# Unleashing the Power of Graph Data Augmentation on Covariate Distribution Shift

Yongduo Sui, Qitian Wu, Jiancan Wu, Qing Cui, Longfei Li, Jun Zhou, Xiang Wang, Xiangnan He
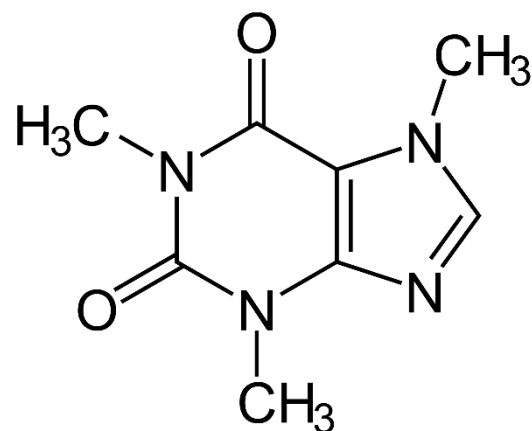
# 1.1 Background

☐ Graph data are everywhere

- Social network
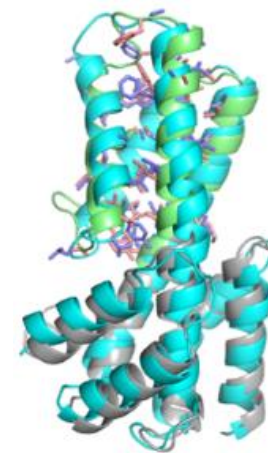- Chemical molecule
- Biological protein



Social network

☐ Graph learning tasks

- Node classification
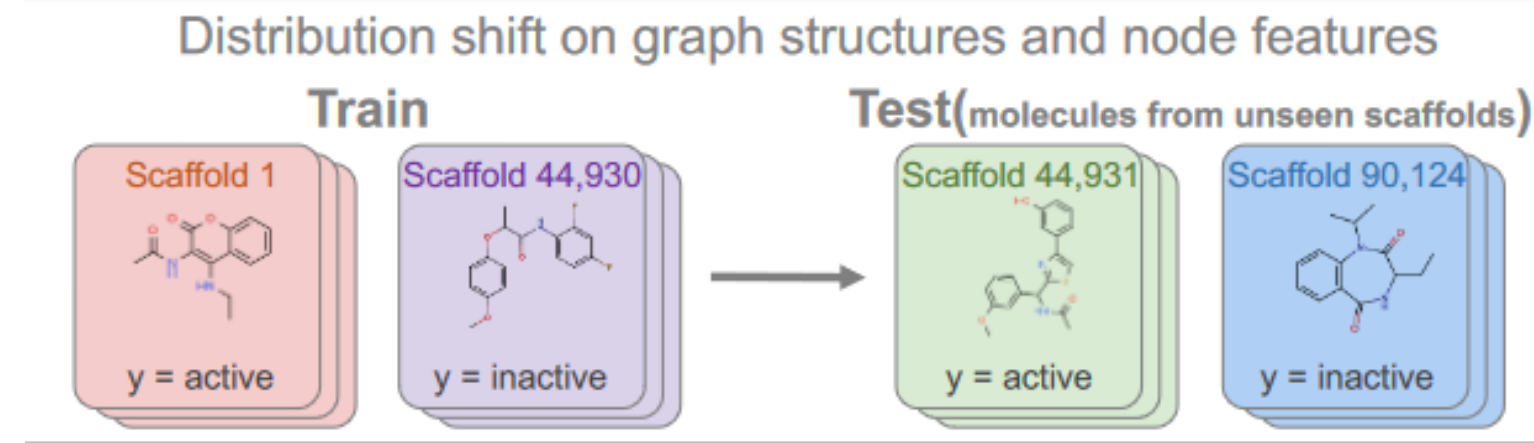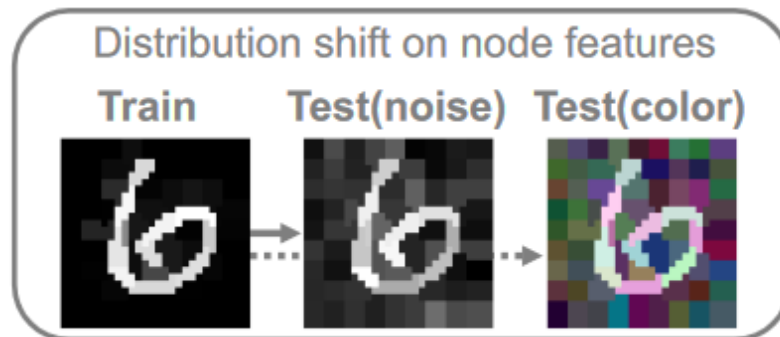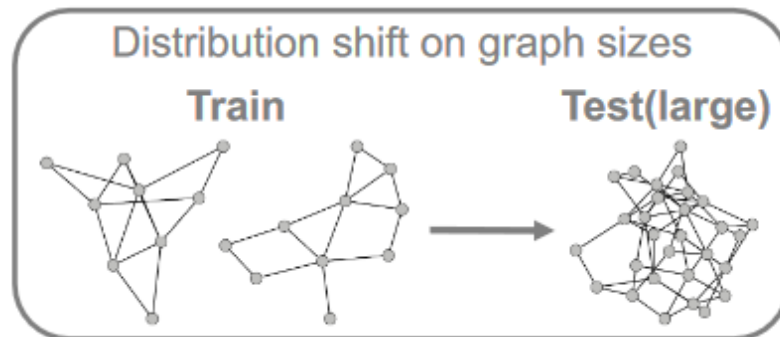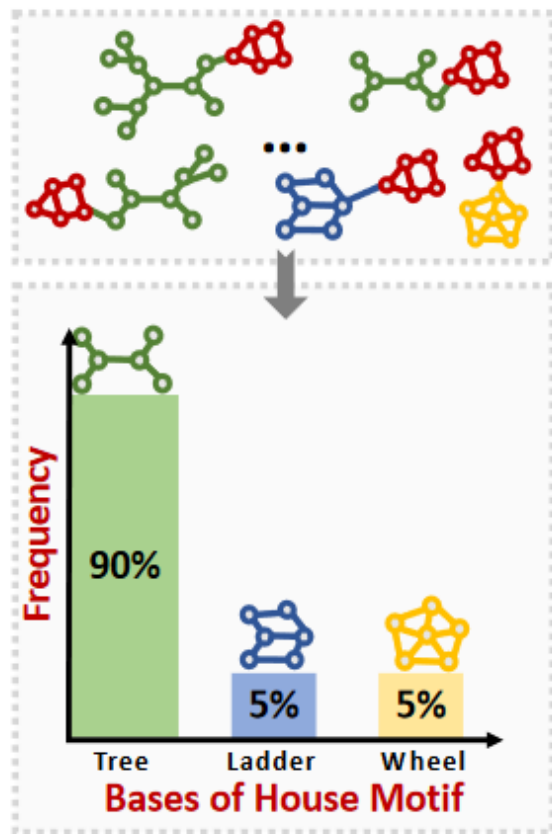- Link prediction
- Graph classification



Molecule



Protein structure

# 1.2 Graph Out-of-distribution Issue

☐ OOD Issue in Graph Classification

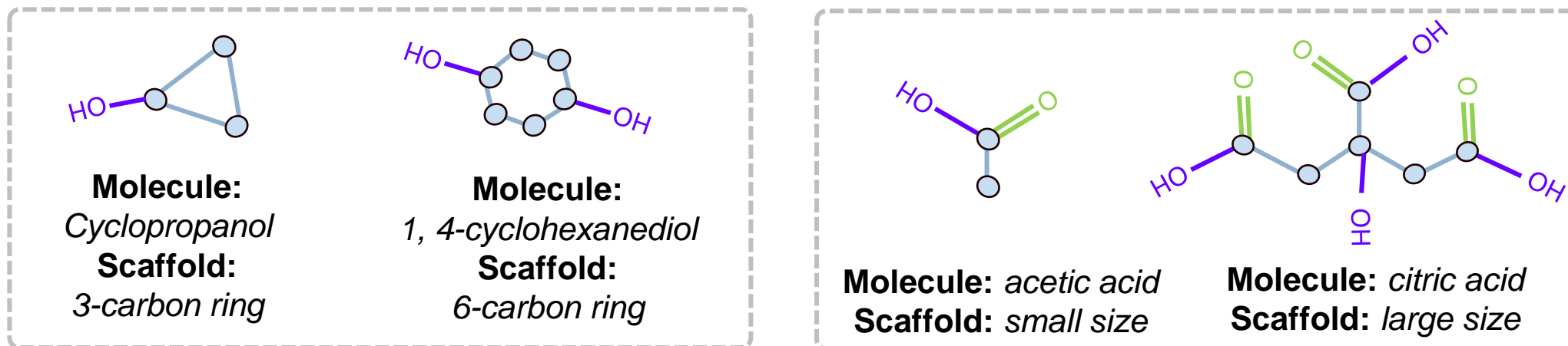[1] Discovering Invariant Rationales for Graph Neural Networks, ICLR 2022

[2] OOD-GNN: Out-of-Distribution Generalized Graph Neural Network, TKDE 2022

# 2.1 Assumption of Graph Generation

☐ Stable (*aka. Causal, Invariant, Rationale*) Feature & Environmental Feature

Stable feature:  functional group, e.g. –OH, -COOH
Environmental feature: scaffold, e.g. carbon ring, carbon chain



**Molecule:**
*Cyclopropanol*
**Scaffold:**
*3-carbon ring*

**Molecule:**
*1, 4-cyclohexanediol*
**Scaffold:**
*6-carbon ring*

**Molecule:** *acetic acid*
**Scaffold:** *small size*

**Molecule:** *citric acid*
**Scaffold:** *large size*

☐ Sufficiency & Invariance Assumption

**Assumption 3.1.** *Given* $G$, *there exists an optimal invariant subgraph generator* $\Phi^*(G)$ *satisfying:*
a. *Invariance property:* $\forall e, e' \in \text{supp}(\mathcal{E})$, $P^e(Y|\Phi^*(G)) = P^{e'}(Y|\Phi^*(G))$.
b. *Sufficiency property:* $Y = w^*(g^*(\Phi^*(G))) + \epsilon$, $\epsilon \perp G$, *where* $g^*(\cdot)$ *denotes a representation learning function,* $w^*$ *is the classifier,* $\perp$ *indicates statistical independence, and* $\epsilon$ *is random noise.*

[4] Learning Invariant Graph Representations for Out-of-Distribution Generalization, NeurIPS 2022
[5] Learning Substructure Invariance for Out-of-Distribution Molecular Representations, NeurIPS 2022

# 3.1 Two Types of Distribution Shifts

☐ Correlation Shift v.s. Covariate Shift

➤ The joint distribution of training and test data as $P_{\text{tr}}(G,Y)$ and $P_{\text{te}}(G,Y)$

Distribution shift means that $P_{\text{tr}}(G,Y) \neq P_{\text{te}}(G,Y)$

$$P_{\text{tr}}(G,Y) = P_{\text{tr}}(Y|G)\,P_{\text{tr}}(G)$$
$$P_{\text{te}}(G,Y) = P_{\text{te}}(Y|G)\,P_{\text{te}}(G)$$

➤ Correlation shift: $P_{\text{tr}}(G) = P_{\text{te}}(G)$ but $P_{\text{tr}}(Y|G) \neq P_{\text{te}}(Y|G)$

➤ Covariate shift: $P_{\text{tr}}(G) \neq P_{\text{te}}(G)$ but $P_{\text{tr}}(Y|G) = P_{\text{te}}(Y|G)$

Graph Data Generation

Our scope: these distribution shifts mainly caused by the environmental features.

[1] OoD-Bench: Quantifying and Understanding Two Dimensions of Out-of-Distribution Generalization, CVPR 2022
[6] GOOD: A Graph Out-of-Distribution Benchmark, NeurIPS 2022

# 3.1 Two Types of Distribution Shifts

☐ Correlation Shift v.s. Covariate Shift



Covariate (diversity) shift

*e.g. Domain Generalization*

Correlation (concept) shift

[1] OoD-Bench: Quantifying and Understanding Two Dimensions of Out-of-Distribution Generalization, CVPR 2022

# 3.1 Two Types of Distribution Shifts

☐ Correlation Shift v.s. Covariate Shift



*Unleashing the Power of Graph Data Augmentation on Covariate Distribution Shift, NeurIPS 2023*

# 3.2 Related Studies

## Methods

**General Generalization Algorithms**

*Zhang, et al, ICLR' 18,*
*Sagawa, et al, ICLR' 20,*
*Arjovsky, et al, Arxiv' 19.*

**Graph Data Augmentation**

*Rong, et al, ICLR' 2020,*
*You, et al, NeurIPS' 2020,*
*Han, et al, ICML' 2022.*

**Graph Invariant Learning**

*Wu, et al, ICLR' 2022(a),*
*Wu, et al, ICLR' 2022(b),*
*Li, et al, NeurIPS' 2022,*
*Sui, et al, KDD' 2022.*

## Limitations

Due to the irregularity of graph data, they are difficult to achieve significant performance improvements.

They are prone to destroy the stable features in the data, resulting in the insensitivity to stable features.

They are difficult to improve the environmental discrepancy.

# 3.3 How to Address Covariate Shift on Graphs?

Existing Issue: Insufficient discrepancy of environmental features

☐ Graph Covariate Shift

$$\text{GCS}(P_{\text{tr}}, P_{\text{te}}) = \frac{1}{2} \int_{\mathcal{S}} |P_{\text{tr}}(g) - P_{\text{te}}(g)| dg$$

$$\mathcal{S} = \{g \in \mathbb{G} | P_{\text{tr}}(g) \cdot P_{\text{te}}(g) = 0\}$$



☐ Our Idea

- Using data augmentation to increase the environmental discrepancy

*Unleashing the Power of Graph Data Augmentation on Covariate Distribution Shift, NeurIPS 2023*

# 3.3 How to Address Covariate Shift on Graphs?

☐ Two Principles for Graph Augmentation

- **Principle 1** (Environmental Feature Discrepancy): Environmental features should remain discrepant during augmentation

> **Principle 3.1 (Environmental Feature Discrepancy)** *Given a graph set $\{g\}$ with distribution function $P$, let $T(\cdot)$ denote an augmentation function that augments graphs $\{T(g)\}$ to distribution $\widetilde{P}$. Then $T(\cdot)$ should meet $\mathrm{GCS}(P, \widetilde{P}) \to 1$.*

- **Principle 2** (Stable Feature Consistency): Stable features should remain consistent during augmentation

> **Principle 3.2 (Stable Feature Consistency)** *Given a set of graphs $\{g\}$ with a corresponding stable feature set $\{g_{\mathrm{sta}} = (\mathbf{A}_{\mathrm{sta}}, \mathbf{X}_{\mathrm{sta}})\}$. Let $T(\cdot)$ denote an augmentation function that augments graphs $\{T(g)\}$ with a corresponding stable feature set $\{\widetilde{g}_{\mathrm{sta}} = (\widetilde{\mathbf{A}}_{\mathrm{sta}}, \widetilde{\mathbf{X}}_{\mathrm{sta}})\}$. Then $T(\cdot)$ should meet $\mathbb{E}[\|\mathbf{A}_{\mathrm{sta}} - \widetilde{\mathbf{A}}_{\mathrm{sta}}\|_F^2] \to 0$ and $\mathbb{E}[\|\mathbf{X}_{\mathrm{sta}} - \widetilde{\mathbf{X}}_{\mathrm{sta}}\|_F^2] \to 0$, where $\|\cdot\|_F$ is the Frobenius norm.*

# 3.4 Adversarial Invariant Augmentation

☐ Distributionally Robust Optimization

$$\min_{\theta} \left\{ \sup_{\widetilde{P}} \{ \mathbb{E}_{\widetilde{P}}[\ell(f(g), y)] : D(\widetilde{P}, P) \le \rho \} \right\},$$

➢ Wasserstein Distance

$$D(\widetilde{P}, P) := \inf_{\mu \in \Gamma(\widetilde{P}, P)} \mathbb{E}_{\mu}[c(\widetilde{g}, g)],$$

➢ Transportation Cost

$$c(\widetilde{g}, g) = \| h(\widetilde{g}) - h(g) \|_2^2.$$

➢ Lagrangian relaxation

$$\min_{\theta} \left\{ \sup_{\widetilde{P}} \{ \mathbb{E}_{\widetilde{P}}[\ell(f(g), y)] - \gamma D(\widetilde{P}, P) \} = \mathbb{E}_{P}[\phi(f(g), y)] \right\},$$

# 3.4 Adversarial Invariant Augmentation

☐ Robust surrogate loss： $\phi(f(g), y)$

$$\phi(f(g), y) := \sup_{\widetilde{g} \in \mathbb{G}} \{\ell(f(\widetilde{g}), y) - \gamma c(\widetilde{g}, g)\}$$

☐ Adversarial Augmentation

$$\nabla_\theta \phi(f(g), y) = \nabla_\theta \ell(f(\widetilde{g}^*), y),$$
$$\text{where } \widetilde{g}^* = \arg\max_{\widetilde{g} \in \mathbb{G}} \{\ell(f(\widetilde{g}), y) - \gamma c(\widetilde{g}, g)\}.$$

# 3.4 Adversarial Invariant Augmentation

☐ Adversarial Augmenter & Stable Feature Generator

$$T_{\theta_1}(g) = (\mathbf{A} \odot \mathbf{M}^a_{\text{adv}}, \mathbf{X} \odot \mathbf{M}^x_{\text{adv}})$$



*Unleashing the Power of Graph Data Augmentation on Covariate Distribution Shift, NeurIPS 2023*

# 3.4 Adversarial Invariant Augmentation

☐ Maximization

$$\max_{\theta_1} \left\{ \mathcal{L}_{\mathrm{adv}} = \mathbb{E}_{P_{\mathrm{tr}}} \left[ \ell(f(T_{\theta_1}(g)), y) - \gamma c(T_{\theta_1}(g), g) \right] \right\}$$

$T_{\theta_2}(g)$



$T_{\theta_1}(g)$

*Unleashing the Power of Graph Data Augmentation on Covariate Distribution Shift, NeurIPS 2023*

# 3.4 Adversarial Invariant Augmentation

☐ Minimization

$$\min_{\theta,\theta_2} \left\{ \mathcal{L}_{\text{cau}} = \mathbb{E}_{P_{\text{tr}}} \left[ \ell(f(T_{\theta_2}(g)), y) + \ell(f(\widetilde{g}), y) \right] \right\} \begin{cases} \widetilde{g} = (\mathbf{A} \odot \widetilde{\mathbf{M}}^a, \mathbf{X} \odot \widetilde{\mathbf{M}}^x) \\ \widetilde{\mathbf{M}}^a = (\mathbf{1}^a - \mathbf{M}^a_{\text{cau}}) \odot \mathbf{M}^a_{\text{adv}} + \mathbf{M}^a_{\text{cau}} \\ \widetilde{\mathbf{M}}^x = (\mathbf{1}^x - \mathbf{M}^x_{\text{cau}}) \odot \mathbf{M}^x_{\text{adv}} + \mathbf{M}^x_{\text{cau}} \end{cases}$$



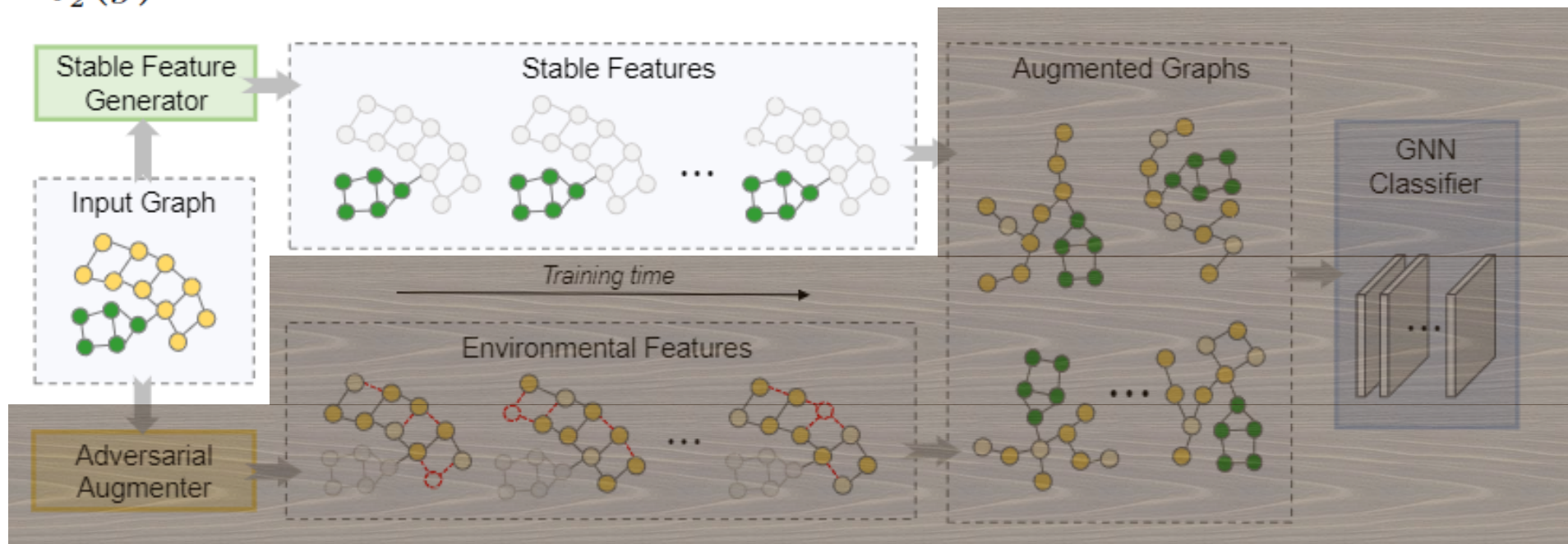*Unleashing the Power of Graph Data Augmentation on Covariate Distribution Shift, NeurIPS 2023*

# 3.4 Adversarial Invariant Augmentation

☐ Mask Combination

$$\min_{\theta,\theta_2}\left\{\mathcal{L}_{\text{cau}} = \mathbb{E}_{P_{\text{tr}}}\left[\ell\big(f(T_{\theta_2}(g)),y\big) + \ell\big(f(\widetilde{g}),y\big)\right]\right\}$$

$$\begin{cases} \widetilde{g} = (\mathbf{A} \odot \widetilde{\mathbf{M}}^a, \mathbf{X} \odot \widetilde{\mathbf{M}}^x) \\ \widetilde{\mathbf{M}}^a = (\mathbf{1}^a - \mathbf{M}^a_{\text{cau}}) \odot \mathbf{M}^a_{\text{adv}} + \mathbf{M}^a_{\text{cau}} \\ \widetilde{\mathbf{M}}^x = (\mathbf{1}^x - \mathbf{M}^x_{\text{cau}}) \odot \mathbf{M}^x_{\text{adv}} + \mathbf{M}^x_{\text{cau}} \end{cases}$$
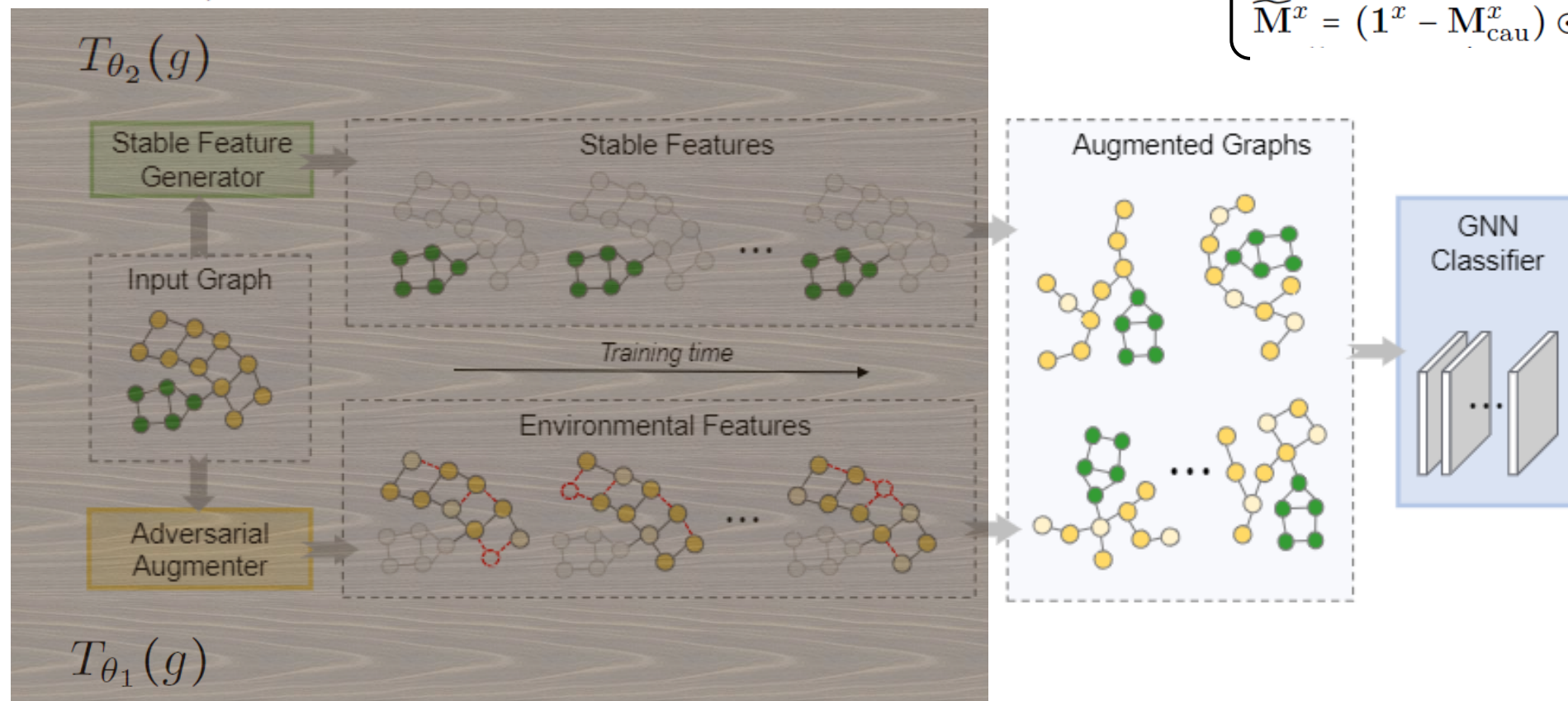


*Unleashing the Power of Graph Data Augmentation on Covariate Distribution Shift, NeurIPS 2023*

# 4.1 Experiments: Main Results

Table 1: Performance on synthetic and real-world datasets. Numbers in **bold** indicate the best performance, while the underlined numbers indicate the second best performance.

| Type | Method | Motif | | CMNIST | Molbbbp | | Molhiv | |
|------|--------|-------|------|--------|---------|------|--------|------|
| | | base | size | color | scaffold | size | scaffold | size |
| General Generalization | ERM | $68.66_{\pm4.25}$ | $51.74_{\pm2.88}$ | $28.60_{\pm1.87}$ | $68.10_{\pm1.68}$ | $78.29_{\pm3.76}$ | $69.58_{\pm2.51}$ | $59.94_{\pm2.37}$ |
| | IRM | $70.65_{\pm4.17}$ | $51.41_{\pm3.78}$ | $27.83_{\pm2.13}$ | $67.22_{\pm1.15}$ | $77.56_{\pm2.48}$ | $67.97_{\pm1.84}$ | $59.00_{\pm2.92}$ |
| | GroupDRO | $68.24_{\pm8.92}$ | $51.95_{\pm5.86}$ | $29.07_{\pm3.14}$ | $66.47_{\pm2.39}$ | $79.27_{\pm2.43}$ | $70.64_{\pm2.57}$ | $58.98_{\pm2.16}$ |
| | VREx | $\underline{71.47}_{\pm6.69}$ | $52.67_{\pm5.54}$ | $28.48_{\pm2.87}$ | $68.74_{\pm1.03}$ | $78.76_{\pm2.37}$ | $70.77_{\pm2.84}$ | $58.53_{\pm2.88}$ |
| Graph Generalization | DIR | $62.07_{\pm8.75}$ | $52.27_{\pm4.56}$ | $\underline{33.20}_{\pm6.17}$ | $66.86_{\pm2.25}$ | $76.40_{\pm4.43}$ | $68.07_{\pm2.29}$ | $58.08_{\pm2.31}$ |
| | CAL | $65.63_{\pm4.29}$ | $51.18_{\pm5.60}$ | $27.99_{\pm3.24}$ | $68.06_{\pm2.60}$ | $\underline{79.50}_{\pm4.81}$ | $67.37_{\pm3.61}$ | $57.95_{\pm2.24}$ |
| | GSAT | $62.80_{\pm11.41}$ | $53.20_{\pm8.35}$ | $28.17_{\pm1.26}$ | $66.78_{\pm1.45}$ | $75.63_{\pm3.83}$ | $68.66_{\pm1.35}$ | $58.06_{\pm1.98}$ |
| | OOD-GNN | $61.10_{\pm7.87}$ | $52.61_{\pm4.67}$ | $26.49_{\pm2.94}$ | $66.72_{\pm1.23}$ | $79.48_{\pm4.19}$ | $70.46_{\pm1.97}$ | $60.60_{\pm3.77}$ |
| | StableGNN | $57.07_{\pm14.10}$ | $46.93_{\pm8.85}$ | $28.38_{\pm3.49}$ | $66.74_{\pm1.30}$ | $77.47_{\pm4.69}$ | $68.44_{\pm1.33}$ | $56.71_{\pm2.79}$ |
| | CIGA | $66.43_{\pm11.31}$ | $49.14_{\pm8.34}$ | $32.22_{\pm2.67}$ | $64.92_{\pm2.09}$ | $65.98_{\pm3.31}$ | $69.40_{\pm2.39}$ | $59.55_{\pm2.56}$ |
| | DisC | $51.08_{\pm3.08}$ | $50.39_{\pm1.15}$ | $24.99_{\pm1.78}$ | $67.12_{\pm2.11}$ | $56.59_{\pm10.09}$ | $68.07_{\pm1.75}$ | $58.76_{\pm0.91}$ |
| Graph Augmentation | DropEdge | $45.08_{\pm4.46}$ | $45.63_{\pm4.61}$ | $22.65_{\pm2.90}$ | $66.49_{\pm1.55}$ | $78.32_{\pm3.44}$ | $\underline{70.78}_{\pm1.38}$ | $58.53_{\pm1.26}$ |
| | GREA | $56.74_{\pm9.23}$ | $\underline{54.13}_{\pm10.02}$ | $29.02_{\pm3.26}$ | $\underline{69.72}_{\pm1.66}$ | $77.34_{\pm3.52}$ | $67.79_{\pm2.56}$ | $\underline{60.71}_{\pm2.20}$ |
| | FLAG | $61.12_{\pm5.39}$ | $51.66_{\pm4.14}$ | $32.30_{\pm2.69}$ | $67.69_{\pm2.36}$ | $79.26_{\pm2.26}$ | $68.45_{\pm2.30}$ | $60.59_{\pm2.95}$ |
| | M-Mixup | $70.08_{\pm3.82}$ | $51.48_{\pm4.91}$ | $26.47_{\pm3.45}$ | $68.75_{\pm0.34}$ | $78.92_{\pm2.43}$ | $68.88_{\pm2.63}$ | $59.03_{\pm3.11}$ |
| | $\mathcal{G}$-Mixup | $59.66_{\pm7.03}$ | $52.81_{\pm6.73}$ | $31.85_{\pm5.82}$ | $67.44_{\pm1.62}$ | $78.55_{\pm4.16}$ | $70.01_{\pm2.52}$ | $59.34_{\pm2.43}$ |
| | AIA (ours) | $\mathbf{73.64}_{\pm5.15}$ | $\mathbf{55.85}_{\pm7.98}$ | $\mathbf{36.37}_{\pm4.44}$ | $\mathbf{70.79}_{\pm1.53}$ | $\mathbf{81.03}_{\pm5.15}$ | $\mathbf{71.15}_{\pm1.81}$ | $\mathbf{61.64}_{\pm3.37}$ |

*Unleashing the Power of Graph Data Augmentation on Covariate Distribution Shift, NeurIPS 2023*

# 5 Conclusion

➢ We aim to address the covariate shift issue in graph learning, which is important yet largely unexplored.

➢ We introduce a novel graph augmentation method, AIA, grounded in two principles: environmental feature discrepancy and stable feature consistency.

➢ We conduct extensive experiments and the results demonstrate the effectiveness of our method.