# Size-Invariant Graph Representations for Graph Classification Extrapolations

Beatrice Bevilacqua*, Yangze Zhou*, Bruno Ribeiro
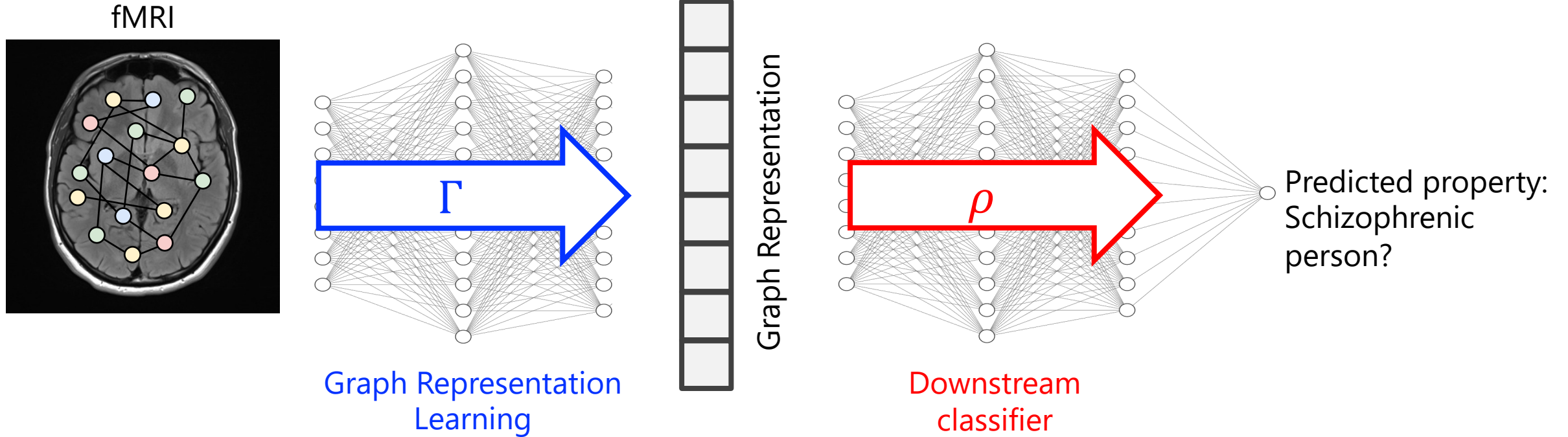
Purdue University

ICML 2021

*Equal contribution

This work focuses on out-of-distribution (OOD) extrapolations in Graph Representation Learning
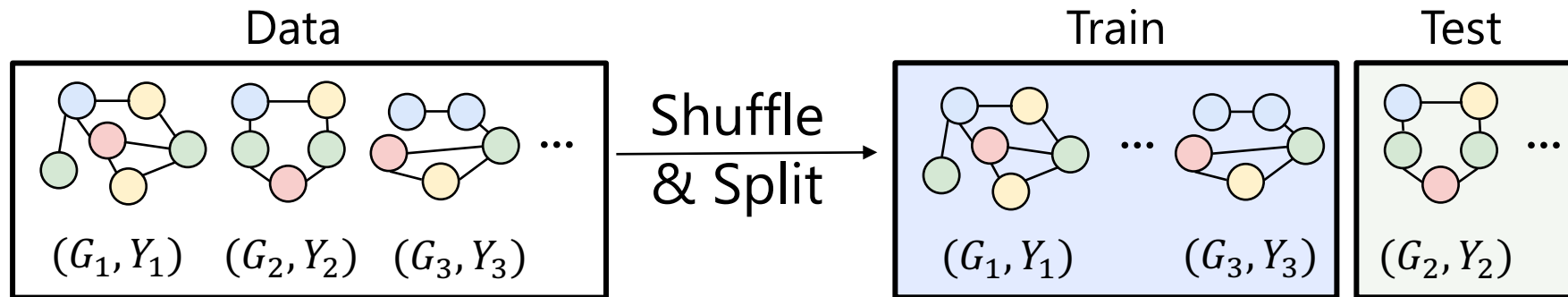
Toolbox:
- Causality
- Graph limits
- Graph Neural Networks

# Graph Classification Tasks

fMRI

$\Gamma$

**Graph Representation Learning**

Graph Representation

$\rho$

**Downstream classifier**

Predicted property: Schizophrenic person?

# Current Graph Classification Approach

Graph Representation Learning generally assumes:

Train distribution = Test distribution



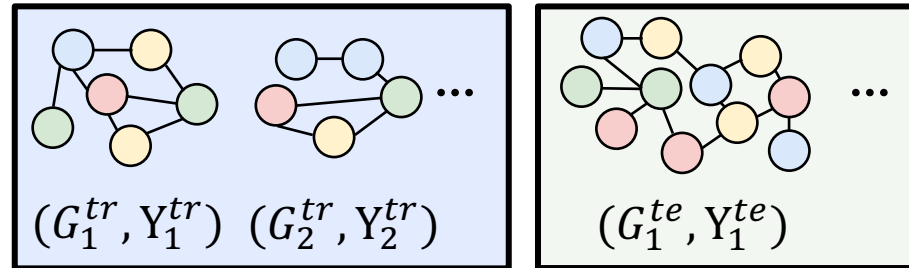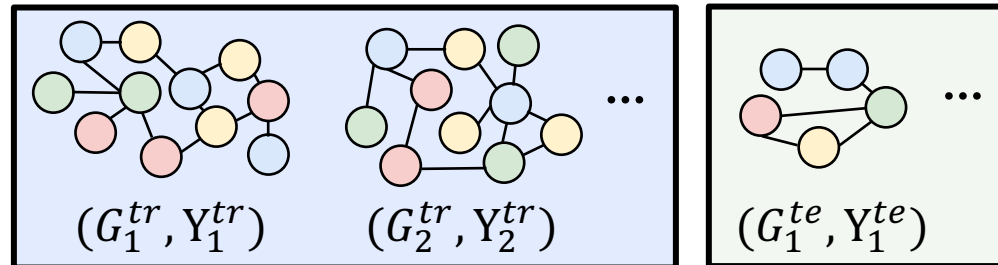**What if test data were out of distribution (OOD)?**

What if train has **small** graphs but test has **large** graphs?

Train (small graphs)    Test (large graphs)

$(G_1^{tr}, Y_1^{tr})$    $(G_2^{tr}, Y_2^{tr})$    $(G_1^{te}, Y_1^{te})$

What if train has **large** graphs but test has **small** graphs?

Train (large graphs)    Test (small graphs)

$(G_1^{tr}, Y_1^{tr})$    $(G_2^{tr}, Y_2^{tr})$    $(G_1^{te}, Y_1^{te})$
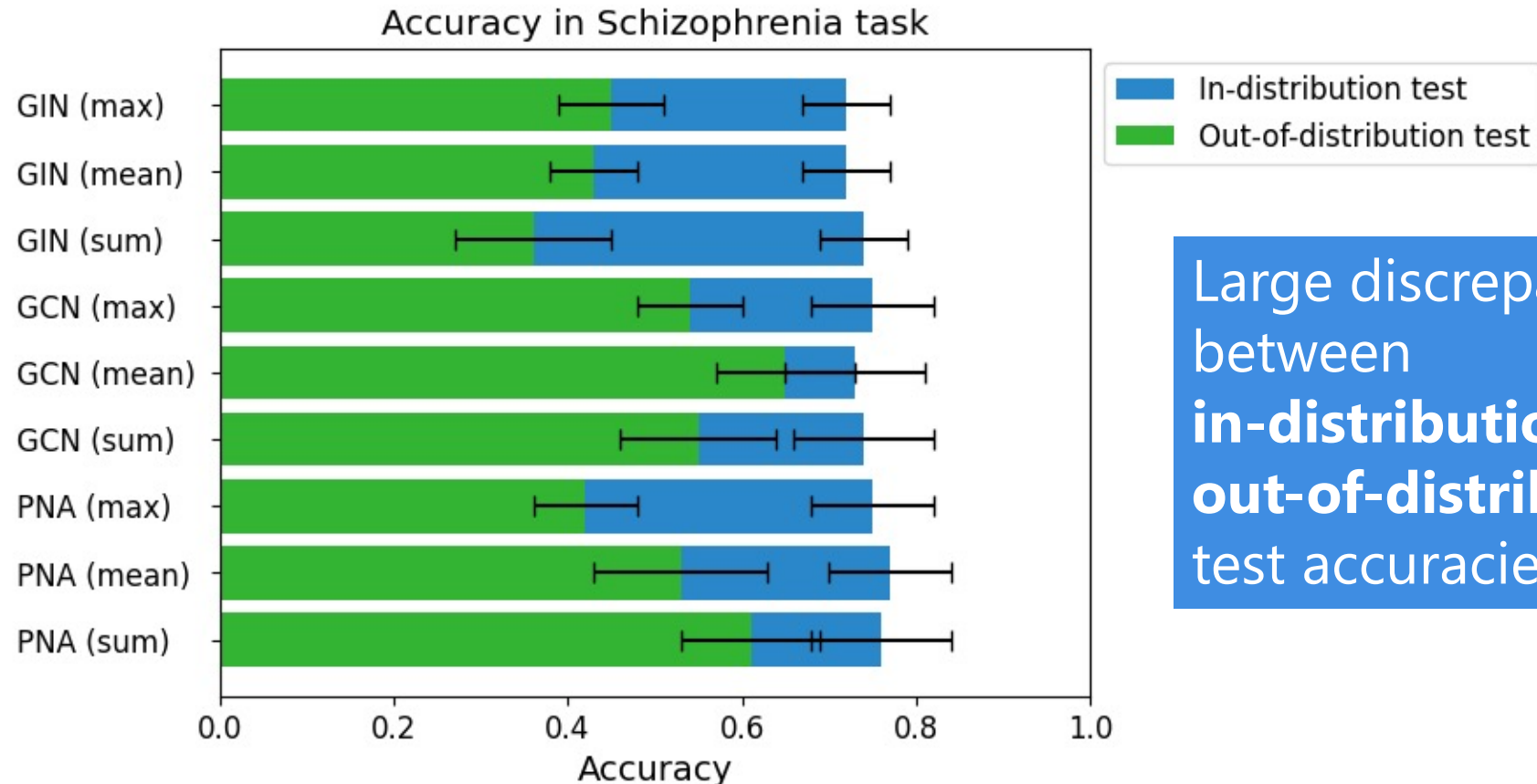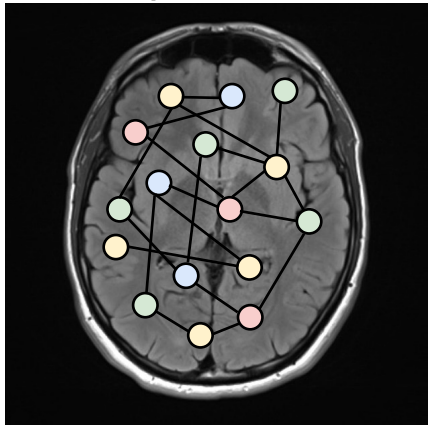
# Size Extrapolation with GNNs?

Do Graph Neural Networks (GNNs) extrapolate?

⇒ GNNs can be applied to graphs of any size

⇒ But may not extrapolate between **small (train)** and **large (test)** graphs:
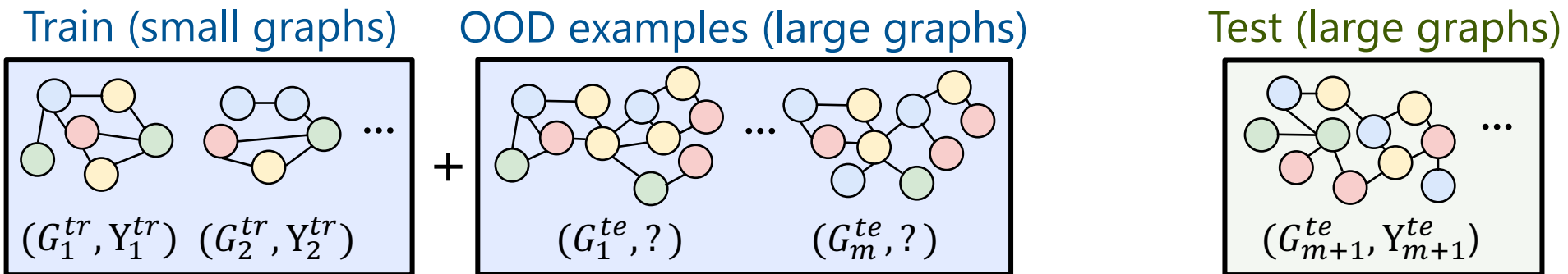
Schizophrenia task



Large discrepancy between **in-distribution** and **out-of-distribution** test accuracies

How do we extrapolate beyond the training distribution?

If OOD examples available, **data-driven methods** work:
▸ Domain Adaptation
▸ Covariate Shift Adaptation
▸ Few-shot Learning
▸ Data Augmentation
▸ Invariant Risk-Minimization (IRM)*

Train (small graphs)    OOD examples (large graphs)    Test (large graphs)



$(G_1^{tr}, Y_1^{tr})$ $(G_2^{tr}, Y_2^{tr})$    +    $(G_1^{te}, ?)$    $(G_m^{te}, ?)$    $(G_{m+1}^{te}, Y_{m+1}^{te})$

# How to Extrapolate in Graph Classification Tasks?

**Data-driven** methods:

<span style="color:green">Pros</span>

- ▸ Can use existing GNN methods

- ▸ Don't assume a mechanism for distribution shift

<span style="color:red">Cons</span>

- ▸ Must have OOD examples during training

<span style="color:red">What if no access to OOD data?</span>

- ▸ Must define a causal mechanism

Next: Observational vs Causal (Counterfactual) modeling

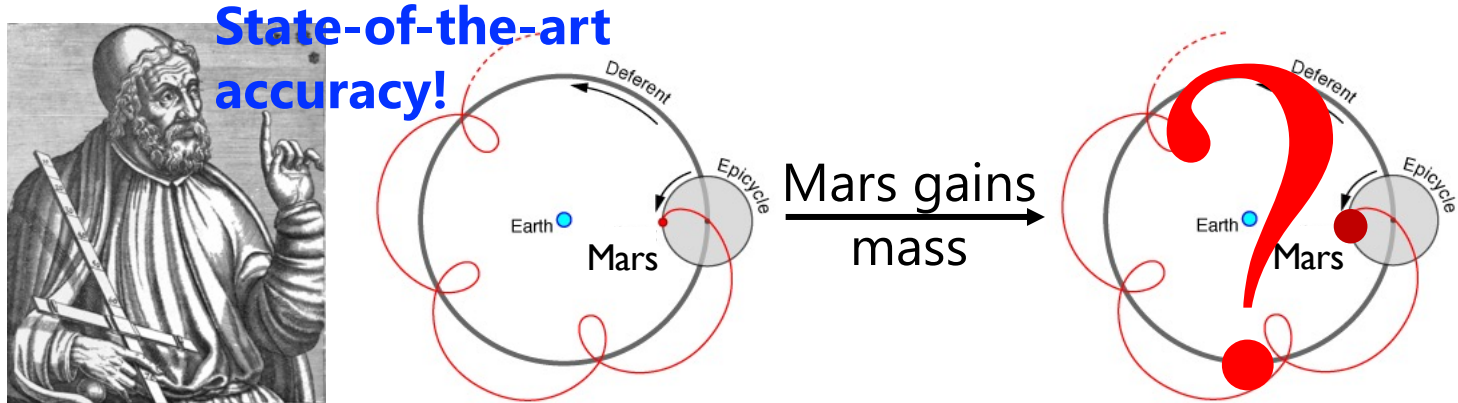# Why are Causal Mechanisms Needed for Extrapolations without OOD Data?

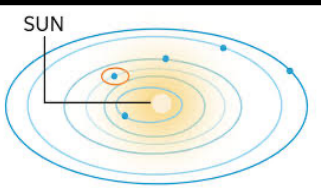**Historical analogy to Graph Representation Learning methods:**

▸ Ptolemaic geocentric model of planetary motion

  ◦ Very **accurate** to predict positions **observationally**

  ◦ **Cannot** predict positions in new **scenarios**

New scenario:
What would happen if
Mars became 10x more massive?

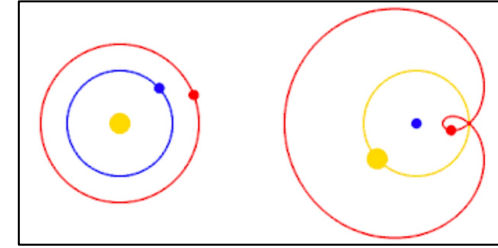**State-of-the-art accuracy!**

Mars gains mass

img credits: wikipedia

Interpretable model
**cannot** predict new scenarios

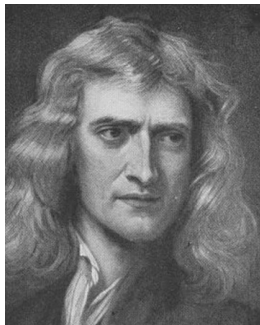**Lesson:**
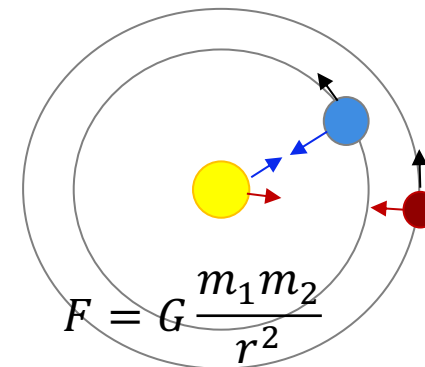Occam's razor & interpretability
≠ out-of-distribution extrapolation

▸ Observational predictions can be purely data-driven



▸ Predicting new scenarios (larger and smaller mass) without OOD examples requires a **mechanism**



New scenario:
What would happen if
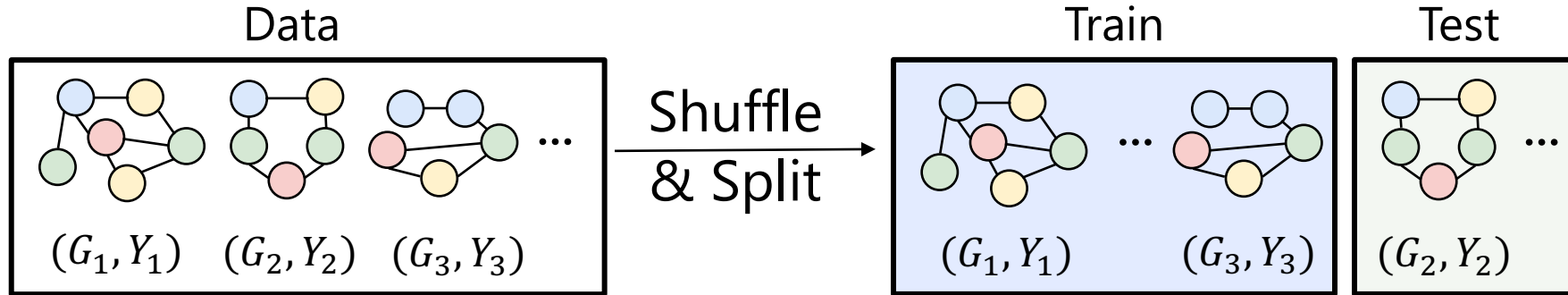Mars became 10x more massive?

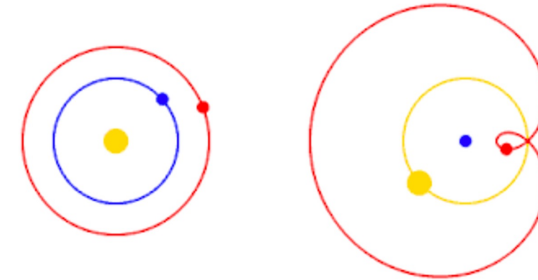$$F = G\frac{m_1 m_2}{r^2}$$

# Size Extrapolations on Graphs

# Differences between Observational and Counterfactual Tasks

**Observational Task:**
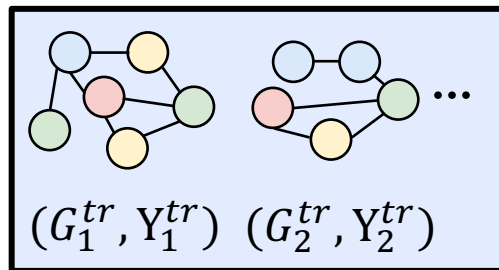Predicting unseen examples of training distribution



Data → Shuffle & Split → Train | Test

$(G_1, Y_1)$  $(G_2, Y_2)$  $(G_3, Y_3)$ → $(G_1, Y_1)$ ... $(G_3, Y_3)$ | $(G_2, Y_2)$

*Planetary Motion Equivalent*

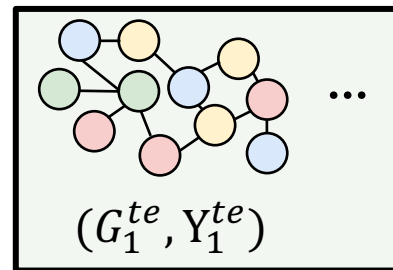**Counterfactual Task** (since we have no access to test data)**:**
What would be the label of a graph if it were larger?

Train (small graphs)

$(G_1^{tr}, Y_1^{tr})$  $(G_2^{tr}, Y_2^{tr})$

Test (large graphs)

$(G_1^{te}, Y_1^{te})$

or vice-versa

$$F = G \frac{m_1 m_2}{r^2}$$

Reminder of talk:

What would be the labels if the graphs were larger?
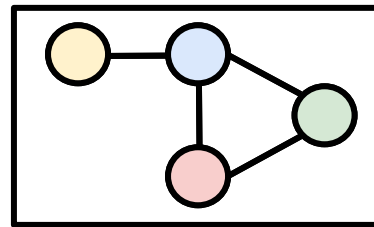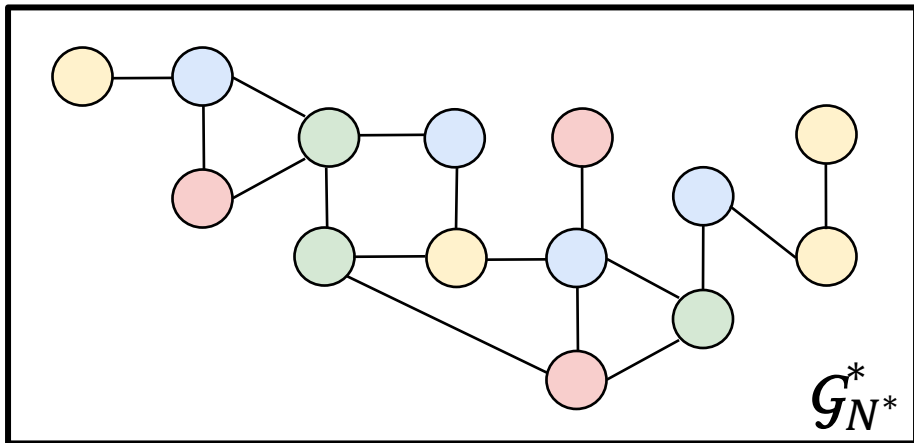
# To Infinity and Beyond...

**Q:** What would be the label if the graph were infinitely large?
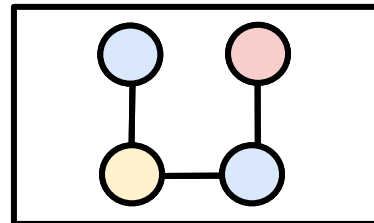
$$N \rightarrow \infty$$

# Lovász Graph Limits

▸ What graph property is invariant as graphs become larger?

　○ Lovász & Szegedy (2006) shows:

　　• Density of induced subgraphs of a dense random graph converges as $N \rightarrow \infty$

Induced k-sized subgraph density

$$t_{\text{ind}}(F_k, \mathcal{G}_{N^*}^*) = \frac{\text{ind}(F_k, \mathcal{G}_{N^*}^*)}{N^*! \, / \, (N^* - k)!}$$

$\mathcal{G}_{N^*}^*$ can be train $\mathcal{G}_{N^{tr}}^{tr}$ or test $\mathcal{G}_{N^{te}}^{te}$ graph



$$t_{\text{ind}}(\,\cdot\!\!\!\!\triangleright, \mathcal{G}_{N^*}^*) = \frac{2}{14! \, / \, (14 - 4)!}$$

Count = 2

$$t_{\text{ind}}(\,\square\,, \mathcal{G}_{N^*}^*) = \frac{1}{14! \, / \, (14 - 4)!}$$

Count = 1

What if we constructed a graph representation from subgraph densities?

# Graph Representation based on Densities

**New graph representation**

**Induced subgraph density**

$$\Gamma_{\text{GNN}}(\mathcal{G}_{N^*}^*) = \sum_{F_{k'} \in \mathcal{F}_{\leq k}} \text{t}_{\text{ind}}(F_{k'}, \mathcal{G}_{N^*}^*) \text{READOUT}_\Gamma(\text{GNN}(F_{k'}))$$

**GNN-representation of $F_{k'}$**

$$\text{GNN}(F_{k'}):$$

$$\text{READOUT}_\Gamma\left(\left\{ \right\}\right)$$

# OOD Error in Schizophrenia Task

▶ Can subgraph density representation $\Gamma_{GNN}$ extrapolate OOD?

### Accuracy in Schizophrenia task



Legend:
- In-distribution test (blue)
- Out-of-distribution test (green)

OOD error same as in-distribution error

# Theory

Understand why $\Gamma_{\mathrm{GNN}}$ can OOD extrapolate

# A Most Expressive Representation

Preliminary: **1-hot encoded graph representations are most-expressive**

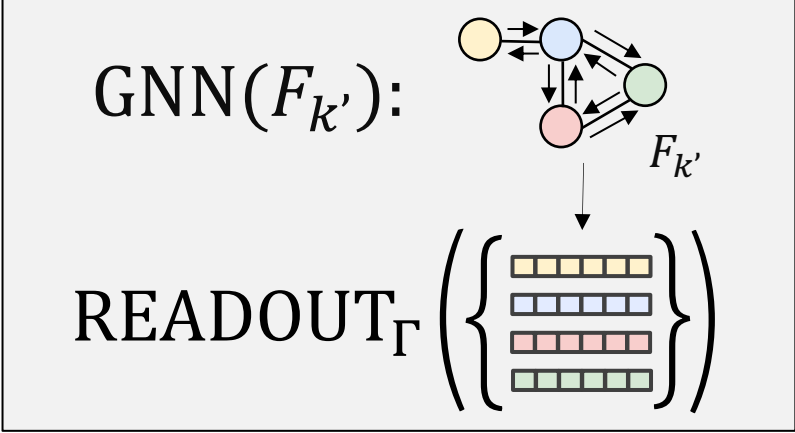**First:** replace the GNN-representation of $F_{k'}$ with a 1-hot encoded representation
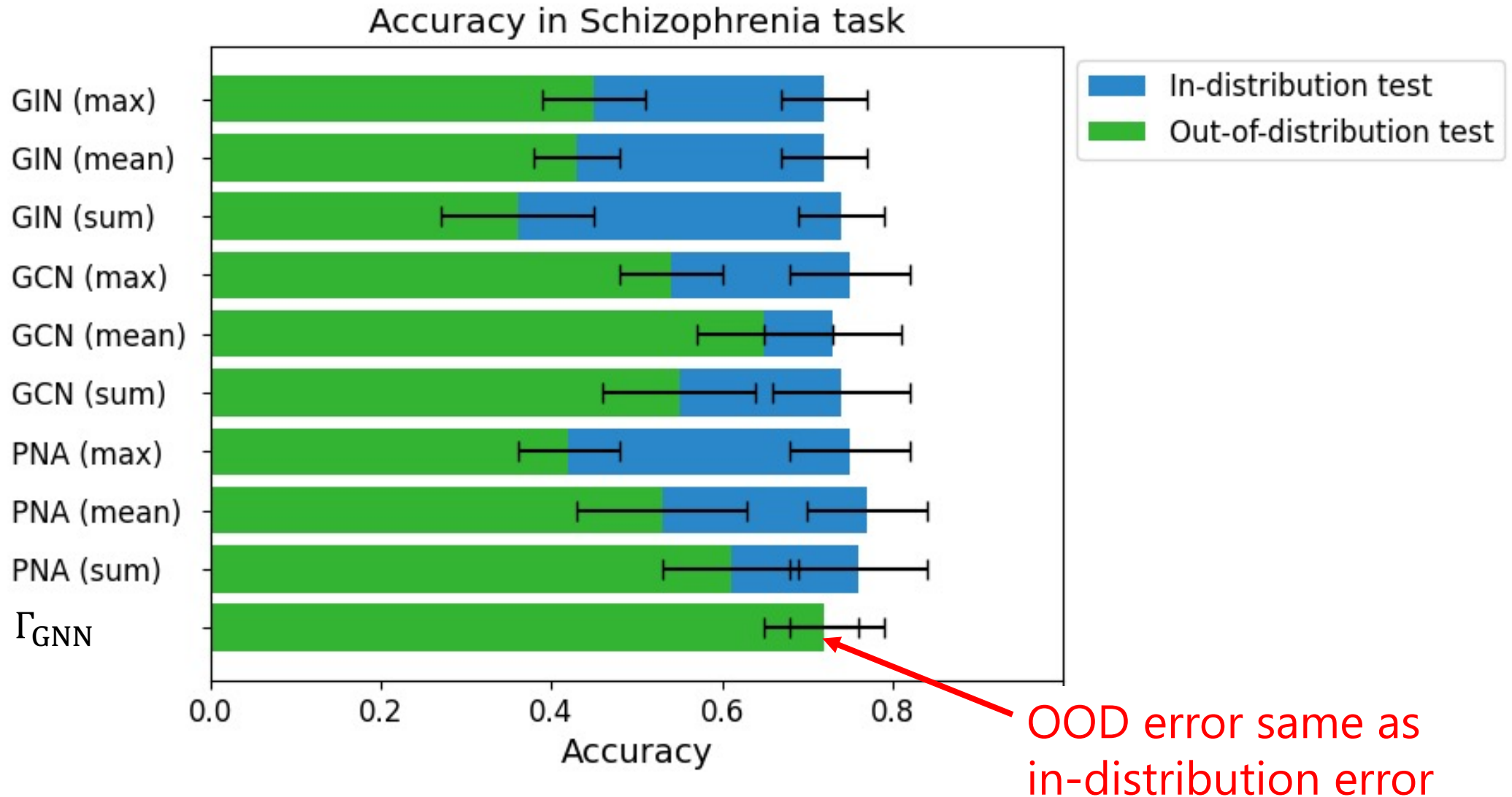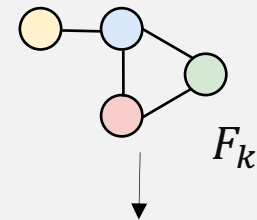
Graph representation

Induced subgraph density

$$\Gamma_{1-\text{hot}}(\mathcal{G}_{N^*}^*) = \sum_{F_{k'} \in \mathcal{F}_{\leq k}} \text{t}_{\text{ind}}(F_{k'}, \mathcal{G}_{N^*}^*) 1_{\text{one}-\text{hot}}\{F_{k'}, \mathcal{F}_{\leq k}\}$$

One-hot encoded
identifier of $F_{k'}$

$F_{k'}$

$(0,0, .., 1, 0, ..)$

# Size-Invariant Representation

**Theorem 1 (informal): Approximately size-invariant representations**
Under certain conditions (explained later), the change in graph representation between train and counterfactual test graph is upper bounded by $k$ and graph sizes (in train and test):

$$P\left(|| \Gamma_{1-\text{hot}}(\mathcal{G}_{N^{tr}}^{tr}) - \Gamma_{1-\text{hot}}(\mathcal{G}_{N^{te}}^{te})||_\infty > \epsilon\right) \leq 2|\mathcal{F}_{\leq k}|(\exp(-\frac{\epsilon^2 N^{tr}}{8k^2}) + \exp(-\frac{\epsilon^2 N^{te}}{8k^2}))$$

Training graph

Counterfactual test graph

▸ Proof relies on Lovász graph limits (formal definition in paper)

Note that $\Gamma_{\text{GNN}}$ is less expressive (more invariant) than $\Gamma_{1-\text{hot}}$

# Effects of Invariant Representations

▸ Why are we interested in invariant representations?

**Proposition 1 (informal): Effect of invariant representations**

Consider:

- $\Gamma$ : A permutation invariant graph representation
- $\rho$ : A downstream classifier

In-distribution generalization error : $\forall y \in Y$, for some $\epsilon, \delta \geq 0$
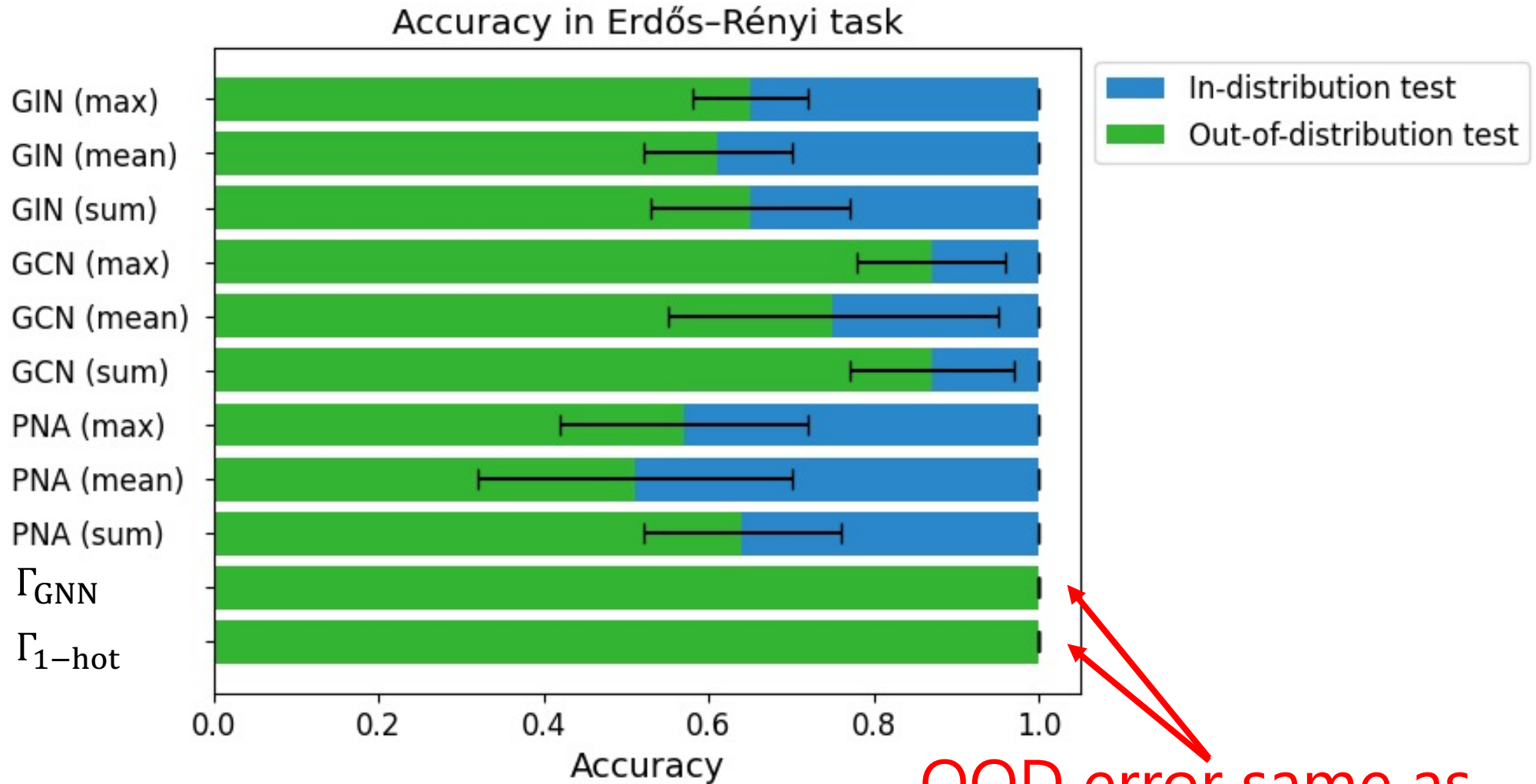
$$\mathrm{P}(|\mathrm{P}(Y = y \mid \mathcal{G}_{N^{tr}}^{tr}) - \rho(y, \Gamma(\mathcal{G}_{N^{tr}}^{tr}))| \leq \epsilon) \geq 1 - \delta$$

If $\Gamma$ is OOD-invariant then test error is the same

$$\mathrm{P}(|\mathrm{P}(Y = y \mid \mathcal{G}_{N^{te}}^{te}) - \rho(y, \Gamma(\mathcal{G}_{N^{te}}^{te}))| \leq \epsilon) \geq 1 - \delta$$

A size-invariant representation has same error
**in-distribution** and **out-of-distribution**

# Erdős–Rényi Task Example

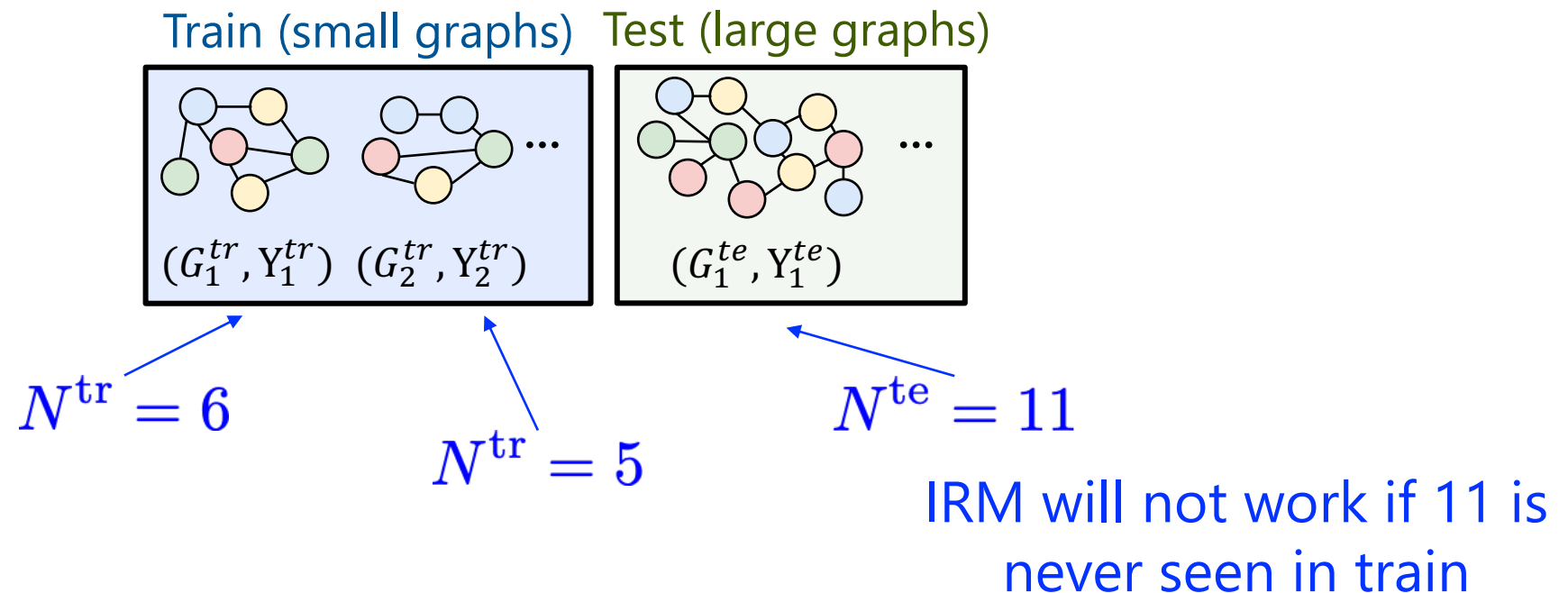Accuracy in Erdős–Rényi task

OOD error same as in-distribution error

# Invariant Risk Minimization (IRM)



IRM (Arjovskj et al., 2019) aims to learn an invariant representation.
However:

▸ 🚫 no guarantees if representation is nonlinear (e.g., GNN)

▸ 🚫 not applicable if training graphs have same size

▸ 🚫 not invariant if OOD support ≠ training support

Train (small graphs)   Test (large graphs)

$(G_1^{tr}, Y_1^{tr})$   $(G_2^{tr}, Y_2^{tr})$   $(G_1^{te}, Y_1^{te})$

$N^{\mathrm{tr}} = 6$

$N^{\mathrm{tr}} = 5$

$N^{\mathrm{te}} = 11$

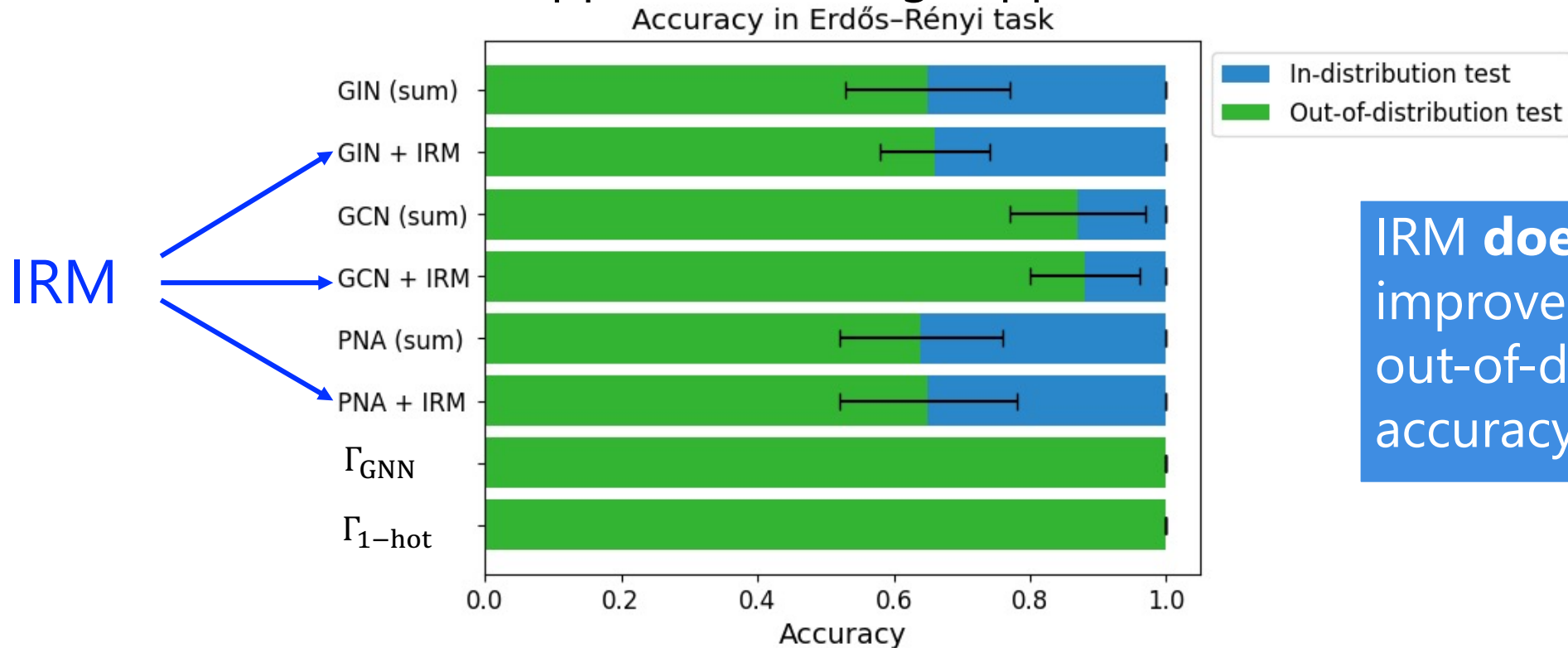IRM will not work if 11 is never seen in train

# Invariant Risk Minimization (IRM)

IRM (Arjovskj et al., 2019) aims to learn an invariant representation.
However:
- 🚫 no guarantees if representation is nonlinear (e.g., GNN)
- 🚫 not applicable if training graphs have same size
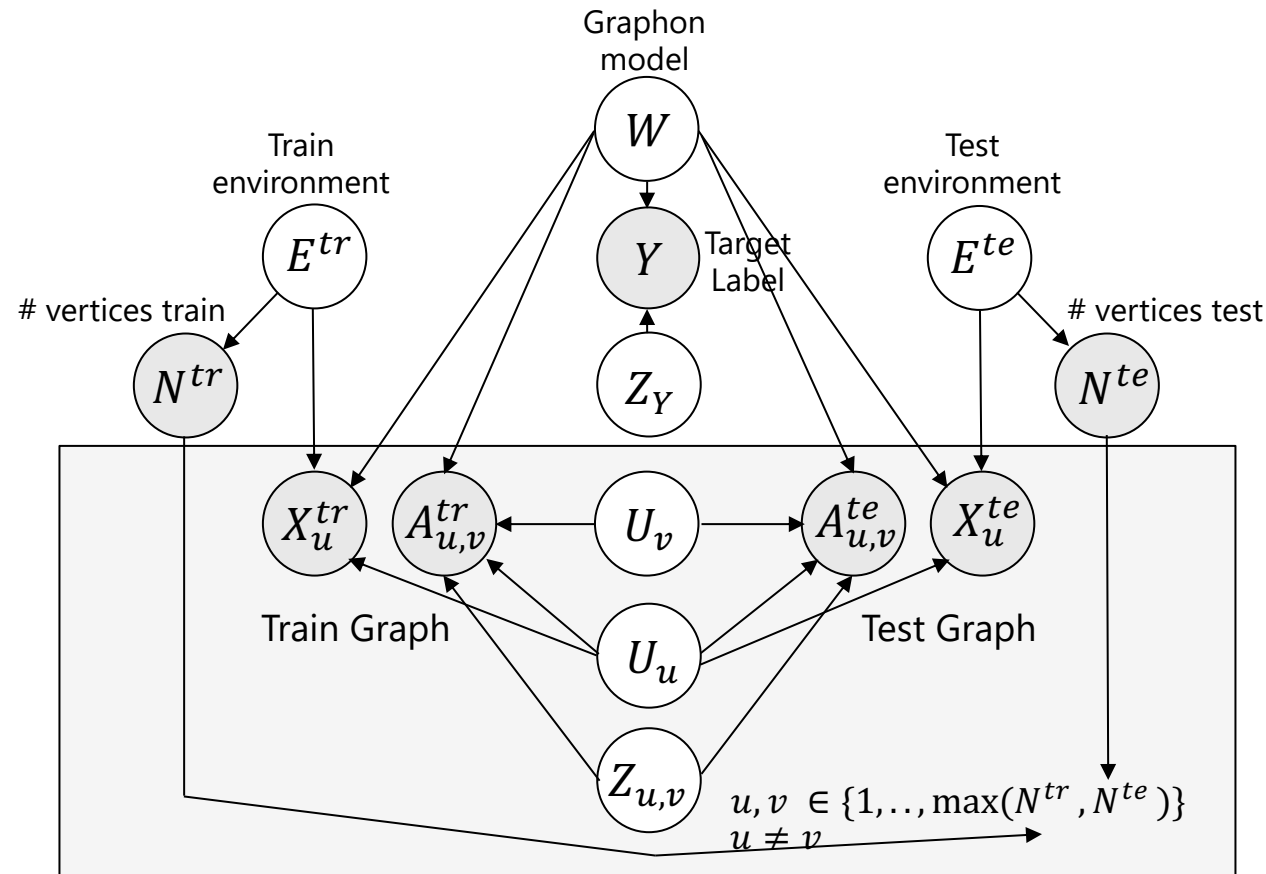- 🚫 not invariant if OOD support ≠ training support

**IRM**



Accuracy in Erdős–Rényi task

IRM **does not** improve out-of-distribution accuracy

# Causal Mechanism Assumed by Theorem 1 & Proposition 1

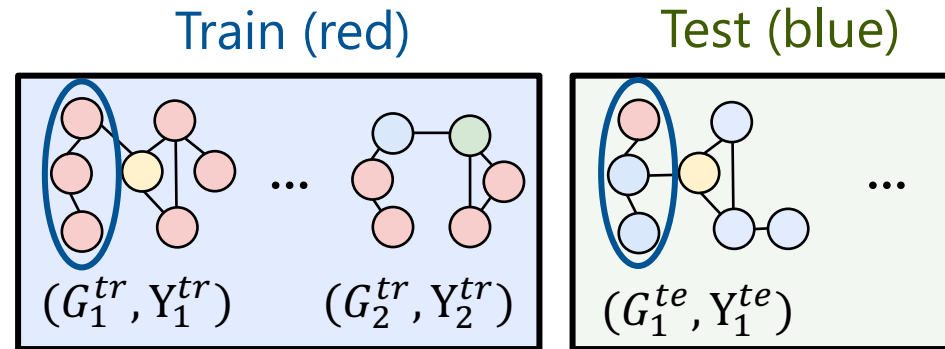# Causal Mechanism Assumed by **Theorem 1 & Proposition 1**

▸ *Structural Causal Model:*

- ◦ Graph label $Y$ is a function of the graph model $W$ + some random noise
- ◦ Graph size $N^{tr}$ ($N^{te}$) is a function of "environment" $E^{tr}$ ($E^{te}$) only
- ◦ Train (test) graphs are generated by $W$ and $E^{tr}$ ($E^{te}$) with same random noises

Improving OOD extrapolation
of vertex attributes

# Symmetry Regularization
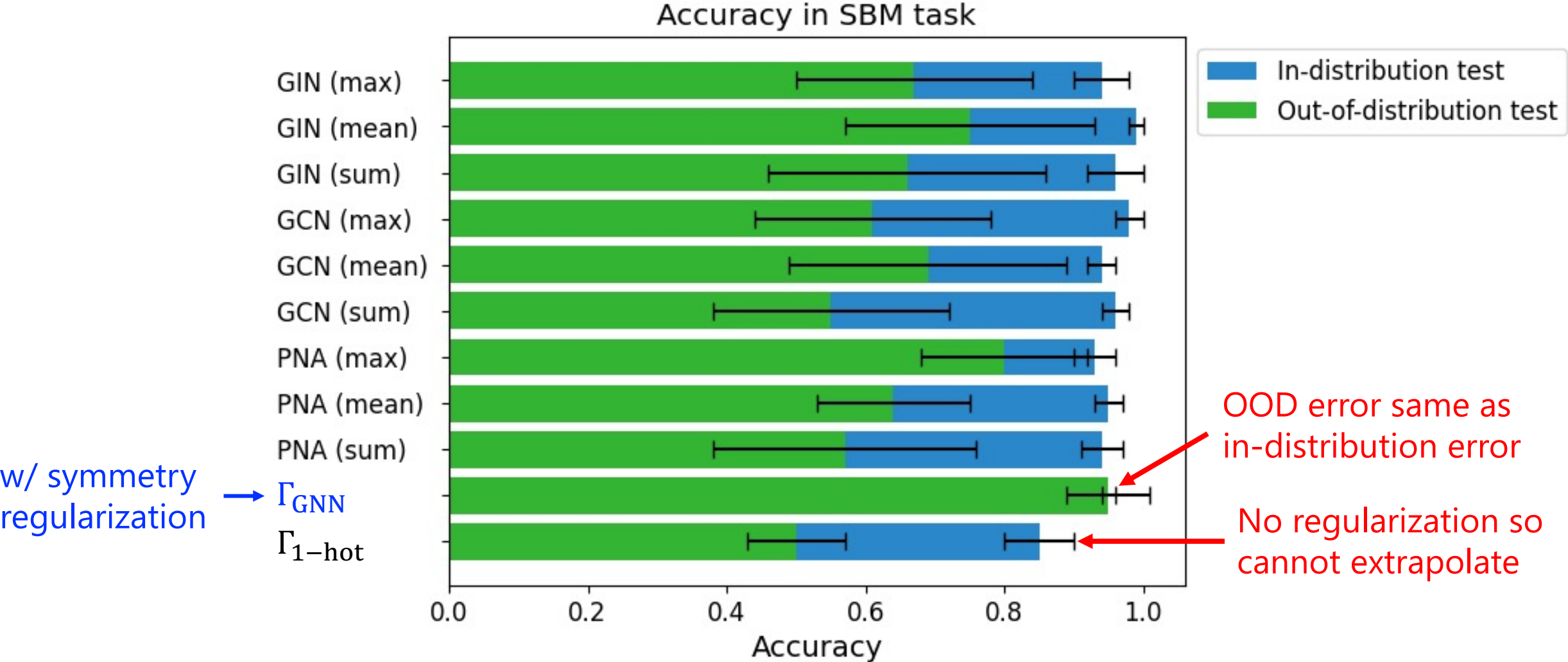
▸ What if OOD shift in attribute distribution?

Train (red)        Test (blue)



$(G_1^{tr}, Y_1^{tr})$        $(G_2^{tr}, Y_2^{tr})$        $(G_1^{te}, Y_1^{te})$

▸ Attribute symmetry regularization for representation $\Gamma_{\mathrm{GNN}}$:

$$\mathrm{Loss} + \lambda\|\mathrm{READOUT}_\Gamma(\mathrm{GNN}(\text{⬡})) - \mathrm{READOUT}_\Gamma(\mathrm{GNN}(\text{⬡}))\|$$

$$+ \lambda\|\mathrm{READOUT}_\Gamma(\mathrm{GNN}(\text{⬡})) - \mathrm{READOUT}_\Gamma(\mathrm{GNN}(\text{⬡}))\|$$

$$+ \lambda \ldots$$

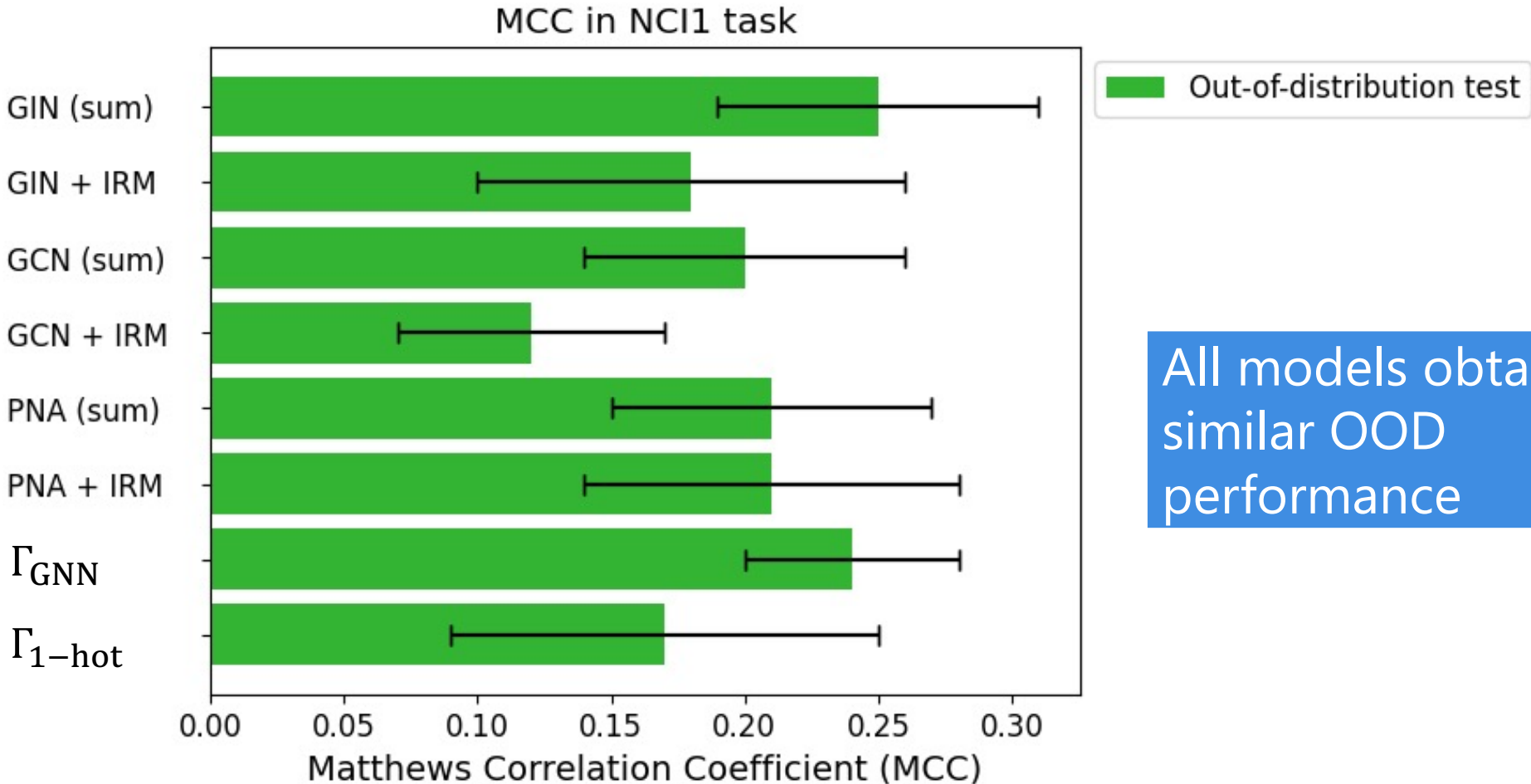Pushes subgraph representations towards topology-only unless hurts training loss

Symmetry regularization helps $\Gamma_{GNN}$ extrapolate to OOD attributes

OOD Extrapolation Depends on
Causal Mechanism Driving Distribution Shift

I.e.: no OOD universal representations!

# No OOD Universal Representations
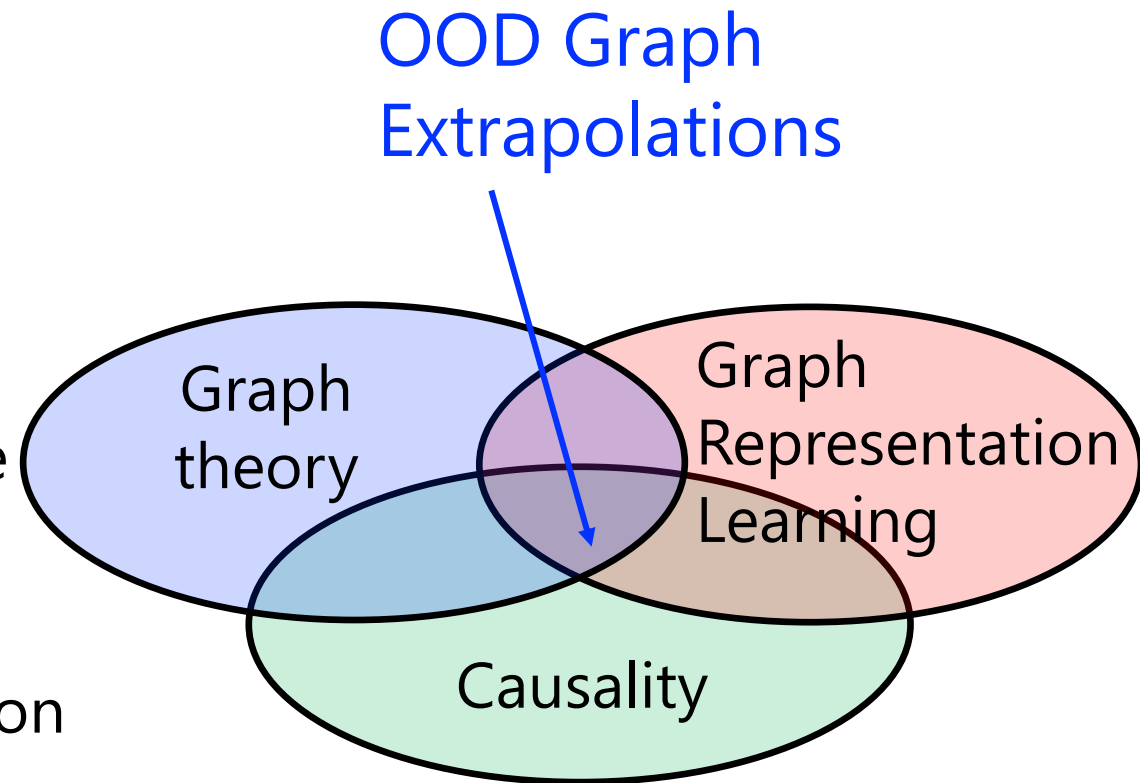
## NCI1 task does not follow our causal mechanism



MCC in NCI1 task

All models obtain similar OOD performance

Exciting new area in Graph Representation Learning:

▸ **OOD extrapolation without examples**

  ◦ Connects counterfactual predictions to stable graph properties

    • E.g., we use subgraph densities as a stable property

▸ There is no universal OOD graph representation

OOD Graph Extrapolations

Graph theory

Graph Representation Learning

Causality

## Thank you!

bbevilac@purdue.edu
zhou950@purdue.edu

# References

▸ Lovász, L. and Szegedy, B. Limits of dense graph sequences. Journal of Combinatorial Theory, Series B, 96(6):933–957, 2006.

▸ Arjovsky, M., Bottou, L., Gulrajani, I., and Lopez-Paz, D. Invariant risk minimization. arXiv preprint arXiv:1907.02893, 2019.

▸ Mouli, S. C. and Ribeiro, B. Neural networks for learning counterfactual g-invariances from single environments. ICLR, 2021.