
Demystifying amortized causal discovery with transformers

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Supervised learning approaches for causal discovery from observational data often
2 achieve competitive performance despite seemingly avoiding explicit assumptions
3 that traditional methods make for identifiability. In this work, we investigate CSIVa
4 [1], a transformer-based model promising to train on synthetic data and transfer
5 to real data. First, we bridge the gap with existing identifiability theory and show
6 that constraints on the training data distribution implicitly define a prior on the test
7 observations. Consistent with classical approaches, good performance is achieved
8 when we have a good prior on the test data, and the underlying model is identifiable.
9 At the same time, we find new trade-offs. Training on datasets generated from
10 different classes of causal models, unambiguously identifiable in isolation,
11 improves the test generalization. Performance is still guaranteed, as the ambiguous
12 cases resulting from the mixture of identifiable causal models are unlikely to occur
13 (which we formally prove). Overall, our study finds that amortized causal discovery
14 still needs to obey identifiability theory, but it also differs from classical methods
15 in how the assumptions are formulated, trading more reliance on assumptions on
16 the noise type for fewer hypotheses on the mechanisms.

17 1 Introduction

18 Causal discovery aims to uncover the underlying causal relationships between variables of a system
19 from pure observations, which is crucial for answering interventional and counterfactual queries when
20 experimentation is impractical or unfeasible [2, 3, 4]. Unfortunately, causal discovery is inherently
21 ill-posed [5]: unique identification of causal directions requires restrictive assumptions on the class
22 of structural causal models (SCMs) that generated the data [6, 7, 8]. These theoretical limitations
23 often render existing methods inapplicable, as the underlying assumptions are usually untestable or
24 difficult to verify in practice [9].

25 Recently, supervised learning algorithms trained on synthetic data have been proposed to overcome
26 the need for specific hypotheses, which restrains the application of classical causal discovery methods
27 to real-world problems [1, 10, 11, 12, 13]. Seminal work from Lopez-Paz et al. [10] argues that
28 this learning-based approach to causal discovery would allow dealing with complex data-generating
29 processes and would greatly reduce the need for explicitly crafting identifiability conditions a-priori:
30 despite this ambitious goal, the output of these methods is generally considered unreliable, as no
31 theoretical guarantee is provided. A pair of non-identifiable structural causal models can be associated
32 with different causal graphs $\mathcal{G} \neq \tilde{\mathcal{G}}$, while entailing the same joint distribution p on the system's
33 variables. It is thus unclear how a learning algorithm presented with observational data generated from
34 p would be able to overcome these theoretical limits and correctly identify a unique causal structure.
35 However, the available empirical evidence seems not to care about impossibility results, as these
36 methods yield surprising generalization results on several synthetic benchmarks. Our work aims to
37 bridge this gap by studying the performance of a transformer architecture for causal discovery through

38 the lens of the theory of identifiability from observational data. Specifically, we analyze the CSIVa
 39 (Causal Structure Induction via Attention) model for causal discovery [1], focusing on bivariate graphs,
 40 as they offer a controlled yet non-trivial setting for the investigation. As our starting point, we provide
 41 closed-form examples that identify the limitations of CSIVa in recovering causal structures of linear
 42 non-Gaussian and nonlinear additive noise models, which are notably identifiable, and demonstrate the
 43 expected failures through empirical evidence. These findings suggest that the class of structural causal
 44 models that can be identified by CSIVa is inherently dependent on the specific class of SCMs observed
 45 during training. Thus, the need for restrictive hypotheses on the data-generating process is intrinsic
 46 to causal discovery, both in the traditional and modern learning-based approaches: assumptions on
 47 the test distribution either are posited when selecting the algorithm (traditional methods) or in the
 48 choice of the training data (learning-based methods). To address this limitation, we theoretically and
 49 empirically analyze *when* training CSIVa on datasets generated by multiple identifiable SCMs with
 50 different structural assumptions improves its generalization at test time. In summary:

- 51 • We show that the class of structural causal models that CSIVa can identify is defined by the
 52 class of SCMs observed through samples during the training. We reinforce the notion that
 53 identifiability in causal discovery inherently requires assumptions, which must be encoded
 54 in the training data in the case of learning-based approaches.
- 55 • To overcome this limitation, we study the benefits of CSIVa training on mixtures of causal
 56 models. We analyze when algorithms learned on multiple models are expected to identify
 57 broad classes of SCMs (unlike many classical methods). Empirically, we show that training
 58 on samples generated by multiple identifiable causal models with different assumptions on
 59 mechanisms and noise distribution results in significantly improved generalization abilities.

60 **Closely related works and their relation with CSIVa.** In this paper, we study *amortized inference*
 61 *of causal graphs*, i.e. optimization of an inference model to directly predict a causal structure from
 62 newly provided data. This is the first work that attempts to understand the connection between
 63 identifiability theory and amortized inference, while several algorithms have been proposed. In the
 64 context of purely observational data, Lopez-Paz et al. [10] defines a distribution regression problem
 65 [14] mapping the kernel mean embedding of the data distribution to a causal graph, while Li et al.
 66 [11] relies on equivariant neural network architectures. More recently, Lippe et al. [12] and Lorch
 67 et al. [13] proposed learning on interventional data, in addition to observations (in the same spirit as
 68 CSIVa). Despite different algorithmic implementations, the target object of estimation of most of
 69 these methods is the distribution over the space of all possible graphs, conditional on the input dataset
 70 (similarly, the ENCO algorithm in Lippe et al. [12] models the conditional distribution of individual
 71 edges). This justifies our choice of restricting our study to the CSIVa architecture (despite this
 72 being a clear limitation), as in the infinite observational sample limit, these methods approximate the
 73 same distribution. Methods necessarily requiring interventional data [15, 16, 17], and learning-based
 74 algorithms unsuitable for amortized inference [18, 19, 20, 21, 22] are out of the scope of this work.

75 2 Background and motivation

76 We start introducing structural causal models (SCMs), an intuitive framework that formalizes causal
 77 relations. Let X be a set of random variables in \mathbb{R} defined according to the set of structural equations:

$$X_i := f_i(X_{\text{PA}_i^{\mathcal{G}}}, N_i), \quad \forall i = 1, \dots, k. \quad (1)$$

78 $N_i \in \mathbb{R}$ are *noise* random variables. The function f_i is the *causal mechanism* mapping the set of *direct*
 79 *causes* $X_{\text{PA}_i^{\mathcal{G}}}$ of X_i and the noise term N_i , to X_i 's value. The *causal graph* \mathcal{G} is a directed acyclic
 80 graph (DAG) with nodes $X = \{X_1, \dots, X_k\}$, and edges $\{X_j \rightarrow X_i : X_j \in X_{\text{PA}_i^{\mathcal{G}}}\}$, with $\text{PA}_i^{\mathcal{G}}$
 81 indices of the parent nodes of X_i in \mathcal{G} . The causal model induces a density p_X over the vector X .

82 2.1 Causal discovery from observational data

83 Causal discovery from observational data is the inference of the causal graph \mathcal{G} from a dataset
 84 of i.i.d. observations of the random vector X . In general, without restrictive assumptions on the
 85 mechanisms and the noise distributions, the direction of edges in the graph \mathcal{G} is not identifiable, i.e.
 86 it can not be found from the population density p_X . In particular, it is possible to identify only a

87 Markov equivalence class, which is the set of graphs encoding the same conditional independencies
 88 as the density p_X . To clarify with an example, consider the causal graph $X_1 \rightarrow X_2$ associated
 89 with a structural causal model inducing a density p_{X_1, X_2} . If the model is not identifiable, there
 90 exists an SCM with causal graph $X_2 \rightarrow X_1$ that entails the same joint density p_{X_1, X_2} . The set
 91 $\{X_1 \rightarrow X_2, X_2 \rightarrow X_1\}$ is the Markov equivalence class of the graph $X_1 \rightarrow X_2$, i.e. the set of all
 92 graphs with X_1, X_2 mutually dependent. Clearly, in this setting, even the exact knowledge of p_{X_1, X_2}
 93 cannot inform us about the correct causal direction.

94 **Definition 1** (Identifiable causal model). Consider a structural causal model with underlying graph \mathcal{G}
 95 and p_X joint density of the causal variables. We say that the model is *identifiable* from observational
 96 data if the density p_X can not be entailed by a structural causal model with graph $\tilde{\mathcal{G}} \neq \mathcal{G}$.

97 We define the *post-additive noise model* (post-ANM) as the causal model with the set of equations:

$$X_i := f_{2,i}(f_{1,i}(X_{\text{PA}_i^{\mathcal{G}}}) + N_i), \quad \forall i = 1, \dots, d, \quad (2)$$

98 with $f_{2,i}$ invertible map and mutually independent noise terms. When $f_{2,i}$ is a nonlinear function,
 99 the post-ANM amounts to the identifiable *post-nonlinear* model (PNL) [8]. When $f_{2,i}$ is the identity
 100 function and $f_{1,i}$ nonlinear, it simplifies to the nonlinear *additive noise model* (ANM)[7, 23], which
 101 is known to be identifiable, and is described by the set of structural equations:

$$X_i := f_{1,i}(X_{\text{PA}_i^{\mathcal{G}}}) + N_i. \quad (3)$$

102 If, additionally, we restrict the mechanisms $f_{1,i}$ to be linear and the noise terms N_i to a non-Gaussian
 103 distribution, we recover the identifiable *linear non-Gaussian additive model* or LiNGAM [6]:

$$X_i = \sum_{j \in \text{PA}_i^{\mathcal{G}}} \alpha_j X_j + N_i, \quad \alpha_j \in \mathbb{R}. \quad (4)$$

104 2.2 Motivation and problem definition

105 Causal discovery from observational data relies on specific assumptions, which can be challenging to
 106 verify in practice [9]. To address this, recent methods leverage supervised learning for the amortized
 107 inference of causal graphs [1, 10, 11, 12, 13, 16, 24], optimizing an inference model to directly
 108 predict a causal structure from a provided dataset. While these approaches aim to reduce reliance on
 109 explicit identifiability assumptions, they often lack a clear connection to the existing causal discovery
 110 theory, making their outputs generally unreliable. We illustrate this limitation through an example.

111 **Example 1.** We consider the CSIVa transformer architecture proposed by Ke et al. [1], which can
 112 learn a map from observational data to a causal graph. The authors of the paper show that, in the
 113 finite sample regime, the CSIVa architecture exactly approximates the conditional distribution $p(\cdot|\mathcal{D})$
 114 over the space of possible graphs, given a dataset \mathcal{D} . Identifiability theory in causal discovery tells us
 115 that if the class of structural causal models that generated the observations is sufficiently constrained,
 116 then there is only one graph that can fit the data within that class. For example, consider the case
 117 of a dataset that is known to be generated by a nonlinear additive noise model, and let $p(\cdot|\mathcal{D}, \text{ANM})$
 118 be the conditional distribution that incorporates this prior knowledge on the SCM: then $p(\cdot|\mathcal{D}, \text{ANM})$
 119 concentrates all the mass on a single point \mathcal{G}^* , the true graph underlying the \mathcal{D} observations. Instead,
 120 in the absence of restrictions on the structural causal model, all the graphs in a Markov equivalence
 121 class are equally likely to be the correct solution given the data. Hence, $p(\cdot|\mathcal{D})$, the distribution
 122 learned by CSIVa, assigns equal probability to each graph in the Markov equivalence class of \mathcal{G}^* .

123 Our arguments of Example 1 are valid for all learning methods that approximate the conditional
 124 distribution over the space of graphs given the input data [1, 10, 11, 12, 13], and suggest that these
 125 algorithms are at most informative about the equivalence class of the causal graph underlying the
 126 observations. However, the available empirical evidence does not seem to highlight these limitations,
 127 as in practice these methods can infer the true causal DAG on several synthetic benchmarks. Thus, fur-
 128 ther investigation is necessary if we want to rely on their output in any meaningful sense. In this work,
 129 we analyze these "black-box" approaches through the lens of established theory of causal discovery
 130 from observational data (causal inference often lacks experimental data, which we do not consider).
 131 We study in detail the CSIVa architecture [1] (see Appendix A), a variation of the transformer neural
 132 network [25] for the supervised learning of algorithms for amortized causal discovery. This model is
 133 optimized via maximum likelihood estimation, i.e. finding Θ that minimizes $-\mathbb{E}_{\mathcal{G}, \mathcal{D}}[\ln \hat{p}(\mathcal{G}|\mathcal{D}; \Theta)]$,

134 where $\hat{p}(\mathcal{G}|\mathcal{D};\Theta)$ is the conditional distribution of a graph \mathcal{G} given a dataset \mathcal{D} parametrized by Θ .
135 We limit the analysis to CSIvA as it is a simple yet competitive end-to-end approach to learning causal
136 models. While this is clearly a limitation of the paper, our theoretical and empirical conclusions
137 exemplify both the role of theoretical identifiability in modern approaches and the new opportunities
138 they provide. Additionally, it fits well within a line of works arguing that specifically transformers
139 can learn causal concepts [26, 27, 28] and identify different assumptions in context [29].

140 3 Experimental results through the lens of theory

141 In this section, we present a comprehensive analysis of causal discovery with transformers and its
142 relation to the theoretical boundaries of causal discovery from observational data. We show that
143 suitable assumptions must be encoded in the training distribution to ensure the identifiability of the
144 test data, and we additionally study the effectiveness of training on mixtures of causal models to
145 overcome these limitations, improving generalization abilities.

146 3.1 Experimental design

147 We concentrate our research on causal models of two variables, causally related according to one of the
148 two graphs $X \rightarrow Y, Y \rightarrow X$. Bivariate models are the simplest non-trivial setting with a well-known
149 theory of causality inference [7, 8, 23], but also amenable to manipulation. This allows for compre-
150 hensive training and analysis of diverse SCMs and facilitates a clear interpretation of the results.

151 **Datasets.** Unless otherwise specified, in our experiments we train CSIvA on a sample of 15000
152 synthetically generated datasets, consisting of 1500 i.i.d. observations. Each dataset is generated ac-
153 cording to a single class of SCMs, defined by the mechanism type and the noise terms distribution. The
154 coefficients of the linear mechanisms are sampled in the range $[-3, -0.5] \cup [0.5, 3]$, removing small co-
155 efficients to avoid *close-to-unfaithful* effects [30]. Nonlinear mechanisms are parametrized according
156 to a neural network with random weights, a strategy commonly adopted in the literature of causal dis-
157 covery [1, 9]. The post-nonlinearity of the PNL model consists of a simple map $z \mapsto z^3$. Noise terms
158 are sampled from common distributions and a randomly generated density that we call *mlp*, previously
159 adopted in Montagna et al. [9], defined by a standard Gaussian transformed by a multilayer perceptron
160 (MLP) (Appendix B.2). We name these datasets *mechanism-noise* to refer to their underlying causal
161 model. For example, data sampled from a nonlinear ANM with Gaussian noise are named *nonlinear-*
162 *gaussian*. More details on the synthetic data generation schema are found in Appendix B.2. All data
163 are standardized by their empirical variance to remove opportunities to learn shortcuts [31, 32, 33].

164 **Metric and random baseline.** As our metric we use the structural Hamming distance (SHD), which
165 is the number of edge removals, insertions or flips to transform one graph to another. In the context
166 of bivariate causal graphs with a single edge, this is simply an error count, so correct inference corre-
167 sponds to $\text{SHD} = 0$, and an incorrect prediction gives $\text{SHD} = 1$. Additionally, we define a reference
168 random baseline, which assigns a causal direction according to a fair coin, achieving $\text{SHD} = 0.5$ in ex-
169 pectation. Each architecture we analyze in the experiments is trained 3 times, with different parameter
170 initialization and training samples: the SHD presented in the plots is the average of each of the 3 mod-
171 els on 1500 distinct test datasets of 1500 points each, and the error bars are 95% confidence intervals.

172 We detail the training hyperparameters in Appendix B.1. Next, we analyze our experimental results,
173 starting by investigating how well CSIvA generalizes on distributions unseen during training.

174 3.2 Warm up: is CSIvA capable of in and out-of-distribution generalization?

175 **In-distribution generalization.** First, we investigate the generalization of CSIvA on datasets
176 sampled from the structural casual model that generates the train distribution, with mechanisms and
177 noise distributions fixed between training and testing. We call this *in-distribution generalization*. As
178 a benchmark, we present the performance of several state-of-the-art approaches from the literature
179 on causal discovery: we consider the DirectLiNGAM, and NoGAM algorithms [34, 35], respectively
180 designed for the inference on LiNGAM and nonlinear ANM generated data¹. The results of Figure 1

¹The causal-learn implementation of the PNL algorithm could not perform better than random on our synthetic post-nonlinear data, and we observed that this was due to the sensitivity of the algorithm to the variance

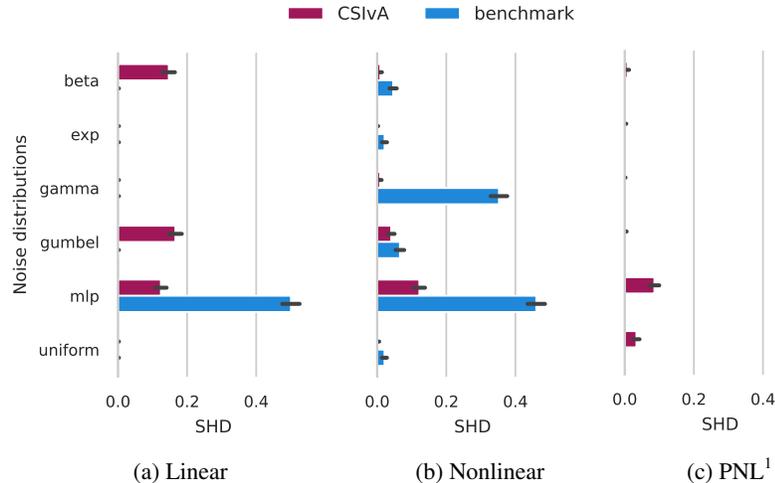


Figure 1: In-distribution generalization of CSivA trained and tested on data generated according to the same structural causal models, fixing mechanisms, and noise distributions between training and testing). As baselines for comparison, we use DirectLiNGAM on linear SCMs and NoGAM on nonlinear ANM (we use their causal-learn and dodiscover implementations). CSivA performance is clearly non-trivial and generalizing well.

181 show that CSivA can properly generalize to unseen samples from the training distribution: the majority
 182 of the trained models present SHD close to zero and comparable to the relative benchmark algorithm.

183 **Out-of-distribution generalization.** In practice, we generally do not know the SCM defining the
 184 test distribution, so we are interested in CSivA’s ability to generalize to data sampled from a class
 185 of causal models that is unobserved during training. We call this *out-of-distribution generalization*
 186 (OOD). We study OOD generalization to different noise terms, analyzing the network performance
 187 on datasets generated from causal models where the mechanisms are fixed with respect to the
 188 training, while the noise distribution varies (e.g., given linear-mlp training samples, testing occurs
 189 on linear-uniform data). Orthogonally to these experiments, we empirically validate CSivA’s OOD
 190 generalization over different mechanism types (linear, nonlinear, post-nonlinear), while leaving the
 191 noise distribution (mlp) fixed across test and training. In Figure 2a, we observe that CSivA cannot
 192 generalize across the different mechanisms, as the SHD of a network tested on unseen causal mech-
 193 anisms approximates that of the random baseline. Further, Figure 2b shows that out-of-distribution
 194 generalization across noise terms does not work reliably, and it is hard to predict when it might occur.

195 **Implications.** CSivA generalizes well to test data generated by the same class of SCMs used
 196 for training, in line with the findings in Ke et al. [1], which validates our implementation and
 197 training procedure. However, it struggles when the test data are out-of-distribution, not generated
 198 by causal models with the *same mechanisms and noise terms* it was trained on. While training on
 199 a wider class of SCMs might overcome this limitation, it requires caution. The identifiability of
 200 causal graphs indeed results from the interplay between the data-generating mechanisms and noise
 201 distribution. However, as we argue in our Example 1, the class of causal models that a supervised
 202 learning algorithm can identify is generally not clear. In what follows, we investigate this point and
 203 its implications for CSivA, showing that the identifiability of the test samples can be ensured by
 204 imposing suitable assumptions on the class of SCMs generating the training distribution.

205 3.3 How does CSivA relate to identifiability theory for causal graphs?

206 The CSivA algorithm does not make structural assumptions about the causal model underlying the
 207 input data. This implies that the output of this method is unclear: as CSivA targets the conditional dis-
 208 tribution $p(\cdot|\mathcal{D})$ over the space of graphs, in the absence of restrictions on the functional mechanisms

scale. So we report the plot of Figure 1c without benchmark comparison. We remark that the point of this experiment is not to make any claims on CSivA being state-of-the-art but to validate that the performance we obtain in our re-implementation is non-trivial. This is clear for PNL, even without comparison.

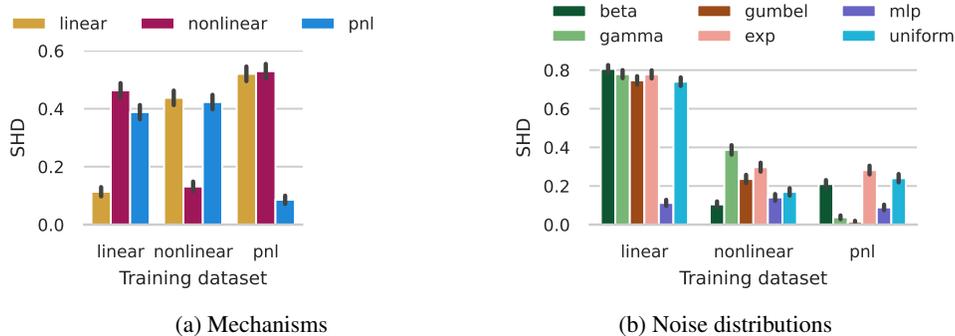


Figure 2: Out-of-distribution generalisation. We train three CSiVA models on data sampled from SCMs with linear, nonlinear additive, and post-nonlinear mechanisms; and noise fixed *mlp* noise distribution. In Figure (a) we test across different noise distributions, with test mechanisms fixed from training. In Figure (b) we test each network on different mechanisms and fixed *mlp* noise. CSiVA struggles to generalize to unseen causal mechanisms and often displays degraded performance over new noise distributions.

209 and the distribution of the noise terms, the causal graph $X \rightarrow Y$ is indistinguishable from $Y \rightarrow X$,
 210 as they are both equally likely to underlie the joint density $p_{X,Y}$ generating the data. As we discuss in
 211 Example 1, the graphical output of the trained architecture could at most identify the equivalence class
 212 of the true causal graph. Yet, our experiments of Section 3.2 show that CSiVA is capable of good in-
 213 distribution generalization, often inferring the correct DAG at test time. We explain this seeming con-
 214 tradiction with the following hypothesis, which motivates the analysis in the remainder of this section.

215 **Hypothesis 1.** *The class of structural causal models that can be identified by CSiVA is defined by the*
 216 *class of structural causal models underlying the generation of the training data.*

217 To support and clarify our statement, we present the following example, adapted from Hoyer et al. [7].

218 **Example 2.** Consider the causal model $Y = f(X) + N$, where $f(X) = -X$ and p_X, p_N are
 219 Gumbel densities $p_X(x) = \exp(-x - \exp(-x))$ and $p_N(n) = \exp(-n - \exp(-n))$. This model
 220 satisfies the assumptions of the LiNGAM, so it is identifiable, in the sense that a backward linear
 221 model with the same distribution does not exist. However, in this special case, we can build a
 222 backward nonlinear additive noise model $X = g(Y) + \tilde{N}$ with independent noise terms: taking
 223 $p_Y(y) = \exp(-y - 2\log(1 + \exp(-y)))$ to be the density of a logistic distribution, $p_{\tilde{N}}(\tilde{n}) =$
 224 $\exp(-2\tilde{n} - \exp(-\tilde{n}))$ and $g(y) = \log(1 + \exp(-y))$; we see that $p_{X,Y}$ can factorize according
 225 to two opposite causal directions, as $p_{X,Y}(x, y) = p_N(y - f(x))p_X(x) = p_{\tilde{N}}(x - g(y))p_Y(y)$.
 226 Given a dataset \mathcal{D} of observations from the forward linear model, causal discovery methods like
 227 DirectLiNGAM [34] can provably identify the correct causal direction $X \rightarrow Y$, assuming that
 228 sufficient samples are provided. Instead, the behavior of CSiVA seems hard to predict: given that
 229 the network approximates the conditional distribution $p(\cdot|\mathcal{D})$ over the possible graphs, for \mathcal{D} with
 230 arbitrary many samples we have $p(X \rightarrow Y|\mathcal{D}) = p(Y \rightarrow X|\mathcal{D}) = 0.5$. On the other hand, given
 231 the prior knowledge that the data-generating SCM is a linear non-gaussian additive noise model, we
 232 have $p(X \rightarrow Y|\mathcal{D}, \text{LiNGAM}) = 1$, because the LiNGAM is identifiable. In this sense, the class
 233 of structural causal models that CSiVA correctly infers appears to be determined by the structural
 234 causal models underlying the generation of the training data. Under our Hypothesis 1, training CSiVA
 235 exclusively on LiNGAM-generated data is equivalent to learning the distribution $p(\cdot|\mathcal{D}, \text{LiNGAM})$,
 236 such that the network should be able to identify the forward linear model, whereas it could only infer
 237 the equivalence class of the causal graph if its training datasets include observations from a nonlinear
 238 additive noise model.

239 The empirical results of Figure 3a show that CSiVA behaves according to our hypothesis: when
 240 training exclusively occurs on datasets $\{\mathcal{D}_{i,\rightarrow}\}_i$ generated by the *forward linear-gumbel model* of
 241 Example 2, the network can identify the causal direction of test data generated according to the same
 242 SCM. Similarly, the transformer trained on datasets $\{\mathcal{D}_{i,\leftarrow}\}_i$ from the *backward nonlinear model*
 243 of the example can generalize to test data coming from the same distribution. According to our claim,
 244 instead, the network that is trained on the union of the training samples $\{\mathcal{D}_{i,\rightarrow}\}_i \cup \{\mathcal{D}_{i,\leftarrow}\}_i$ from
 245 the forward and backward models (50:50 ratio in Figure 3a) displays the same test SHD (around
 246 0.5) as a random classifier assigning the causal direction with equal probability.

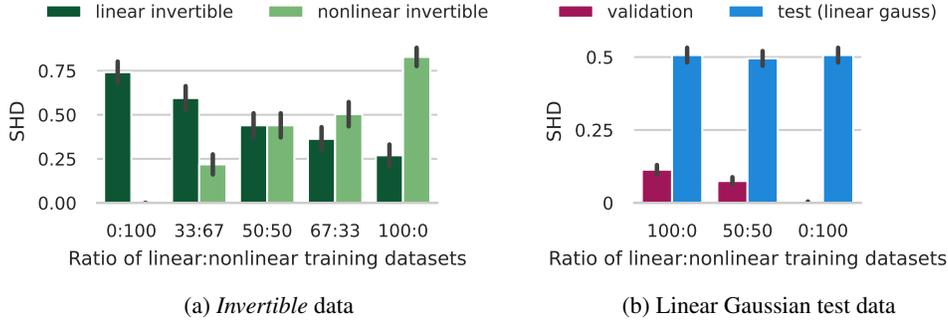


Figure 3: Experiments on identifiability theory. In Figure (a) we test the performance on linear-Gaussian data. Models are trained with different ratios of samples from linear and nonlinear SCMs with Gaussian noise terms. The validation results showcase that the networks were trained successfully. Figure (b) shows the SHD of models trained on different ratios of *linear* and *nonlinear invertible* data of Example 2. CSIVa behaves according to identifiability theory, failing to predict on linear Gaussian models and *invertible* data (50:50 ratio).

247 Further, we investigate CSIVa’s relation with known identifiability theory by training and testing the
 248 architecture on data from a linear Gaussian model, which is well-known to be unidentifiable. Not
 249 surprisingly, the results of Figure 3b show that none of the algorithms that we learn can infer the
 250 causal order of linear Gaussian models with test SHD any better than a random baseline.

251 **Implications.** Our experiments show that CSIVa learns algorithms that closely follow identifiability
 252 theory for causal discovery. In particular, while the method itself does not require explicit assumptions
 253 on the data-generating process, the chosen training data ultimately determines the class of causal
 254 models identifiable during inference. Notably, previous work has argued that supervised learning
 255 approaches in causal discovery would help with "dealing with complex data-generating processes and
 256 greatly reduce the need of explicitly crafting identifiability conditions a-priori", Lopez-Paz et al. [10].
 257 In the case of CSIVa, this expectation does not appear to be fulfilled, as the assumptions still need
 258 to be encoded explicitly in the training data. However, this observation opens two new important
 259 questions: (1) Can we train a single network to encompass multiple (or even all) identifiable causal
 260 structures? (2) How much ambiguity might exist between these identifiable models?

261 3.4 A *low-dimensions* argument in favor of learning from multiple causal models

262 Example 2 of the previous section shows that elements of distinct classes of identifiable structural
 263 causal models, such as LiNGAM and nonlinear ANM, may become non-identifiable when we
 264 consider their union. In this section, we show that in the class of post-additive noise models given
 265 by equation (2) (obtained as the union of the LiNGAM, the nonlinear ANM, and the post-nonlinear
 266 model), the set of distributions that is non-identifiable is negligible. Our proposition extends the
 267 results of Hoyer et al. [7], which are limited to the case of linear and nonlinear additive noise models,
 268 and Zhang and Hyvärinen [8], which provides the conditions of identifiability of the post-ANM
 269 without bounding the set of non-identifiable distributions.

270 Let X, Y be a pair of random variables generated according to the causal direction $X \rightarrow Y$ and the
 271 post-additive noise model structural equation:

$$Y = f_2(f_1(X) + N_Y), \quad (5)$$

272 where N_Y and X are independent random variables, and f_2 is invertible. If the SCM is non-
 273 identifiable, the data-generating process can be described by a *backward* model with the structural
 274 equation:

$$X = g_2(g_1(Y) + N_X), \quad (6)$$

275 N_X independent from Y , and g_2 invertible. We introduce the random variables \tilde{X}, \tilde{Y} , such that the
 276 forward and backward equations can be rewritten as

$$\begin{aligned} Y &= f_2(\tilde{Y}), & \tilde{Y} &:= f_1(X) + N_Y, \\ X &= g_2(\tilde{X}), & \tilde{X} &:= g_1(Y) + N_X. \end{aligned}$$

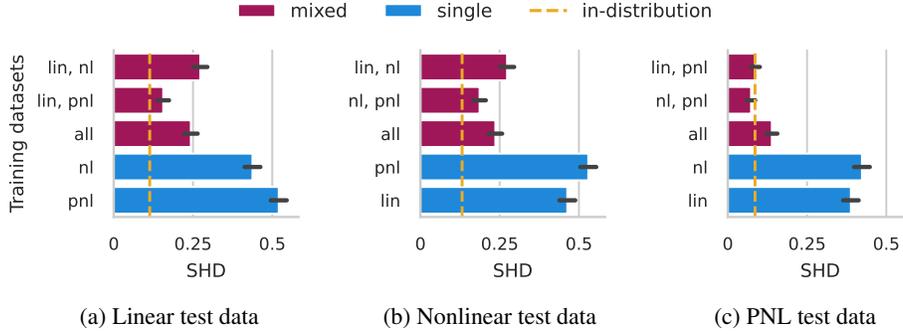


Figure 4: Mixture of causal mechanisms. We train four models on samples from structural casual models with different mechanism types. We compare their test SHD (the lower, the better) against networks trained on datasets generated according to a single type of mechanism. The dashed line indicates the test SHD of a model trained on samples with the same mechanisms as test SCM. Training on multiple causal models with different mechanisms (*mixed* bars) always improves performance compared to training on single SCMs.

277 We note that this implies that the following invertible additive noise models on \tilde{X}, \tilde{Y} hold:

$$\tilde{Y} = h_Y(\tilde{X}) + N_Y, \quad h_Y := f_1 \circ g_2, \quad (7)$$

$$\tilde{X} = h_X(\tilde{Y}) + N_X, \quad h_X := g_1 \circ f_2. \quad (8)$$

278 **Proposition 1** (Adapted from Hoyer et al. [7]). *Let p_{N_Y}, h_X, h_Y be fixed, and define $\nu_Y := \log p_{N_Y},$*
 279 *$\xi := \log p_{\tilde{X}}$. Suppose that p_{N_Y} and $p_{\tilde{X}}$ are strictly positive densities, and that $\nu_Y, \xi, f_1, f_2, g_1,$ and*
 280 *g_2 are three times differentiable. Further, assume that for a fixed pair h_Y, ν_Y exists $\tilde{y} \in \mathbb{R}$ s.t. $\nu_Y''(\tilde{y} -$*
 281 *$h_Y(\tilde{x}))h_Y'(\tilde{x}) \neq 0$ is satisfied for all but a countable set of points $\tilde{x} \in \mathbb{R}$. Then, the set of all densities*
 282 *$p_{\tilde{X}}$ of \tilde{X} such that both equations (5) and (6) are satisfied is contained in a 2-dimensional space.*

283 **Implications.** Our result is closely related to Theorem 1 of Hoyer et al. [7], which we simply
 284 generalize to the post-ANM. Intuitively, it says that the space of all continuous distributions such that
 285 the bivariate post-ANM is non-identifiable is contained in a 2-dimensional space. As the space of
 286 continuous distributions of random variables is infinite-dimensional, we conclude that the post-ANM
 287 is generally identifiable, which suggests that the setting of Example 2 is rather artificial. Our results
 288 provide a theoretical ground for training causal discovery algorithms on datasets generated from
 289 multiple identifiable SCMs. This is particularly appealing in the case of CSiVA, given the poor OOD
 290 generalization ability observed in our experiments of Section 3.2.

291 3.5 Can we train CSiVA on multiple causal models for better generalization?

292 In this section, we investigate the benefits of training over multiple causal models, i.e. on samples
 293 generated by a combination of classes of identifiable SCMs characterized by different mechanisms
 294 and noise terms distribution. Our motivation is as follows: given that our empirical evidence
 295 shows that CSiVA is capable of in-distribution generalization, whereas dramatically degrades the
 296 performance when testing occurs out-of-distribution, it is thus desirable to increase the class of
 297 causal models represented in the training datasets. We separately study the effects of training over
 298 multiple mechanisms and multiple noise distributions and compare the testing performance against
 299 architectures trained on samples of a single SCM.

300 **Mixture of causal mechanisms.** We consider four networks optimized by training of CSiVA on
 301 datasets generated from pairs (or triples) of distinct SCMs, with fixed *mlp* noise and which differ in
 302 terms of their mechanisms type: linear and nonlinear; nonlinear and post-nonlinear; linear and post-
 303 nonlinear; linear, nonlinear and post-nonlinear. The number of training datasets for each architecture is
 304 fixed (15000) and equally split between the causal models with different mechanism types. The results
 305 of Figure 4 show that the networks trained on mixtures of mechanisms all present significantly better
 306 test SHD compared to CSiVA models trained on a single mechanism type. We find that learning on
 307 multiple SCMs improves the SHD from ~ 0.5 to ~ 0.2 both on linear and nonlinear test data (Figures
 308 4a and 4b), and even better accuracy is achieved on post-nonlinear samples, as shown in Figure 4c.

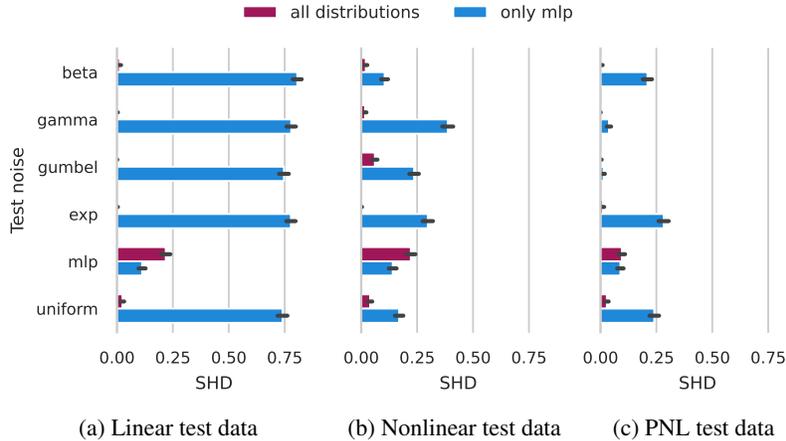


Figure 5: Mixture of noise distributions. We train three networks on samples from SCMs with different noise terms distributions and fixed mechanism types: linear, nonlinear, and post-nonlinear. We present their test SHD (the lower, the better) on data from SCMs with the mechanisms fixed with respect to training, and noise terms changing between each dataset. Training on multiple causal models with different noises (*all distributions* bars) always improves performance compared to training on single SCMs with fixed mlp noise (*only mlp* bars).

309 **Mixture of noise distributions.** Next, we analyze the test performance of three CSIVa networks
 310 optimized on samples from structural causal models that have different distributions for their noise
 311 terms, while keeping the mechanism types fixed. Figure 5 shows that training over different noises
 312 (beta, gamma, gumbel, exponential, mlp, uniform) always results in a network that is agnostic with
 313 respect to the noise distributions of the SCM generating the test samples, always achieving $\text{SHD} < 0.1$,
 314 with the exception of datasets with mlp error terms (0.2 average SHD on nonlinear and pnl data).

315 **Implications.** We have shown that learning on mixtures of SCMs with different noise term dis-
 316 tributions and mechanism types leads to models generalizing to a much broader class of structural
 317 causal models during testing. Hence, combining datasets generated from multiple models looks
 318 like a promising framework to overcome the limited out-of-distribution generalization abilities of
 319 CSIVa observed in Section 3.2. However, it is easier to incorporate prior assumptions on the class of
 320 causal mechanisms (linear, non-linear, post-non-linear) compared to the noise distributions (which are
 321 potentially infinite). This introduces a trade-off between amortized inference and classical methods
 322 for causal discovery: for example, RESIT, NoGAM, and CAM [23, 35, 36] algorithms require no
 323 assumptions on the noise type, but only work for a limited class of mechanisms (nonlinear).

324 4 Conclusion

325 In this work, we investigate the interplay between identifiability theory and supervised learning
 326 for amortized inference of causal graphs, using CSIVa as the ground of our study. Consistent
 327 with classical algorithms, we demonstrate that good performance can be achieved if (i) we have
 328 a good prior on the structural causal model generating the test data (ii) the setting is identifiable.
 329 In particular, prior knowledge of the test distribution is encoded in the training data in the form
 330 of constraints on the structural causal model underlying their generation. With these results, we
 331 highlight the need for identifiability theory in modern learning-based approaches to causality, while
 332 past works have mostly disregarded this connection. Further, our findings provide the theoretical
 333 ground for training on observations sampled from multiple classes of identifiable SCMs, a strategy
 334 that improves test generalization to a broad class of causal models. Finally, we highlight an interesting
 335 new trade-off regarding identifiability: traditional methods like LiNGAM, RESIT, and PNL require
 336 strong restrictions on the structural mechanisms underlying the data generation (linear, nonlinear
 337 or post-nonlinear) while generally being agnostic relative to the noise terms distribution. Training
 338 on mixtures of causal models instead offers an alternative that is less reliant on assumptions on the
 339 mechanisms, while incorporating knowledge about all possible noise distributions in the training data
 340 is practically impossible to achieve. We leave it to future work to reproduce our analysis on a wider
 341 class of architectures, as well as extending our study to interventional data with more than two nodes.

References

- 342
- 343 [1] Nan Rosemary Ke, Silvia Chiappa, Jane X. Wang, Jorg Bornschein, Anirudh Goyal, Melanie
344 Rey, Theophane Weber, Matthew Botvinick, Michael Curtis Mozer, and Danilo Jimenez
345 Rezende. Learning to Induce Causal Structure. In *International Conference on Learning Representations*,
346 September 2022. URL https://openreview.net/forum?id=hp_RwhKDJ5.
- 347 [2] Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. *Elements of Causal Inference:
348 Foundations and Learning Algorithms*. Adaptive Computation and Machine Learning. The MIT
349 Press, Cambridge, Mass, 2017. ISBN 978-0-262-03731-0.
- 350 [3] Judea Pearl. *Causality*. Cambridge University Press, Cambridge, 2nd edition, 2009.
- 351 [4] Peter Spirtes. Introduction to causal inference. *Journal of Machine Learning Research*, 11(54):
352 1643–1662, 2010. URL <http://jmlr.org/papers/v11/spirtes10a.html>.
- 353 [5] Clark Glymour, Kun Zhang, and Peter Spirtes. Review of causal discovery methods based on
354 graphical models. *Frontiers in Genetics*, 10, 2019. ISSN 1664-8021. doi: 10.3389/fgene.2019.
355 00524. URL <https://www.frontiersin.org/articles/10.3389/fgene.2019.00524>.
- 356 [6] Shohei Shimizu, Patrik O. Hoyer, Aapo Hyvärinen, and Antti Kerminen. A linear non-gaussian
357 acyclic model for causal discovery. *Journal of Machine Learning Research*, 7:2003–2030, dec
358 2006. ISSN 1532-4435.
- 359 [7] Patrik Hoyer, Dominik Janzing, Joris M Mooij, Jonas Peters, and Bernhard Schölkopf. Non-
360 linear causal discovery with additive noise models. In D. Koller, D. Schuurmans, Y. Bengio,
361 and L. Bottou, editors, *Advances in Neural Information Processing Systems*, volume 21. Cur-
362 ran Associates, Inc., 2008. URL [https://proceedings.neurips.cc/paper/2008/file/
363 f7664060cc52bc6f3d620bcdec94a4b6-Paper.pdf](https://proceedings.neurips.cc/paper/2008/file/f7664060cc52bc6f3d620bcdec94a4b6-Paper.pdf).
- 364 [8] Kun Zhang and Aapo Hyvärinen. On the identifiability of the post-nonlinear causal model. In
365 *Proceedings of the Twenty-Fifth Conference on Uncertainty in Artificial Intelligence*, UAI '09,
366 page 647–655, Arlington, Virginia, USA, 2009. AUAI Press. ISBN 9780974903958.
- 367 [9] Francesco Montagna, Atalanti Mastakouri, Elias Eulig, Nicoletta Noceti, Lorenzo Rosasco,
368 Dominik Janzing, Bryon Aragam, and Francesco Locatello. Assumption violations
369 in causal discovery and the robustness of score matching. In A. Oh, T. Neumann,
370 A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural
371 Information Processing Systems*, volume 36, pages 47339–47378. Curran Associates,
372 Inc., 2023. URL [https://proceedings.neurips.cc/paper_files/paper/2023/file/
373 93ed74938a54a73b5e4c52bbaf42ca8e-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/93ed74938a54a73b5e4c52bbaf42ca8e-Paper-Conference.pdf).
- 374 [10] David Lopez-Paz, Krikamol Muandet, Bernhard Schölkopf, and Ilya Tolstikhin. Towards a
375 learning theory of cause-effect inference. In *Proceedings of the 32nd International Conference
376 on International Conference on Machine Learning - Volume 37*, ICML'15, page 1452–1461.
377 JMLR.org, 2015.
- 378 [11] Hebi Li, Qi Xiao, and Jin Tian. Supervised Whole DAG Causal Discovery, June 2020.
- 379 [12] Phillip Lippe, Taco Cohen, and Efstratios Gavves. Efficient neural causal discovery without
380 acyclicity constraints. In *International Conference on Learning Representations*, 2022. URL
381 <https://openreview.net/forum?id=eYciPrLuUhG>.
- 382 [13] Lars Lorch, Scott Sussex, Jonas Rothfuss, Andreas Krause, and Bernhard Schölkopf. Amortized
383 inference for causal structure learning. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and
384 Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022. URL
385 <https://openreview.net/forum?id=eV4JI-MMeX>.
- 386 [14] Zoltan Szabo, Bharath Sriperumbudur, Barnabas Poczos, and Arthur Gretton. Learning theory
387 for distribution regression. *Journal of Machine Learning Research*, 17:1–40, 09 2016.

- 388 [15] Philippe Brouillard, Sébastien Lachapelle, Alexandre Lacoste, Simon Lacoste-Julien,
389 and Alexandre Drouin. Differentiable causal discovery from interventional data. In
390 H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neu-*
391 *ral Information Processing Systems*, volume 33, pages 21865–21877. Curran Associates,
392 Inc., 2020. URL https://proceedings.neurips.cc/paper_files/paper/2020/file/f8b7aa3a0d349d9562b424160ad18612-Paper.pdf.
393
- 394 [16] Nan Rosemary Ke, Olexa Bilaniuk, Anirudh Goyal, Stefan Bauer, Hugo Larochelle, Bernhard
395 Schölkopf, Michael Curtis Mozer, Christopher Pal, and Yoshua Bengio. Neural causal structure
396 discovery from interventions. *Transactions on Machine Learning Research*, 2023. ISSN
397 2835-8856. URL <https://openreview.net/forum?id=rdHVPPVuXa>. Expert Certification.
- 398 [17] Nino Scherrer, Olexa Bilaniuk, Yashas Annadani, Anirudh Goyal, Patrick Schwab, Bernhard
399 Schölkopf, Michael C. Mozer, Yoshua Bengio, Stefan Bauer, and Nan Rosemary Ke. Learning
400 neural causal models with active interventions, 2022.
- 401 [18] Sébastien Lachapelle, Philippe Brouillard, Tristan Deleu, and Simon Lacoste-Julien. Gradient-
402 based neural dag learning. In *International Conference on Learning Representations*, 2020.
403 URL <https://openreview.net/forum?id=rk1bKA4YDS>.
- 404 [19] Ignavier Ng, AmirEmad Ghassami, and Kun Zhang. On the role of sparsity and dag constraints
405 for learning linear dags. In *Proceedings of the 34th International Conference on Neural*
406 *Information Processing Systems*, NIPS ’20, Red Hook, NY, USA, 2020. Curran Associates Inc.
407 ISBN 9781713829546.
- 408 [20] Xun Zheng, Bryon Aragam, Pradeep Ravikumar, and Eric P. Xing. Dags with no tears:
409 Continuous optimization for structure learning. In *Neural Information Processing Systems*,
410 2018. URL <https://api.semanticscholar.org/CorpusID:53217974>.
- 411 [21] Zhen Zhang, Ignavier Ng, Dong Gong, Yuhang Liu, Ehsan M Abbasnejad, Mingming Gong,
412 Kun Zhang, and Javen Qinfeng Shi. Truncated matrix power iteration for differentiable DAG
413 learning. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors,
414 *Advances in Neural Information Processing Systems*, 2022. URL [https://openreview.net/](https://openreview.net/forum?id=I4aSjFR7j0m)
415 [forum?id=I4aSjFR7j0m](https://openreview.net/forum?id=I4aSjFR7j0m).
- 416 [22] Kevin Bello, Bryon Aragam, and Pradeep Kumar Ravikumar. DAGMA: Learning DAGs via
417 m-matrices and a log-determinant acyclicity characterization. In Alice H. Oh, Alekh Agarwal,
418 Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing*
419 *Systems*, 2022. URL <https://openreview.net/forum?id=8rZYMpFUGk>.
- 420 [23] Jonas Peters, Joris M. Mooij, Dominik Janzing, and Bernhard Schölkopf. Causal discovery
421 with continuous additive noise models. *J. Mach. Learn. Res.*, 15(1):2009–2053, jan 2014. ISSN
422 1532-4435.
- 423 [24] Sindy Löwe, David Madras, Richard S. Zemel, and Max Welling. Amortized causal discovery:
424 Learning to infer causal graphs from time-series data. In *CLEaR*, 2020. URL [https://api.](https://api.semanticscholar.org/CorpusID:219955853)
425 [semanticscholar.org/CorpusID:219955853](https://api.semanticscholar.org/CorpusID:219955853).
- 426 [25] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N
427 Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon,
428 U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, ed-
429 itors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates,
430 Inc., 2017. URL [https://proceedings.neurips.cc/paper_files/paper/2017/file/](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf)
431 [3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf).
- 432 [26] Zhijing Jin, Yuen Chen, Felix Leeb, Luigi Gresele, Ojasv Kamal, Zhiheng Lyu, Kevin Blin,
433 Fernando Gonzalez Adauto, Max Kleiman-Weiner, Mrinmaya Sachan, et al. Cladder: A
434 benchmark to assess causal reasoning capabilities of language models. *Advances in Neural*
435 *Information Processing Systems*, 36, 2024.
- 436 [27] Jiaqi Zhang, Joel Jennings, Agrin Hilmkil, Nick Pawlowski, Cheng Zhang, and Chao Ma.
437 Towards causal foundation model: on duality between causal inference and attention, 2024.

- 438 [28] Meyer Scetbon, Joel Jennings, Agrin Hilmkil, Cheng Zhang, and Chao Ma. Fip: a fixed-point
439 approach for causal generative modeling, 2024.
- 440 [29] Shantanu Gupta, Cheng Zhang, and Agrin Hilmkil. Learned causal method prediction, 2023.
- 441 [30] Caroline Uhler, G. Raskutti, Peter Bühlmann, and B. Yu. Geometry of the faithfulness assump-
442 tion in causal inference. *The Annals of Statistics*, 41, 07 2012. doi: 10.1214/12-AOS1080.
- 443 [31] Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel,
444 Matthias Bethge, and Felix Wichmann. Shortcut learning in deep neural networks. *Nature*
445 *Machine Intelligence*, 2:665–673, 11 2020. doi: 10.1038/s42256-020-00257-z.
- 446 [32] Alexander G. Reisach, Christof Seiler, and Sebastian Weichwald. Beware of the simulated dag!
447 causal discovery benchmarks may be easy to game. In *Neural Information Processing Systems*,
448 2021. URL <https://api.semanticscholar.org/CorpusID:239998404>.
- 449 [33] Francesco Montagna, Nicoletta Noceti, Lorenzo Rosasco, and Francesco Locatello. Shortcuts
450 for causal discovery of nonlinear models by score matching, 2023.
- 451 [34] Shohei Shimizu, Takanori Inazumi, Yasuhiro Sogawa, Aapo Hyvarinen, Yoshinobu Kawahara,
452 Takashi Washio, Patrik Hoyer, and Kenneth Bollen. DirectLiNGAM: A direct method for
453 learning a linear non-gaussian structural equation model. *Journal of Machine Learning Research*,
454 12, 01 2011.
- 455 [35] Francesco Montagna, Nicoletta Noceti, Lorenzo Rosasco, Kun Zhang, and Francesco Locatello.
456 Causal discovery with score matching on additive models with arbitrary noise. In *2nd Conference*
457 *on Causal Learning and Reasoning*, 2023. URL <https://openreview.net/forum?id=rV00Bx90deu>.
- 459 [36] Peter Bühlmann, Jonas Peters, and Jan Ernest. CAM: Causal additive models, high-dimensional
460 order search and penalized regression. *The Annals of Statistics*, 42(6), dec 2014. URL
461 <https://doi.org/10.1214/14-aos1260>.
- 462 [37] Jannik Kossen, Neil Band, Clare Lyle, Aidan Gomez, Tom Rainforth, and Yarin Gal. Self-
463 attention between datapoints: Going beyond individual input-output pairs in deep learning.
464 In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors, *Advances in*
465 *Neural Information Processing Systems*, 2021. URL <https://openreview.net/forum?id=wRxz0a2z5T>.
- 467 [38] Juan Lin. Factorizing multivariate function classes. In M. Jordan, M. Kearns, and
468 S. Solla, editors, *Advances in Neural Information Processing Systems*, volume 10. MIT
469 Press, 1997. URL [https://proceedings.neurips.cc/paper_files/paper/1997/
470 file/8fb21ee7a2207526da55a679f0332de2-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/1997/file/8fb21ee7a2207526da55a679f0332de2-Paper.pdf).

471 A Learning to induce: causal discovery with transformers

472 A.1 A supervised learning approach to causal discovery

473 First, we describe the training procedure for the CSIvA architecture, which aims to learn the dis-
474 tribution of causal graphs conditioned on observational and/or interventional datasets. We omit
475 interventional datasets from the discussion as they are not of interest to our work. Training data are
476 generated from the joint distribution $p_{\mathcal{G}, \mathcal{D}}$ between a graph \mathcal{G} and a dataset \mathcal{D} . First, we sample a set
477 of directed acyclic graphs $\{\mathcal{G}^i\}_{i=1}^n$ with nodes X_1, \dots, X_d , from a distribution $p_{\mathcal{G}}$. Then, for each
478 graph we sample a dataset of m observations of the graph nodes $\mathcal{D}^i = \{x_1^j, \dots, x_d^j\}_{j=1}^m$, $i = 1, \dots, n$.
479 Hence, we build a training dataset $\{\mathcal{G}^i, \mathcal{D}^i\}_{i=1}^n$.

The CSIvA model defines a distribution $\hat{p}_{\mathcal{G}|\mathcal{D}}(\cdot; \Theta)$ of graphs conditioned on the observational data
and parametrized by Θ . Given an invertible map $\mathcal{G} \mapsto A$ from a graph to its binary adjacency matrix

representation of $d \times d$ entries (where $A_{ij} = 1$ iff $X_i \rightarrow X_j$ in \mathcal{G}), we consider an equivalent estimated distribution $\hat{p}_{A|\mathcal{D}}(\cdot; \Theta)$, which has the following autoregressive form:

$$\hat{p}_{A,\mathcal{D}}(A|\mathcal{D}; \Theta) = \prod_{l=1}^{d^2} \sigma(A_l; \rho = f_{\Theta}(A_1, \dots, A_{l-1}, \mathcal{D})),$$

where $\sigma(\cdot; \rho)$ is a Bernoulli distribution parametrized by ρ . ρ itself is a function of f_{Θ} defined by the encoder-decoder transformer architecture, taking as input previous elements of the matrix A (here represented as a vector of d^2 entries) and the dataset \mathcal{D} . Θ is optimized via maximum likelihood estimation, i.e. $\Theta^* = \operatorname{argmin}_{\Theta} -\mathbf{E}_{\mathcal{G},\mathcal{D}}[\ln \hat{p}(\mathcal{G}|\mathcal{D}; \Theta)]$, which corresponds to the usual cross-entropy loss for the Bernoulli distribution. Training is achieved using stochastic gradient descent, in which each gradient update is performed using a pair (\mathcal{D}^i, A^i) , $i = 1 \dots, d$. In the infinite sample limit, we have $\hat{p}_{\mathcal{G}|\mathcal{D}}(\cdot; \Theta^*) = p_{\mathcal{G}|\mathcal{D}}(\cdot)$, while in the finite-capacity case, it is only an approximation of the target distribution.

A.2 CSIVa architecture

In this section, we summarize the architecture of CSIVa, a transformer neural network that can learn a map from data to causally interpreted graphs, under supervised training.

Transformer neural network. Transformers [25] are a popular neural network architecture for modeling structured, sequential data. They consist of an *encoder*, a stack of layers that learns a representation of each element in the input sequence based on its relation with all the other sequence’s elements, through the mechanism of self-attention, and a decoder, which maps the learned representation to the target of interest. Note that data for causal discovery are not sequential in their nature, which motivates the adaptations introduced by Ke et al. [1] in their CSIVa architecture.

CSIVa embeddings. Each element x_i^j of an input dataset is embedded into a vector of dimensionality E . Half of this vector is allocated to embed the value x_i^j itself, while the other half is allocated to embed the unique identity for the node X_i . We use a node-specific embedding because the values of each node may have very different interpretations and meanings. The node identity embedding is obtained using a standard 1D transformer positional embedding over node indices. The value embedding is obtained by passing x_i^j , through a multi-layer perceptron (MLP).

CSIVa alternating attention. Similarly to the transformer’s encoder, CSIVa stacks a number of identical layers, performing self-attention followed by a nonlinear mapping, most commonly an MLP layer. The main difference relative to the standard encoder is in the implementation of the self-attention layer: as transformers are in their nature suitable for the representation of sequences, given an input sample of D elements, self-attention is usually run across all elements of the sequence. However, data for causal discovery are tabular, rather than sequential: one option would be to unravel the $n \times d$ matrix of the data, where n is the number of observations and d the number of variables, into a vector of $n \cdot d$ elements, and let this be the input sequence of the encoder. CSIVa adopts a different strategy: the self-attention in each encoder layer consists of alternate passes over the attribute and the sample dimensions, known as *alternating attention* [37]. As a clarifying example, consider a dataset $\{(x_1^i, x_2^i)\}_{i=1}^n$ of n i.i.d. samples from the joint distribution of the pair of random variables X_1, X_2 . For each layer of the encoder, in the first step (known as *attention between attributes*), attention operates across all nodes of a single sample (x_1^i, x_2^i) to encode the relationships between the two nodes. In the second step (*attention between samples*), attention operates across all samples (x_k^1, \dots, x_k^n) , $k \in \{1, 2\}$ of a given node, to encode information about the distribution of single node values.

CSIVa encoder summary. The encoder produces a summary vector s_i with H elements for each node X_i , which captures essential information about the node’s behavior and its interactions with other nodes. The summary representation is formed independently for each node and involves combining information across the n samples. This is achieved with a method often used with transformers that involves a weighted average based on how informative each sample is. The weighting is obtained using the embeddings of a summary "sample" $n + 1$ to form queries, and embeddings of node’s samples $\{x_i^j\}_{j=1}^n$ to provide keys and values, and then using standard key-value attention.

Hyperparameter	Value
Hidden state dimension	64
Encoder transformer layers	8
Decoder transformer layers	8
Num. attention heads	8
Optimizer	Adam
Learning rate	10^{-4}
Samples per dataset (n)	1500
Num. training datasets	15000
Num. iterations	< 150000
Batch size	5

Table 1: Hyperparameters for the training of the CSIvA models of the experiments in Section 3.

526 **CSIvA decoder.** The decoder uses the summary information from the encoder to generate a
527 prediction of the adjacency matrix A of the underlying \mathcal{G} . It operates sequentially, at each step
528 producing a binary output indicating the prediction $\hat{A}_{i,j}$ of $A_{i,j}$, proceeding row by row. The decoder
529 is an autoregressive transformer, meaning that each prediction $\hat{A}_{i,j}$ is obtained based on all elements
530 of A previously predicted, as well as the summary produced by the encoder. The method does not
531 enforce acyclicity, although Ke et al. [1] shows that in cyclic outputs generally don’t occur, in
532 practice.

533 B Training details

534 B.1 Hyperparameters

535 In Table 1 we detail the hyperparameters of the training of the network of the experiments. We define
536 an iteration as a gradient update over a batch of 5 datasets. Models are trained until convergence,
537 using a patience of 5 (training until five consecutive epochs without improvement) on the validation
538 loss - this always occurs before the 25-th epoch (corresponding to ≈ 150000 iterations). The batch
539 size is limited to 5 due to memory constraints.

540 B.2 Synthetic data

541 In this section, we provide additional details on the synthetic data generation, which was performed
542 with the `causally`² Python library [9]. Our data-generating framework follows that of Montagna
543 et al. [9], an extensive benchmark of causal discovery methods on different classes of SCMs.

544 **Causal mechanisms.** The *nonlinear mechanisms* of the PNL model and the nonlinear ANM model
545 are generated by a neural network with one hidden layer with 10 hidden units, with a parametric
546 ReLU activation function. The network weights are randomly sampled according to a standard
547 Gaussian distribution. The *linear mechanisms* are generated by sampling the regression coefficients
548 in the range $[-3, -0.5] \cup [0.5, 3]$.

549 **Distribution of the noise terms.** We generated datasets from structural causal models with the
550 following distribution of the noise terms: Beta, Gamma, Gaussian (for nonlinear data), Gumbel,
551 Exponential, and Uniform. Additionally, we define the *mlp* distribution by nonlinear transformations
552 of gaussian samples from a gaussian distribution centered at zero and with standard deviation σ
553 uniformly sampled between 0.5 and 1. The nonlinear transformation is parametrized by a neural
554 network with one hidden layer with 100 units, and sigmoid activation function. The weights of the
555 network are uniformly sampled in the range $[-1.5, 1.5]$. We additionally standardized the output of
556 each *mlp* sample by the empirical variance computed over all samples.

557 Data are standardized with their empirical variance, which removes the presence of shortcuts which
558 could be learned by the network, notably *varsortability* [32] and *score-sortability* [33].

²<https://causally.readthedocs.io/en/latest/>

559 **B.3 Computer resources**

560 Our experiments were run on a local computing cluster, using any and all available GPUs (all
 561 NVIDIA). For replication purposes, GTX 1080 Ti's are entirely suitable, as the batch size was set
 562 to match their memory capacity, when working with bivariate graphs. All jobs ran with 10GB of
 563 RAM and 4 CPU cores. The results presented in this paper were produced after 145 days of GPU
 564 time, of which 68 were on GTX 1080 Ti's, 13 on RTX 2080 Ti's, 11 on A10s, 19 on A40s, and 35
 565 on RTX 3090s. Together with previous experiments, while developing our code and experimental
 566 design, we used 376 days of GPU time (for reference, at a total cost of 492.14 Euros), similarly split
 567 across whichever GPUs were available at the time: 219 on GTX 1080 Ti's, 38 on RTX 2080 Ti's, 18
 568 on A10s, 63 on RTX 3090s, 31 on A40s, and 6 on A100s.

569 **C Further experiments**

570 We present our experimental results on one further question, to help clarify the results in the main text
 571 of the paper. Our aim is to understand when to make tradeoffs between computational resources, and
 572 having models that have been trained on a wider variety of SCMs. We compare training on multiple
 573 SCMs to single-SCM training, when all models see the same amount of training data from each SCM
 574 type as a non-mixed model (i.e. a mixed network trains on 15,000 linear datasets and 15,000 PNL
 575 datasets, instead of 15,000 divided between the two SCM types).

576 In the main text of this paper, we compare neural networks trained on a mix of structural causal
 577 models (e.g. noise distributions, or mechanism types), to models trained on a single mechanism-noise
 578 combination, where all models have the same amount of training data, 15,000 datasets. In mixed
 579 training, we split these evenly, so a "lin, nl" model is trained on 7,500 datasets from linear SCMs, and
 580 7,500 from nonlinear SCMs. Our results in this framework are promising, and show that for many
 581 combinations of SCM types, we can train one model instead of two, and achieve good progress, while
 582 making a 50% savings on training costs. However, if our training budget is high/unlimited, we should
 583 also ask whether there is a downside to mixed training - can we achieve the same performance as a
 584 model trained on a single SCM type? Fig. 6 shows good results in this direction - the models trained
 585 with the same number of datasets per SCM type as an unmixed model had similar (or even better,
 586 for PNL data) performance as the un-mixed model trained on the same SCM type as the test data.
 587 These mixed models are also significantly more useful than having 2 or 3 separate models per SCM
 588 type, as they have good across-the-board performance. However, if we used the same computational
 589 resources to train 3 separate networks (one for each mechanism type) and wanted to use them for
 590 causal discovery on a dataset with unknown assumptions, we would be left with the rather difficult
 task of deciding which model to trust.

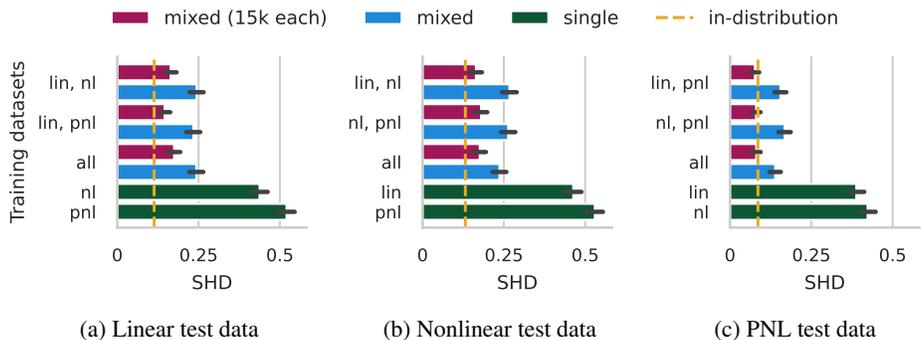


Figure 6: Mixtures of causal mechanisms, with varying amounts of training data. We train eight models on samples from structural casual models with different mechanisms. Four (in purple), were trained on 15,000 samples for each SCM type (so the "lin,nl" model saw 30,000 samples in total, and the "all" model saw 45,000), and the other four (blue) are the same as in Fig. 4, and were trained on 15,000 samples in total, evenly split between the SCM types they were trained on. We compare their test SHD (the lower, the better) against networks trained on datasets generated according to a single type of mechanism. The dashed line indicates the test SHD of a model trained on samples with the same causal mechanisms as the test SCM. Training on multiple causal models with different mechanisms (mixed bars) always improves performance compared to training on single SCMs.

592 **D Theoretical results and proofs**

593 Before stating the proof of Proposition 1, we show under which condition the pair of random
 594 variables X, Y satisfies the forward and backward models of equations (5), (6): this is relevant for
 595 our discussion, as the proof of Proposition 1 consists of showing that this condition is *almost* never
 596 satisfied.

597 **Notation.** We adopt the following notation: $\nu_X := \log p_{N_X}$, $\nu_Y := \log p_{N_Y}$, $\xi := \log p_{\tilde{X}}$, $\eta :=$
 598 $\log p_{\tilde{Y}}$, and $\pi := \log p_{\tilde{X}, \tilde{Y}}$.

599 **Theorem 1** (Theorem 1 of Zhang and Hyvärinen [8]). *Assume that X, Y satisfies both causal*
 600 *relations of equations (5) and (6). Further, suppose that p_{N_Y} and $p_{\tilde{X}}$ are positive densities on the*
 601 *support of N_Y and \tilde{X} respectively, and that $\nu_Y, \xi, f_1, f_2, g_1$, and g_2 are third order differentiable.*
 602 *Then, for each pair (\tilde{x}, \tilde{y}) satisfying $\nu_Y''(\tilde{y} - h_Y(\tilde{x}))h_Y'(\tilde{x}) \neq 0$, the following differential equation*
 603 *holds:*

$$\xi''' = \xi'' \left(\frac{h_Y''}{h_Y'} - \frac{\nu_Y'' h_Y'}{\nu_Y''} \right) + \frac{\nu_Y''' \nu_Y' h_Y'' h_Y'}{\nu_Y''} - \frac{\nu_Y' (h_Y'')^2}{h_Y'} - 2\nu_Y'' h_Y'' h_Y' + \nu_Y' h_Y''',$$

604 and h_X is constrained in the following way:

$$\frac{1}{h_X'} = \frac{\xi'' + \nu_Y'' (h_Y')^2 - \nu_Y' h_Y''}{\nu_Y'' h_Y'}, \quad (9)$$

605 where the arguments of the functions have been left out for clarity.

606 *Proof of Theorem 1.* We demonstrate separately the two statements of the theorem.

607 **Part 1.** Given that equations (5) and (6) hold, this implies that the forward and backward models
 608 on \tilde{X}, \tilde{Y} of equations (7) and (8) are also valid, namely that:

$$\begin{aligned} \tilde{Y} &= h_Y(\tilde{X}) + N_Y, \\ \tilde{X} &= h_X(\tilde{Y}) + N_X. \end{aligned}$$

These are the structural equations of two causal models, associated with the *forward* $\tilde{X} \rightarrow \tilde{Y}$ and
backward $\tilde{Y} \rightarrow \tilde{X}$ graphs, respectively. Applying the Markov factorization of the distribution
 according to the forward direction, we get:

$$p_{\tilde{X}, \tilde{Y}}(\tilde{x}, \tilde{y}) = p_{\tilde{Y}|\tilde{X}}(\tilde{y}|\tilde{x})p_{\tilde{X}}(\tilde{x}) = p_{N_Y}(\tilde{y} - h_Y(\tilde{x}))p_{\tilde{X}}(\tilde{x}),$$

609 which implies

$$\pi(\tilde{x}, \tilde{y}) = \nu_Y(\tilde{y} - h_Y(\tilde{x})) + \xi(\tilde{x}), \quad (10)$$

610 for any \tilde{x}, \tilde{y} . Similarly, the Markov factorization on the backward model implies:

$$\pi(\tilde{x}, \tilde{y}) = \nu_X(\tilde{x} - h_X(\tilde{y})) + \eta(\tilde{y}). \quad (11)$$

611 From (11), we have that:

$$\begin{aligned} \frac{\partial^2}{\partial \tilde{x}^2} \pi(\tilde{x}, \tilde{y}) &= \nu_X''(\tilde{x} - h_X(\tilde{y})) \\ \frac{\partial^2}{\partial \tilde{x} \partial \tilde{y}} \pi(\tilde{x}, \tilde{y}) &= -\nu_X''(\tilde{x} - h_X(\tilde{y}))h_X'(\tilde{y}), \end{aligned}$$

612 which implies

$$\frac{\partial}{\partial \tilde{x}} \left(\frac{\frac{\partial^2}{\partial \tilde{x}^2} \pi(\tilde{x}, \tilde{y})}{\frac{\partial^2}{\partial \tilde{x} \partial \tilde{y}} \pi(\tilde{x}, \tilde{y})} \right) = 0. \quad (12)$$

613 Computing the same set of partial derivatives from (10), we find:

$$\begin{aligned} \frac{\partial^2}{\partial \tilde{x}^2} \pi(\tilde{x}, \tilde{y}) &= \nu_Y''(\tilde{y} - h_Y(\tilde{x}))(h_Y'(\tilde{x}))^2 - \nu_Y'(\tilde{y} - h_Y(\tilde{x}))h_Y''(\tilde{x}) + \xi''(\tilde{x}) \\ \frac{\partial^2}{\partial \tilde{x} \partial \tilde{y}} \pi(\tilde{x}, \tilde{y}) &= -\nu_Y''(\tilde{y} - h_Y(\tilde{x}))h_Y'(\tilde{x}). \end{aligned}$$

614 from which follows:

$$\begin{aligned} \frac{\partial}{\partial \tilde{x}} \left(\frac{\frac{\partial^2}{\partial \tilde{x}^2} \pi(\tilde{x}, \tilde{y})}{\frac{\partial^2}{\partial \tilde{x} \partial \tilde{y}} \pi(\tilde{x}, \tilde{y})} \right) &= -2h_Y'' + \frac{\nu_Y' h_Y'''}{\nu_Y'' h_Y'} - \frac{\xi'''}{\nu_Y'' h_Y'} + \frac{\nu_Y''' \nu_Y' h_Y''}{(\nu_Y'')^2} - \frac{\nu_Y' (h_Y'')^2}{\nu_Y'' (h_Y')^2} + \frac{\xi'' \nu_Y''' h_Y''}{(\nu_Y'')^2 \nu_Y'' (h_Y')^2} \\ &= 0. \end{aligned}$$

615 where we drop the input arguments for conciseness. The equality with 0 is given by the equality with
616 (12). Manipulating the above expression, the first claim follows.

617 **Part 2.** Next, we prove the constraint derived on h_X . To do this, we exploit the fact that \tilde{Y} is
618 independent of N_X , which implies the following condition [38]:

$$\frac{\partial^2}{\partial \tilde{y} \partial n_x} \log p(\tilde{y}, n_x) = 0, \quad (13)$$

619 for any (\tilde{y}, n_x) . According to equations (7), (8), we have that:

$$\begin{aligned} \tilde{Y} &= h_Y(\tilde{X}) + N_Y, \\ N_X &= \tilde{X} - h_X(\tilde{Y}), \end{aligned}$$

such that we can define an invertible map $\Phi : (\tilde{y}, n_x) \mapsto (\tilde{x}, n_Y)$. It is easy to show that the Jacobian of the transformation has determinant $|J_\Phi| = 1$, such that

$$p(\tilde{y}, n_Y) = p(\tilde{x}, n_Y),$$

where $(\tilde{x}, n_Y) = \Phi^{-1}(\tilde{y}, n_x)$. Thus, being \tilde{X}, N_Y independent random variables, we have that:

$$\log p(\tilde{y}, n_x) = \log p(\tilde{x}) + \log p(n_Y) = \xi(\tilde{x}) + \nu_Y(n_Y).$$

Given that $\tilde{X} = h_X(\tilde{Y}) + N_X$, we have that

$$\frac{\partial^2}{\partial \tilde{y} \partial \tilde{n}_X} \log p(\tilde{x}) = \xi'' h_X',$$

while $N_Y = \tilde{Y} - h_Y(\tilde{X})$ implies

$$\frac{\partial^2}{\partial \tilde{y} \partial \tilde{n}_X} \log p(n_Y) = -\nu_Y'' h_Y' + \nu_Y'' h_X' (h_Y')^2 - \nu_Y' h_X' h_Y'',$$

such that

$$\log p(\tilde{x}, n_Y) = \xi'' h_X' + -\nu_Y'' h_Y' + \nu_Y'' h_X' (h_Y')^2 - \nu_Y' h_X' h_Y'',$$

which must be equal to zero, being equal to the LHS of (13). Thus, we conclude that

$$\frac{1}{h_X'} = \frac{\xi'' + \nu_Y'' (h_Y')^2 - \nu_Y' h_Y''}{\nu_Y'' h_Y'},$$

620 proving the claim. □

621 D.1 Proof of Proposition 1

622 *Proof.* Under the hypothesis that equations (5), (6) hold, i.e. when the data generating process satisfy
623 both a forward and a backward model, by Theorem 1 we have that:

$$\xi'''(\tilde{x}) = \xi''(\tilde{x})G(\tilde{x}, \tilde{y}) + H(\tilde{x}, \tilde{y}), \quad (14)$$

624 where

$$\begin{aligned} G(\tilde{x}, \tilde{y}) &= \left(\frac{h_Y''}{h_Y'} - \frac{\nu_Y''' h_Y'}{\nu_Y''} \right), \\ H(\tilde{x}, \tilde{y}) &= \frac{\nu_Y''' \nu_Y' h_Y'' h_Y'}{\nu_Y''} - \frac{\nu_Y' (h_Y'')^2}{h_Y'} - 2\nu_Y'' h_Y'' h_Y' + \nu_Y' h_Y'''. \end{aligned}$$

625 Define $z := \xi'''$, such that the above equation can be written as $z'(\tilde{x}) = z(\tilde{x})G(\tilde{x}, \tilde{y}) + H(\tilde{x}, \tilde{y})$.
 626 given that such function z exists, it is given by:

$$z(\tilde{x}) = z(\tilde{x}_0)e^{\int_{\tilde{x}_0}^{\tilde{x}} G(t,y)dt} + \int_{\tilde{x}_0}^{\tilde{x}} e^{\int_t^{\tilde{x}} G(t,y)dt} H(\hat{t}, y)d\hat{t}. \quad (15)$$

Let \tilde{y} such that $\nu_Y''(\tilde{y} - h_Y(\tilde{x}))h_Y'(\tilde{x}) \neq 0$ holds for all but countable values of \tilde{x} . Then, z is determined by $z(\tilde{x}_0)$, as we can extend equation (15) to all the remaining points. The set of all functions ξ satisfying the differential equation (14) is a 3-dimensional affine space, as fixing $\xi(\tilde{x}_0), \xi''(\tilde{x}_0), \xi'''(\tilde{x}_0)$ for some point \tilde{x}_0 completely determines the solution ξ . Moreover, given ν_Y, h_X, h_Y fixed, ξ'' is specified by (9) of theorem 1, which implies:

$$\xi'' = \frac{\nu_Y'' h_Y'}{h_X'} + \nu_Y' h_Y'' - \nu_Y'' (h_Y')^2,$$

627 which confines ξ solutions of (14) to a 2-dimensional affine space. □

628 **NeurIPS Paper Checklist**

629 **1. Claims**

630 Question: Do the main claims made in the abstract and introduction accurately reflect the
631 paper's contributions and scope?

632 Answer: [\[Yes\]](#)

633 Justification: Supervised learning models in causal discovery do not provide connections
634 with the known identifiability theory. In the abstract, we present this open problem, and
635 highlight our main empirical findings and how they connect to the theory of identifiability in
636 causality. The content of the paper (mostly Section 3) unravels the abstract claims in all of
637 their details.

638 Guidelines:

- 639 • The answer NA means that the abstract and introduction do not include the claims
640 made in the paper.
- 641 • The abstract and/or introduction should clearly state the claims made, including the
642 contributions made in the paper and important assumptions and limitations. A No or
643 NA answer to this question will not be perceived well by the reviewers.
- 644 • The claims made should match theoretical and experimental results, and reflect how
645 much the results can be expected to generalize to other settings.
- 646 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
647 are not attained by the paper.

648 **2. Limitations**

649 Question: Does the paper discuss the limitations of the work performed by the authors?

650 Answer: [\[Yes\]](#)

651 Justification: We discuss the limitations of our work in Section 1, paragraph "Closely related
652 works and their relation with CSIVa", regarding the use of CSIVa as our only architecture
653 for the experiments. Additionally, in the same paragraph, we remark that the scope of this
654 study is limited to the context of causal discovery on observational data. Finally, in Section
655 2.2, we discuss our choice of limiting the empirical study to the case of bivariate graphs.

656 Guidelines:

- 657 • The answer NA means that the paper has no limitation while the answer No means that
658 the paper has limitations, but those are not discussed in the paper.
- 659 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 660 • The paper should point out any strong assumptions and how robust the results are to
661 violations of these assumptions (e.g., independence assumptions, noiseless settings,
662 model well-specification, asymptotic approximations only holding locally). The authors
663 should reflect on how these assumptions might be violated in practice and what the
664 implications would be.
- 665 • The authors should reflect on the scope of the claims made, e.g., if the approach was
666 only tested on a few datasets or with a few runs. In general, empirical results often
667 depend on implicit assumptions, which should be articulated.
- 668 • The authors should reflect on the factors that influence the performance of the approach.
669 For example, a facial recognition algorithm may perform poorly when image resolution
670 is low or images are taken in low lighting. Or a speech-to-text system might not be
671 used reliably to provide closed captions for online lectures because it fails to handle
672 technical jargon.
- 673 • The authors should discuss the computational efficiency of the proposed algorithms
674 and how they scale with dataset size.
- 675 • If applicable, the authors should discuss possible limitations of their approach to
676 address problems of privacy and fairness.
- 677 • While the authors might fear that complete honesty about limitations might be used by
678 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
679 limitations that aren't acknowledged in the paper. The authors should use their best

680 judgment and recognize that individual actions in favor of transparency play an impor-
681 tant role in developing norms that preserve the integrity of the community. Reviewers
682 will be specifically instructed to not penalize honesty concerning limitations.

683 3. Theory Assumptions and Proofs

684 Question: For each theoretical result, does the paper provide the full set of assumptions and
685 a complete (and correct) proof?

686 Answer: [Yes]

687 Justification: Proposition 1 is proved in detail in Appendix D.1, which is based on Theorem
688 1 of Zhang and Hyvärinen [8], which we report in the Appendix together with its proof. We
689 do not provide an explicit sketch of the proof of our Proposition 1 in the main text, as we
690 already detail the intuition behind it in the content of Section 3.3.

691 Guidelines:

- 692 • The answer NA means that the paper does not include theoretical results.
- 693 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
694 referenced.
- 695 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 696 • The proofs can either appear in the main paper or the supplemental material, but if
697 they appear in the supplemental material, the authors are encouraged to provide a short
698 proof sketch to provide intuition.
- 699 • Inversely, any informal proof provided in the core of the paper should be complemented
700 by formal proofs provided in appendix or supplemental material.
- 701 • Theorems and Lemmas that the proof relies upon should be properly referenced.

702 4. Experimental Result Reproducibility

703 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
704 perimental results of the paper to the extent that it affects the main claims and/or conclusions
705 of the paper (regardless of whether the code and data are provided or not)?

706 Answer: [Yes]

707 Justification: We have specified our data generation methods in Appendix B.2, as well
708 as the CSIVa method (which is a previously published model) in Appendix A, and our
709 hyperparameters for training in Appendix B.1. We will also release our implementation of
710 CSIVa, our data generation code (which is a thin wrapper around the `causally` <https://causally.readthedocs.io/en/latest/> Python library), and our experimental code.

712 Guidelines:

- 713 • The answer NA means that the paper does not include experiments.
- 714 • If the paper includes experiments, a No answer to this question will not be perceived
715 well by the reviewers: Making the paper reproducible is important, regardless of
716 whether the code and data are provided or not.
- 717 • If the contribution is a dataset and/or model, the authors should describe the steps taken
718 to make their results reproducible or verifiable.
- 719 • Depending on the contribution, reproducibility can be accomplished in various ways.
720 For example, if the contribution is a novel architecture, describing the architecture fully
721 might suffice, or if the contribution is a specific model and empirical evaluation, it may
722 be necessary to either make it possible for others to replicate the model with the same
723 dataset, or provide access to the model. In general, releasing code and data is often
724 one good way to accomplish this, but reproducibility can also be provided via detailed
725 instructions for how to replicate the results, access to a hosted model (e.g., in the case
726 of a large language model), releasing of a model checkpoint, or other means that are
727 appropriate to the research performed.
- 728 • While NeurIPS does not require releasing code, the conference does require all submis-
729 sions to provide some reasonable avenue for reproducibility, which may depend on the
730 nature of the contribution. For example
731 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
732 to reproduce that algorithm.

- 733 (b) If the contribution is primarily a new model architecture, the paper should describe
734 the architecture clearly and fully.
- 735 (c) If the contribution is a new model (e.g., a large language model), then there should
736 either be a way to access this model for reproducing the results or a way to reproduce
737 the model (e.g., with an open-source dataset or instructions for how to construct
738 the dataset).
- 739 (d) We recognize that reproducibility may be tricky in some cases, in which case
740 authors are welcome to describe the particular way they provide for reproducibility.
741 In the case of closed-source models, it may be that access to the model is limited in
742 some way (e.g., to registered users), but it should be possible for other researchers
743 to have some path to reproducing or verifying the results.

744 5. Open access to data and code

745 Question: Does the paper provide open access to the data and code, with sufficient instruc-
746 tions to faithfully reproduce the main experimental results, as described in supplemental
747 material?

748 Answer: [Yes]

749 Justification: We will release our implementation of CSIVa, our data generation code
750 (which is a thin wrapper around the causally [https://causally.readthedocs.io/](https://causally.readthedocs.io/en/latest/)
751 [en/latest/](https://causally.readthedocs.io/en/latest/) Python library), and our experimental code.

752 Guidelines:

- 753 • The answer NA means that paper does not include experiments requiring code.
- 754 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/](https://nips.cc/public/guides/CodeSubmissionPolicy)
755 [public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 756 • While we encourage the release of code and data, we understand that this might not be
757 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
758 including code, unless this is central to the contribution (e.g., for a new open-source
759 benchmark).
- 760 • The instructions should contain the exact command and environment needed to run to
761 reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 762 • The authors should provide instructions on data access and preparation, including how
763 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 764 • The authors should provide scripts to reproduce all experimental results for the new
765 proposed method and baselines. If only a subset of experiments are reproducible, they
766 should state which ones are omitted from the script and why.
- 767 • At submission time, to preserve anonymity, the authors should release anonymized
768 versions (if applicable).
- 769 • Providing as much information as possible in supplemental material (appended to the
770 paper) is recommended, but including URLs to data and code is permitted.

772 6. Experimental Setting/Details

773 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
774 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
775 results?

776 Answer: [Yes]

777 Justification: Yes, we provide these details in Section 3.1 and Appendix B.

778 Guidelines:

- 779 • The answer NA means that the paper does not include experiments.
- 780 • The experimental setting should be presented in the core of the paper to a level of detail
781 that is necessary to appreciate the results and make sense of them.
- 782 • The full details can be provided either with the code, in appendix, or as supplemental
783 material.

784 7. Experiment Statistical Significance

785 Question: Does the paper report error bars suitably and correctly defined or other appropriate
786 information about the statistical significance of the experiments?

787 Answer: [Yes]

788 Justification: For each plot, we provide error bars in the form of 95% confidence intervals
789 computed on 1.5k points (hence, it's reasonable to apply the central limit theorem to argue
790 that the confidence intervals are valid).

791 Guidelines:

- 792 • The answer NA means that the paper does not include experiments.
- 793 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
794 dence intervals, or statistical significance tests, at least for the experiments that support
795 the main claims of the paper.
- 796 • The factors of variability that the error bars are capturing should be clearly stated (for
797 example, train/test split, initialization, random drawing of some parameter, or overall
798 run with given experimental conditions).
- 799 • The method for calculating the error bars should be explained (closed form formula,
800 call to a library function, bootstrap, etc.)
- 801 • The assumptions made should be given (e.g., Normally distributed errors).
- 802 • It should be clear whether the error bar is the standard deviation or the standard error
803 of the mean.
- 804 • It is OK to report 1-sigma error bars, but one should state it. The authors should
805 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
806 of Normality of errors is not verified.
- 807 • For asymmetric distributions, the authors should be careful not to show in tables or
808 figures symmetric error bars that would yield results that are out of range (e.g. negative
809 error rates).
- 810 • If error bars are reported in tables or plots, The authors should explain in the text how
811 they were calculated and reference the corresponding figures or tables in the text.

812 8. Experiments Compute Resources

813 Question: For each experiment, does the paper provide sufficient information on the com-
814 puter resources (type of compute workers, memory, time of execution) needed to reproduce
815 the experiments?

816 Answer: [Yes]

817 Justification: We provide all details on our computer resources in Appendix B.3.

818 Guidelines:

- 819 • The answer NA means that the paper does not include experiments.
- 820 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
821 or cloud provider, including relevant memory and storage.
- 822 • The paper should provide the amount of compute required for each of the individual
823 experimental runs as well as estimate the total compute.
- 824 • The paper should disclose whether the full research project required more compute
825 than the experiments reported in the paper (e.g., preliminary or failed experiments that
826 didn't make it into the paper).

827 9. Code Of Ethics

828 Question: Does the research conducted in the paper conform, in every respect, with the
829 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

830 Answer: [Yes]

831 Justification: We do not believe any of the concerns in the Code of Ethics apply to our work.

832 Guidelines:

- 833 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 834 • If the authors answer No, they should explain the special circumstances that require a
835 deviation from the Code of Ethics.

- 836 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
837 eration due to laws or regulations in their jurisdiction).

838 10. Broader Impacts

839 Question: Does the paper discuss both potential positive societal impacts and negative
840 societal impacts of the work performed?

841 Answer: [NA]

842 Justification: Our work is about assessing and studying pre-existing causal discovery models.
843 As we release no new model, there is no societal impact that could be caused by our work.

844 Guidelines:

- 845 • The answer NA means that there is no societal impact of the work performed.
- 846 • If the authors answer NA or No, they should explain why their work has no societal
847 impact or why the paper does not address societal impact.
- 848 • Examples of negative societal impacts include potential malicious or unintended uses
849 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
850 (e.g., deployment of technologies that could make decisions that unfairly impact specific
851 groups), privacy considerations, and security considerations.
- 852 • The conference expects that many papers will be foundational research and not tied
853 to particular applications, let alone deployments. However, if there is a direct path to
854 any negative applications, the authors should point it out. For example, it is legitimate
855 to point out that an improvement in the quality of generative models could be used to
856 generate deepfakes for disinformation. On the other hand, it is not needed to point out
857 that a generic algorithm for optimizing neural networks could enable people to train
858 models that generate Deepfakes faster.
- 859 • The authors should consider possible harms that could arise when the technology is
860 being used as intended and functioning correctly, harms that could arise when the
861 technology is being used as intended but gives incorrect results, and harms following
862 from (intentional or unintentional) misuse of the technology.
- 863 • If there are negative societal impacts, the authors could also discuss possible mitigation
864 strategies (e.g., gated release of models, providing defenses in addition to attacks,
865 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
866 feedback over time, improving the efficiency and accessibility of ML).

867 11. Safeguards

868 Question: Does the paper describe safeguards that have been put in place for responsible
869 release of data or models that have a high risk for misuse (e.g., pretrained language models,
870 image generators, or scraped datasets)?

871 Answer: [NA]

872 Justification: The data and models in this paper do not have high risk for misuse.

873 Guidelines:

- 874 • The answer NA means that the paper poses no such risks.
- 875 • Released models that have a high risk for misuse or dual-use should be released with
876 necessary safeguards to allow for controlled use of the model, for example by requiring
877 that users adhere to usage guidelines or restrictions to access the model or implementing
878 safety filters.
- 879 • Datasets that have been scraped from the Internet could pose safety risks. The authors
880 should describe how they avoided releasing unsafe images.
- 881 • We recognize that providing effective safeguards is challenging, and many papers do
882 not require this, but we encourage authors to take this into account and make a best
883 faith effort.

884 12. Licenses for existing assets

885 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
886 the paper, properly credited and are the license and terms of use explicitly mentioned and
887 properly respected?

888 Answer: [Yes]

889 Justification: We cite the authors of all papers we build our work on. Additionally, we
890 provide the URL to all previously existing code we rely on, which is available in the form of
891 public GitHub repository under MIT license.

892 Guidelines:

- 893 • The answer NA means that the paper does not use existing assets.
- 894 • The authors should cite the original paper that produced the code package or dataset.
- 895 • The authors should state which version of the asset is used and, if possible, include a
896 URL.
- 897 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 898 • For scraped data from a particular source (e.g., website), the copyright and terms of
899 service of that source should be provided.
- 900 • If assets are released, the license, copyright information, and terms of use in the
901 package should be provided. For popular datasets, `paperswithcode.com/datasets`
902 has curated licenses for some datasets. Their licensing guide can help determine the
903 license of a dataset.
- 904 • For existing datasets that are re-packaged, both the original license and the license of
905 the derived asset (if it has changed) should be provided.
- 906 • If this information is not available online, the authors are encouraged to reach out to
907 the asset's creators.

908 13. **New Assets**

909 Question: Are new assets introduced in the paper well documented and is the documentation
910 provided alongside the assets?

911 Answer: [\[Yes\]](#)

912 Justification: As our work is an analysis of pre-existing methods of causal discovery, we do
913 not release new assets other than the code strictly needed for reproducing our experimental
914 results. This code is attached to this submission to facilitate the reproducibility of our
915 results. All the documentation necessary for reproducing our results is provided in the main
916 manuscript.

917 Guidelines:

- 918 • The answer NA means that the paper does not release new assets.
- 919 • Researchers should communicate the details of the dataset/code/model as part of their
920 submissions via structured templates. This includes details about training, license,
921 limitations, etc.
- 922 • The paper should discuss whether and how consent was obtained from people whose
923 asset is used.
- 924 • At submission time, remember to anonymize your assets (if applicable). You can either
925 create an anonymized URL or include an anonymized zip file.

926 14. **Crowdsourcing and Research with Human Subjects**

927 Question: For crowdsourcing experiments and research with human subjects, does the paper
928 include the full text of instructions given to participants and screenshots, if applicable, as
929 well as details about compensation (if any)?

930 Answer: [\[NA\]](#)

931 Justification: We do not work with human subjects or crowdsourcing.

932 Guidelines:

- 933 • The answer NA means that the paper does not involve crowdsourcing nor research with
934 human subjects.
- 935 • Including this information in the supplemental material is fine, but if the main contribu-
936 tion of the paper involves human subjects, then as much detail as possible should be
937 included in the main paper.
- 938 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
939 or other labor should be paid at least the minimum wage in the country of the data
940 collector.

941 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human**
942 **Subjects**

943 Question: Does the paper describe potential risks incurred by study participants, whether
944 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
945 approvals (or an equivalent approval/review based on the requirements of your country or
946 institution) were obtained?

947 Answer: [NA]

948 Justification: We do not work with human subjects or crowdsourcing.

949 Guidelines:

- 950 • The answer NA means that the paper does not involve crowdsourcing nor research with
951 human subjects.
- 952 • Depending on the country in which research is conducted, IRB approval (or equivalent)
953 may be required for any human subjects research. If you obtained IRB approval, you
954 should clearly state this in the paper.
- 955 • We recognize that the procedures for this may vary significantly between institutions
956 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
957 guidelines for their institution.
- 958 • For initial submissions, do not include any information that would break anonymity (if
959 applicable), such as the institution conducting the review.