

SINCE FAITHFULNESS FAILS: THE PERFORMANCE LIMITS OF NEURAL CAUSAL DISCOVERY

Anonymous authors

Paper under double-blind review

ABSTRACT

Neural causal discovery methods have recently improved in terms of scalability and computational efficiency. However, there are still opportunities for improving their accuracy in uncovering causal structures. We argue that the key obstacle in unlocking this potential is the *faithfulness assumption*, commonly used by contemporary neural approaches. We show that this assumption, which is often not satisfied in real-world or synthetic datasets, limits the effectiveness of existing methods. We evaluate the impact of faithfulness violations both qualitatively and quantitatively and provide a unified evaluation framework to facilitate further research.

1 INTRODUCTION

Causal discovery is essential to scientific research, driving a growing demand for machine learning methods to support this process. Despite the development of several neural-based causal discovery methods in recent years (Brouillard et al., 2020; Lorch et al., 2021; Annadani et al., 2023; Nazaret et al., 2024), their performance remains insufficient for real-world applications, particularly in fields like medicine and biology (de Castro et al., 2019; Peters et al., 2016). Furthermore, these methods are usually evaluated using synthetic datasets, which vary between studies, obscuring the overall picture and making assessment of advancements difficult.

To address this challenge, we introduce a unified benchmark for evaluating neural causal discovery methods. Specifically, we use identical datasets, tune hyperparameters consistently, and use a standardized functional approximation across all methods. Our systematic evaluation reveals that, although there has been progress in computational efficiency over the past few years, significant gains in causal discovery accuracy have yet to emerge. Further underscoring the challenges, we discover that the existing methods can not take advantage of the increasing amount of data, countering the universally held assumption that more data leads to better learning.

The key claim of this work is that progress in causal discovery requires moving beyond the faithfulness assumption. Although it is widely known that real-world and synthetic data rarely satisfy this assumption (Hoover, 2001; Andersen, 2013), most neural-based methods overlook its impact. We develop techniques to measure how faithfulness violations degrade performance and set an upper bound for current benchmarks. Our results show a clear correlation: faithfulness violations significantly hinder performance, and improvements within the current paradigm are limited.

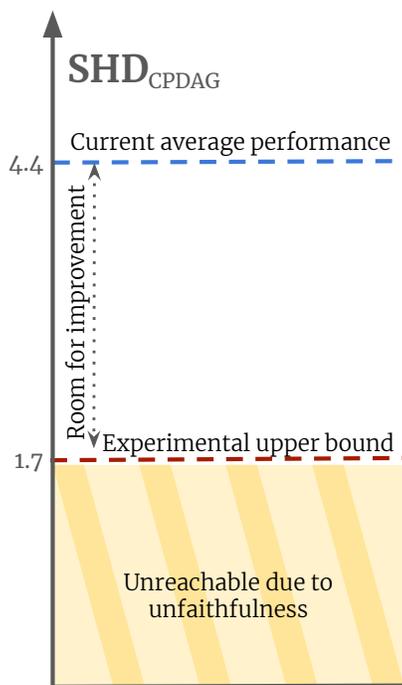


Figure 1: Neural causal discovery methods suffer from inherent performance limit due to violation of faithfulness assumption, but there is still room for improvement. Values computed for ER(5,1) class of graphs. See Sections 3, 5.

We believe that our work establishes a solid foundation that will propel future research in ML methods for causal discovery. Our original contributions are as follows:

- We identify violations of faithfulness as the core challenge and analyze its consequences both qualitatively and quantitatively.
- We develop an open unified benchmark for causal discovery evaluation.
- We present a soft upper bound on the performance of neural causal discovery methods for synthetic benchmarks.

2 BACKGROUND

Structural Causal Models (SCMs) and graph representation Causal relationships are commonly formalized using SCMs, which represent causal dependencies through a set of structural equations. For a directed acyclic graph (DAG) $G = (V, E)$, an SCM is defined by a set of equations

$$X_i = f_i(Pa_i, U_i), \quad (1)$$

where $i \in V$, X_i is a random variable, $f_i: \mathbb{R}^{|Pa_i|+1} \rightarrow \mathbb{R}$ is a function, Pa_i denotes the set of parents of the vertex i in the graph G , and U_i is an independent noise term associated with X_i . In this paper, we assume *additive noise* SCMs, also referred to as *additive noise models* (ANM), where:

$$f_i(Pa_i, U_i) = g_i(Pa_i) + U_i \quad (2)$$

for some $g_i: \mathbb{R}^{|Pa_i|} \rightarrow \mathbb{R}$.

Causal discovery Causal structure discovery aims to recover the ground truth DAG representing causal relationships among variables. However, the unique solution cannot be identified from the observational data only; instead, one can only identify the structure up to a Markov Equivalence Class (MEC), the set of DAGs that encode the same conditional independencies. This can be uniquely represented by a Complete Partially Directed Acyclic Graph (CPDAG), which is a sum of DAGs from the same class. This results in a graph that includes both directed and undirected edges, reflecting consistent and uncertain causal directions within the MEC (Verma & Pearl, 1990).

Faithfulness assumption A probability distribution P is said to be *faithful* to a DAG $G = (V, E)$ if all the conditional independence relations present in the data correspond to those implied by the d -separation criteria of the DAG (for more on d -separation, see Appendix A.1 or Pearl (2009)). Formally, this can be written as:

$$X_a \perp\!\!\!\perp X_b \mid X_S \Rightarrow a \text{ is } d\text{-separated from } b \mid S, \quad (3)$$

where $\perp\!\!\!\perp$ denotes conditional independence of the variables, $a, b \in V$ are nodes of the graph, and $S \subseteq V \setminus \{a, b\}$ is a set of nodes. Intuitively, the faithfulness assumption can be understood as the statement that all statistical independencies in the observed data are the result of the underlying causal structure. Faithfulness assumption can be violated, for example, in a situation when paths cancel each other effects out, leading to statistical independence despite an existing causal relationship. **An example of this kind of violation is shown in Appendix A.2.**

While the faithfulness is a useful and powerful assumption in causal discovery, it is rarely satisfied in the practical scenarios (Cartwright, 2001; Andersen, 2013).

Score-based neural causal discovery To allow for scalable causal discovery on graphs with hundreds of nodes, recent approaches focus on heuristics employing continuous optimization techniques that use neural networks as functional approximators to model the underlying probability distribution of the data (Nazaret et al., 2024). These approaches use a continuous representation of the graph structure, enforcing a differentiable acyclicity constraint to ensure the result is a valid DAG. The primary objective is to maximize $\log p_\theta(X|G)$, that is the log-likelihood of the data given the graph while incorporating regularization terms to control graph complexity. The training procedure comprises two parts: fitting functional approximators and structure search. They are usually done in parallel to maximize compute efficiency. Methods of this class are guaranteed to recover a DAG from the MEC class of ground true graph when the faithfulness assumption is fulfilled (see Brouillard et al. (2020)).

We benchmark four differentiable causal discovery methods DCDI (Brouillard et al., 2020), SDCD (Nazaret et al., 2024), BayesDAG (Annadani et al., 2023), and DiBS (Lorch et al., 2021), as they summarize various research directions and improvements explored in neural causal discovery over the last four years (see Appendix E). DCDI and SDCD represent the graph using an adjacency matrix, and optimize using the Augmented Lagrangian method (Zheng et al., 2018), aiming to find a single graph that maximizes the likelihood, with regularization added to penalize complex structures. In contrast, BayesDAG and DiBS take a Bayesian approach, approximating the posterior distribution over graphs rather than finding a single solution, with regularization introduced via prior distributions on graph structures. All four methods assume that the distribution is faithful to the ground truth DAG.

Structure evaluation We evaluate graph discovery within the MEC using $\text{ESHD}_{\text{CPDAG}}$ and $\text{F1-Score}_{\text{CPDAG}}$, where $\text{ESHD}_{\text{CPDAG}} = 0$ and $\text{F1-Score}_{\text{CPDAG}} = 1$ when the predicted graph is in the same MEC as the ground truth. For Bayesian methods, we compute the average by sampling 100 graphs from the posterior; for non-Bayesian methods, we use a single graph.

The Structural Hamming Distance (SHD) (Tsamardinos et al., 2006) counts edge insertions, deletions, and reversals needed to match the predicted graph to the true graph. We define **Expected SHD between CPDAGs** as:

$$\text{ESHD}_{\text{CPDAG}}(\mathcal{G}, \mathbb{G}) = \mathbb{E}_{\mathcal{G}^* \sim \mathbb{G}}[\text{SHD}(\text{CPDAG}(\mathcal{G}), \text{CPDAG}(\mathcal{G}^*))], \quad (4)$$

where \mathbb{G} is the resulting distribution of graphs, \mathcal{G}^* is a graph sampled from \mathbb{G} and \mathcal{G} is the ground true graph. **The F1-Score measures the harmonic mean of precision and recall for edge predictions.** We compute the **Expected F1-Score between the CPDAGs** as follows:

$$\text{F1-Score}_{\text{CPDAG}}(\mathcal{G}, \mathbb{G}) = \mathbb{E}_{\mathcal{G}^* \sim \mathbb{G}}[\text{F1-Score}(\text{CPDAG}(\mathcal{G}), \text{CPDAG}(\mathcal{G}^*))]. \quad (5)$$

For more details and justification on the selection of metrics please refer to Appendix D.

3 UNIFIED BENCHMARK FOR SCORE-BASED NEURAL CAUSAL DISCOVERY METHODS ON SYNTHETIC DATA

In this section, we present a unified benchmark that exposes both the strengths and limitations of neural-based causal discovery methods. We evaluate methods DiBS, DCDI, BayesDAG, and SDCD introduced in Section 2 on identical datasets, tune hyperparameters consistently, and use a common functional approximation.

Our analysis spans several key dimensions of performance. In Section 3.2, we show that despite advancements in causal discovery over the past few years, $\text{ESHD}_{\text{CPDAG}}$ and $\text{F1-Score}_{\text{CPDAG}}$ metrics do not improve significantly. In Section 3.3, we demonstrate that structure discovery accuracy does not scale with the amount of data. Finally, in Section 3.4, we confirm that variations in MLP architecture have minimal impact on performance. **In Appendix F we provide additional results on real-world structures which align with the conclusions presented in this section.**

3.1 EXPERIMENTAL SETUP

Dataset generation We sample three types of graphs from the Erdős-Rényi (ER) distribution (Erdős & Rényi, 1959): one with 5 nodes and the expected degree of 1, another with 10 nodes and the expected degree of 2, and the third with 30 nodes and the expected degree of 2. These datasets are referred to as ER(5, 1), ER(10, 2), and ER(30, 2), respectively. These parameter choices align with commonly studied medium-sized graphs in causal discovery research (Brouillard et al., 2020; Nazaret et al., 2024). Data generation follows the SCM formalism introduced in Section 2, with functional relationships modeled by two-layer neural networks (hidden dimension 8, ReLU activation) and additive Gaussian noise. The noise has zero mean, and its variance is sampled independently for each node. This setup is known to be challenging (Geffner et al., 2024; Nazaret et al., 2024). For more details refer to Appendix C.1.

Method	ER(5, 1)		ER(10, 2)		ER(30, 2)	
	ESHDCPDAG	F1-ScoreCPDAG	ESHDCPDAG	F1-ScoreCPDAG	ESHDCPDAG	F1-ScoreCPDAG
DCDI	5.7 (3.7, 8.1)	0.60 (0.46, 0.74)	16.9 (15.7, 18.1)	0.52 (0.50, 0.56)	45.9 (42.0, 49.9)	0.73 (0.69, 0.77)
BayesDAG	3.9 (3.6, 4.3)	0.78 (0.77, 0.81)	18.3 (16.9, 19.8)	0.56 (0.54, 0.59)	51.7 (48.2, 55.9)	0.59 (0.57, 0.61)
DiBS	2.6 (1.7, 3.7)	0.85 (0.80, 0.90)	16.9 (14.2, 20.1)	0.61 (0.57, 0.68)	68.0 (65.3, 70.9)	0.23 (0.22, 0.24)
SDCD	5.4 (3.8, 6.7)	0.60 (0.35, 0.69)	20.9 (19.5, 22.2)	0.54 (0.46, .62)	62.8 (58.8, 67.7)	0.55 (0.53, 0.58)

Table 1: Comparison of ESHDCPDAG and F1-ScoreCPDAG for different methods on ER(10, 2) (left) and ER(30, 2) (right) dataset. The numbers in the subscripts correspond to 95% confidence intervals. The statistics were computed based on 30 graphs.

Hyperparameter tuning To ensure a fair comparison across all methods, we perform systematic hyperparameter tuning, selecting the best-performing parameters for each model. We employ a grid search approach based on the parameter ranges suggested by the original authors. This process optimizes key variables such as regularization coefficients, sparsity controls, and kernel configurations. Details can be found in Appendix C.2.

Functional approximators We standardize the choice of functional approximators across all experiments, using a two-layer MLP with a hidden dimension of 4. This model size is consistent with previous work (Brouillard et al., 2020; Nazaret et al., 2024) and has proven to perform well across all the benchmarked methods, as discussed in Section 3.4. Additionally, we use trainable variance to allow the model to adapt to varying noise levels, in line with our dataset generation setup.

3.2 PERFORMANCE COMPARISON

Table 1, summarizes the benchmark results of neural-based causal discovery methods on graphs from ER(5, 1), ER(10, 2), and ER(30, 2) classes. We tune hyperparameters to optimize the ESHDCPDAG metric. For all classes of graphs, metrics were computed based on 30 graphs.

The results show that DiBS is particularly effective for smaller graphs (ER(5, 1) and ER(10, 2)), while DCDI is able to achieve the best results for moderate-size graphs (ER(10, 2) and ER(30, 2)). The ranking of the methods changes with the size of the graphs but SDCD consistently exhibits the worst performance in terms of ESHDCPDAG. Nevertheless, the performance of all the methods remains unsatisfactory with all methods predicting more than half of the edges incorrectly.

3.3 IMPACT OF SAMPLE SIZE

We investigate whether the number of observations affects the performance of causal discovery methods. One could expect that neural based models, similarly to independence testing ones, will improve when more data is supplied (Kalisch & Bühlmann, 2007). We compare benchmarked methods on dataset with varying number of observational samples, ranging from 20 to 8,000 observations.

The results, presented in Figure 13, reveal no consistent pattern of improvement in the ESHDCPDAG metric as observational sample size increases, despite extensive hyperparameter tuning (as described in Section 3.1). For example, DiBS shows the best performance on larger datasets, but its improvements plateau after around 800 samples. Similarly, BayesDAG shows only marginal improvements with larger sample sizes and is unable to outperform DiBS. DCDI improves up to 250 samples and then maintains consistent performance regardless of the sample size, similar to DiBS. Interestingly, SDCD’s performance is poor on datasets with small number of observations but begins to improve once sample sizes exceed 250, though is unable to reach DCDI’s performance, for larger sample sizes the rate of improvement decreases.

Further analysis of the effect of sample size on smaller graphs ER(5, 1) is presented in Figure 12 in Appendix C.4. Overall, the results on smaller graphs align with the trends observed on larger graphs. Specifically, while some methods improve with increasing sample size, others show inconsistent or even degraded performance.

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

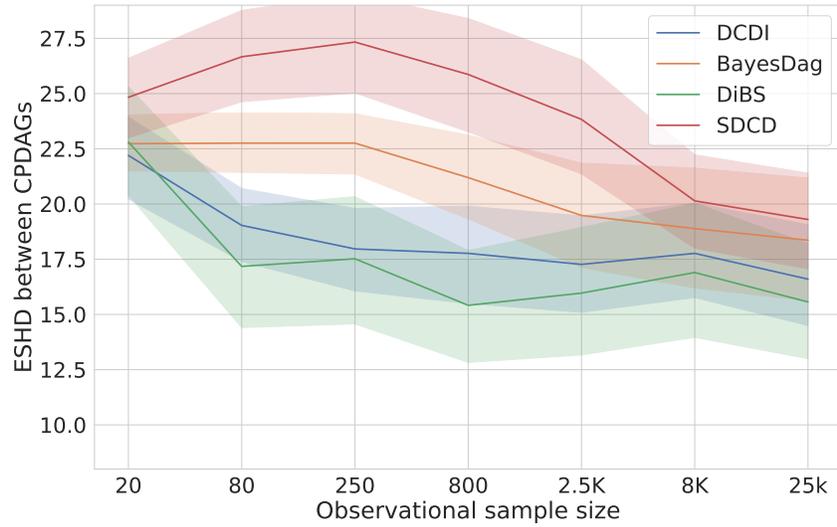


Figure 2: Comparison of $ESH D_{CPDAG}$ for different methods using the [4, 4] architecture, for ER(10, 2) dataset, averaged over 30 samples.

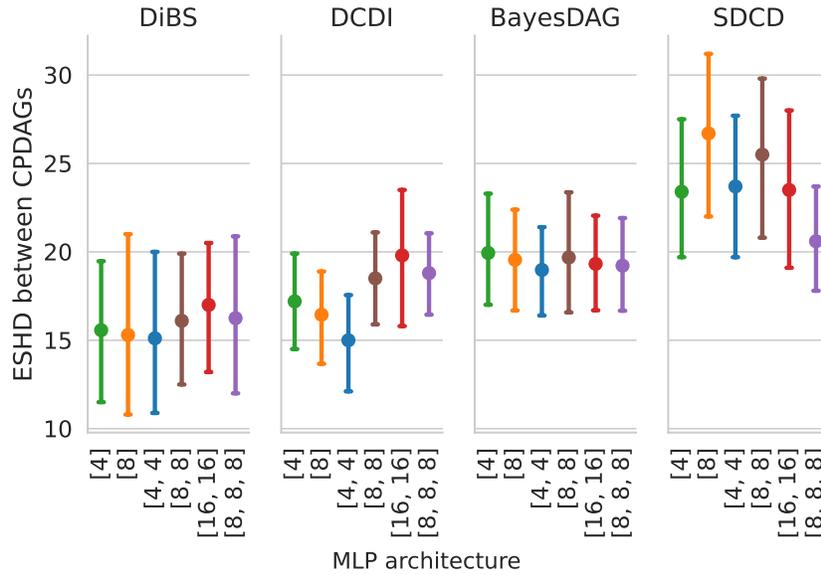


Figure 3: Comparison of $ESH D_{CPDAG}$ using different MLP architectures as functional approximator for ER(10, 2) dataset and 800 observational samples, averaged over 30 samples.

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

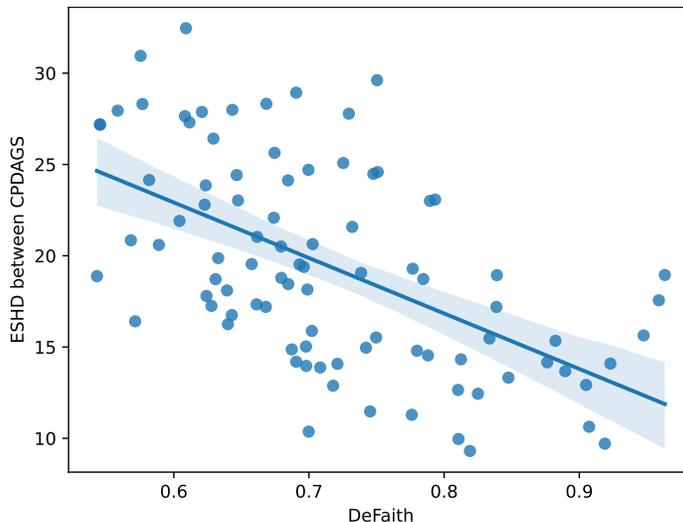


Figure 4: Linear regression fit between the average performance of neural causal discovery methods and faithfulness accuracy measure.

3.4 COMPARISON OF NEURAL MODEL ARCHITECTURES

Finally, we investigate the impact of the neural model architecture, used as the functional approximator, on the performance of the benchmarked methods. Specifically, we assess how the capacity of different architectures influences the ability to uncover causal relationships from synthetic data. To provide a comprehensive evaluation, we explored architectures with 1, 2, and 3 layers, configured with 4, 8, and 16 hidden units.

Results, presented in Figure 3 show the comparison of $\text{ESHD}_{\text{CPDAG}}$ metric for the benchmarked architectures across all methods on dataset with 800 samples. We find that the choice of neural architecture has no significant impact on performance across methods. We conclude that any of the tested MLP architectures provides sufficient capacity to model the underlying distribution effectively. **Additionally for BayesDAG and SDCD we implemented layer normalization and residual connections. We investigated the impact of this changes in architectures and did not find any significant differences, see Figure 11. The details and additional experimental results are in Appendix C.3.**

4 MEASURING IMPACT OF FAITHFULNESS VIOLATION

In this section we explore how violations of the faithfulness assumption impact the performance of neural causal discovery methods. In Section 3, we showed that despite various attempts to scale up data and model complexity, the performance of these methods remains stagnant, possibly due to deeper challenges related to the underlying data properties and the limitations inherent to the algorithms. This leads us to investigate whether violations of the faithfulness assumption, common in synthetic non-linear data, might be the key factor limiting performance improvements.

The faithfulness assumption translates into the set of conditional independence statements that all need to be satisfied. As mentioned in Section 2, synthetic non-linear data rarely adheres to faithfulness assumption, rendering binary criterion not practical. To address this, we introduce a degree of faithfulness metric, denoted *DeFaith*, **which captures the faithfulness violations on a continuous scale.**

Inspired by Zhang & Spirtes (2003), we use Spearman’s rank correlation coefficient to quantify the conditional dependencies in the dataset. **We define a predictor that classifies nodes as independent if conditional Spearman’s rank correlation coefficient computed based on a dataset D is smaller than a certain threshold.**

DeFaith is the quality of this predictor measured by Area Under Receiver Operator Curve **computed over all possible pairs of variables a, b and separation sets $S \subseteq V \setminus \{a, b\}$.** Formally,

Algorithm 1 Overview of NN-opt method

1: **Input:** Set of nodes V , training data $\{D_i\}_{i \in V}$, regularization coefficient λ , \mathbb{G} the space of DAGs with nodes V

2: # Part 1: Network fitting

3: **for** $i \in V$ and $\pi \subseteq V \setminus \{i\}$ **do** ▷ For each variable and each possible parent set

4: $\theta_{i,\pi} \leftarrow \text{TRAINNETWORK}(i, D, \pi)$ ▷ Train ensembles of 3 networks

5: **end for**

6: # Part 2: Exhaustive graph search

7: **for** $G \in \mathbb{G}$ **do** ▷ Evaluate all possible DAGs

8: $\text{score}_G \leftarrow \sum_{i \in V} \text{COMPUTENLL}(D_i, D_{Pa_i^G}, \theta_{i, Pa_i^G})$ ▷ Compute NLL using ensemble

9: $\text{score}_G \leftarrow \text{score}_G + \lambda \cdot |G|$ ▷ Add regularizing term

10: **end for**

11: **Output:** $\arg \max\{\text{score}_G : G \in \mathbb{G}\}$

$$DeFaith(D, G) = \underset{a, b \in V, S \subseteq V \setminus \{a, b\}}{\text{AUROC}} (1 - \text{abs}(\rho_s^D(a, b|S)), \mathbf{1}[a \perp_G b|S])$$

where V is set of nodes in graph G , $a \perp_G b|S$ denotes d -separation between nodes a and b given S , and $\rho_s^D(a, b|S)$ denotes conditional Spearman’s rank correlation coefficient computed based on dataset D . The measure attains a value of 1.0 for faithful distributions.

In this experiment, we generate 30 graphs from the ER(10, 2) class, introduced in Section 3.1. Based on each graph, we define three different SCMs, resulting in 90 distinct distributions. Each dataset consists of 8,000 observational samples. We then evaluate the *DeFaith* of each distribution and compute the performance of the selected neural-based causal discovery methods.

In Figure 4 we present the relationship between average performance of all methods and the degree of faithfulness for all 90 distributions in the dataset. The performance is better (lower SHD) for distributions with higher degree of faithfulness. The Spearman’s rank correlation coefficient is $\rho = -0.58$. This result proves the strong anti-monotonicity between the faithfulness accuracy and methods’ performance.

5 ESTIMATING UPPER BOUND ON PERFORMANCE

In this section, we investigate the limits of the performance of score-based neural causal discovery methods. To do this we develop a method dubbed as NN-opt method, to compute an experimental upper bound on the performance. As for the benchmarked methods, the goal of NN-opt method is to find a structure that minimizes the regularized log-likelihood of data, therefore it is expected to recover a graph from the correct MEC class when the faithfulness assumption holds (see Section 2).

The method overview is in Algorithm 1. It is based on the common approach used by score-based neural causal discovery methods described in Section 2. The procedure consists of two steps. First, we train neural networks to approximate functional relationships between variables. Contrary to benchmarked methods we train a separate network for each parent set instead of training one for all. This renders functional approximation fitting procedure completely independent from structure search. Therefore, it simplifies the training task and allows for strict control of the training procedure via validation loss monitoring. Second, we conduct an exhaustive search over the space of DAGs to find the structure that minimizes the log-likelihood loss. For increased stability of this step, we use an ensemble of 3 neural networks to compute the log-likelihood of the data under various structures.

The approach is exhaustive both in the sense of structure search and in neural network training, trading computational efficiency for additional precision. NN-opt is a brute force technique intended to be able to reach the limits of score-based neural causal discovery approaches. NN-opt is helpful as an upper-bound benchmark but is not practical to use.

We expect the method to improve with the number of samples and stabilize when the data becomes sufficiently large. Therefore, we applied NN-opt method to datasets of various sizes. The results are presented in Figure 5 on the left. For very small datasets we observe rapid improvement in

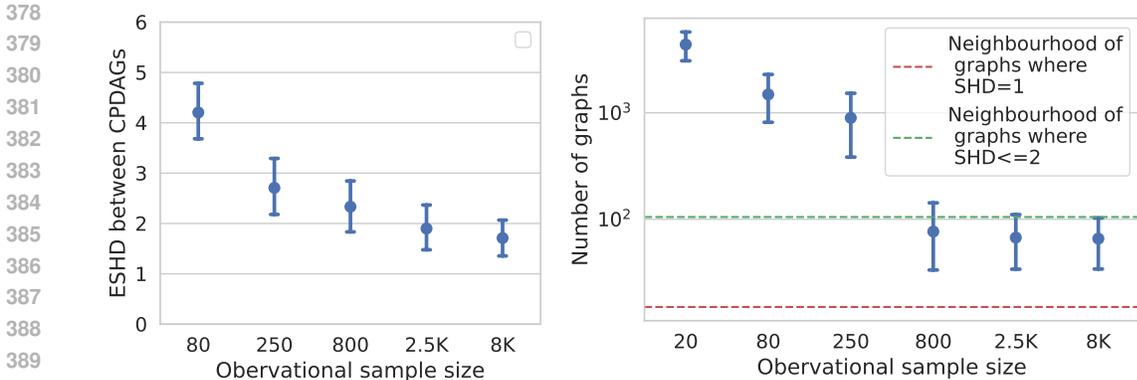


Figure 5: Comparison of the performance of NN-opt method depending on data size (left), and comparison of number of DAGs with score higher than true graph (right). Averaged over 90 samples

terms of $\text{ESHD}_{\text{CPDAG}}$, but as the sample size grows, the structure discovery accuracy stabilizes. For sample sizes of 2,500 and 8,000, the value of $\text{ESHD}_{\text{CPDAG}}$ is just below 2. In the dataset used for this experiment, the average number of edges in CPDAG is around 8.4, meaning that on average almost 25% of the edges are predicted incorrectly.

Furthermore, to show that the problem is systematic, we present the number of graphs with a higher score than the ground true DAG in Figure 5 on the right. For smaller datasets (with no more than 250 samples) there are around 1000 graphs or more with scores higher than the ground true graph. The number stabilizes around 65 structures, that scored higher than the ground true graph, for bigger datasets. This number is close to the number of graphs with SHD distance ≤ 2 from the ground truth, depicted by the green line in the figure. These findings demonstrate the methods' consistent inability to identify correct structures.

We argue that this result shows the inherent limitations of the score-based neural causal discovery algorithms due to the violation of the faithfulness assumption. Our NN-opt method controls errors raised from both functional approximations fitting and structure search. Thus violation of faithfulness is the only probable source of errors.

To ensure the validity of the result we performed an extensive hyperparameter search, including models with various architectures. Details of described experiments can be found in Appendix B.

6 RELATED WORK

Causal discovery without the faithfulness assumption While many causal discovery methods rely on the faithfulness assumption, alternative conditions have been proposed. One notable approach is the adjacency-faithfulness assumption, introduced by Ramsey et al. (2006) in the conservative PC algorithm. This assumption, which is less restrictive than full faithfulness, leads to more robust with minimal computational overhead. In the context of linear structural causal models (SCMs), Van de Geer & Bühlmann (2013) demonstrated that a sparsity-based assumption can effectively reveal the underlying causal structure. Similarly, Isozaki (2014) proposed a method to reduce unnecessary independence tests during structure discovery, offering greater robustness against violations of faithfulness due to statistical errors. More recently, Ng et al. (2021) suggested another causal discovery method, based on relaxed faithfulness assumption that requires less independence tests to be fulfilled. Marx et al. (2021) explores a weaker alternative to the faithfulness assumption, called the 2-faithfulness assumption, and suggests how to construct a causal discovery algorithm based on it. Moreover, Lippe et al. (2022) introduced a neural-based approach that uses interventional data, avoiding the faithfulness assumption altogether.

Describing faithfulness violations Faithfulness violation has been extensively explored in the linear setting by (Uhler et al., 2013). They showed that the conditions that would allow for discovering

432 the true independencies in a finite sample regime are rarely met when making use of linear synthetic
 433 data. Additionally, they proved that the bigger the graph the more difficult it is to find a faithful
 434 distribution. Zhang & Spirtes (2003) provided theoretical conditions for violation of faithfulness being
 435 detectable during training. More generally, Andersen (2013) described reasons, why faithfulness is
 436 likely violated in complex, evolved real-world systems. To the best of our knowledge, we are the first
 437 to estimate the limits of the score-based neural causal discovery methods on unfaithful data.

438
 439 **Benchmarking** There is a multitude of recent benchmarks that use real-world data to assess the
 440 performance of causal discovery methods (Chevalley et al., 2022; Mehrjou et al., 2022). However,
 441 these datasets lack the ground truth structure rendering structure discovery accuracy assessment
 442 impossible. Additionally, these works usually focus on classical, not neural, causal discovery methods.
 443 Some recent work is concerned with the quality of evaluations and performance under assumptions
 444 violations. Karimi-Mamaghan et al. (2024) investigates metrics for Bayesian causal discovery in
 445 a linear setting. Their finding suggests that the standard structure-based metrics do not align well
 446 with downstream task performance when structure uncertainty is high (especially for bigger graphs),
 447 Montagna et al. (2023) evaluates classical causal discovery methods under different assumption
 448 violations. In our work, we focus on a unified, synthetic, and challenging setup to thoroughly evaluate
 449 neural causal discovery claims of being general and accurate. Most recently, Zhou et al. (2024)
 450 introduced a comprehensive benchmark, but they did not compare neural-based methods in their work.

451 7 LIMITATIONS & FUTURE WORK

- 453 • Work of Lippe et al. (2022) suggests that interventional data can replace the need for
 454 faithfulness assumption. A valuable extension of our research would be to evaluate the
 455 performance of the benchmarked methods on interventional datasets to understand their
 456 limitations and potential improvements in this context.
- 457 • Our work provides experimental evidence for the scale of the impact of violation of faith-
 458 fulness on performance in a challenging non-linear setting. It would be beneficial for the
 459 community if some theoretical results (akin Uhler et al. (2013); Zhang & Spirtes (2003))
 460 were derived in a non-linear setting.
- 461 • While our, experimental upper bound, NN-opt method is based on common, with bench-
 462 marked methods, theoretical principles. We leave strict theoretical justification of its
 463 optimality for future work.
- 464 • In this work we present the method that allows to estimate the upper bound on performance
 465 of score-based neural causal discovery methods on any dataset and provide numerical results
 466 for the Erodos-Renyi class of graphs. The results could be computed for more classes and
 467 even some small real-world or real-world inspired graphs, see Elidan (2001).

469 8 CONCLUSIONS

471 In this work, we present compelling evidence that the faithfulness assumption is a major limiting
 472 factor in advancing causal discovery. Our findings demonstrate that the accuracy of structure recovery
 473 is correlated with the degree of faithfulness violation. Additionally, we introduce a novel method to
 474 calculate the upper bound of performance for score-based neural causal discovery methods, revealing
 475 serious limitations. Our results highlight the need for a paradigm shift. We argue that further progress
 476 in causal discovery requires moving beyond the faithfulness assumption and encourage researchers to
 477 explore alternative conditions. The implications of our work extend beyond theoretical advancements.
 478 By challenging the faithfulness assumption, we open up avenues for more robust and generalizable
 479 methods in causal discovery, which could have far-reaching consequences in fields like healthcare,
 480 economics, and policy-making.

482 9 REPRODUCIBILITY STATMENT

484 We put effort and resources to ensure that presented experiments can be reproduced by the re-
 485 search community. Specifically, we provide detailed descriptions of the data generation process,
 benchmarking score-based neural causal discovery methods and proposed NN-opt method.

Our dataset generation process is based on code included in DCDI code repository (Brouillard et al., 2020) and the details of this processed can be found in Section 3.1 and Appendix C.1. The performance of the selected causal discovery methods, for the benchmark, was compute using official repositories released by authors: DCDI (Brouillard et al., 2020), SDCD (Nazaret et al., 2024), BayesDag (Annadani et al., 2023) and DiBS (Lorch et al., 2021). The range of tested hyperparameters and the selected values can be found in Section 3.1 and Appendix C.2.

The description on NN-opt method is provided in Section 5 moreover the high level overview of method is in Algorithm 1. Hyperparameter selection is described in Appendix B.

REFERENCES

- Holly Andersen. When to expect violations of causal faithfulness and why it matters. *Philosophy of Science*, 80, 12 2013. doi: 10.1086/673937.
- Yashas Annadani, Nick Pawlowski, Joel Jennings, Stefan Bauer, Cheng Zhang, and Wenbo Gong. Bayesdag: Gradient-based posterior inference for causal discovery. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023. URL http://papers.nips.cc/paper_files/paper/2023/hash/05cf28e3d3c9a179d789c55270fe6f72-Abstract-Conference.html.
- Kevin Bello, Bryon Aragam, and Pradeep Ravikumar. Dagma: Learning dags via m-matrices and a log-determinant acyclicity characterization. *Advances in Neural Information Processing Systems*, 35:8226–8239, 2022.
- Philippe Brouillard, Sébastien Lachapelle, Alexandre Lacoste, Simon Lacoste-Julien, and Alexandre Drouin. Differentiable causal discovery from interventional data. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. URL <https://proceedings.neurips.cc/paper/2020/hash/f8b7aa3a0d349d9562b424160ad18612-Abstract.html>.
- Nancy Cartwright. What is wrong with bayes nets? *Monist*, 84, 04 2001. doi: 10.2307/27903726.
- Mathieu Chevalley, Yusuf Roohani, Arash Mehrjou, Jure Leskovec, and Patrick Schwab. Causal-bench: A large-scale benchmark for network inference from single-cell perturbation data. *CoRR*, abs/2210.17283, 2022. doi: 10.48550/ARXIV.2210.17283. URL <https://doi.org/10.48550/arXiv.2210.17283>.
- Daniel Coelho de Castro, Ian Walker, and Ben Glocker. Causality matters in medical imaging. *CoRR*, abs/1912.08142, 2019. URL <http://arxiv.org/abs/1912.08142>.
- G. Elidan. Bayesian Network Repository, 2001. <https://www.cse.huji.ac.il/galel/Repository/>.
- P Erdős and A Rényi. On random graphs i. *Publicationes Mathematicae Debrecen*, 6:290–297, 1959.
- Tomas Geffner, Javier Antorán, Adam Foster, Wenbo Gong, Chao Ma, Emre Kiciman, Amit Sharma, Angus Lamb, Martin Kukla, Nick Pawlowski, Agrin Hilmkil, Joel Jennings, Meyer Scetbon, Miltiadis Allamanis, and Cheng Zhang. Deep end-to-end causal inference. *Trans. Mach. Learn. Res.*, 2024, 2024. URL <https://openreview.net/forum?id=e6sqttxEGX>.
- Kevin D. Hoover. *Causality in Macroeconomics*, pp. 89–134. Cambridge University Press, 2001.
- Takashi Isozaki. A robust causal discovery algorithm against faithfulness violation. *Inf. Media Technol.*, 9(1):121–131, 2014. doi: 10.11185/IMT.9.121. URL <https://doi.org/10.11185/imt.9.121>.
- Markus Kalisch and Peter Bühlmann. Estimating high-dimensional directed acyclic graphs with the pc-algorithm. *J. Mach. Learn. Res.*, 8:613–636, 2007. doi: 10.5555/1314498.1314520. URL <https://dl.acm.org/doi/10.5555/1314498.1314520>.

- 540 Amir Mohammad Karimi-Mamaghan, Panagiotis Tigas, Karl Henrik Johansson, Yarin Gal, Yashas
541 Annadani, and Stefan Bauer. Challenges and considerations in the evaluation of bayesian causal
542 discovery. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net, 2024. URL <https://openreview.net/forum?id=bqgtkBDkNs>.
- 543
544
- 545 Sébastien Lachapelle, Philippe Brouillard, Tristan Deleu, and Simon Lacoste-Julien. Gradient-based
546 neural dag learning. *arXiv preprint arXiv:1906.02226*, 2019.
- 547
- 548 Hao-Chih Lee, Matteo Danieletto, Riccardo Miotto, Sarah T Cherng, and Joel T Dudley. Scaling
549 structural learning with no-bears to infer causal transcriptome networks. In *Pacific Symposium on*
550 *Biocomputing 2020*, pp. 391–402. World Scientific, 2019.
- 551
- 552 Phillip Lippe, Taco Cohen, and Efstratios Gavves. Efficient neural causal discovery without acyclicity
553 constraints. In *The Tenth International Conference on Learning Representations, ICLR 2022,*
554 *Virtual Event, April 25-29, 2022*. OpenReview.net, 2022. URL <https://openreview.net/forum?id=eYciPrLuUhG>.
- 555
- 556 Romain Lopez, Jan-Christian Hütter, Jonathan Pritchard, and Aviv Regev. Large-scale differentiable
557 causal discovery of factor graphs. *Advances in Neural Information Processing Systems*, 35:
558 19290–19303, 2022.
- 559
- 560 Lars Lorch, Jonas Rothfuss, Bernhard Schölkopf, and Andreas Krause. Dibs: Differentiable bayesian
561 structure learning. In Marc’Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang,
562 and Jennifer Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems 34:*
563 *Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December*
564 *6-14, 2021, virtual*, pp. 24111–24123, 2021. URL <https://proceedings.neurips.cc/paper/2021/hash/ca6ab34959489659f8c3776aaf1f8efd-Abstract.html>.
- 565
- 566 Alexander Marx, Arthur Gretton, and Joris M. Mooij. A weaker faithfulness assumption based
567 on triple interactions. In Cassio P. de Campos, Marloes H. Maathuis, and Erik Quaeghebeur
568 (eds.), *Proceedings of the Thirty-Seventh Conference on Uncertainty in Artificial Intelligence,*
569 *UAI 2021, Virtual Event, 27-30 July 2021*, volume 161 of *Proceedings of Machine Learning*
570 *Research*, pp. 451–460. AUAI Press, 2021. URL <https://proceedings.mlr.press/v161/marx21a.html>.
- 571
- 572
- 573 Arash Mehrjou, Ashkan Soleymani, Andrew Jesson, Pascal Notin, Yarin Gal, Stefan Bauer, and
574 Patrick Schwab. Genedisco: A benchmark for experimental design in drug discovery. In
575 *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event,*
576 *April 25-29, 2022*. OpenReview.net, 2022. URL <https://openreview.net/forum?id=-w2oom06qgc>.
- 577
- 578 Francesco Montagna, Atalanti-Anastasia Mastakouri, Elias Eulig, Nicoletta Noceti, Lorenzo
579 Rosasco, Dominik Janzing, Bryon Aragam, and Francesco Locatello. Assumption viola-
580 tions in causal discovery and the robustness of score matching. In Alice Oh, Tristan
581 Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Ad-*
582 *vances in Neural Information Processing Systems 36: Annual Conference on Neural Infor-*
583 *mation Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 -*
584 *16, 2023*, 2023. URL http://papers.nips.cc/paper_files/paper/2023/hash/93ed74938a54a73b5e4c52bbaf42ca8e-Abstract-Conference.html.
- 585
- 586 Vinod Nair and Geoffrey E. Hinton. Rectified linear units improve restricted boltzmann machines.
587 In Johannes Fürnkranz and Thorsten Joachims (eds.), *Proceedings of the 27th International*
588 *Conference on Machine Learning (ICML-10), June 21-24, 2010, Haifa, Israel*, pp. 807–814.
589 Omnipress, 2010. URL <https://icml.cc/Conferences/2010/papers/432.pdf>.
- 590
- 591 Achille Nazaret, Justin Hong, Elham Azizi, and David M. Blei. Stable differentiable causal discov-
592 ery. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria,*
593 *July 21-27, 2024*. OpenReview.net, 2024. URL <https://openreview.net/forum?id=JJZBZW28Gn>.

- 594 Ignavier Ng, Yujia Zheng, Jiji Zhang, and Kun Zhang. Reliable causal discovery with im-
595 proved exact search and weaker assumptions. In Marc’Aurelio Ranzato, Alina Beygelz-
596 imer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (eds.), *Advances*
597 *in Neural Information Processing Systems 34: Annual Conference on Neural Informa-*
598 *tion Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp. 20308–
599 20320, 2021. URL [https://proceedings.neurips.cc/paper/2021/hash/](https://proceedings.neurips.cc/paper/2021/hash/a9b4ec2eb4ab7b1b9c3392bb5388119d-Abstract.html)
600 [a9b4ec2eb4ab7b1b9c3392bb5388119d-Abstract.html](https://proceedings.neurips.cc/paper/2021/hash/a9b4ec2eb4ab7b1b9c3392bb5388119d-Abstract.html).
- 601 Judea Pearl. *Causality*. Cambridge University Press, 2 edition, 2009. doi: 10.1017/
602 CBO9780511803161.
- 603 Jonas Peters, Peter Bühlmann, and Nicolai Meinshausen. Causal inference by using invariant
604 prediction: identification and confidence intervals. *Journal of the Royal Statistical Society Series*
605 *B: Statistical Methodology*, 78(5):947–1012, 2016.
- 606 Joseph D. Ramsey, Jiji Zhang, and Peter Spirtes. Adjacency-faithfulness and conserva-
607 tive causal inference. In *UAI ’06, Proceedings of the 22nd Conference in Uncer-*
608 *tainty in Artificial Intelligence, Cambridge, MA, USA, July 13-16, 2006*. AUAI Press,
609 2006. URL [https://dslpitt.org/uai/displayArticleDetails.jsp?mmnu=](https://dslpitt.org/uai/displayArticleDetails.jsp?mmnu=1&smnu=2&article_id=1259&proceeding_id=22)
610 [1&smnu=2&article_id=1259&proceeding_id=22](https://dslpitt.org/uai/displayArticleDetails.jsp?mmnu=1&smnu=2&article_id=1259&proceeding_id=22).
- 611 Paul Rolland, Volkan Cevher, Matthäus Kleindessner, Chris Russell, Dominik Janzing, Bernhard
612 Schölkopf, and Francesco Locatello. Score matching enables causal discovery of nonlinear additive
613 noise models. In *International Conference on Machine Learning*, pp. 18741–18753. PMLR, 2022.
- 614 Ioannis Tsamardinos, Laura E. Brown, and Constantin F. Aliferis. The max-min hill-climbing
615 bayesian network structure learning algorithm. *Mach. Learn.*, 65(1):31–78, 2006. doi: 10.1007/
616 S10994-006-6889-7. URL <https://doi.org/10.1007/s10994-006-6889-7>.
- 617 Caroline Uhler, Garvesh Raskutti, Peter Bühlmann, and Bin Yu. Geometry of the faithfulness
618 assumption in causal inference. *The Annals of Statistics*, 41(2), April 2013. ISSN 0090-5364. doi:
619 10.1214/12-aos1080. URL <http://dx.doi.org/10.1214/12-AOS1080>.
- 620 Sara Van de Geer and Peter Bühlmann. 0-penalized maximum likelihood for sparse directed acyclic
621 graphs. 2013.
- 622 Thomas Verma and Judea Pearl. Equivalence and synthesis of causal models. In
623 Piero P. Bonissone, Max Henrion, Laveen N. Kanal, and John F. Lemmer (eds.),
624 *UAI ’90: Proceedings of the Sixth Annual Conference on Uncertainty in Artificial*
625 *Intelligence, MIT, Cambridge, MA, USA, July 27-29, 1990*, pp. 255–270. Elsevier,
626 1990. URL [https://dslpitt.org/uai/displayArticleDetails.jsp?mmnu=1&](https://dslpitt.org/uai/displayArticleDetails.jsp?mmnu=1&smnu=2&article_id=1918&proceeding_id=1006)
627 [smnu=2&article_id=1918&proceeding_id=1006](https://dslpitt.org/uai/displayArticleDetails.jsp?mmnu=1&smnu=2&article_id=1918&proceeding_id=1006).
- 628 Yue Yu, Tian Gao, Naiyu Yin, and Qiang Ji. Dags with no curl: An efficient dag structure learning
629 approach. In *International Conference on Machine Learning*, pp. 12156–12166. Pmlr, 2021.
- 630 Jiji Zhang and Peter Spirtes. Strong faithfulness and uniform consistency in causal inference. In
631 Christopher Meek and Uffe Kjærulff (eds.), *UAI ’03, Proceedings of the 19th Conference in*
632 *Uncertainty in Artificial Intelligence, Acapulco, Mexico, August 7-10 2003*, pp. 632–639. Morgan
633 Kaufmann, 2003. URL [https://dslpitt.org/uai/displayArticleDetails.](https://dslpitt.org/uai/displayArticleDetails.jsp?mmnu=1&smnu=2&article_id=983&proceeding_id=19)
634 [jsp?mmnu=1&smnu=2&article_id=983&proceeding_id=19](https://dslpitt.org/uai/displayArticleDetails.jsp?mmnu=1&smnu=2&article_id=983&proceeding_id=19).
- 635 Xun Zheng, Bryon Aragam, Pradeep Ravikumar, and Eric P. Xing. Dags with NO TEARS:
636 continuous optimization for structure learning. In Samy Bengio, Hanna M. Wallach, Hugo
637 Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett (eds.), *Advances*
638 *in Neural Information Processing Systems 31: Annual Conference on Neural Information*
639 *Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pp.
640 9492–9503, 2018. URL [https://proceedings.neurips.cc/paper/2018/hash/](https://proceedings.neurips.cc/paper/2018/hash/e347c51419fffb23ca3fd5050202f9c3d-Abstract.html)
641 [e347c51419fffb23ca3fd5050202f9c3d-Abstract.html](https://proceedings.neurips.cc/paper/2018/hash/e347c51419fffb23ca3fd5050202f9c3d-Abstract.html).
- 642 Wei Zhou, Hong Huang, Guowen Zhang, Ruize Shi, Kehan Yin, Yuanyuan Lin, and Bang Liu.
643 OCDB: revisiting causal discovery with a comprehensive benchmark and evaluation framework.
644 *CoRR*, abs/2406.04598, 2024. doi: 10.48550/ARXIV.2406.04598. URL [https://doi.org/](https://doi.org/10.48550/arXiv.2406.04598)
645 [10.48550/arXiv.2406.04598](https://doi.org/10.48550/arXiv.2406.04598).

A ADDITIONAL BACKGROUND INFORMATION

A.1 *d*-SEPARATION

Two nodes A and B in a DAG are said to be ***d*-separated** by a set of nodes Z if all paths between A and B are blocked when conditioning on Z . A path is considered blocked under the following conditions:

- If a path includes a non-collider node (a node where arrows do not converge, i.e., a chain or fork), conditioning on that node blocks the path. For example, if $A \rightarrow C \rightarrow B$, or $A \leftarrow C \rightarrow B$, conditioning on C makes A and B independent.
- If the path includes a collider (a node where arrows converge, i.e., $A \rightarrow C \leftarrow B$), the path is blocked unless either the collider itself or one of its descendants is conditioned on. For instance, in the path $A \rightarrow C \leftarrow B$, conditioning on C or its descendants would unblock the path, making A and B dependent.
- If there are multiple paths connecting A and B , all paths must be blocked for A and B to be considered *d*-separated. Even if one path remains unblocked, A and B are *d*-connected, meaning they are dependent.

In causal discovery, we are interested in making statements about the relationship between the causal graph and the data distribution. Given a causal graph G and the data distribution P , the **Markov assumption** states that if variables A and B are *d*-separated in the graph G by some conditioning set C , then A and B are conditionally independent in the distribution P when conditioned on the same conditioning set C . Formally, this can be written as:

$$A \perp_G B | C \Rightarrow A \perp_P B | C \quad (6)$$

A.2 EXAMPLE OF FAITHFULNESS VIOLATION

In this subsection we will illustrate a faithfulness violation for a simple 3 nodes structural causal model with linear functions and additive Gaussian noise. Such a setup is aimed at showing example of faithfulness violation while maintaining simplicity. The example and graphics is from (Uhler et al., 2013).

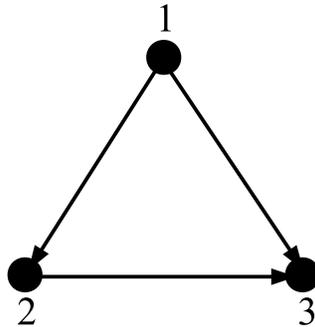


Figure 6: Simple 3 nodes graph G .

First lets define a structural causal model on a graph G shown in graph 6.

$$\begin{aligned} X_1 &= \varepsilon_1, \\ X_2 &= a_{12}X_1 + \varepsilon_2, \\ X_3 &= a_{13}X_1 + a_{23}X_2 + \varepsilon_3, \\ (\varepsilon_1, \varepsilon_2, \varepsilon_3) &\sim \mathcal{N}(0, I), \end{aligned}$$

Since data is linear we can use covariance to measure dependency of variables. Using defined structural causal model, we can write:

702
703
704
705
706
707
708
709
710

$$\text{cov}(X_1, X_2) = a_{12}, \tag{7}$$

$$\text{cov}(X_1, X_3) = a_{13} + a_{12}a_{23}, \tag{8}$$

$$\text{cov}(X_2, X_3) = a_{12}^2a_{23} + a_{12}a_{13} + a_{23}, \tag{9}$$

$$\text{cov}(X_1, X_2 | X_3) = a_{13}a_{23} - a_{12}, \tag{10}$$

$$\text{cov}(X_1, X_3 | X_2) = -a_{13}, \tag{11}$$

$$\text{cov}(X_2, X_3 | X_1) = -a_{23}. \tag{12}$$

711 If we define a_{13}, a_{23}, a_{12} in such a way that:

$$a_{13} * a_{23} - a_{1,2} = 0$$

712 then we get a situation where: nodes 1 and 2 are not d-separated given node 3 in a graph G and $X_1 \perp\!\!\!\perp X_2 | X_3$ which is a violation of faithfulness.

717 B NN-OPT METHOD DETAILS

718
719 **Details of experiments with NN-opt method** In order to test which architecture perform best, we conducted an experiment, training NN-opt method with different sizes of neural networks. The trained models were judged in terms of negative log likelihood and their performance on the task of causal discovery measured as $\text{ESHD}_{\text{CPDAG}}$. For each tested architecture, we performed the search for the best regularization coefficient, the tested coefficients were: $[0.1, 0.3, 1.0]$. Among all models, the best results were consistently obtained for regularization coefficient = 0.3. The learning rate was set to 0.0003. The results of the experiments are shown in Table 2. As we can see, the best, both in case of NLL and $\text{ESHD}_{\text{CPDAG}}$ was model with two layers and hidden dimension of size 8. Notably this is the same architecture, as was used to generate data.

720
721
722
723
724
725
726
727 **Selected hyperparameters:** Number of layers = 2, hidden dimension = 8, regularization coefficient = 0.3.

Model architecture	NLL	$\text{ESHD}_{\text{CPDAG}}$
[4]	0.33 _(0.22, 0.43)	3.63 _(2.83, 4.67)
[4, 4]	0.2 _(0.1, 0.3)	3.15 _(2.0, 4.65)
[4, 4, 4]	0.23 _(0.14, 0.34)	3.03 _(2.33, 4.07)
[8]	0.18 _(0.06, 0.29)	2.13 _(1.43, 3.07)
[8, 8]	0.13 _(0.02, 0.24)	1.23 _(0.77, 1.87)
[8, 8, 8]	0.22 _(0.12, 0.32)	2.77 _(1.97, 3.67)
[16]	0.14 _(0.03, 0.26)	1.77 _(1.1, 2.73)
[16, 16]	0.33 _(0.24, 0.42)	2.4 _(1.0, 4.32)
[16, 16, 16]	0.88 _(0.8, 1.0)	4.0 _(3.07, 4.97)

728
729
730
731
732
733
734
735
736
737
738
739
740
741 Table 2: The performance of NN-opt method models with different architectures. The numbers in the subscripts, correspond to 0.95 confidence intervals. The experiments were performed on 30 graphs.

745 C DETAILS ABOUT BENCHMARK AND EXTENSIONS

746 C.1 DATASET GENERATION DETAILS

747
748
749 The data is generated using a fully connected MLP with two hidden layers of 8 units each, initialized with random weights drawn from a uniform distribution and use the ReLU (Nair & Hinton, 2010) activation function to introduce non-linearity. The neural network models the relationships between variables in the underlying DAG, where each node represents a variable and the edges capture dependencies between these variables. The input variables, which serve as the initial causes in the graph, are sampled from normal distributions. The noise added to the system is sampled from a Gaussian distribution $\mathcal{N}(0, 0.1^2)$, simulating uncertainty in the model. The dataset consists of 100,000 data points, and the data is rescaled to maintain consistency across samples.

C.2 MODEL HYPERPARAMETERS

We performed extensive hyperparameter tuning for all methods. In addition to the MLP architecture grids described in Appendix C.3, the following hyperparameter grids were explored:

DCDI Grid search: Regularization coefficients tested: [0.1, 0.3, 1, 2]. Values below 0.001 or above 5 led to poor performance. **Selected:** Regularization coefficient = 1, learning rate = 0.001, Augmented Lagrangian tolerance = 10^{-8} .

DiBS Grid search: Alpha linear: [0.01, 0.02, 0.05], kernel parameters: h latent: [0.5, 1.0, 2.0], h theta: [20.0, 50.0, 200.0], **step size:**[0.05, 0.03, 0.01, 0.005, 0.003]. **Selected:** Alpha linear = 0.02, h latent = 1.0, h theta = 50.0, **step size = 0.03**.

BayesDAG Grid search: Scale noise: [0.1, 0.01], scale noise p: [0.1, 0.01, 1.0], lambda sparse: [50.0, 100.0, 300.0, 500.0]. **Selected:** Scale noise = 0.1, scale noise p = 0.01, lambda sparse = 500.0.

SDCD Grid search: Constraint modes: ["exp", "spectral radius", "matrix power"]. The $ESHD_{CPDAG}$ metric showed similar results across modes. **Selected:** Spectral radius was chosen for faster computation, with a learning rate of 0.0003.

For each of these method, all other parameters were retained from the original paper or code.

C.3 MODEL ARCHITECTURE COMPARISON WITHIN METHOD

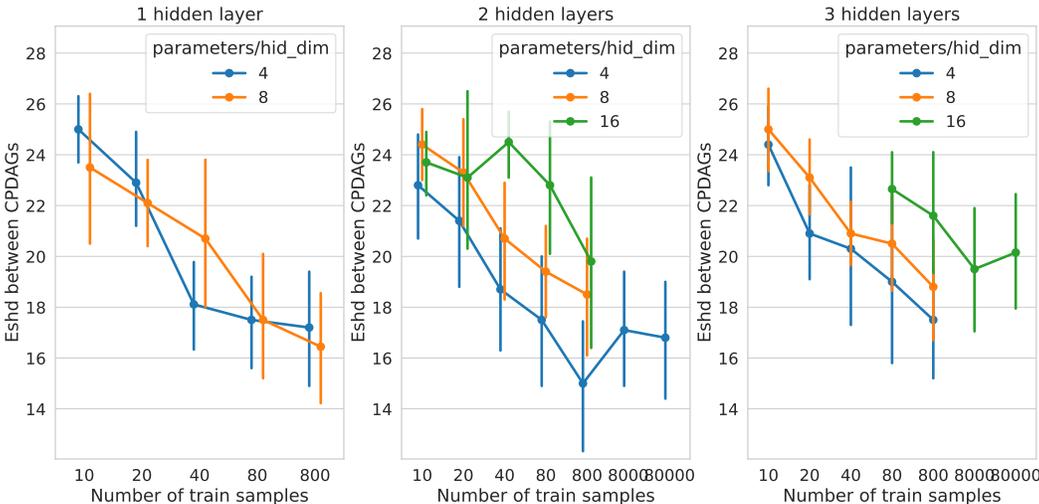
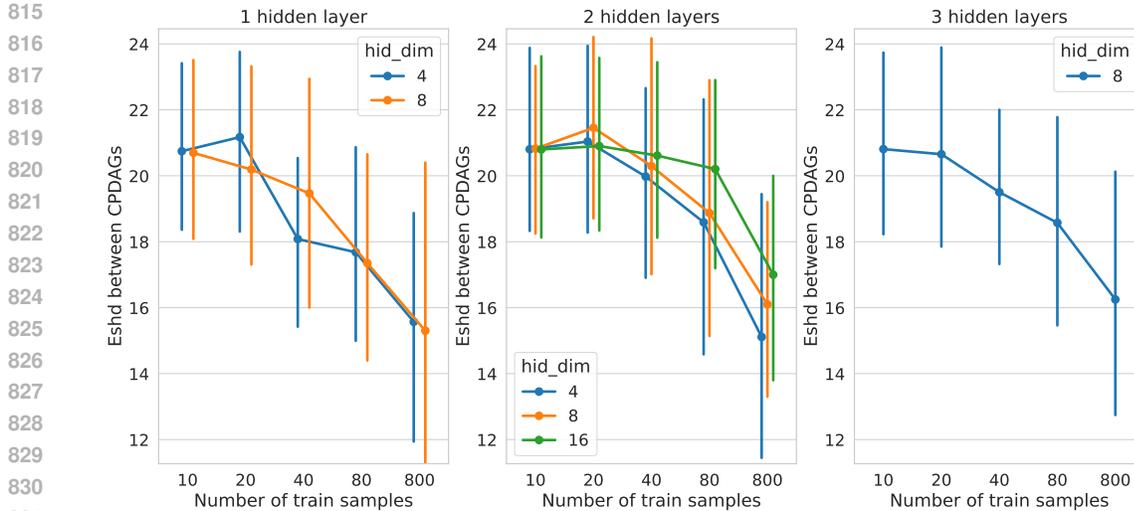


Figure 7: Comparison of the $ESHD_{CPDAG}$ of DCDI for datasets with different observational sample size. The result is based on 10 graphs.

DCDI In Figure 7, we present the performance analysis of the DCDI across various neural network configurations. Our results reveal that the optimal performance is generally achieved by a two-layer model with a hidden dimension of 4. Interestingly, we observe that more expressive models exhibit diminished performance relative to the smaller models.

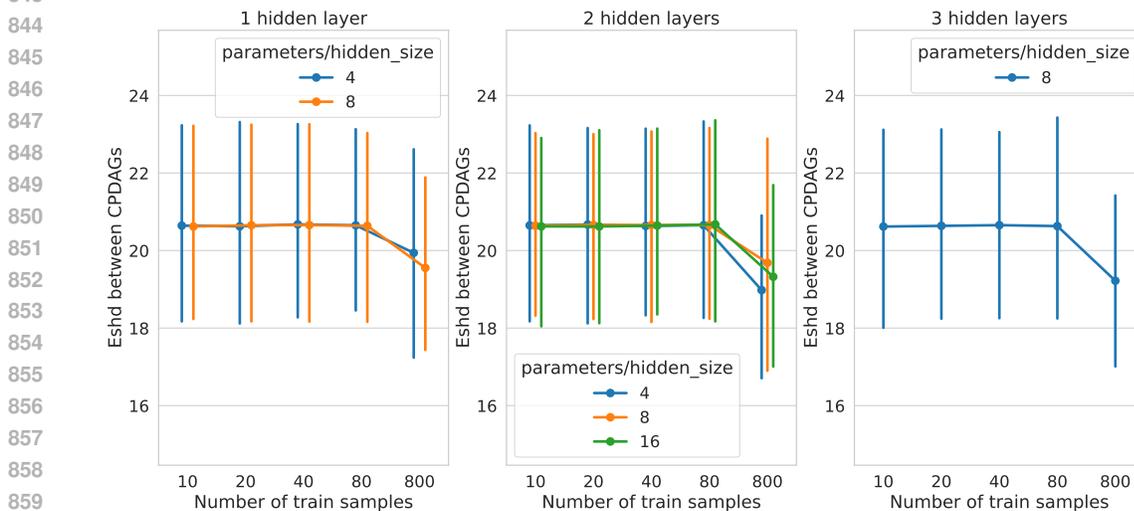
DiBS Figure 8 presents the performance analysis of the DiBS method across various neural network configurations. As with the DCDI method, we evaluate models with different numbers of layers and hidden dimension sizes. Consistent with DCDI, we find that the optimal performance for DiBS is achieved by a two-layer model with a hidden dimension of 4. However, the performance landscape for DiBS exhibits less variability across different model configurations. Single-layer models perform nearly as well as the optimal two-layer model.

810 Furthermore, we observe that more expressive models do not show a significant degradation in
 811 performance as was seen with DCDI. The overall differences in metric across all tested configurations
 812 are relatively small for DiBS, indicating a more consistent performance across varying levels of
 813 model complexity.
 814



833 Figure 8: Comparison of the performance of DiBS depending on the model architecture and number
 834 of samples.
 835

837 **BayesDAG** Figure 9 compares the performance of BayesDAG across different model architectures
 838 and sample sizes. For smaller sample sizes, BayesDAG’s performance remains consistent, with
 839 noticeable differences emerging only at a sample size of 800. This suggests that BayesDAG requires
 840 more data to fully leverage its model capacity, unlike what we observed for DCDI and DiBS, where
 841 performance varied more significantly across sample sizes. Notably, the best-performing architecture
 842 for DiBS is a two-layer MLP with a hidden dimension of 4.
 843



858
 859
 860
 861
 862 Figure 9: Comparison of the performance of DiBS depending on the model architecture and number
 863 of samples.

SDCD Figure 10 presents a similar comparison of SDCD performance across different MLP architectures and sample sizes. Interestingly, the three-layer architectures show stagnant performance regardless of sample size, while the one-layer models exhibit significant improvement as the sample size increases. Overall, the best performance is achieved with a one-layer MLP with 8 hidden units, although it remains comparable to the one-layer MLP with 4 hidden units and the two-layer MLP with 4 hidden units.

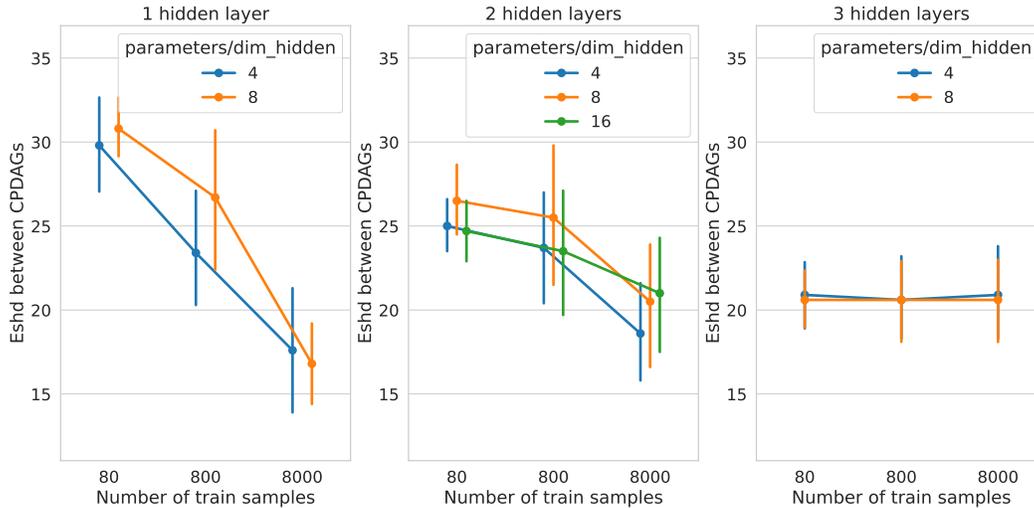


Figure 10: Comparison of the performance of SDCD depending on the model architecture and number of samples.

Model architecture Inspired by BayesDAG, we also implemented layer normalization and residual connections to assess their impact. We conducted additional experiments on both the best-performing model ([4, 4]) and the largest model ([8, 8, 8]). The size of networks was similar to the one proposed in articles introducing tested methods: in DCI it was [16, 16], for SDCD it was [10, 10], for DiBS [5, 5] and for BayesDAG it was a two layer network with a hidden size varying with dimensionality. The results of these tests are presented Figure 11. We show, there is no significant and consistent improvement across all networks, supporting our initial conclusion that variations in MLP architecture have minimal impact on performance.

C.4 INFLUENCE OF SAMPLE SAMPLES ON PERFORMANCE ON THE GRAPH WITH ER(5, 1)

Figure 12 shows the ESH_{CPDAG} of benchmarked methods for different sample sizes. For all observational sample sizes, SDCD and DCI have a large confidence interval. For datasets with 2,500 and 8,000 samples, BayesDAG performs better than other benchmarked methods, getting small confidence interval for 8,000 samples.

C.5 ADDITIONAL RESULTS FOR SDCD AND DiBS

D JUSTIFICATION OF EVALUATION METRICS

We design metrics based on popular SHD, F1-score metrics, which we explain shortly below.

The Structural Hamming Distance. SHD (Tsamardinos et al., 2006) quantifies the difference between the predicted graph and the ground truth graph by counting the number of edge insertions, deletions, and reversals required to transform one into the other. SHD values indicate the degree of error in recovering the true causal structure: lower SHD values signify better predictions, while higher values indicate more significant discrepancies.

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

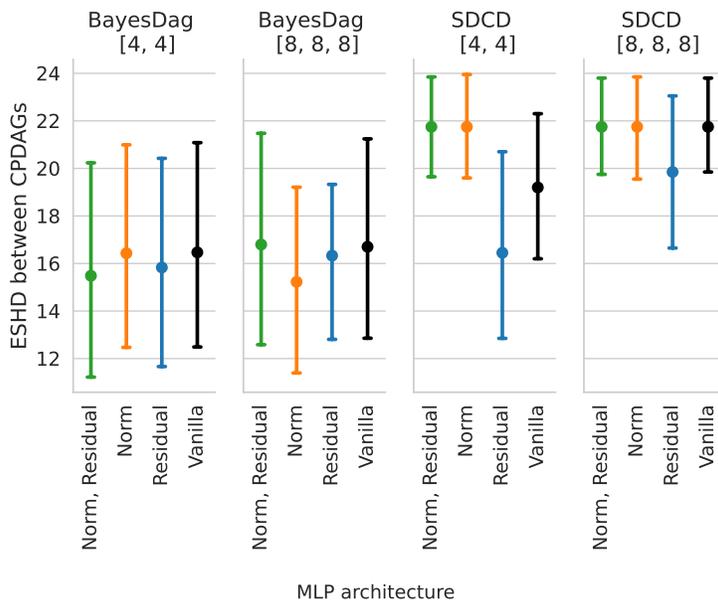


Figure 11: Comparison of the performance of SDCD depending on the model architecture and number of samples.

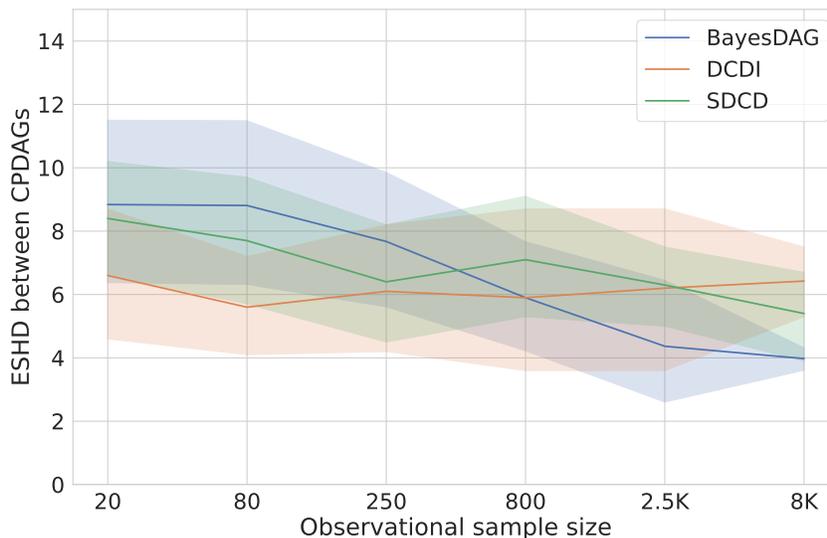


Figure 12: Comparison of $ESH D_{CPDAG}$ for benchmarked methods on ER(5, 1) dataset, averaged over 10 graphs.

The F1-score. The F1-Score measures the harmonic mean of precision and recall for edge predictions, where precision reflects the fraction of correctly predicted edges among all predicted edges, and recall reflects the fraction of correctly predicted edges among the true edges.

We evaluate causal discovery methods based on observational data. In general, in this setup, it is only possible to recover true DAG up to a Markov Equivalence Class, a class of graphs with the same conditional independence relationships, due to identifiability issues TODO cite pearl?. If we were to compare the predicted and ground true graphs using standard metrics like SHD or F1-score we would obtain distorted results — graphs from the MEC class do not generally receive these metrics’ optimal values.

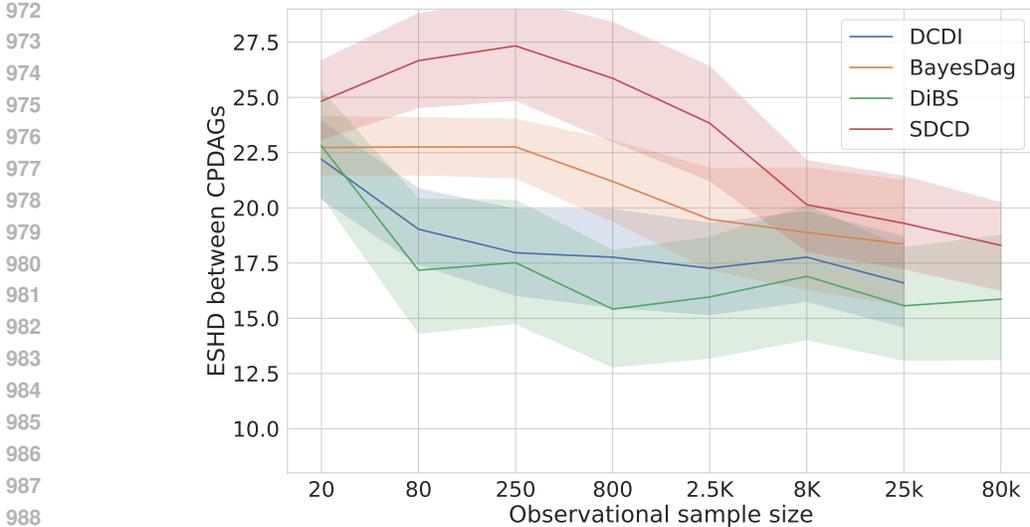


Figure 13: Comparison of ESHDCPDAG for different methods using the [4, 4] architecture, for ER(10, 2) dataset, averaged over 30 samples.

Therefore, we modify the formulation of the metrics to account for the limitations of causal discovery from observational data. We define ESHDCPDAG and $\text{F1-Score}_{\text{CPDAG}}$. These metrics attain their optimal values, 0 and 1 correspondingly, for all DAG from ground truth MEC. Additionally, some of the benchmarked methods are Bayesian thus return the posterior over possible solutions. For those methods, we design metrics that compute the expected value over the posterior and approximate it with the Montecarlo estimator based on a sample of size 100.

We define **Expected SHD between CPDAGs** as:

$$\text{ESHDCPDAG}(\mathcal{G}, \mathbb{G}) = \mathbb{E}_{\mathcal{G}^* \sim \mathbb{G}}[\text{SHD}(\text{CPDAG}(\mathcal{G}), \text{CPDAG}(\mathcal{G}^*))], \tag{13}$$

where \mathbb{G} is the resulting distribution of graphs, \mathcal{G}^* is a graph sampled from \mathbb{G} and \mathcal{G} is the ground true graph. Similarly, we compute the **Expected F1-Score between the CPDAGs**:

$$\text{F1-Score}_{\text{CPDAG}}(\mathcal{G}, \mathbb{G}) = \mathbb{E}_{\mathcal{G}^* \sim \mathbb{G}}[\text{F1-Score}(\text{CPDAG}(\mathcal{G}), \text{CPDAG}(\mathcal{G}^*))]. \tag{14}$$

E JUSTIFICATION OF THE SELECTION OF METHODS

During the preliminary phase, we considered the following methods NO-TEARS (Zheng et al., 2018), NO-BEARS (Lee et al., 2019), NO-CURL (Yu et al., 2021), GRAN-DAG (Lachapelle et al., 2019), SCORE (Rolland et al., 2022), DAGMA (Bello et al., 2022), DCDFG (Lopez et al., 2022), DCDI (Brouillard et al., 2020), DiBS (Lorch et al., 2021), BayesDAG (Annadani et al., 2023), SDCD (Nazaret et al., 2024), from which we chose DCDI, SDCD, DiBS and BayesDAG. Below we explain why the included ones cover non-included methods.

NO-TEARS is the first method to use augmented Lagrangian and differentiable constraints to enforce DAGness. However, the suggested formulation entangles functional and structural parameters, making NO-TEARS applicable only to linear models or restricted neural networks. The NO-TEARS method was improved in GRAN-DAG (introduces separate adjacency matrix and sampling based on Gumbel softmax) and then in DCDI (accounts for interventional data). We chose to use DCDI as it is the most developed method in this line of work and has clean implementation.

An interesting line of work shows articles introducing methods such as NO-BEARS and DAGMA, that were focused on improving the acyclicity constraint introduced in NO-TEARS, all proposed constraints were unified in the SDCD paper, and a new constraint was proposed, that was shown to

perform the best. Additionally, SDCD is compared against SCORE and DCDFG again presenting better performance.

The two other methods are from the class of Bayesian approaches. DiBS method is selected as a Bayesian approach that uses classic NO-TEARS-based regularization embedded in its prior. The BayesDAG is based on the NO-CURL parametrization of DAGs and provides improvements to the optimization pipeline (uses MCMC instead of SVGD).

We argue that this selection of four methods summarizes various research directions and improvements explored in neural causal discovery over the last four years and well represents the spectrum of existing approaches.

F EXPERIMENTS ON REAL-WORLD STRUCTURES

To further substantiate our findings, we conducted additional experiments using Bayesian network structures sourced from the bnlearn repository (Elidan, 2001). This repository provides networks that represent real-world systems. However, the functional relationships in these networks are often limited to simple models, such as linear Gaussian or discrete distributions. To address this limitation, we utilized the graph structures from bnlearn but generated functional relationships consistent with those used in the synthetic benchmark.

We selected `cancer` and `sachs` structures. `cancer` has 5 nodes and 4 edges and we employed the set of hyperparameters that yielded the best performance for the ER(5, 1) for each method. `sachs` has 11 nodes and 17 edges and we employed the set of hyperparameters that yielded the best performance for the ER(10, 2) for each method.

The results, presented in Figure 14, align with the observations detailed in Section 3. Across all methods, we observed either consistently poor performance regardless of sample size or very slow improvements, exhibiting diminishing returns with increasing data.

Most methods significantly underperformed compared to the NNOpt approach. An exception is DiBS, which achieved results comparable to the upper bound on the `cancer` graph, a behavior similar to its performance on the ER(5, 1) class of graphs.

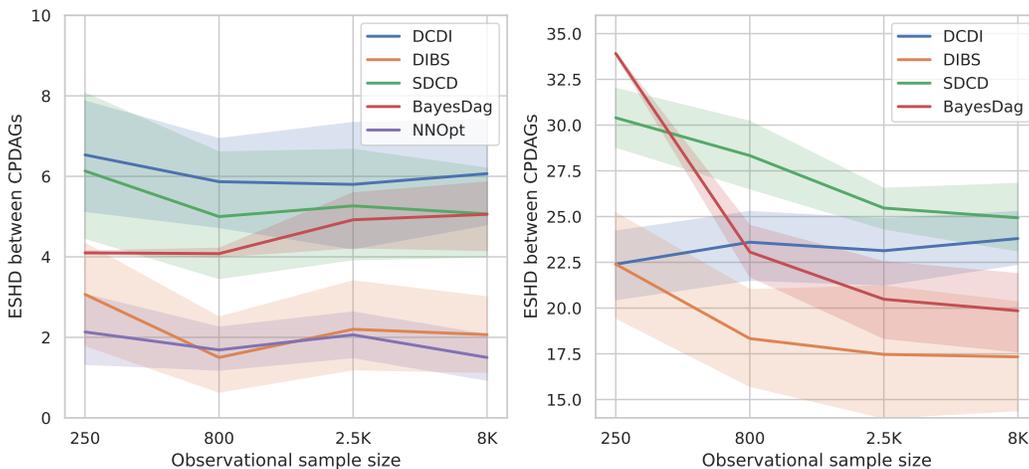


Figure 14: *On the left:* Results of the benchmark on `cancer` structure. *On the right:* Results of the benchmark on `sachs` structure. In both plots, the 95% bootstrap confidence interval is provided as a shaded area. The results are computed on 15 distributions.