I(Y; \mathcal{G}_{CIB}^1, \mathcal{G}^2) + I(Y; \mathcal{G}^2)$$

$$-I(Y; \mathcal{G}_{CIB}^1, \mathcal{G}^2) \leq \mathbb{E}_{\mathcal{G}_{CIB}^1, \mathcal{G}^2, Y} [-\log p_{\theta}(Y | \mathcal{G}_{CIB}^1, \mathcal{G}^2)]$$

Prediction Loss

$$I(\mathcal{G}^1; \mathcal{G}_{CIB}^1 | \mathcal{G}^2) = I(\mathcal{G}_{CIB}^1; \mathcal{G}^1, \mathcal{G}^2) - I(\mathcal{G}_{CIB}^1; \mathcal{G}^2)$$

$$I(\mathcal{G}_{CIB}^1; \mathcal{G}^1, \mathcal{G}^2) \leq \mathbb{E}_{\mathcal{G}^1, \mathcal{G}^2} \left[-\frac{1}{2} \log A + \frac{1}{2N^1} A + \frac{1}{2N^1} B^2 \right]$$

$$:= \mathcal{L}_{MI^1}(\mathcal{G}_{CIB}^1, \mathcal{G}^1, \mathcal{G}^2)$$

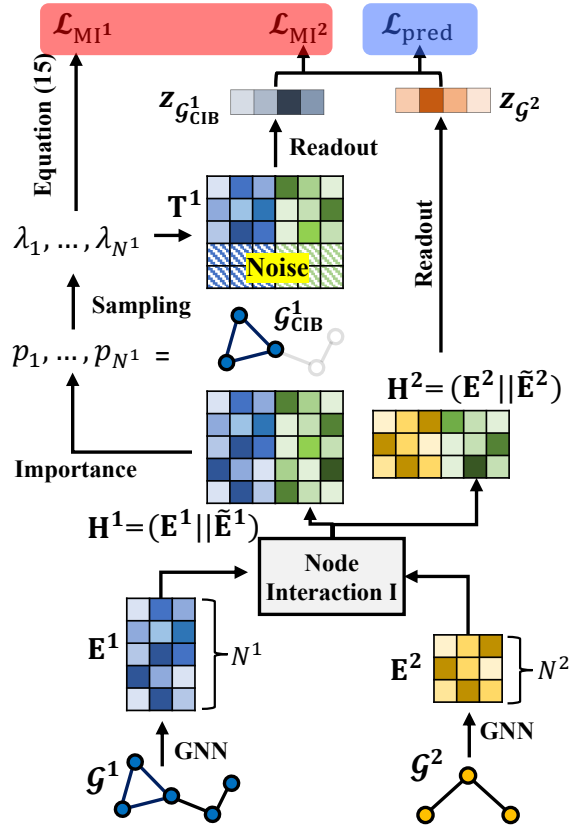
$$-I(\mathcal{G}_{CIB}^1; \mathcal{G}^2) \leq \mathbb{E}_{\mathcal{G}_{CIB}^1, \mathcal{G}^2} [-\log p_{\xi}(\mathcal{G}^2 | \mathcal{G}_{CIB}^1)]$$

$$:= \mathcal{L}_{MI^2}(\mathcal{G}_{CIB}^1, \mathcal{G}^2)$$

Compression Loss

$$\min -I(Y; \mathcal{G}_{\text{CIB}}^1 | \mathcal{G}^2) + \beta I(\mathcal{G}^1; \mathcal{G}_{\text{CIB}}^1 | \mathcal{G}^2)$$

METHODOLOGY CONDITIONAL GRAPH INFORMATION BOTTLENECK



$$-I(Y; \mathcal{G}_{\text{CIB}}^1 | \mathcal{G}^2) = -I(Y; \mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2) + I(Y; \mathcal{G}^2) \quad \because \text{Chain rule of mutual information}$$

Direct calculation of mutual information is intractable;
Instead, we minimize the upper bound

Proposition. (Upper bound of $-I(Y; \mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2)$) Given a pair of graph $(\mathcal{G}^1, \mathcal{G}^2)$, its label information Y , and the learned CIB-graph $\mathcal{G}_{\text{CIB}}^1$, we have:

$$-I(Y; \mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2) \leq \mathbb{E}_{Y, \mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2} [\log p_{\theta}(Y | \mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2)]$$

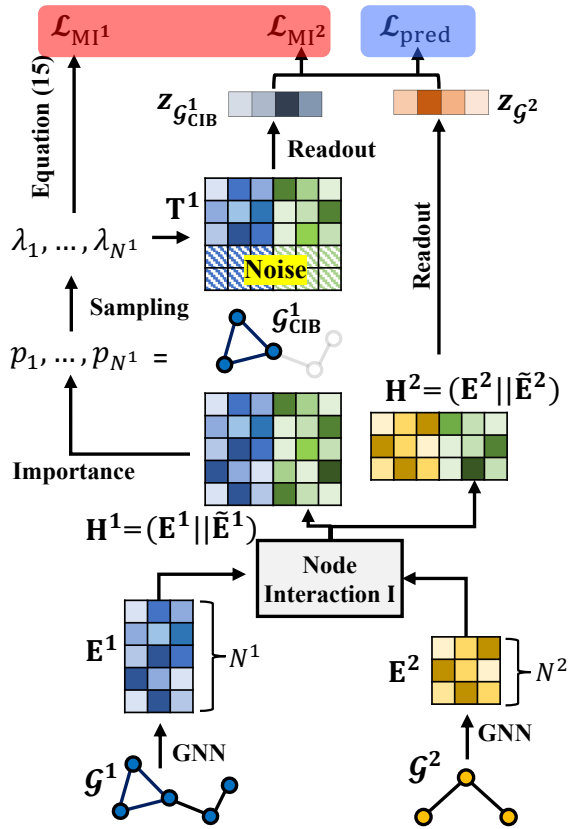
where $p_{\theta}(Y | \mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2)$ is variational approximation of $p(Y | \mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2)$.

Proof. By the definition of mutual information and introducing variational approximation $p_{\theta}(Y | \mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2)$ of intractable distribution $p(Y | \mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2)$, we have:

$$\begin{aligned} I(Y; \mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2) &= \mathbb{E}_{Y, \mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2} \left[\log \frac{p(Y | \mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2)}{p(Y)} \right] \\ &= \mathbb{E}_{Y, \mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2} \left[\log \frac{p_{\theta}(Y | \mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2)}{p(Y)} \right] + \mathbb{E}_{\mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2} [p(Y | \mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2) || p_{\theta}(Y | \mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2)] \\ &\geq \mathbb{E}_{Y, \mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2} \left[\log \frac{p_{\theta}(Y | \mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2)}{p(Y)} \right] \quad \because \text{Non-negativity of KL divergence} \\ &= \mathbb{E}_{Y, \mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2} [\log p_{\theta}(Y | \mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2)] + H(Y) \end{aligned}$$

$$\min -I(Y; \mathcal{G}_{\text{CIB}}^1 | \mathcal{G}^2) + \beta I(\mathcal{G}^1; \mathcal{G}_{\text{CIB}}^1 | \mathcal{G}^2)$$

METHODOLOGY CONDITIONAL GRAPH INFORMATION BOTTLENECK



$$-I(Y; \mathcal{G}_{\text{CIB}}^1 | \mathcal{G}^2) = -I(Y; \mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2) + I(Y; \mathcal{G}^2) \quad \because \text{Chain rule of mutual information}$$

Direct calculation of mutual information is intractable;
 Instead, we minimize the upper bound

Proposition. (Upper bound of $-I(Y; \mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2)$) Given a pair of graph $(\mathcal{G}^1, \mathcal{G}^2)$, its label information Y , and the learned CIB-graph $\mathcal{G}_{\text{CIB}}^1$, we have:

$$-I(Y; \mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2) \leq \mathbb{E}_{Y, \mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2} [\log p_\theta(Y | \mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2)]$$

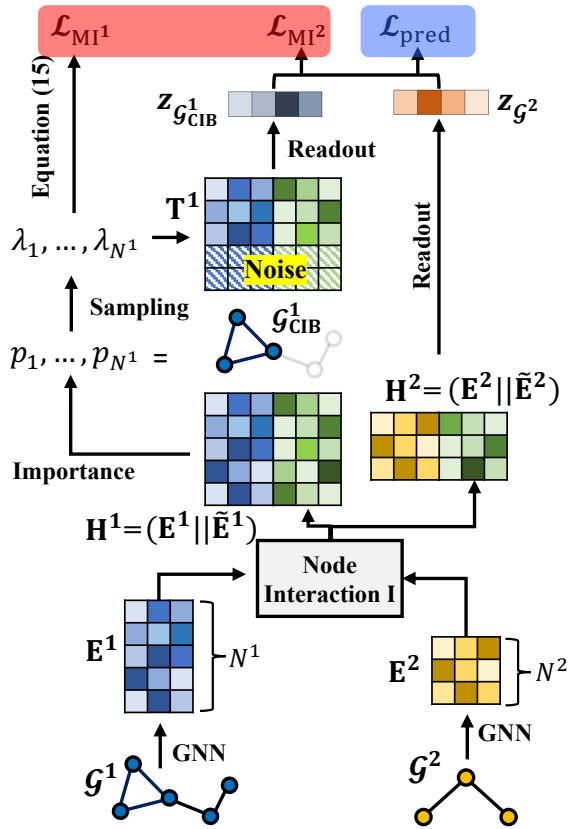
where $p_\theta(Y | \mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2)$ is variational approximation of $p(Y | \mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2)$.

Implementation

- Consider $p_\theta(Y | \mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2)$ as a predictor parameterized by θ , which outputs the model prediction Y based on the input pair $(\mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2)$
- The upper bound is minimized by minimizing the prediction loss $\mathcal{L}_{\text{pred}}(Y, \mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2)$

$$\min -I(Y; \mathcal{G}_{\text{CIB}}^1 | \mathcal{G}^2) + \beta I(\mathcal{G}^1; \mathcal{G}_{\text{CIB}}^1 | \mathcal{G}^2)$$

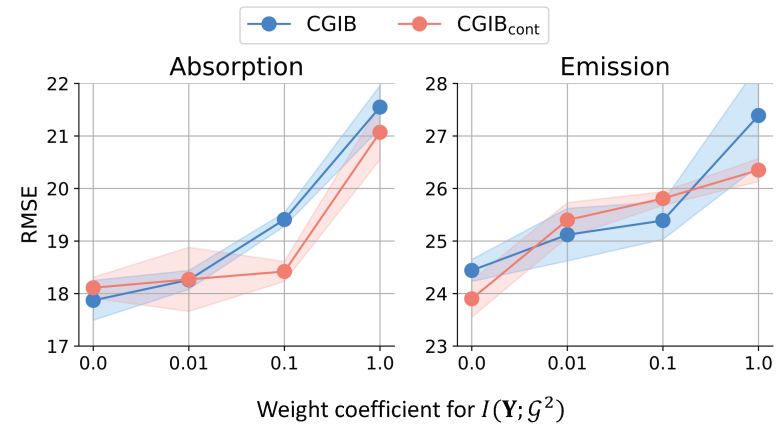
METHODOLOGY CONDITIONAL GRAPH INFORMATION BOTTLENECK



$$-I(Y; \mathcal{G}_{\text{CIB}}^1 | \mathcal{G}^2) = -I(Y; \mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2) + I(Y; \mathcal{G}^2)$$

∴ Chain rule of mutual information

The 2nd term is empirically found to be not helpful



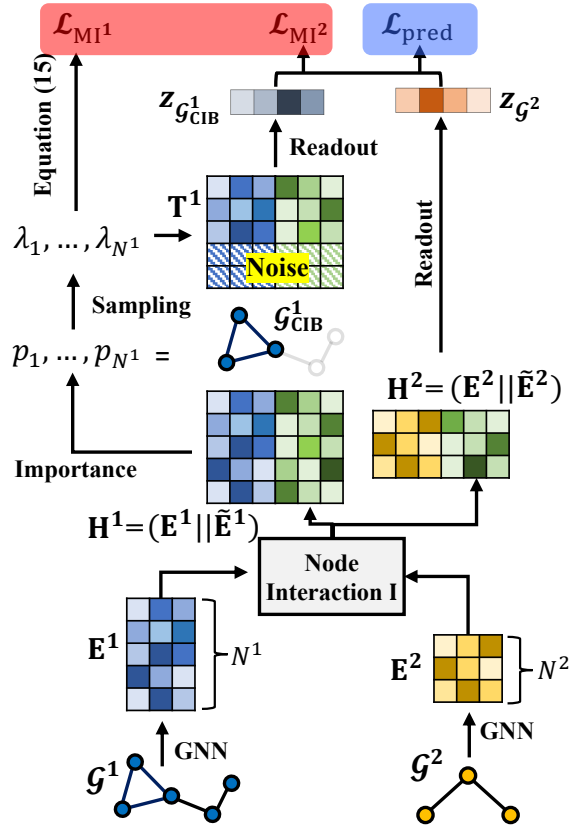
Following VIB, we treat $r(Y)$ as fixed spherical Gaussian,
 $I(Y; \mathcal{G}^2) \leq \mathbb{E}_{\mathcal{G}^2} [KL(p_{\xi}(Y|\mathcal{G}^2) || r(Y))]$
 where $r(Y) \sim N(Y|0, 1)$

Increasing the contribution of this term deteriorates the model performance

We remove the term $I(Y, \mathcal{G}^2)$ from the model

$$\min -I(Y; \mathcal{G}_{\text{CIB}}^1 | \mathcal{G}^2) + \beta I(\mathcal{G}^1; \mathcal{G}_{\text{CIB}}^1 | \mathcal{G}^2)$$

METHODOLOGY CONDITIONAL GRAPH INFORMATION BOTTLENECK



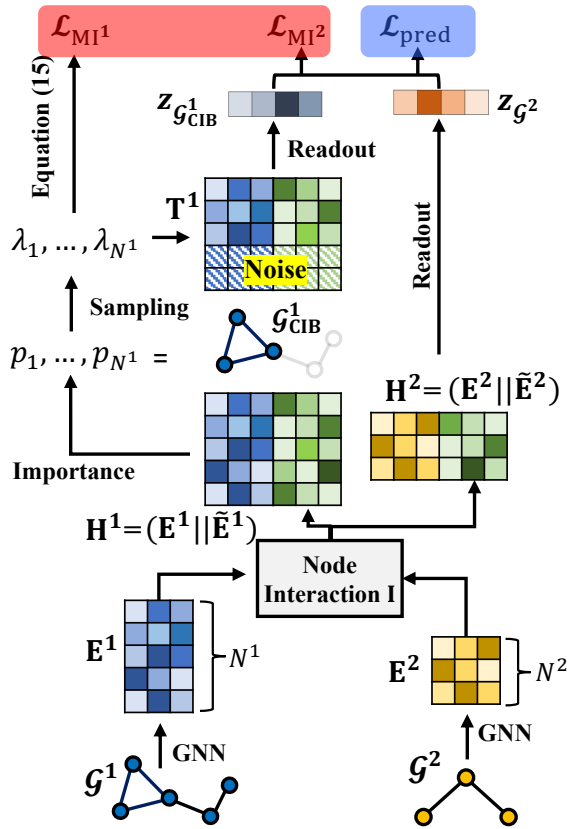
$$I(\mathcal{G}^1; \mathcal{G}_{\text{CIB}}^1 | \mathcal{G}^2) = \underbrace{I(\mathcal{G}_{\text{CIB}}^1; \mathcal{G}^1, \mathcal{G}^2)}_{\mathcal{L}_{\text{MI}^1}} - \underbrace{I(\mathcal{G}_{\text{CIB}}^1; \mathcal{G}^2)}_{\mathcal{L}_{\text{MI}^2}} \quad \because \text{Chain rule of mutual information}$$

$\mathcal{L}_{\text{MI}^1}$: **Compression through Noise injection**
 → Injecting noise into unimportant nodes

$\mathcal{L}_{\text{MI}^2}$: **Solute Prediction**
 → Encourage $\mathcal{G}_{\text{CIB}}^1$ to contain as much information about \mathcal{G}^2 as possible
 → The term that arises from the Conditional Mutual Information
 → **Key to the success of CGIB!** Enables the conditional information compression of CGIB

$$\min -I(Y; \mathcal{G}_{\text{CIB}}^1 | \mathcal{G}^2) + \beta I(\mathcal{G}^1; \mathcal{G}_{\text{CIB}}^1 | \mathcal{G}^2)$$

METHODOLOGY CONDITIONAL GRAPH INFORMATION BOTTLENECK



$$I(\mathcal{G}^1; \mathcal{G}_{\text{CIB}}^1 | \mathcal{G}^2) = \underbrace{I(\mathcal{G}_{\text{CIB}}^1; \mathcal{G}^1, \mathcal{G}^2)}_{\mathcal{L}_{\text{MI}^1}} - I(\mathcal{G}_{\text{CIB}}^1; \mathcal{G}^2) \quad \because \text{Chain rule of mutual information}$$

Compression through Noise Injection

* Injecting noise into unimportant nodes

H_i^1 : Representation of node i of \mathcal{G}^1 that contains information about both $\mathcal{G}^1, \mathcal{G}^2$

$p_i = \text{MLP}(H_i^1)$: Important of node i of \mathcal{G}^1

$T_i^1 = \lambda_i H_i^1 + (1 - \lambda_i) \epsilon$ where $\lambda_i \sim \text{Bernoulli}(p_i)$ and $\epsilon \sim N(\mu_{H^1}, \sigma_{H^1}^2)$

Intuition) Unimportant nodes would not affect the model performance even if they are replaced with noise

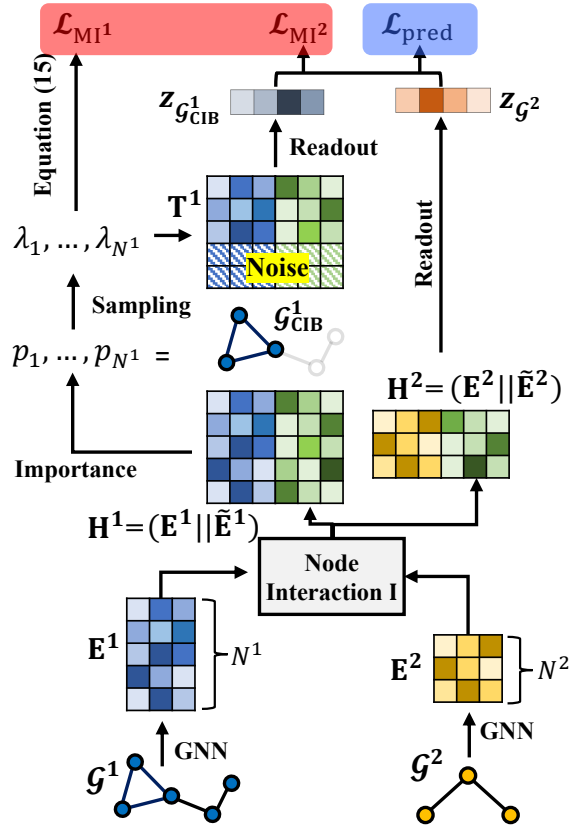
Upper bound of $I(\mathcal{G}_{\text{CIB}}^1; \mathcal{G}^1, \mathcal{G}^2)$

$$I(\mathcal{G}_{\text{CIB}}^1; \mathcal{G}^1, \mathcal{G}^2) \leq \mathbb{E}_{\mathcal{G}^1, \mathcal{G}^2} \left[-\frac{1}{2} \log A + \frac{1}{2N^1} A + \frac{1}{2N^1} B^2 \right] \quad \text{where } A = \sum_{j=1}^{N^1} (1 - \lambda_j)^2 \text{ and } B = \frac{\sum_{j=1}^{N^1} \lambda_j (H_j^1 - \mu_{H^1})^2}{\sigma_{H^1}}$$

$$:= \mathcal{L}_{\text{MI}^1}(\mathcal{G}_{\text{CIB}}^1, \mathcal{G}^1, \mathcal{G}^2)$$

$$\min -I(Y; \mathcal{G}_{\text{CIB}}^1 | \mathcal{G}^2) + \beta I(\mathcal{G}^1; \mathcal{G}_{\text{CIB}}^1 | \mathcal{G}^2)$$

METHODOLOGY CONDITIONAL GRAPH INFORMATION BOTTLENECK



Upper bound of $I(\mathcal{G}_{\text{CIB}}^1; \mathcal{G}^1, \mathcal{G}^2)$

$$I(\mathcal{G}_{\text{CIB}}^1; \mathcal{G}^1, \mathcal{G}^2) \leq \mathbb{E}_{\mathcal{G}^1, \mathcal{G}^2} \left[-\frac{1}{2} \log A + \frac{1}{2N^1} A + \frac{1}{2N^1} B^2 \right] \quad \text{where } A = \sum_{j=1}^{N^1} (1 - \lambda_j)^2 \text{ and } B = \frac{\sum_{j=1}^{N^1} \lambda_j (H_j^1 - \mu_{H^1})^2}{\sigma_{H^1}}$$

$$:= \mathcal{L}_{MI^1}(\mathcal{G}_{\text{CIB}}^1, \mathcal{G}^1, \mathcal{G}^2)$$

Proof. Given the perturbed graph $\mathcal{G}_{\text{CIB}}^1$ and its representation $z_{\mathcal{G}_{\text{CIB}}^1}$, we assume there is no information loss during the readout process, i.e., $I(z_{\mathcal{G}_{\text{CIB}}^1}; \mathcal{G}^1, \mathcal{G}^2) = I(\mathcal{G}_{\text{CIB}}^1; \mathcal{G}^1, \mathcal{G}^2)$.

$$I(z_{\mathcal{G}_{\text{CIB}}^1}; \mathcal{G}^1, \mathcal{G}^2) = \mathbb{E}_{z_{\mathcal{G}_{\text{CIB}}^1}, \mathcal{G}^1, \mathcal{G}^2} \left[-\log \frac{p_{\phi}(z_{\mathcal{G}_{\text{CIB}}^1} | \mathcal{G}^1, \mathcal{G}^2)}{p(z_{\mathcal{G}_{\text{CIB}}^1})} \right]$$

$$= \mathbb{E}_{\mathcal{G}^1, \mathcal{G}^2} \left[-\log \frac{p_{\phi}(z_{\mathcal{G}_{\text{CIB}}^1} | \mathcal{G}^1, \mathcal{G}^2)}{q(z_{\mathcal{G}_{\text{CIB}}^1})} \right] - \mathbb{E}_{z_{\mathcal{G}_{\text{CIB}}^1}, \mathcal{G}^1, \mathcal{G}^2} \left[KL(p(z_{\mathcal{G}_{\text{CIB}}^1}) || q(z_{\mathcal{G}_{\text{CIB}}^1})) \right]$$

$$\leq \mathbb{E}_{z_{\mathcal{G}_{\text{CIB}}^1}, \mathcal{G}^1, \mathcal{G}^2} \left[KL(p_{\phi}(z_{\mathcal{G}_{\text{CIB}}^1} | \mathcal{G}^1, \mathcal{G}^2) || q(z_{\mathcal{G}_{\text{CIB}}^1})) \right] \quad (1) \quad \because \text{Non-negativity of KL divergence}$$

Assuming that $q(z_{\mathcal{G}_{\text{CIB}}^1})$ is Gaussian distribution.

The noise $\varepsilon \sim N(\mu_{H^1}, \sigma_{H^1})$ is sampled from Gaussian distribution where μ_{H^1} and σ_{H^1} are mean and variance of H^1 .

Thus, $q(z_{\mathcal{G}_{\text{CIB}}^1}) = N(N^1 \mu_{H^1}, N^1 \sigma_{H^1})$ (2) \because Summation of Gaussian is Gaussian

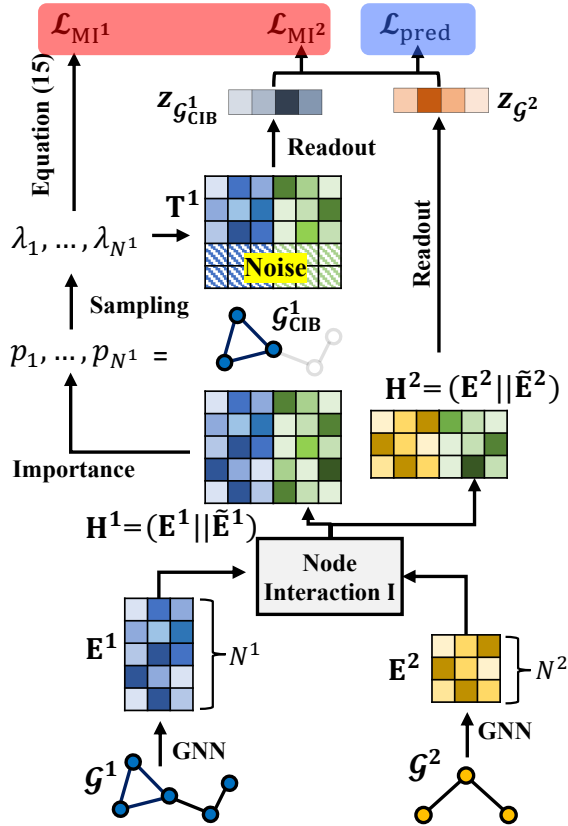
And, $p(z_{\mathcal{G}_{\text{CIB}}^1} | \mathcal{G}^1, \mathcal{G}^2) = N(N^1 \mu_{H^1} + \sum_{j=1}^{N^1} \lambda_j H_j^1 - \sum_{j=1}^{N^1} \lambda_j \mu_{H^1}, \sum_{j=1}^{N^1} (1 - \lambda_j)^2 \sigma_{H^1}^2)$ (3)

By plugging Equation (2) and (3) into (1), we have:

$$-I(\mathcal{G}_{\text{CIB}}^1; \mathcal{G}^1, \mathcal{G}^2) \leq \mathbb{E}_{\mathcal{G}^1, \mathcal{G}^2} \left[-\frac{1}{2} \log A + \frac{1}{2N^1} A + \frac{1}{2N^1} B^2 \right] + C \quad \text{where } A = \sum_{j=1}^{N^1} (1 - \lambda_j)^2 \text{ and } B = \frac{\sum_{j=1}^{N^1} \lambda_j (H_j^1 - \mu_{H^1})^2}{\sigma_{H^1}}$$

$$\min -I(Y; \mathcal{G}_{\text{CIB}}^1 | \mathcal{G}^2) + \beta I(\mathcal{G}^1; \mathcal{G}_{\text{CIB}}^1 | \mathcal{G}^2)$$

METHODOLOGY CONDITIONAL GRAPH INFORMATION BOTTLENECK



$$I(\mathcal{G}^1; \mathcal{G}_{\text{CIB}}^1 | \mathcal{G}^2) = I(\mathcal{G}_{\text{CIB}}^1; \mathcal{G}^1, \mathcal{G}^2) - \underbrace{I(\mathcal{G}_{\text{CIB}}^1; \mathcal{G}^2)}_{\text{Chain rule of mutual information}}$$

Solute Prediction

Encourage $\mathcal{G}_{\text{CIB}}^1$, which is compressed conditioned on \mathcal{G}^2 , to contain as much information about \mathcal{G}^2 as possible

Intuition) Make use of \mathcal{G}^2 when detecting $\mathcal{G}_{\text{CIB}}^1$

1) Variational IB-based approach

Derive upper bound similar to the prediction loss

$$-I(\mathcal{G}_{\text{CIB}}^1; \mathcal{G}^2) \leq \mathbb{E}_{\mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2} [-\log p_{\xi}(\mathcal{G}^2 | \mathcal{G}_{\text{CIB}}^1)] := \mathcal{L}_{\text{MI}^2}(\mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2)$$

2) Contrastive Learning-based approach

- Minimizing the contrastive loss is proven to be equivalent to maximizing the mutual information

- Hence, minimize $-I(\mathcal{G}_{\text{CIB}}^1; \mathcal{G}^2)$ by minimizing the contrastive loss $\rightarrow \text{CGIB}_{\text{cont}}$

$$\mathcal{L}_{\text{MI}^2} = -\frac{1}{K} \sum_{i=1}^K \log \frac{\exp(\text{sim}(\mathbf{z}_{\mathcal{G}_{\text{CIB},i}^1}, \mathbf{z}_{\mathcal{G}_i^2})/\tau)}{\sum_{j=1, j \neq i}^K \exp(\text{sim}(\mathbf{z}_{\mathcal{G}_{\text{CIB},i}^1}, \mathbf{z}_{\mathcal{G}_j^2})/\tau)}$$

EXPERIMENTS DATASET

Dataset		\mathcal{G}^1	\mathcal{G}^2	# \mathcal{G}^1	# \mathcal{G}^2	# Pairs	Task
Chromophore ¹	Absorption	Chrom.	Solvent	6416	725	17276	reg.
	Emission	Chrom.	Solvent	6412	1021	18141	reg.
	Lifetime	Chrom.	Solvent	2755	247	6960	reg.
MNSol ²		Solute	Solvent	372	86	2275	reg.
FreeSolv ³		Solute	Solvent	560	1	560	reg.
CompSol ⁴		Solute	Solvent	442	259	3548	reg.
Abraham ⁵		Solute	Solvent	1038	122	6091	reg.
CombiSolv ⁶		Solute	Solvent	1495	326	10145	reg.
ZhangDDI ⁷		Drug	Drug	544	544	40255	cls.
ChChMiner ⁸		Drug	Drug	949	949	21082	cls.

Dataset statistics

Chromophore dataset

Absorption max, Emission max, Lifetime

Solvation Free Energy dataset

- MNSol
- FreeSolv
- CompSol
- Abraham
- CombiSolv

Drug-Drug Interaction dataset

- ZhangDDI
- ChChMiner

EXPERIMENTS MAIN TABLE

	Chromophore			MNSol	FreeSolv	CompSol	Abraham	CombiSolv
	Absorption	Emission	Lifetime					
GCN	25.75 (1.48)	31.87 (1.70)	0.866 (0.015)	0.675 (0.021)	1.192 (0.042)	0.389 (0.009)	0.738 (0.041)	0.672 (0.022)
GAT	26.19 (1.44)	30.90 (1.01)	0.859 (0.016)	0.731 (0.007)	1.280 (0.049)	0.387 (0.010)	0.798 (0.038)	0.662 (0.021)
MPNN	24.43 (1.55)	30.17 (0.99)	0.802 (0.024)	0.682 (0.017)	1.159 (0.032)	0.359 (0.011)	0.601 (0.035)	0.568 (0.005)
GIN	24.92 (1.67)	32.31 (0.26)	0.829 (0.027)	0.669 (0.017)	1.015 (0.041)	0.331 (0.016)	0.648 (0.024)	0.595 (0.014)
CIGIN	19.32 (0.35)	25.09 (0.32)	0.804 (0.010)	0.607 (0.024)	0.905 (0.014)	0.308 (0.018)	0.411 (0.008)	0.451 (0.009)
CGIB	17.87 (0.38)	24.44 (0.21)	0.796 (0.010)	0.568 (0.013)	0.831 (0.012)	0.277 (0.008)	0.396 (0.009)	0.428 (0.009)
CGIB _{cont}	18.11 (0.20)	23.90 (0.35)	0.771 (0.005)	0.538 (0.007)	0.852 (0.022)	0.276 (0.017)	0.390 (0.006)	0.422 (0.005)

Performance on Molecular Interaction

Observations

Outperforms baselines on both Molecular Interaction / Drug-Drug Interaction tasks

Evaluation on drugs unseen during training

	(a) Transductive				(b) Inductive			
	ZhangDDI		ChChMiner		ZhangDDI		ChChMiner	
	AUROC	Accuracy	AUROC	Accuracy	AUROC	Accuracy	AUROC	Accuracy
GCN	91.64 (0.31)	83.31 (0.61)	94.71 (0.33)	87.36 (0.24)	68.39 (1.85)	63.78 (1.55)	73.63 (0.44)	67.07 (0.66)
GAT	92.10 (0.28)	84.14 (0.38)	96.15 (0.53)	89.49 (0.88)	69.99 (2.95)	64.41 (1.39)	75.72 (1.66)	68.77 (1.48)
MPNN	92.34 (0.35)	84.56 (0.31)	96.25 (0.53)	90.02 (0.42)	71.54 (1.24)	65.12 (1.14)	75.45 (0.32)	68.24 (1.42)
GIN	93.16 (0.04)	85.59 (0.05)	97.52 (0.05)	91.89 (0.66)	72.74 (1.32)	66.16 (1.21)	74.63 (0.48)	67.80 (0.46)
SSI-DDI	92.74 (0.12)	84.61 (0.18)	98.44 (0.08)	93.50 (0.16)	73.29 (2.23)	66.53 (1.31)	78.24 (1.29)	70.69 (1.47)
MIRACLE	93.05 (0.07)	84.90 (0.36)	88.66 (0.37)	84.29 (0.14)	73.23 (3.32)	50.00 (0.00)	60.25 (0.56)	50.09 (0.11)
CIGIN	93.28 (0.13)	85.54 (0.30)	98.51 (0.10)	93.77 (0.25)	74.02 (0.10)	66.81 (0.09)	79.23 (0.51)	71.56 (0.38)
CGIB	94.27 (0.47)	86.88 (0.56)	98.80 (0.04)	94.69 (0.16)	74.59 (0.88)	67.65 (1.07)	81.14 (1.20)	72.47 (0.16)
CGIB _{cont}	93.78 (0.62)	86.36 (0.75)	98.84 (0.31)	94.52 (0.38)	75.08 (0.34)	67.31 (0.82)	81.51 (0.67)	74.29 (0.14)

Performance on Drug-Drug Interaction

Observations

Improvement gap is larger in inductive setting

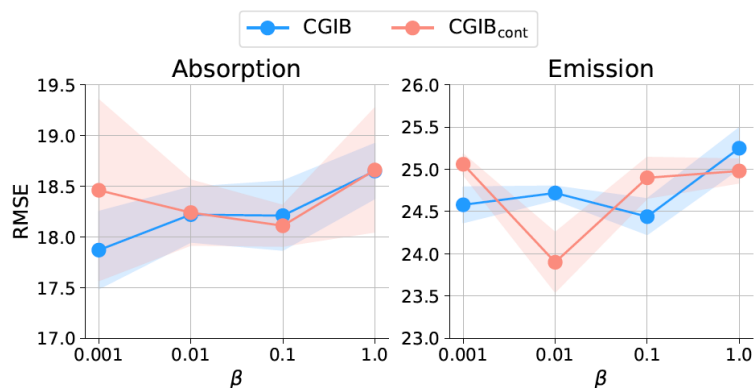
∴ By detecting function group that is basic in nature → helps generalization

EXPERIMENTS SENSITIVITY ANALYSIS

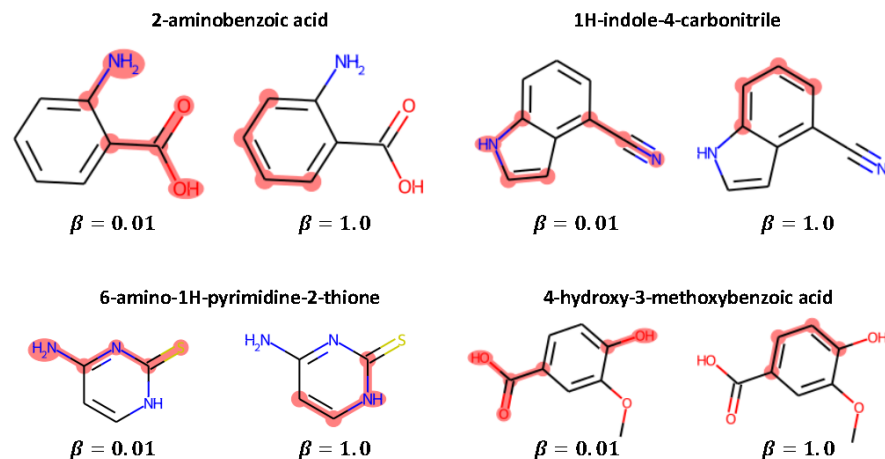
$$\min -I(Y; \mathcal{G}_{\text{CIB}}^1 | \mathcal{G}^2) + \beta I(\mathcal{G}^1; \mathcal{G}_{\text{CIB}}^1 | \mathcal{G}^2)$$

- β Controls Trade-off btw prediction and compression

As β increases, Compression > Prediction



Sensitivity Analysis on beta



Qualitative Analysis on beta

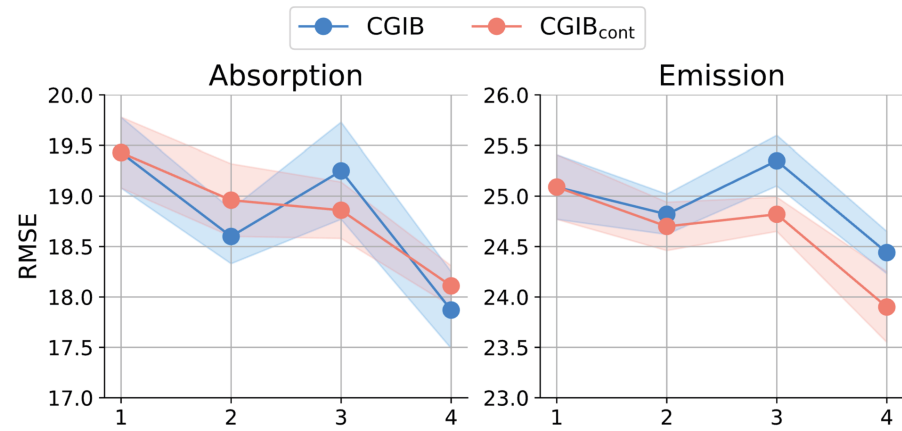
Observations

- $\beta = 1.0$: Poor performance in general (focus on compression)
- However, the model fails to detect functional group when β is too small
- poor generalization
- Hence, finding an appropriate β is crucial

Observations

- $\beta = 1.0$ → CGIB focuses on compression
e.g., CGIB focuses an aromatic ring, which is not relevant to chemical reactions
- $\beta = 0.01$ → CGIB focuses on prediction
e.g., CGIB focuses on external part, which generally more relevant to chemical reactions

EXPERIMENTS ABLATION STUDIES



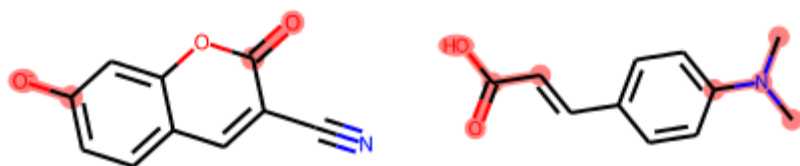
Observations

- Considering conditional MI is the key for success in relational learning
- A naïve consideration of \mathcal{G}^1 and \mathcal{G}^2 rather performs worse than considering \mathcal{G}^1 only

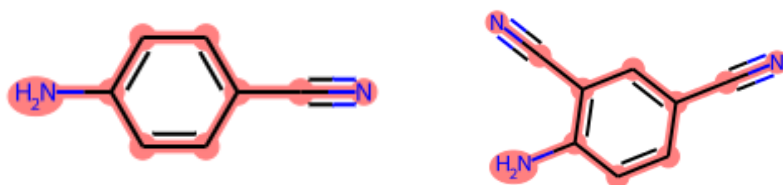
1. Without IB $\rightarrow \min - I(Y; \mathcal{G}^1, \mathcal{G}^2)$ (Same as CIGIN)
2. $I(\mathcal{G}_{\text{CIB}}^1; \mathcal{G}^1)$ $\rightarrow \min - I(Y; \mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2) + I(\mathcal{G}_{\text{CIB}}^1; \mathcal{G}^1)$ (Same as VGIB)
3. $I(\mathcal{G}_{\text{CIB}}^1; \mathcal{G}^1, \mathcal{G}^2)$ $\rightarrow \min - I(Y; \mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2) + I(\mathcal{G}_{\text{CIB}}^1; \mathcal{G}^1, \mathcal{G}^2)$
4. $I(\mathcal{G}^1; \mathcal{G}_{\text{CIB}}^1 | \mathcal{G}^2)$ $\rightarrow \min - I(Y; \mathcal{G}_{\text{CIB}}^1, \mathcal{G}^2) + I(\mathcal{G}^1; \mathcal{G}_{\text{CIB}}^1 | \mathcal{G}^2)$ (Same as CGIB)

EXPERIMENTS QUALITATIVE ANALYSIS

(a) Ordinary solvents



(b) Liquid oxygen solvent



Observations

(a) Chromophore (\mathcal{G}^1) interact with **ordinary solvents** (\mathcal{G}^2)

Focus on external parts → Aligns with domain knowledge

(b) Chromophore (\mathcal{G}^1) interact with **liquid oxygen solvents** (\mathcal{G}^2)

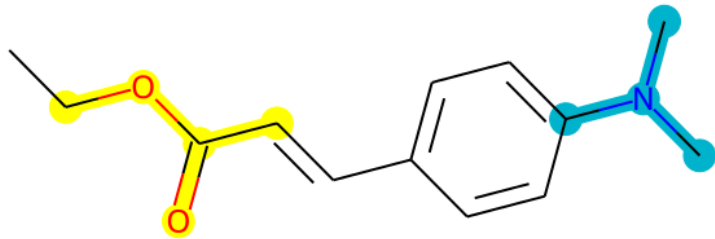
Focus on all parts → Aligns with domain knowledge

EXPERIMENTS QUALITATIVE ANALYSIS

(c)

Ethanol, THF, 1-Hexanol, 1-Butanol

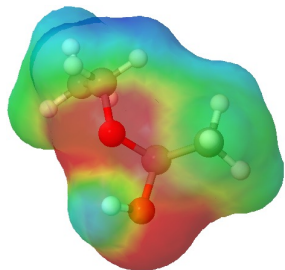
Benzene



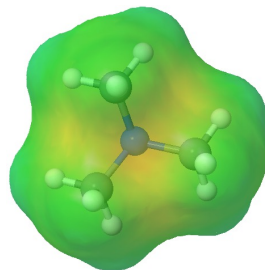
Oxygen-Carbon

Nitrogen-Carbon

(d) Polarity of the structure



Oxygen-Carbon
(High Polarity)



Nitrogen-Carbon
(Low Polarity)

Observations

(c) Chromophore (G^1) interacts with various solvents (G^2)
(e.g., Trans-ethyl p-(dimethylamino) cinnamate (EDAC))
Detected parts in chromophore depend on the polarity of solvent

- Case 1: High polarity solvent (Ethanol, THF, 1-hexanol, 1-butanol)

Structure with high polarity is detected (e.g., Oxygen-carbon)

→ Interact with high polarity solvent

- Case 2: Low polarity solvent (Benzene solvent)

Structure with low polarity is detected (e.g., Nitrogen-Carbon)

→ Interact with low polarity solvent

Detected structure of Chromophore (G^1) depends on the paired solvents (G^2)

CONCLUSION

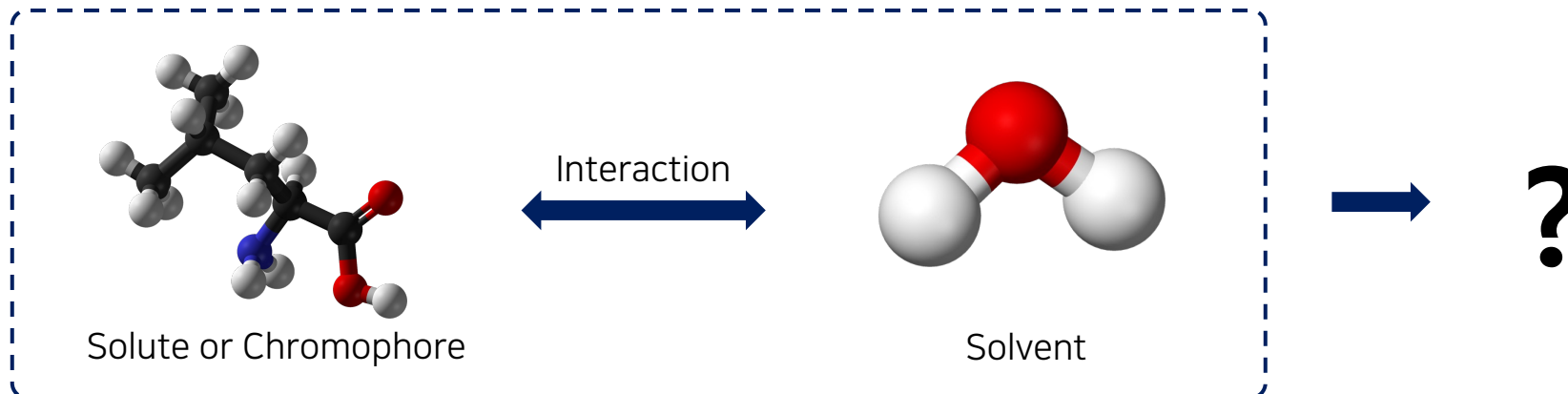
Proposed a method for tackling relation learning tasks, which are crucial for scientific discovery

- Based on Conditional Information Bottleneck

It is crucial to consider Graph 2 (Solvent) when detecting the important subgraph from Graph 1 (Chromophore)

- i.e., Make use of \mathcal{G}^2 when detecting $\mathcal{G}_{\text{CIB}}^1$ of \mathcal{G}^1

CGIB has interpretability, which makes it highly practical



THANK YOU!

[Full Paper] <https://openreview.net/forum?id=5hz3GV4IPq>

[Source Code] <https://github.com/Namkyeong/CGIB>

[Author Email] namkyeong96@kaist.ac.kr