

Meyer Scetbon, Joel Jennings, Agrin Hilmkil, Cheng Zhang, Chao Ma

## Why Causality?

**Generative Modeling**

Easily sampleable Noise  $N \sim P_N$  is mapped by  $H$  to Data  $X \sim P_X$  (Samples from  $P_X$ ).

**Goal:** Learn a function  $H$  that maps  $P_N$  to  $P_X$

**Examples:** VAEs, GANs, Normalizing Flows, Diffusions, ...

**Question:** what if  $P_X$  admits an (unknown) structure?

**Causal Generative Modeling**

$X := [X_1, X_2, X_3, \dots, X_T] \sim P_X$  |  $Y := [Y_1, Y_2, Y_3, \dots, Y_T] \sim P_Y$

**Goal:** Learn a function  $H$  that maps  $P_N$  to  $P_X$  and that satisfies the unknown structure of  $P_X$

**Advantage:** we can now simulate the effects of changes on the variables  $X_i$  w.r.t the structured generative process  $H$

➔ Causality offers a reliable tool for **decision-making** in generative modeling

## From DAGs to Causal Order

A DAG (always) induces a causal order

Ex:  $\pi = (X_1, X_2, X_3)$

**Equivalent Parameterization of SCMs**

- Noise Distribution:  $P_N$
- Topological Ordering:  $\pi$
- Map:  $H: (x, n) \in \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^d$  s.t.

**Fixed-Point SCM:**  $X = P_\pi^T H(P_\pi X, P_\pi N)$

Permutation matrix Associated to  $\pi$

**Partial Identifiability of Fixed-Point SCMs**

**Problem (Strong Identifiability):** Given  $\pi$  and  $P_X$ , can we recover uniquely the fixed-point SCM?

**Theorem (Monotonic Model)**

If  $H$  is **monotonic increasing** w.r.t the noise, and  $P_N$  is fixed, then there exists a **unique SCM** generating  $P_X$  with order  $\pi$  and noise  $P_N$

**Problem (Weak Identifiability):** Given  $\pi$  and  $P_X$ , can we recover uniquely the **interventional** and **counterfactual** distributions?

**Theorem (Invertible Model)**

- If  $H$  is **invertible** and  $C^1$ , then the **interventional** and **counterfactual** distributions are uniquely identifiable.
- Also, one can choose **arbitrarily** the noise distribution  $P_N$ , e.g. standard Gaussian.

## Learning Fixed-Point SCMs

**Recovering Causal Orders from Observations**

**Goal:** Infer in a zero-shot manner the causal order of variables from observations

**Dataset Generative Process**

- $P_N$  is a generated noise distribution on  $\mathbb{R}^d$
- $G$  is a generated DAG
- $H$  is a generated sequence of  $d$  functions  $H_i: \mathbb{R}^d \rightarrow \mathbb{R}$
- $D \in \mathbb{R}^{m \times d}$  is obtained by sampling  $m$  i.i.d samples  $N^{(k)} \sim P_N$  and then by solving for each  $k: X_i^{(k)} = H_i(PA_G(X_i^{(k)}), N_i^{(k)})$ ,  $\forall i \in [d]$

**Training:** Given  $(D_1, G_1), \dots, (D_n, G_n)$  i.i.d, we train a model  $\mathcal{M}$  that predicts the leaves of the graphs in a sequential manner.

**Loss:**  $\mathcal{M}$  is learned by minimizing d-TOE where at each step we remove a leaf.

**Algorithm 1 d-TOE( $\mathcal{M}, (D_n, G_n)$ )**

- Input:  $\mathcal{M}, (D_n, G_n)$
- Initialize d-TOE = 0.
- for  $q = 1$  to  $d$  do
- $p \leftarrow \mathcal{M}(D_n)$ ,  $y \leftarrow \mathcal{L}(G_n)$
- d-TOE  $\leftarrow$  d-TOE + BN( $p, y$ )
- $\hat{\ell} \leftarrow \text{argmax}_j |p_j|$ ,  $\ell \leftarrow \mathcal{B}(y, \hat{\ell})$
- $D_n \leftarrow \mathcal{R}_1(D_n, \ell)$ ,  $G_n \leftarrow \mathcal{R}_2(G_n, \ell)$
- end for
- Return d-TOE

**Learning SCMs on the Ordered Variables**

**Goal:** Learn the fixed-point SCM associated to a single dataset given the causal order

**Proposed Architecture**

- Causal Embedding:  $\mathcal{E}: X \in \mathbb{R}^d \rightarrow [X_1 * \theta_1, \dots, X_d * \theta_d] \in \mathbb{R}^{d \times d}$
- DAG Attention:  $DA(Q, K) = \frac{\exp((QK^T - M)/\sqrt{d})}{\mathfrak{S}(\exp(\frac{QK^T - M}{\sqrt{d}}))^{1/d}}$ , where  $M = \begin{bmatrix} +\infty & +\infty & +\infty \\ 0 & +\infty & +\infty \\ 0 & 0 & +\infty \end{bmatrix}$  and  $[\mathfrak{S}(v)]_i = v_i$  if  $v_i \geq 1$ , and  $[\mathfrak{S}(v)]_i = 1$  otherwise
- Causal Encoder  $C_i: h_{\ell+1} = h(DA(h_\ell, X)X + h_\ell) \in \mathbb{R}^{d \times d}$
- Causal Decoder  $\mathcal{F}: X \in \mathbb{R}^{d \times d} \rightarrow [(X_1, \omega_1), \dots, (X_d, \omega_d)] \in \mathbb{R}^d$

**Model:**  $\mathcal{T}: X \in \mathbb{R}^d \rightarrow \mathcal{F} \circ \mathcal{C}_L(\mathcal{E}(X)) \in \mathbb{R}^d$  satisfies the simple structure of fixed-point SCM.

**Training:** minimize the MSE  $\mathbb{E}_{Z \sim P_{P_\pi X}} \|T(Z) - Z\|^2$  where  $P_{P_\pi X}$  is the causally ordered distribution of observations.

**End-to-End Pipeline**

**Problem Set Up:** Let  $X^{(1)}, \dots, X^{(m)}$  i.i.d samples from  $P_X$ . Our goal is to learn an SCM generating  $P_X$ .

**Training:** We learn the SCM by training  $\mathcal{T}$  on  $D$  ordered according to the causal order predicted by  $\mathcal{M}$  on  $D$ .

**Inference:** Once trained, we show that under ANM, and correctness of  $\hat{\pi}$ ,  $\hat{\mathcal{T}}$  learned converges (in the limit of infinite sample) to the true SCM generating the data.

Diagram:  $D$  is input to Leaf Predictor  $\mathcal{M}$  and SCM Learner  $\mathcal{T}$ . Leaf Predictor  $\mathcal{M}$  outputs Predicted Causal Order  $\hat{\pi}$ . SCM Learner  $\mathcal{T}$  outputs Predicted Dataset  $\hat{D}$ . Loss to backpropagate is  $\|\hat{D} - D\|^2$ .

## Structural Causal Models

**Parameterization of SCMs**

- Noise distribution:  $P_N$
- Directed Acyclic Graph:  $G$
- System of Equations:  $X_i = H_i(PA_G(X_i), N_i), \forall i \in [d]$

**Generate observations**

- Sample  $(N_1, N_2, N_3) \sim P_N$
- Solve  $X_i = H_i(PA_G(X_i), N_i)$

**Generate effects of changes**

- Intervene on  $i: X_i \leftarrow a$
- Generate observations with the intervened system

**Learning SCMs from observations is hard**

- Computationally **Challenging**

Finding the DAG from data is an **NP-hard combinatorial** problem

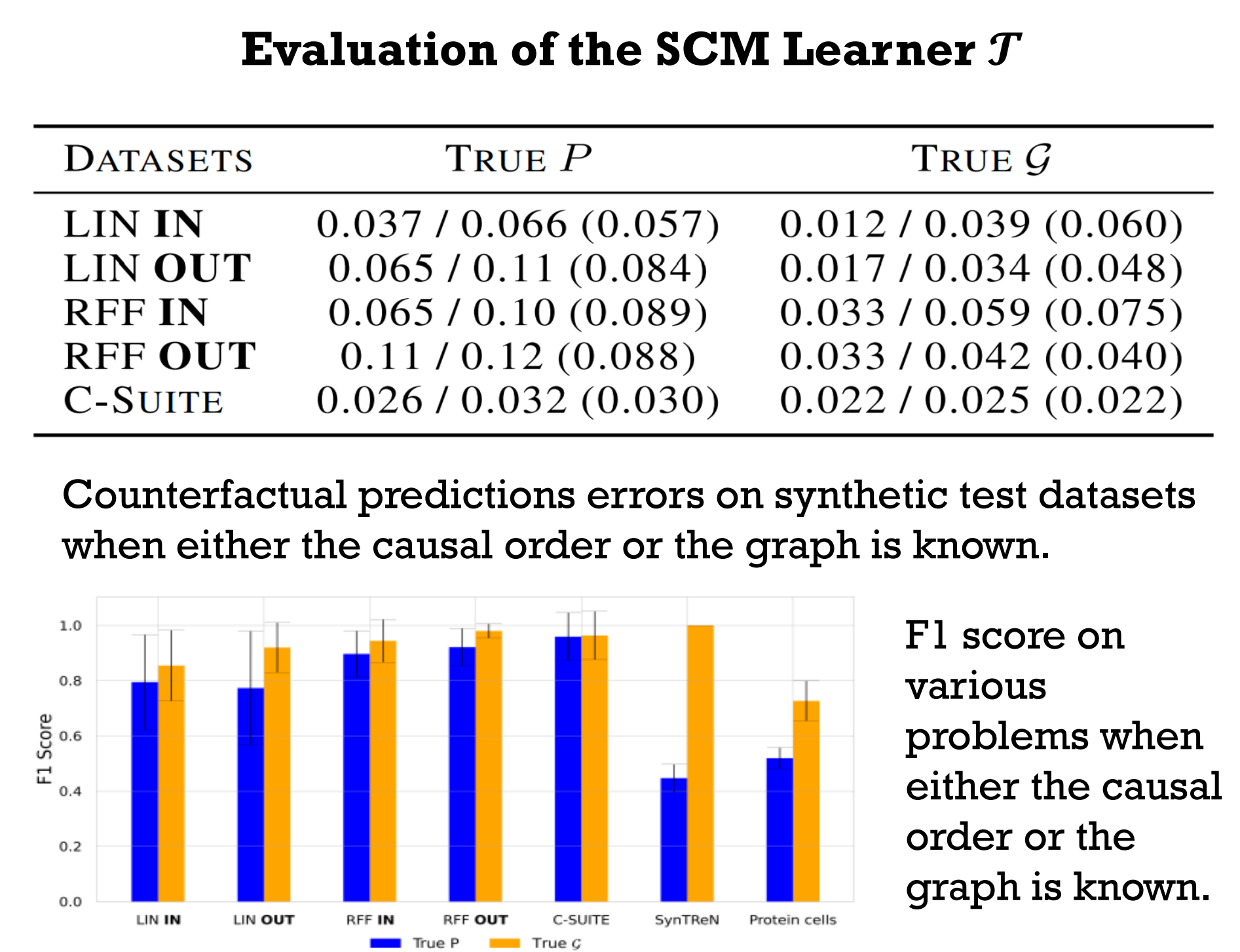
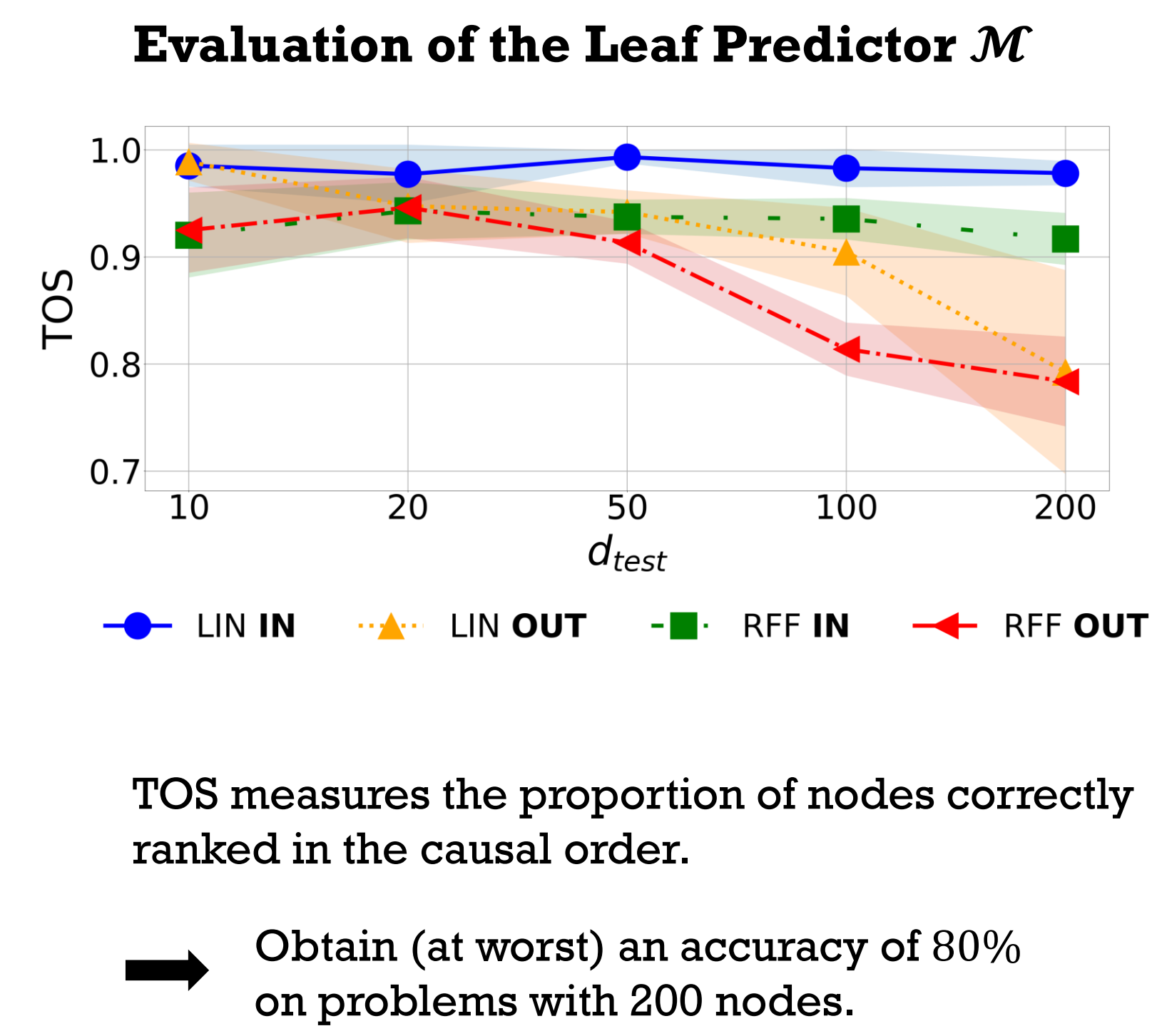
- Ill-posed** Inverse Problem

The **uniqueness** of the SCM is not always **guaranteed**

**Question:** can we improve on both limitations?

## Empirical Evaluations

- We train  $\mathcal{M}$  on  $\approx 200k$  datasets with  $n = 200$  samples and  $d = 50$  dimensions.
- We evaluate the models on  $\approx 400$  test datasets of various sizes, newly sampled from either the distribution used during training (**in-distribution**), or from a variant with a substantial shift (**out-of-distribution**).



### Benchmarking of End-to-End Pipeline

DATASETS	LIN OUT	RFF OUT
PC	0.47 (0.14)	0.40 (0.12)
GES	0.56 (0.12)	0.37 (0.060)
GOLEM	0.73 (0.29)	0.31 (0.13)
DECI	0.36 (0.13)	0.74 (0.14)
GRAN-DAG	0.29 (0.19)	0.50 (0.26)
DAG-GNN	0.61 (0.19)	0.44 (0.15)
DP-DAG	0.17 (0.074)	0.16 (0.067)
AVICI	0.73 (0.16)	0.74 (0.17)
<b>FiP (OURS)</b>	<b>0.76(0.20)</b>	<b>0.81(0.15)</b>

DATASETS	LIN OUT	RFF OUT
DECI	0.39 (0.29)	0.18 (0.12)
DoWHY - AVICI	0.20 (0.18)	0.16 (0.096)
<b>FiP (OURS)</b>	<b>0.13 (0.10)</b>	<b>0.13 (0.096)</b>
FiP w. $\mathcal{G}$	0.034 (0.048)	0.042 (0.040)
DoWHY w. $\mathcal{G}$	0.0017 (0.0017)	0.088 (0.072)

Comparison of the counterfactual predictions against various baselines on O.O.D datasets.

DATASET	MODEL	CF ERROR (RMSE)
TRIANGLE	CAUSAL NF	0.13 (0.02)
	CAREFL	0.17 (0.03)
	VACA	4.19 (0.04)
	<b>FiP</b>	<b>0.094(0.021)</b>
SIMPSON	CAUSAL NF	0.12(0.02)
	CAREFL	0.17 (0.04)
	VACA	1.50 (0.04)
	<b>FiP</b>	<b>0.12(0.0089)</b>

Comparison of the counterfactual predictions on synthetic datasets when the order is known.

**We obtain SoTA results for causal discovery and causal inference tasks on O.O.D test datasets.**