

---

# BayesDAG: Gradient-Based Posterior Sampling for Causal Discovery

---

Yashas Annadani<sup>\*1,2</sup> Nick Pawlowski<sup>3</sup> Joel Jennings<sup>3</sup> Stefan Bauer<sup>1,4,5</sup> Cheng Zhang<sup>3</sup> Wenbo Gong<sup>\*3</sup>

## Abstract

Bayesian causal discovery aims to infer the posterior distribution over causal models from observed data, quantifying epistemic uncertainty and benefiting downstream tasks. However, computational challenges arise due to joint inference over combinatorial space of Directed Acyclic Graphs (DAGs) and nonlinear functions. In this work, we introduce a scalable Bayesian causal discovery framework based on stochastic gradient Markov Chain Monte Carlo (SG-MCMC) that directly samples DAGs from the posterior without any DAG regularization, simultaneously draws function parameter samples and is applicable to both linear and nonlinear causal models. To enable our approach, we derive a novel equivalence to the permutation-based DAG learning, which opens up possibilities of using any relaxed gradient estimator defined over permutations. To our knowledge, this is the first framework applying gradient-based MCMC sampling for causal discovery. Empirical evaluations on synthetic and real-world datasets demonstrate our approach’s effectiveness compared to state-of-the-art baselines.

## 1. Introduction

The quest for discovering causal relationships in data-generating processes lies at the heart of empirical sciences and decision-making (Pe’er et al., 2001; Sachs et al., 2005; Van Koten & Gray, 2006). Structural Causal Models (SCMs) (Pearl, 2009) and their associated Directed Acyclic Graphs (DAGs) provide a mathematical framework for modeling such relationships. Knowledge of the underlying SCM permits predictions of unseen interventions and causal reasoning, thus making causal discovery – learning an unknown SCM and its associated DAG from observed data

<sup>\*</sup>Equal contribution. This work was done during YA’s internship at Microsoft Research. <sup>1</sup>Helmholtz AI <sup>2</sup>KTH Stockholm <sup>3</sup>Microsoft Research <sup>4</sup>TU Munich <sup>5</sup>CIFAR Azrieli Global Scholar. Correspondence to: Wenbo Gong <wenbogong@microsoft.com>.

– a subject of extensive research (Peters et al., 2017).

In contrast to traditional methods that infer a single graph or its Markov equivalence class (MEC) (Chickering, 2002; Spirtes et al., 2000), Bayesian causal discovery (Friedman & Koller, 2003; Heckerman et al., 2006; Tong & Koller, 2001) aims to infer a posterior distribution over SCMs and their DAGs from observed data. This approach encapsulates the epistemic uncertainty, degree of confidence in every causal hypothesis, which is particularly valuable for real-world applications when data is scarce. It is also beneficial for downstream tasks such as experimental design (Annadani et al., 2023; Murphy, 2001; Tigas et al., 2022).

The central challenge in Bayesian causal discovery lies in inferring the posterior distribution over the union of the exponentially growing (discrete) DAGs space and (continuous) causal model parameters. Prior works have used Markov Chain Monte Carlo (MCMC) to directly sample DAGs or bootstrap traditional discovery methods (Chickering, 2002; Murphy, 2001; Tong & Koller, 2001), but these methods are typically limited to linear models which admit closed-form marginalization over continuous parameters. Recent approaches have begun to utilize gradient information for more efficient inference. These approaches are based on either the DAG regularizer (Zheng et al., 2018) (for e.g. DIBS (Lorch et al., 2021)) with continuous relaxation of adjacency matrices, or permutation based methods which directly infer permutation matrices over nodes through Variational Inference (VI) (Cundy et al., 2021). DAG regularizer based methods are usually computationally expensive due to the DAG constraint and cannot guarantee DAG generation, while existing permutation based methods are usually restricted to only linear causal models.

In this work, we propose BayesDAG, which overcomes the above limitations. Our contributions are:

1. We prove that an augmented space of edge beliefs and node potentials  $(\mathbf{W}, \mathbf{p})$ , similar to NoCurl (Yu et al., 2021), permits equivalent Bayesian inference in DAG space without the need for any regularizer.
2. We derive an equivalence relation from this augmented space to permutation-based DAG learning which provides a general framework for gradient-based posterior inference.

- Based on this general framework, we propose a scalable Bayesian causal discovery approach with SG-MCMC for DAGs and causal model parameters. Our approach is model-agnostic for linear and non-linear cases and also offers improved inference quality with sampling approach (Ma et al., 2015).

## 2. Background

**Structural Causal Model** The causal relationships among  $d$  variables  $\mathbf{X} \in \mathbb{R}^d$  can be represented by a Structural Causal Model (SCM) which consists of a set of structural equations (Peters et al., 2017) where each variable  $X_i$  is a function of its direct causes  $\mathbf{X}_{\mathbf{pa}^i}$  and an exogenous noise variable  $\epsilon_i$  (for e.g. Gaussian):

$$X_i := f_i(\mathbf{X}_{\mathbf{pa}^i}) + \epsilon_i \quad \text{with } \epsilon_i \sim \mathcal{N}(0, \sigma_i) \quad (1)$$

These equations induce a causal graph  $\mathbf{G} = (\mathbf{V}, \mathbf{E})$ , typically assumed to be a DAG (Pearl, 2009). A directed edge between a node pair  $v_i, v_j \in \mathbf{V}$  (i.e.,  $v_i \rightarrow v_j$ ) represents that  $X_i$  causes  $X_j$ . We use the binary adjacency matrix  $\mathbf{G} \in \{0, 1\}^{d \times d}$  to represent the causal graph. Given a dataset  $\mathbf{D} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}$ , DAG  $\mathbf{G}$  and SCM functional parameters  $\Theta$ , they induce a unique joint distribution  $p(\mathbf{D}, \Theta, \mathbf{G}) = p(\mathbf{D}|\mathbf{G}, \Theta)p(\mathbf{G}, \Theta)$  with the prior  $p(\mathbf{G}, \Theta)$  and likelihood  $p(\mathbf{D}|\mathbf{G}, \Theta)$  (Friedman & Koller, 2003). Bayesian causal discovery aims to infer the posterior  $p(\mathbf{G}, \Theta|\mathbf{D}) = p(\mathbf{D}, \Theta, \mathbf{G})/p(\mathbf{D})$ . However, this posterior is intractable due to super-exponential growth of the possible DAGs  $\mathbf{G}$  (Robinson, 1973) and continuously valued model parameters  $\Theta$  of nonlinear functions. VI (Zhang et al., 2018) or SG-MCMC (Gong, 2022; Ma et al., 2015) are two types of methods developed to tackle general Bayesian inference problems, but adaptations are required for Bayesian causal discovery.

**NoCurl Characterization** Inferring causal graphs is challenging due to the DAG constraint. Recently, (Yu et al., 2021) introduced NoCurl, a novel characterization of the **weighted DAG** space. They define a potential  $p_i \in \mathbb{R}$  for each node  $i$ , grouped as potential vector  $\mathbf{p} \in \mathbb{R}^d$ . They also introduce a gradient operator on  $\mathbf{p}$ , i.e.  $(\text{grad } \mathbf{p})(i, j) = p_i - p_j$ . Based on the above operation, a mapping that directly maps from an augmented space  $(\mathbf{W}, \mathbf{p})$  to the DAG space,  $\gamma(\cdot, \cdot) : \mathbb{R}^{d \times d} \times \mathbb{R}^d \rightarrow \mathbb{R}^{d \times d}$ , was proposed:

$$\gamma(\mathbf{W}, \mathbf{p}) = \mathbf{W} \odot \text{ReLU}(\text{grad } \mathbf{p}) \quad (2)$$

where  $\text{ReLU}(\cdot)$  is the ReLU activation function and  $\mathbf{W}$  is a skew-symmetric **continuously weighted** matrix. This formulation is complete (Theorem 2.1 in (Yu et al., 2021)), as any continuously weighted DAG can be represented by a  $(\mathbf{W}, \mathbf{p})$  pair and vice versa. NoCurl translates the learning of a single weighted DAG to a corresponding  $(\mathbf{W}, \mathbf{p})$  pair.

However, direct gradient-based optimization is challenging due to a highly non-convex loss landscape, which leads to the reported failure in (Yu et al., 2021). Although NoCurl appears suitable for our purpose, the failure in directly learning suggests non-trivial optimizations. We hypothesize that this arises from the continuously weighted matrix  $\mathbf{W}$ . In the following, we introduce our proposed parametrization inspired by NoCurl to characterize the **binary** DAG adjacency matrix.

## 3. Sampling the DAGs

### 3.1. Bayesian Inference in $(\mathbf{W}, \mathbf{p})$ Space

The NoCurl formulation (Equation (2)) focuses on learning a *single weighted* DAG, which is not directly useful for our purpose. We need to address two key questions: (1) considering only binary adjacency matrices without weights; (2) ensuring Bayesian inference in  $(\mathbf{W}, \mathbf{p})$  is valid.

Consider the mapping  $\tau : \{0, 1\}^{d \times d} \times \mathbb{R}^d \rightarrow \{0, 1\}^{d \times d}$ :

$$\tau(\mathbf{W}, \mathbf{p}) = \mathbf{W} \odot \text{Step}(\text{grad } \mathbf{p}) \quad (3)$$

where we abuse the term  $\mathbf{W}$  for binary matrices, and replace  $\text{ReLU}(\cdot)$  with  $\text{Step}(\cdot)$ .  $\mathbf{W}$  acts as mask to disable the edge existence. Thus, due to the  $\text{Step}$ ,  $\tau$  can only output a binary adjacency matrix. We show that performing Bayesian inference in this modified augmented  $(\mathbf{W}, \mathbf{p})$  space is valid, i.e., using the posterior  $p(\mathbf{W}, \mathbf{p}|\mathbf{D})$  to replace  $p(\mathbf{G}|\mathbf{D})$ . This differs from NoCurl, which focuses on a single graph rather than the validity for Bayesian inference, requiring a new theory for soundness. We provide the main results and relegate all proofs and details to the appendix.

**Theorem 3.1** (Equivalence of inference in  $(\mathbf{W}, \mathbf{p})$  and binary DAG space). *Assume graph  $\mathbf{G}$  is a binary adjacency matrix representing a DAG and node potential  $\mathbf{p}$  does not contain the same values, i.e.  $p_i \neq p_j \forall i, j$ . Then, with the induced joint observational distribution  $p(\mathbf{D}, \mathbf{G})$ , dataset  $\mathbf{D}$ , and a corresponding prior  $p(\mathbf{G})$ , we have*

$$p(\mathbf{G}|\mathbf{D}) = \int p_\tau(\mathbf{p}, \mathbf{W}|\mathbf{D}) \mathbb{1}(\mathbf{G} = \tau(\mathbf{W}, \mathbf{p})) d\mathbf{W} d\mathbf{p} \quad (4)$$

if  $p(\mathbf{G}) = \int p_\tau(\mathbf{p}, \mathbf{W}) \mathbb{1}(\mathbf{G} = \tau(\mathbf{W}, \mathbf{p})) d\mathbf{W} d\mathbf{p}$ , where  $p_\tau(\mathbf{W}, \mathbf{p})$  is the prior,  $\mathbb{1}(\cdot)$  is the indicator function, and  $p_\tau(\mathbf{p}, \mathbf{W}|\mathbf{D})$  is the posterior distribution over  $\mathbf{p}, \mathbf{W}$ .

This theorem guarantees that instead of performing inference directly in the constrained space (i.e. DAG space), we can apply Bayesian inference in a less complex  $(\mathbf{W}, \mathbf{p})$  space where  $\mathbf{W} \in \{0, 1\}^{d \times d}$  and  $\mathbf{p} \in \mathbb{R}^d$ .

For inference of  $\mathbf{p}$ , we adopt a sampling-based approach, which is asymptotically accurate (Ma et al., 2015) and provides better inference quality compared to VI (Gong et al.,

2019; Springenberg et al., 2016; Trippe & Turner, 2018). In particular, we consider SG-MCMC (refer to Section 4), which avoids the expensive Metropolis-Hastings acceptance step and scales to large datasets. We emphasize that any other suitable sampling algorithms can be directly plugged in, thanks to the generality of the framework.

### 3.2. Equivalent Formulation

The mapping  $\tau$  does not provide meaningful gradient information for  $\mathbf{p}$  due to the piecewise constant  $\text{Step}(\cdot)$  function. We address this issue by deriving an equivalence to a permutation learning problem that enables approximating the gradient of  $\mathbf{p}$ .

**Intuition** The node potential  $\mathbf{p}$  implicitly defines a topological ordering through the mapping  $\text{Step}(\text{grad}(\cdot))$ . In particular,  $\text{grad}(\cdot)$  outputs a skew-symmetric adjacency matrix, where each entry specifies the potential difference between nodes.  $\text{Step}(\text{grad}(\cdot))$  zeros out the negative potential differences (i.e.  $p_i \leq p_j$ ), and only permits the edge direction from higher potential to the lower one (i.e.  $p_i > p_j$ ). This implicitly defines a sorting operation based on the node potentials, which can be cast as a particular  $\arg \max$  problem (Blondel et al., 2020; Kuhn, 1955; Mena et al., 2018; Niculae et al., 2018) involving a permutation matrix.

**Alternative formulation** We define  $\mathbf{L} \in \{0, 1\}^{d \times d}$  as a matrix with lower triangular part to be 1, and vector  $\mathbf{o} = [1, \dots, d]$ . We propose the following formulation:

$$\mathbf{G} = \mathbf{W} \odot [\boldsymbol{\sigma}(\mathbf{p})\mathbf{L}\boldsymbol{\sigma}(\mathbf{p})^T], \boldsymbol{\sigma}(\mathbf{p}) = \arg \max_{\boldsymbol{\sigma}' \in \Sigma_d} \mathbf{p}^T(\boldsymbol{\sigma}'\mathbf{o}) \quad (5)$$

Here,  $\Sigma_d$  represents the space of all  $d$  dimensional permutation matrices. The following theorem states the equivalence of this formulation to Equation (3).

**Theorem 3.2** (Equivalence to NoCurl formulation). *Assuming the conditions in Theorem 3.1 are satisfied. Then, for a given  $(\mathbf{W}, \mathbf{p})$ , we have*

$$\mathbf{G} = \mathbf{W} \odot \text{Step}(\text{grad} \mathbf{p}) = \mathbf{W} \odot [\boldsymbol{\sigma}(\mathbf{p})\mathbf{L}\boldsymbol{\sigma}(\mathbf{p})^T]$$

where  $\mathbf{G}$  is a DAG and  $\boldsymbol{\sigma}(\mathbf{p})$  is defined in Equation (5).

This theorem translates our proposed operator  $\text{Step}(\text{grad}(\mathbf{p}))$  into finding a corresponding permutation matrix  $\boldsymbol{\sigma}(\mathbf{p})$ . Although this does not directly solve the uninformative gradient, it opens the door for approximating this gradient with the tools from the differentiable permutation literature (Blondel et al., 2020; Mena et al., 2018; Niculae et al., 2018). For simplicity, we adopt the relaxed Gumbel-Sinkhorn approach (Mena et al., 2018), but we emphasize that this equivalence is general enough that any past or future approximation methods can be

easily applied. To get the binary permutation matrix, we apply the Hungarian algorithm (Munkres, 1957) and use a straight-through estimator for  $\mathbf{p}$ .

Some of the previous works (Charpentier et al., 2022; Cundy et al., 2021) have leveraged the Sinkhorn operator to model variational distributions over permutation matrices. However, they start with a full rank  $\mathbf{M}$ , which has been reported to require over **1000** Sinkhorn iterations to converge (Cundy et al., 2021). However, our formulation, based on explicit node potential  $\mathbf{p}\mathbf{o}^T$ , generates a rank-1 matrix, requiring much fewer Sinkhorn steps (around **300**) in practice, saving two-thirds of the computational cost.

## 4. Bayesian Causal Discovery via Sampling

In this section, we delve into two inference algorithms that are derived from the proposed framework. The first one, which will be our main focus, combines SG-MCMC and VI in a Gibbs sampling manner. The second one, which is based entirely on SG-MCMC with continuous relaxation, is also derived, but we include its details in Appendix A due to its inferior empirical performance.

**Model Formulation** We build upon the model formulation of (Geffner et al., 2022), which combines the additive noise model with neural networks. Specifically,  $X_i := f_i(\mathbf{X}_{\mathbf{Pa}^i}) + \epsilon_i$ , where  $f_i$  adheres to the adjacency relation specified by  $\mathbf{G}$ , i.e.  $\partial f_i(\mathbf{x})/\partial x_j = 0$  if no edge exists between nodes  $i$  and  $j$ . We define  $f_i$  as

$$f_i(\mathbf{x}) = \zeta_i \left( \sum_{j=1}^d G_{ji} l_j(x_j) \right) \quad (6)$$

where  $\zeta_i$  and  $l_i$  are neural networks with parameters  $\Theta$ , and  $\mathbf{G}$  serves as a mask disabling non-parent values. To reduce the number of neural networks, we adopt a weight-sharing mechanism:  $\zeta_i(\cdot) = \zeta(\mathbf{u}_i, \cdot)$  and  $l_i(\cdot) = l(\mathbf{u}_i, \cdot)$ , with trainable node embeddings  $\mathbf{u}_i$ .

**Likelihood of SCM** The likelihood can be evaluated through the noise variable (Geffner et al., 2022):

$$p(\mathbf{x}|\mathbf{G}) = \prod_{i=1}^d p_{\epsilon_i}(x_i - f_i(\mathbf{x}_{\mathbf{Pa}_G^i})) \quad (7)$$

**Prior design** We implicitly define the prior  $p(\mathbf{G})$  via  $p(\mathbf{p}, \mathbf{W})$ . We propose the following for the joint prior:

$$p(\mathbf{W}, \mathbf{p}, \Theta) \propto \mathcal{N}(\Theta; \mathbf{0}, \mathbf{I})\mathcal{N}(\mathbf{p}; \mathbf{0}, \alpha\mathbf{I})\mathcal{N}(\mathbf{W}; \mathbf{0}, \mathbf{I}) \exp(-\lambda_s \|\tau(\mathbf{W}, \mathbf{p})\|_F^2)$$

where  $\alpha$  controls the initialization scale of  $\mathbf{p}$  and  $\lambda_s$  controls the sparseness of  $\mathbf{G}$ .

#### 4.1. Bayesian Inference of $W, p, \Theta$

The main challenge lies in the binary nature of  $W \in \{0, 1\}^{d \times d}$ . We propose a combination of SG-MCMC for  $p, \Theta$  and VI for  $W$ . It should be noted that our framework can incorporate any suitable discrete sampler if needed. We employ a Gibbs sampling procedure (Casella & George, 1992), which iteratively applies (1) sampling  $p, \Theta \sim p(p, \Theta | D, W)$  with SG-MCMC; (2) updating the variational posterior  $q_\phi(W | p, D) \approx p(W | p, \Theta, D)$ .

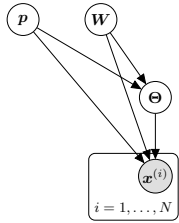


Figure 1: Graphical model of the inference problem.

We define the posterior  $p(p, \Theta | D, W) \propto \exp(-U(p, W, \Theta))$ , where  $U(p, W, \Theta) = -\log p(p, D, W, \Theta)$ . SG-MCMC in continuous time defines a specific form of Itô diffusion that maintains the target distribution invariant (Ma et al., 2015) without the expensive computation of the MH step. We adopt the Euler-Maruyama discretization for simplicity and employ the sampler based on (Gong et al., 2019), which is inspired by Adam (Kingma & Ba, 2014). Detailed update equations and gradient derivations can be found in Appendix C.

**Variational inference for  $W$**  We use the variational posterior  $q_\phi(W | p)$  to approximate the true posterior  $p(W | p, \Theta, D)$ . Specifically, we select an independent Bernoulli distribution with logits defined by the output of a neural network  $\mu_\phi(p)$ :

$$q_\phi(W | p) = \prod_{ij} \text{Ber}(\mu_\phi(p)_{ij}) \quad (8)$$

To train  $q_\phi$ , we derive the corresponding *evidence lower bound* (ELBO):

$$\text{ELBO}(\phi) = \mathbb{E}_{q_\phi(W | p)} [\log p(D, p, \Theta | W)] - D_{\text{KL}} [q_\phi(W | p) \| p(W)] \quad (9)$$

where  $D_{\text{KL}}$  is the Kullback-Leibler divergence. The derivation is in Appendix B.6. Algorithm 2 summarizes this inference procedure.

**SG-MCMC with continuous relaxation** Furthermore, we explore an alternative formulation that circumvents the need for variational inference. Instead, we employ SG-MCMC to sample  $\tilde{W}$ , a continuous relaxation of  $W$ , facilitating a fully sampling-based approach. For a detailed formulation, please refer to Appendix A. We report its performance in Appendix G.2, which surprisingly is inferior to SG-MCMC+VI. We hypothesize that coupling  $W, p$  through

$\mu_\phi$  is important since changes in  $p$  results in changes of the permutation matrix  $\sigma(p)$ , which should also influence  $W$  accordingly during posterior inference. However, through sampling  $\tilde{W}$  with few SG-MCMC steps, this change cannot be immediately reflected, resulting in inferior performance. Thus, we focus only on the performance of SG-MCMC+VI for our experiments.

**Computational complexity** Our proposed SG-MCMC+VI offers a notable improvement in computational cost compared to existing approaches, such as DIBS (Lorch et al., 2021). The computational complexity of our method is  $O(BN_p + N_p d^3)$ , where  $B$  represents the batch size and  $N_p$  is the number of parallel SG-MCMC chains. This former term stems from the forward and backward passes, and the latter comes from the Hungarian algorithm, which can be parallelized to further reduce computational cost. In comparison, DIBS has a complexity of  $O(N_p^2 N + N_p d^3)$  with  $N \gg B$  being the full dataset size. This is due to the kernel computation involving the entire dataset and the evaluation of the matrix exponential in the DAG regularizer (Zheng et al., 2018). As a result, our approach provides linear scalability w.r.t.  $N_p$  with substantially smaller batch size  $B$ . Conversely, DIBS exhibits quadratic scaling in terms of  $N_p$  and lacks support for mini-batch gradients.

## 5. Experiments

We aim to study empirically the following aspects: (1) posterior inference quality of BayesDAG in high dimensional nonlinear causal models with synthetic datasets and (2) performance in semi-synthetic and real world applications.

**Baselines and Metrics.** We mainly compare BayesDAG with the following baselines: Bootstrap GES (BGES) (Chickering, 2002; Friedman et al., 2013), BCD Nets (Cundy et al., 2021), Differentiable DAG Sampling (DDS (Charpentier et al., 2022)) and DIBS (Lorch et al., 2021). For evaluation, we compute the expected SHD (E-SHD), expected orientation F1 score (Edge F1) and negative log-likelihood of the held-out data (NLL). Our synthetic data generation and evaluation protocol follows prior work (Annadani et al., 2021; Geffner et al., 2022; Lorch et al., 2021). All the experimental details, including hyperparameters are provided in Appendix F.

### 5.1. Evaluation on Synthetic Data

We evaluate our method on synthetic data with nonlinear functional relations, where ground truth graphs are known. We generate data by randomly sampling DAGs from Erdos-Rényi (ER) (Erdős et al., 1960) or Scale-Free (SF) (Barabási & Albert, 1999) graphs and drawing at random ground truth MLP parameters for nonlinear functions. We assess perfor-



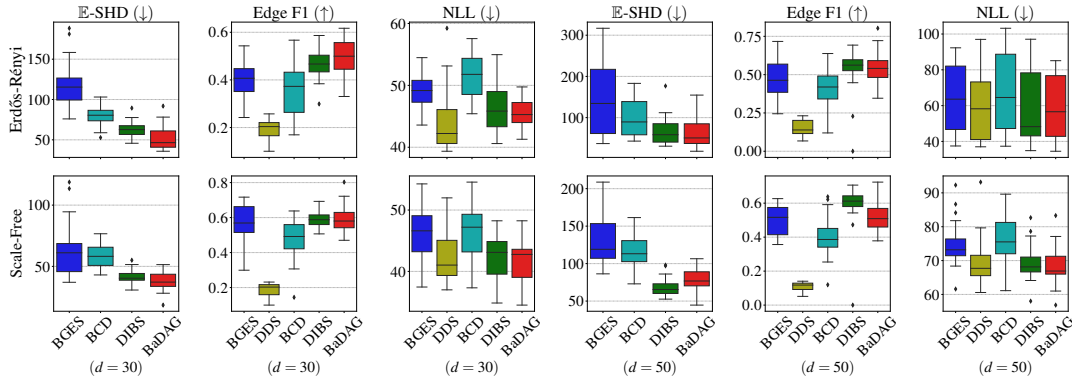


Figure 2: Posterior inference of both graph and nonlinear functional parameters on synthetic datasets of nonlinear causal models with  $d = 30$  and  $d = 50$  variables. BayesDAG gives best results across most metrics and outperforms other permutation based approaches (BCD and DDS).  $\downarrow$  denotes lower is better and  $\uparrow$  denotes higher is better.

mance on 30 random datasets.

Our approach is the first to attempt full posterior inference in nonlinear models using permutation-based inference. Results for  $d = 30$  variables in Figure 2 demonstrate that

Table 1:  $\mathbb{E}$ -SHD (with 95% CI) for ER graphs in higher dimensional nonlinear causal models. DIBS becomes computationally prohibitive for  $d > 50$ .

	$d = 70$	$d = 100$
BGES	$355.77 \pm 18.02$	$563.02 \pm 27.21$
BCD	$217.05 \pm 9.58$	$362.66 \pm 29.18$
DIBS	N/A	N/A
BaDAG	<b><math>143.70 \pm 11.61</math></b>	<b><math>295.92 \pm 24.67</math></b>

BayesDAG significantly outperforms other *permutation-based approaches* and DIBS in most of the metrics. For  $d = 50$ , BayesDAG performs comparably to DIBS in ER but a little worse in SF. However, our method achieves better NLL on held-out data compared to most baselines including DIBS for  $d = 30, 50$ , ER and SF settings.

We additionally evaluate on  $d \in \{70, 100\}$  variables. We find that our method consistently outperforms the baselines with  $d = 70$  and in terms of  $\mathbb{E}$ -SHD with  $d = 100$ . Full results are presented in Appendix G.1. Competitive performance for  $d > 50$  in nonlinear settings further demonstrates the applicability and computational efficiency of the proposed approach. In contrast, the only fully Bayesian nonlinear method, DIBS, is not computationally efficient to run for  $d > 50$ .

## 5.2. Applications

**Evaluation on Semi-Synthetic Data** We evaluate our method on the SynTReN simulator (Van den Bulcke et al., 2006). This simulator creates synthetic transcriptional regulatory networks and produces simulated gene expression data that approximates real experimental data. Table 2

Table 2: Results (with 95% confidence intervals) on Syntren (semi-synthetic) and Sachs Protein Cells (real-world) datasets.  $\downarrow$  denotes lower is better and  $\uparrow$  denotes higher is better.

	Syntren ( $d = 20$ )		Sachs Protein Cells ( $d = 11$ )	
	$\mathbb{E}$ -SHD ( $\downarrow$ )	Edge F1 ( $\uparrow$ )	$\mathbb{E}$ -SHD ( $\downarrow$ )	Edge F1 ( $\uparrow$ )
BGES	$66.18 \pm 9.47$	<b><math>0.21 \pm 0.05</math></b>	<b><math>16.61 \pm 0.44</math></b>	$0.22 \pm 0.02$
DDS	$134.37 \pm 4.58$	$0.13 \pm 0.02$	$34.90 \pm 0.73$	$0.21 \pm 0.02$
BCD	$38.38 \pm 7.12$	$0.15 \pm 0.07$	$17.05 \pm 1.93$	$0.20 \pm 0.08$
DIBS	$46.43 \pm 4.12$	$0.16 \pm 0.02$	$22.3 \pm 0.31$	$0.20 \pm 0.01$
BaDAG	<b><math>34.21 \pm 2.82</math></b>	<b><math>0.20 \pm 0.02</math></b>	$18.92 \pm 1.0$	<b><math>0.26 \pm 0.04</math></b>

presents the results of all the methods. We find that our method recovers the true network much better in terms of  $\mathbb{E}$ -SHD as well as Edge F1 compared to baselines.

**Evaluation on Real Data** We also evaluate on a real dataset which measures the expression level of different proteins and phospholipids in human cells (called the Sachs Protein Cells Dataset) (Sachs et al., 2005). The data corresponds to a network of protein-protein interactions of 11 different proteins with 17 edges in total among them. There are 853 observational samples in total, from which we bootstrap 800 samples of 5 different datasets. Results in Table 2 demonstrate that our method performs well as compared to the baselines, proving the suitability of the proposed sampling framework to real-world settings.

## 6. Discussion

In this work, we propose BayesDAG, a novel, scalable Bayesian causal discovery framework that employs SG-MCMC (and VI) to infer causal models. We establish the validity of performing Bayesian inference in the augmented  $(W, p)$  space and demonstrate its connection to permutation-based DAG learning. Our method offers direct DAG sampling in nonlinear models while demonstrating superior inference accuracy.

## References

- Adams, R. P. and Zemel, R. S. Ranking via sinkhorn propagation. *arXiv preprint arXiv:1106.1925*, 2011.
- Agrawal, R., Squires, C., Yang, K., Shanmugam, K., and Uhler, C. Abcd-strategy: Budgeted experimental design for targeted causal structure discovery. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 3400–3409. PMLR, 2019.
- Annadani, Y., Rothfuss, J., Lacoste, A., Scherrer, N., Goyal, A., Bengio, Y., and Bauer, S. Variational causal networks: Approximate bayesian inference over causal structures. *arXiv preprint arXiv:2106.07635*, 2021.
- Annadani, Y., Tigas, P., Ivanova, D. R., Jesson, A., Gal, Y., Foster, A., and Bauer, S. Differentiable multi-target causal bayesian experimental design. *arXiv preprint arXiv:2302.10607*, 2023.
- Barabási, A.-L. and Albert, R. Emergence of scaling in random networks. *science*, 286(5439):509–512, 1999.
- Bengio, Y., Lahlou, S., Deleu, T., Hu, E. J., Tiwari, M., and Bengio, E. Gflownet foundations. *arXiv preprint arXiv:2111.09266*, 2021.
- Blondel, M., Teboul, O., Berthet, Q., and Djolonga, J. Fast differentiable sorting and ranking. In *International Conference on Machine Learning*, pp. 950–959. PMLR, 2020.
- Casella, G. and George, E. I. Explaining the gibbs sampler. *The American Statistician*, 46(3):167–174, 1992.
- Charpentier, B., Kibler, S., and Günnemann, S. Differentiable dag sampling. *arXiv preprint arXiv:2203.08509*, 2022.
- Chickering, D. M. Optimal structure identification with greedy search. *Journal of machine learning research*, 3 (Nov):507–554, 2002.
- Cooper, G. F. and Herskovits, E. A bayesian method for the induction of probabilistic networks from data. *Machine learning*, 9:309–347, 1992.
- Cundy, C., Grover, A., and Ermon, S. Bcd nets: Scalable variational approaches for bayesian causal discovery. *Advances in Neural Information Processing Systems*, 34: 7095–7110, 2021.
- Deleu, T., Góis, A., Emezue, C., Rankawat, M., Lacoste-Julien, S., Bauer, S., and Bengio, Y. Bayesian structure learning with generative flow networks. In *Uncertainty in Artificial Intelligence*, pp. 518–528. PMLR, 2022.
- Eaton, D. and Murphy, K. Bayesian structure learning using dynamic programming and mcmc. *arXiv preprint arXiv:1206.5247*, 2012.
- Erdős, P., Rényi, A., et al. On the evolution of random graphs. *Publ. Math. Inst. Hung. Acad. Sci.*, 5(1):17–60, 1960.
- Friedman, N. and Koller, D. Being bayesian about network structure. a bayesian approach to structure discovery in bayesian networks. *Machine learning*, 50(1):95–125, 2003.
- Friedman, N., Goldszmidt, M., and Wyner, A. Data analysis with bayesian networks: A bootstrap approach. *arXiv preprint arXiv:1301.6695*, 2013.
- Geffner, T., Antoran, J., Foster, A., Gong, W., Ma, C., Kiciman, E., Sharma, A., Lamb, A., Kukla, M., Pawlowski, N., et al. Deep end-to-end causal inference. *arXiv preprint arXiv:2202.02195*, 2022.
- Gong, W. *Advances in approximate inference: combining VI and MCMC and improving on Stein discrepancy*. PhD thesis, University of Cambridge, 2022.
- Gong, W., Tschitschek, S., Nowozin, S., Turner, R. E., Hernández-Lobato, J. M., and Zhang, C. Icebreaker: Element-wise efficient information acquisition with a bayesian deep latent gaussian model. *Advances in neural information processing systems*, 32, 2019.
- Grzegorzczak, M. and Husmeier, D. Improving the structure mcmc sampler for bayesian networks by introducing a new edge reversal move. *Machine Learning*, 71(2-3):265, 2008.
- Heckerman, D., Meek, C., and Cooper, G. A bayesian approach to causal discovery. *Innovations in Machine Learning: Theory and Applications*, pp. 1–28, 2006.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Kuhn, H. W. The hungarian method for the assignment problem. *Naval research logistics quarterly*, 2(1-2):83–97, 1955.
- Kuipers, J. and Moffa, G. Partition mcmc for inference on acyclic digraphs. *Journal of the American Statistical Association*, 112(517):282–299, 2017.
- Li, C., Chen, C., Carlson, D., and Carin, L. Preconditioned stochastic gradient langevin dynamics for deep neural networks. In *Thirtieth AAAI Conference on Artificial Intelligence*, 2016.
- Liu, Q. and Wang, D. Stein variational gradient descent: A general purpose bayesian inference algorithm. *Advances in neural information processing systems*, 29, 2016.

- Lorch, L., Rothfuss, J., Schölkopf, B., and Krause, A. Dibs: Differentiable bayesian structure learning. *Advances in Neural Information Processing Systems*, 34:24111–24123, 2021.
- Ma, Y.-A., Chen, T., and Fox, E. A complete recipe for stochastic gradient mcmc. *Advances in neural information processing systems*, 28, 2015.
- Madigan, D., York, J., and Allard, D. Bayesian graphical models for discrete data. *International Statistical Review/Revue Internationale de Statistique*, pp. 215–232, 1995.
- Mena, G., Belanger, D., Linderman, S., and Snoek, J. Learning latent permutations with gumbel-sinkhorn networks. *arXiv preprint arXiv:1802.08665*, 2018.
- Munkres, J. Algorithms for the assignment and transportation problems. *Journal of the society for industrial and applied mathematics*, 5(1):32–38, 1957.
- Murphy, K. P. Active learning of causal bayes net structure. Technical report, technical report, UC Berkeley, 2001.
- Niculae, V., Martins, A., Blondel, M., and Cardie, C. Sparsemap: Differentiable sparse structured inference. In *International Conference on Machine Learning*, pp. 3799–3808. PMLR, 2018.
- Nishikawa-Toomey, M., Deleu, T., Subramanian, J., Bengio, Y., and Charlin, L. Bayesian learning of causal structure and mechanisms with gflownets and variational bayes. *arXiv preprint arXiv:2211.02763*, 2022.
- Pearl, J. *Causality*. Cambridge university press, 2009.
- Peters, J., Janzing, D., and Schölkopf, B. *Elements of causal inference: foundations and learning algorithms*. The MIT Press, 2017.
- Pe’er, D., Regev, A., Elidan, G., and Friedman, N. Inferring subnetworks from perturbed expression profiles. *Bioinformatics*, 17(suppl\_1):S215–S224, 2001.
- Robinson, R. W. Counting labeled acyclic digraphs. *New directions in the theory of graphs*, pp. 239–273, 1973.
- Sachs, K., Perez, O., Pe’er, D., Lauffenburger, D. A., and Nolan, G. P. Causal protein-signaling networks derived from multiparameter single-cell data. *Science*, 308(5721): 523–529, 2005.
- Spirtes, P., Glymour, C. N., Scheines, R., and Heckerman, D. *Causation, prediction, and search*. MIT press, 2000.
- Springenberg, J. T., Klein, A., Falkner, S., and Hutter, F. Bayesian optimization with robust bayesian neural networks. *Advances in neural information processing systems*, 29, 2016.
- Tigas, P., Annadani, Y., Jesson, A., Schölkopf, B., Gal, Y., and Bauer, S. Interventions, where and how? experimental design for causal models at scale. *arXiv preprint arXiv:2203.02016*, 2022.
- Tong, S. and Koller, D. Active learning for structure in bayesian networks. In *International joint conference on artificial intelligence*, volume 17, pp. 863–869. Citeseer, 2001.
- Trippe, B. and Turner, R. Overpruning in variational bayesian neural networks. *arXiv preprint arXiv:1801.06230*, 2018.
- Van den Bulcke, T., Van Leemput, K., Naudts, B., van Remortel, P., Ma, H., Verschoren, A., De Moor, B., and Marchal, K. Syntren: a generator of synthetic gene expression data for design and analysis of structure learning algorithms. *BMC bioinformatics*, 7:1–12, 2006.
- Van Koten, C. and Gray, A. An application of bayesian network for predicting object-oriented software maintainability. *Information and Software Technology*, 48(1): 59–67, 2006.
- Yu, Y., Gao, T., Yin, N., and Ji, Q. Dags with no curl: An efficient dag structure learning approach. In *International Conference on Machine Learning*, pp. 12156–12166. PMLR, 2021.
- Zhang, C., Bütepage, J., Kjellström, H., and Mandt, S. Advances in variational inference. *IEEE transactions on pattern analysis and machine intelligence*, 41(8):2008–2026, 2018.
- Zheng, X., Aragam, B., Ravikumar, P. K., and Xing, E. P. Dags with no tears: Continuous optimization for structure learning. *Advances in Neural Information Processing Systems*, 31, 2018.

## Appendix – BayesDAG: Gradient-Based Posterior Sampling for Causal Discovery

### A. Joint Inference with SG-MCMC

In this section, we propose an alternative formulation that enables a joint inference framework for  $\mathbf{p}, \mathbf{W}, \Theta$  using SG-MCMC, thereby avoiding the need for variational inference for  $\mathbf{W}$ .

We adopt a continuous relaxation of  $\mathbf{W}$ , similar to (Lorch et al., 2021), by introducing a latent variable  $\tilde{\mathbf{W}}$ . The graphical model is illustrated in Figure 3. We can define

$$p(\mathbf{W}|\tilde{\mathbf{W}}) = \prod_{i,j} p(W_{ij}|\tilde{W}_{ij}) \quad (10)$$

with  $p(W_{ij} = 1|\tilde{W}_{ij}) = \sigma(\tilde{W}_{ij})$  where  $\sigma(\cdot)$  is the sigmoid function. In other words,  $\tilde{W}_{ij}$  defines the existence logits of  $W_{ij}$ .

With the introduction of  $\tilde{\mathbf{W}}$ , the original posterior expectations of  $p(\mathbf{p}, \mathbf{W}, \Theta|\mathbf{D})$ , e.g. during evaluation, can be translated using the following proposition.

**Proposition A.1** (Equivalence of posterior expectation). *Under the generative model Figure 3, we have*

$$\mathbb{E}_{p(\mathbf{p}, \mathbf{W}, \Theta|\mathbf{D})} [f(\mathbf{G} = \tau(\mathbf{p}, \mathbf{W}), \Theta)] = \mathbb{E}_{p(\mathbf{p}, \tilde{\mathbf{W}}, \Theta)} \left[ \frac{\mathbb{E}_{p(\mathbf{W}|\tilde{\mathbf{W}})} [f(\mathbf{G}, \Theta)p(\mathbf{D}, \Theta|\mathbf{p}, \mathbf{W})]}{\mathbb{E}_{p(\mathbf{W}|\tilde{\mathbf{W}})} [p(\mathbf{D}, \Theta|\mathbf{p}, \mathbf{W})]} \right] \quad (11)$$

where  $f$  is the target quantity.

This proof is in Appendix B.3.

With this proposition, instead of sampling  $\mathbf{W}$ , use SG-MCMC to draw  $\tilde{\mathbf{W}}$  samples. Similar to Section 4.1, to use SG-MCMC for  $\mathbf{p}, \tilde{\mathbf{W}}, \Theta$ , we need their gradient information. The following proposition specifies the required gradients.

**Proposition A.2.** *With the generative model defined as Figure 3, we have*

$$\begin{aligned} \nabla_{\mathbf{p}, \Theta, \tilde{\mathbf{W}}} U(\mathbf{p}, \tilde{\mathbf{W}}, \Theta) &= -\nabla_{\mathbf{p}} \log p(\mathbf{p}) - \nabla_{\Theta} \log p(\Theta) \\ &\quad - \nabla_{\tilde{\mathbf{W}}} \log p(\tilde{\mathbf{W}}) - \nabla_{\mathbf{p}, \Theta, \tilde{\mathbf{W}}} \log \mathbb{E}_{p(\mathbf{W}|\tilde{\mathbf{W}})} [p(\mathbf{D}|\mathbf{W}, \mathbf{p}, \Theta)] \end{aligned} \quad (12)$$

The proof is in Appendix B.4.

With these gradients, we can directly plug in existing SG-MCMC samplers to draw samples for  $\mathbf{p}, \tilde{\mathbf{W}}$ , and  $\Theta$  in joint inference (Algorithm 1).

## B. Theory

### B.1. Proof of Theorem 3.1

For completeness, we recite the theorem here.

*Theorem 3.1* (Equivalence of inference in  $(\mathbf{W}, \mathbf{p})$  and binary DAG space). Assume graph  $\mathbf{G}$  is a binary adjacency matrix representing a DAG and node potential  $\mathbf{p}$  does not contain the same values, i.e.  $p_i \neq p_j \forall i, j$ . Then, with the induced joint observational distribution  $p(\mathbf{D}, \mathbf{G})$ , dataset  $\mathbf{D}$  and a corresponding prior  $p(\mathbf{G})$ , we have

$$p(\mathbf{G}|\mathbf{D}) = \int p_{\tau}(\mathbf{p}, \mathbf{W}|\mathbf{D}) \mathbb{1}(\mathbf{G} = \tau(\mathbf{W}, \mathbf{p})) d\mathbf{W} d\mathbf{p} \quad (13)$$

if  $p(\mathbf{G}) = \int p_{\tau}(\mathbf{p}, \mathbf{W}) \mathbb{1}(\mathbf{G} = \tau(\mathbf{W}, \mathbf{p})) d\mathbf{W} d\mathbf{p}$ , where  $p_{\tau}(\mathbf{W}, \mathbf{p})$  is the prior,  $\mathbb{1}(\cdot)$  is the indicator function and  $p_{\tau}(\mathbf{p}, \mathbf{W}|\mathbf{D})$  is the posterior distribution over  $\mathbf{p}, \mathbf{W}$ .

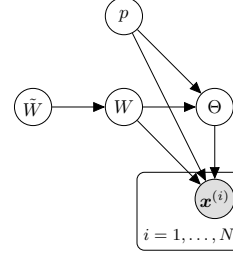


Figure 3: Graphical model with latent variable  $\tilde{\mathbf{W}}$ .



---

**Algorithm 1** Joint inference
 

---

**Input:** dataset  $D$ , prior  $p(\mathbf{p}, \tilde{\mathbf{W}}, \Theta)$ , SG-MCMC sampler update  $\text{Sampler}(\cdot)$ ; sampler hyperparameter  $\Psi$ ; training steps  $T$ .

**Output:** posterior samples  $\{\mathbf{p}, \tilde{\mathbf{W}}, \Theta\}$

Initialize  $\mathbf{p}_0, \tilde{\mathbf{W}}_0, \Theta_0$

**for**  $t = 1, \dots, T$  **do**

Evaluate gradient  $\nabla_{\mathbf{p}_{t-1}, \tilde{\mathbf{W}}_{t-1}, \Theta_{t-1}} U$  based on Equation (12).

Update samples  $\mathbf{p}_t, \tilde{\mathbf{W}}_t, \Theta_t = \text{Sampler}(\nabla_{\mathbf{p}_{t-1}, \tilde{\mathbf{W}}_{t-1}, \Theta_{t-1}} U; \Psi)$

**if** storing condition met **then**

$\{\mathbf{p}, \tilde{\mathbf{W}}, \Theta\} \leftarrow \mathbf{p}_t, \tilde{\mathbf{W}}_t, \Theta_t$

**end if**

**end for**

---



---

**Algorithm 2** BayesDAG SG-MCMC+VI Inference
 

---

**Input:** dataset  $D$ ; prior  $p(\mathbf{p}, \mathbf{W}), p(\Theta)$ ; SG-MCMC sampler  $\text{Sampler}$ ; sampler hyperparameters  $\Psi$ ; network  $\mu_\phi(\cdot)$ ; training iteration  $T$ .

**Output:** samples  $\{\Theta, \mathbf{p}\}$  and variational posterior  $q_\phi$

Initialize  $\Theta^{(0)}, \mathbf{p}^{(0)}, \phi$

**for**  $t = 1 \dots, T$  **do**

Sample  $\mathbf{W}^{(t-1)} \sim q_\phi(\mathbf{W}|\mathbf{p}^{(t-1)})$

Evaluate  $\nabla_{\mathbf{p}, \Theta} U$  (Equations (24) and (25)) with  $\Theta^{(t-1)}, \mathbf{p}^{(t-1)}, \mathbf{W}^{(t-1)}$

$\Theta^{(t)}, \mathbf{p}^{(t)} = \text{Sampler}(\nabla_{\mathbf{p}, \Theta} U; \Psi)$

**if** storing condition met **then**

$\{\mathbf{p}, \Theta\} \leftarrow \mathbf{p}^{(t)}, \Theta^{(t)}$

**end if**

Maximize ELBO (Equation (9)) w.r.t.  $\phi$  with  $\mathbf{p}^{(t)}, \Theta^{(t)}$

**end for**

---

To prove this theorem, we first prove the following lemma stating the equivalence of  $\tau$  (Equation (3)) to binary DAG space.

**Lemma B.1** (Equivalence of  $\tau$  to DAG space). *Consider  $d$  random variables, a node potential vector  $\mathbf{p} \in \mathbb{R}^d$  and a binary matrix  $\mathbf{W} \in \{0, 1\}^{d \times d}$ . Then the following holds:*

- (a) For any  $\mathbf{W} \in \{0, 1\}^{d \times d}, \mathbf{p} \in \mathbb{R}^d, \mathbf{G} = \tau(\mathbf{W}, \mathbf{p})$  is a DAG.
- (b) For any DAG  $\mathbf{G} \in \mathbb{D}$ , where  $\mathbb{D}$  is the space of all DAGs, there exists a corresponding  $(\mathbf{W}, \mathbf{p})$  such that  $\tau(\mathbf{W}, \mathbf{p}) = \mathbf{G}$ .

*Proof.* The main proof directly follows the theorem 2.1 in (Yu et al., 2021). For (a), we show the output from  $\tau(\mathbf{W}, \mathbf{p})$  must be a DAG. By leveraging the Lemma 3.4 in (Yu et al., 2021), we can easily obtain that  $\text{Step}(\text{grad } \mathbf{p})$  emits a binary adjacency matrix representing a DAG. The only difference is that we replace the  $\text{ReLU}(\cdot)$  with  $\text{Step}(\cdot)$  but the conclusion can be directly generalized.

For (b), we show that for any DAG  $\mathbf{G}$ , there exists a  $(\mathbf{W}, \mathbf{p})$  pair s.t.  $\tau(\mathbf{W}, \mathbf{p}) = \mathbf{G}$ . To see this, we can observe that  $\mathbf{p}$  implicitly defines a topological order in the mapping  $\tau$ . For any  $p_i > p_j$ , we have  $j \rightarrow i$  after the mapping  $\text{Step}(\text{grad } \mathbf{p})$ . Thus, by leveraging Theorem 3.7 in (Yu et al., 2021), we obtain that there exists a potential vector  $\mathbf{p} \in \mathbb{R}^d$  for any DAG  $\mathbf{G}$  such that

$$(\text{grad } \mathbf{p})(i, j) > 0 \quad \text{when } G_{ij} = 1$$

Thus, we can choose  $\mathbf{W}$  in the following way:

$$\mathbf{W} = \begin{cases} W_{ij} = 0 & \text{if } G_{ij} = 0 \\ W_{ij} = 1 & \text{if } G_{ij} = 1 \end{cases}$$

□

Next, let's prove the [Theorem 3.1](#).

*Proof of Theorem 3.1.* From [Lemma B.1](#), we see that the mapping is complete. Namely, the  $(\mathbf{W}, \mathbf{p})$  space can represent the entire DAG space. Next, we show that performing Bayesian inference in  $(\mathbf{W}, \mathbf{p})$  space can also correspond to the inference in DAG space.

Assume we have the prior  $p_\tau(\mathbf{W}, \mathbf{p})$ . Then through mapping  $\tau$ , we implicitly define a prior over the DAG  $\mathbf{G}$  in the following:

$$p_\tau(\mathbf{G}) = \int p_\tau(\mathbf{W}, \mathbf{p}) \mathbb{1}(\mathbf{G} = \tau(\mathbf{W}, \mathbf{p})) d\mathbf{W} d\mathbf{p} \quad (14)$$

This basically states that the corresponding prior over  $\mathbf{G}$  is an accumulation of the corresponding probability associated with  $(\mathbf{W}, \mathbf{p})$  pairs.

Similarly, we can define a corresponding posterior  $p_\tau(\mathbf{G}|\mathbf{D})$ :

$$p_\tau(\mathbf{G}|\mathbf{D}) = \int p_\tau(\mathbf{W}, \mathbf{p}|\mathbf{D}) \mathbb{1}(\mathbf{G} = \tau(\mathbf{W}, \mathbf{p})) d\mathbf{W} d\mathbf{p} \quad (15)$$

Now, let's show that this posterior  $p_\tau(\mathbf{G}|\mathbf{D}) = p(\mathbf{G}|\mathbf{D})$  if prior matches, i.e.  $p(\mathbf{G}) = p_\tau(\mathbf{G})$ . From Bayes's rule, we can easily write down

$$p_\tau(\mathbf{W}, \mathbf{p}|\mathbf{D}) = \frac{p(\mathbf{D}|\mathbf{G} = \tau(\mathbf{W}, \mathbf{p}))p(\mathbf{p}, \mathbf{W})}{\sum_{\mathbf{G}' \in \mathbb{D}} p(\mathbf{D}, \mathbf{G}')} \quad (16)$$

Then, by substituting [Equation \(16\)](#) into [Equation \(15\)](#), we have

$$\begin{aligned} p_\tau(\mathbf{G}|\mathbf{D}) &= \int \frac{p(\mathbf{D}|\mathbf{G})p_\tau(\mathbf{W}, \mathbf{p})}{\sum_{\mathbf{G}' \in \mathbb{D}} p(\mathbf{D}, \mathbf{G}')} \mathbb{1}(\mathbf{G} = \tau(\mathbf{W}, \mathbf{p})) d\mathbf{W} d\mathbf{p} \\ &= \frac{\int p(\mathbf{D}|\mathbf{G})p_\tau(\mathbf{W}, \mathbf{p}) \mathbb{1}(\mathbf{G} = \tau) d\mathbf{W} d\mathbf{p}}{\sum_{\mathbf{G}' \in \mathbb{D}} p(\mathbf{D}, \mathbf{G}')} \end{aligned} \quad (17)$$

$$= \frac{p(\mathbf{D}|\mathbf{G}) \int p_\tau(\mathbf{W}, \mathbf{p}) \mathbb{1}(\mathbf{G} = \tau) d\mathbf{W} d\mathbf{p}}{\sum_{\mathbf{G}' \in \mathbb{D}} p(\mathbf{D}, \mathbf{G}')} \quad (18)$$

$$\begin{aligned} &= \frac{p(\mathbf{D}|\mathbf{G})p_\tau(\mathbf{G})}{\sum_{\mathbf{G}' \in \mathbb{D}} p(\mathbf{D}|\mathbf{G}')p_\tau(\mathbf{G}')} \\ &= p(\mathbf{G}|\mathbf{D}) \end{aligned} \quad (19)$$

where [Equation \(17\)](#) is from the fact that  $\sum_{\mathbf{G}' \in \mathbb{D}} p(\mathbf{D}, \mathbf{G}')$  is independent of  $(\mathbf{W}, \mathbf{p})$  due to marginalization. [Equation \(18\)](#) is obtained because  $p(\mathbf{D}|\mathbf{G})$  is also independent of  $(\mathbf{W}, \mathbf{p})$  due to (1)  $\mathbb{1}(\mathbf{G} = \tau(\mathbf{W}, \mathbf{p}))$  and (2)  $p(\mathbf{D}|\mathbf{G})$  is a constant when fixing  $\mathbf{G}$ . [Equation \(19\)](#) is obtained by applying Bayes's rule and  $p_\tau(\mathbf{G}) = p(\mathbf{G})$ .  $\square$

## B.2. Proof of [Theorem 3.2](#)

*Theorem 3.2* (Equivalence of NoCurl formulation). Assuming the conditions in [Theorem 3.1](#) are satisfied. Then, for a given  $(\mathbf{W}, \mathbf{p})$ , we have

$$\mathbf{G} = \mathbf{W} \odot \text{Step}(\text{grad } \mathbf{p}) = \mathbf{W} \odot [\boldsymbol{\sigma}^*(\mathbf{p}) \mathbf{L} \boldsymbol{\sigma}^*(\mathbf{p})^T]$$

where  $\mathbf{G}$  is a DAG and  $\boldsymbol{\sigma}^*(\mathbf{p})$  is defined in [Equation \(5\)](#).

To prove this theorem, we need to first prove the following lemma.

**Lemma B.2.** For any permutation matrix  $\mathbf{M} \in \Sigma_d$ , we have

$$\text{grad}(\mathbf{M}\mathbf{p}) = \mathbf{M}^T \text{grad}(\mathbf{p})\mathbf{M}$$

where  $\text{grad}$  is the operator defined in ??.

*Proof.* By definition of  $\text{grad}(\cdot)$ , we have

$$\begin{aligned}\text{grad}(\mathbf{M}\mathbf{p}) &= (\mathbf{M}\mathbf{p})_i - (\mathbf{M}\mathbf{p})_j \\ &= \mathbf{1}(i)^T \mathbf{M}\mathbf{p} - \mathbf{1}(j)^T \mathbf{M}\mathbf{p} \\ &= \mathbf{M}_{i,:}\mathbf{p} - \mathbf{M}_{j,:}\mathbf{p}\end{aligned}$$

where  $\mathbf{1}(i)$  is a one-hot vector with  $i^{\text{th}}$  entry 1, and  $\mathbf{M}_{i,:}$  is the  $i^{\text{th}}$  row of matrix  $\mathbf{M}$ . The above is equivalent to computing the grad with new labels obtained by permuting  $\mathbf{p}$  with  $\mathbf{M}$ . Therefore, we can see that  $\text{grad}(\mathbf{M}\mathbf{p})$  can be computed by permuting the original  $\text{grad}(\mathbf{p})$  by matrix  $\mathbf{M}$ .

$$\text{grad}(\mathbf{M}\mathbf{p}) = \mathbf{M}^T \text{grad}(\mathbf{p})\mathbf{M}$$

□

*Proof of Theorem 3.2.* Since  $\mathbf{W}$  plays the same role in both formulations, we focus on the equivalence of  $\text{Step}(\text{grad}(\cdot))$ .

Define a sorted  $\tilde{\mathbf{p}} = \boldsymbol{\sigma}\mathbf{p}$ , where  $\boldsymbol{\sigma} \in \Sigma_d$ , such that for  $i < j$ , we have  $\tilde{p}_i > \tilde{p}_j$ . Namely,  $\boldsymbol{\sigma}$  is a permutation matrix. Thus, we have

$$\text{grad}(\mathbf{p}) = \text{grad}(\boldsymbol{\sigma}^T \tilde{\mathbf{p}}).$$

By Lemma B.2, we have

$$\text{grad}(\boldsymbol{\sigma}^T \tilde{\mathbf{p}}) = \boldsymbol{\sigma} \text{grad}(\tilde{\mathbf{p}})\boldsymbol{\sigma}^T.$$

Since  $\tilde{\mathbf{p}}$  is an ordered vector. Therefore,  $\text{grad}(\tilde{\mathbf{p}})$  is a skew-symmetric matrix with a positive lower half part.

Therefore, we have

$$\text{Step}(\text{grad}(\mathbf{p})) = \text{Step}(\boldsymbol{\sigma} \text{grad}(\tilde{\mathbf{p}})\boldsymbol{\sigma}^T) = \boldsymbol{\sigma} \text{Step}(\text{grad}(\tilde{\mathbf{p}}))\boldsymbol{\sigma}^T = \boldsymbol{\sigma} \mathbf{L}\boldsymbol{\sigma}^T$$

This is true because  $\boldsymbol{\sigma}$  is just a permutation matrix that does not alter the sign of  $\text{grad}(\tilde{\mathbf{p}})$ .

Since  $\boldsymbol{\sigma}$  is a permutation matrix that sort  $\mathbf{p}$  value in a ascending order, from Lemma 1 in (Blondel et al., 2020), we have

$$\boldsymbol{\sigma} = \arg \max_{\boldsymbol{\sigma}' \in \Sigma_d} \mathbf{p}^T (\boldsymbol{\sigma}' \mathbf{o})$$

□

### B.3. Proof of Proposition A.1

*Proof.*

$$\begin{aligned}& \mathbb{E}_{p(\mathbf{p}, \mathbf{W}, \boldsymbol{\Theta} | \mathbf{D})} [f(\mathbf{G} = \tau(\mathbf{p}, \mathbf{W}), \boldsymbol{\Theta})] \\ &= \int p(\mathbf{p}, \mathbf{W}, \boldsymbol{\Theta}, \tilde{\mathbf{W}} | \mathbf{D}) f(\mathbf{G}, \boldsymbol{\Theta}) d\mathbf{p} d\mathbf{W} d\boldsymbol{\Theta} d\tilde{\mathbf{W}} \\ &= \int p(\mathbf{p}, \tilde{\mathbf{W}}, \boldsymbol{\Theta} | \mathbf{D}) p(\mathbf{W} | \mathbf{p}, \boldsymbol{\Theta}, \tilde{\mathbf{W}}, \mathbf{D}) f(\mathbf{G}, \boldsymbol{\Theta}) d\mathbf{p} d\mathbf{W} d\boldsymbol{\Theta} d\tilde{\mathbf{W}} \\ &= \mathbb{E}_{p(\mathbf{p}, \tilde{\mathbf{W}}, \boldsymbol{\Theta} | \mathbf{D})} \left[ \frac{\int p(\mathbf{D} | \mathbf{p}, \boldsymbol{\Theta}, \mathbf{W}) p(\mathbf{p}) p(\tilde{\mathbf{W}}) p(\mathbf{W} | \tilde{\mathbf{W}}) p(\boldsymbol{\Theta} | \mathbf{p}, \mathbf{W}) f(\mathbf{G}, \boldsymbol{\Theta}) d\mathbf{W}}{\int p(\mathbf{D} | \mathbf{p}, \boldsymbol{\Theta}, \mathbf{W}) p(\mathbf{p}) p(\tilde{\mathbf{W}}) p(\mathbf{W} | \tilde{\mathbf{W}}) p(\boldsymbol{\Theta} | \mathbf{p}, \mathbf{W}) d\mathbf{W}} \right] \\ &= \mathbb{E}_{p(\mathbf{p}, \tilde{\mathbf{W}}, \boldsymbol{\Theta})} \left[ \frac{\mathbb{E}_{p(\mathbf{W} | \tilde{\mathbf{W}})} [f(\mathbf{G}, \boldsymbol{\Theta}) p(\mathbf{D}, \boldsymbol{\Theta} | \mathbf{p}, \mathbf{W})]}{\mathbb{E}_{p(\mathbf{W} | \tilde{\mathbf{W}})} [p(\mathbf{D}, \boldsymbol{\Theta} | \mathbf{p}, \mathbf{W})]} \right]\end{aligned}$$

□

#### B.4. Proof of Proposition A.2

*Proof.*

$$\begin{aligned}
 \nabla_{\mathbf{p}} U(\mathbf{p}, \tilde{\mathbf{W}}, \Theta) &= -\nabla_{\mathbf{p}} \log p(\mathbf{p}, \tilde{\mathbf{W}}, \Theta, D) \\
 &= -\nabla_{\mathbf{p}} \log p(\mathbf{p}) - \nabla_{\mathbf{p}} \log p(\tilde{\mathbf{W}}, \Theta, D | \mathbf{p}) \\
 &= -\nabla_{\mathbf{p}} \log p(\mathbf{p}) - \frac{\nabla_{\mathbf{p}} \int p(D | \mathbf{W}, \mathbf{p}, \Theta) p(\Theta | \mathbf{p}, \mathbf{W}) p(\mathbf{W} | \tilde{\mathbf{W}}) p(\tilde{\mathbf{W}}) d\mathbf{W}}{\int p(D | \mathbf{W}, \mathbf{p}, \Theta) p(\Theta | \mathbf{p}, \mathbf{W}) p(\mathbf{W} | \tilde{\mathbf{W}}) p(\tilde{\mathbf{W}}) d\mathbf{W}} \\
 &= -\nabla_{\mathbf{p}} \log p(\mathbf{p}) - \frac{\nabla_{\mathbf{p}} \mathbb{E}_{p(\mathbf{W} | \tilde{\mathbf{W}})} [p(D | \mathbf{W}, \mathbf{p}, \Theta)]}{\mathbb{E}_{p(\mathbf{W} | \tilde{\mathbf{W}})} [p(D | \mathbf{W}, \mathbf{p}, \Theta)]} \\
 &= -\nabla_{\mathbf{p}} \log p(\mathbf{p}) - \nabla_{\mathbf{p}} \log \mathbb{E}_{p(\mathbf{W} | \tilde{\mathbf{W}})} [p(D | \mathbf{W}, \mathbf{p}, \Theta)]
 \end{aligned}$$

Other gradient  $\nabla_{\tilde{\mathbf{W}}} U$  and  $\nabla_{\Theta} U$  can be derived using the similar approach, which concludes the proof.  $\square$

#### B.5. Proof of Proposition C.1

*Proof of Proposition C.1.* By definition, we have easily have

$$\begin{aligned}
 \nabla_{\mathbf{p}} U &= -\nabla_{\mathbf{p}} \log p(\mathbf{p}, \mathbf{W}, \Theta, D) \\
 &= -\nabla_{\mathbf{p}} \log p(\mathbf{p}, \mathbf{W}) - \nabla_{\mathbf{p}} \log p(D, \Theta | \tau(\mathbf{W}, \mathbf{p})) \\
 &= -\nabla_{\mathbf{p}} \log p(\mathbf{p}, \mathbf{W}) - \nabla_{\mathbf{p}} \log p(D | \Theta, \tau(\mathbf{W}, \mathbf{p})) + \underbrace{\nabla_{\mathbf{p}} \log p(\Theta | \tau(\mathbf{p}, \mathbf{W}))}_0
 \end{aligned}$$

Similarly, we have

$$\begin{aligned}
 \nabla_{\Theta} U &= -\nabla_{\Theta} \log p(\mathbf{p}, \mathbf{W}, \Theta, D) \\
 &= -\nabla_{\Theta} \log p(D | \Theta, \tau(\mathbf{W}, \mathbf{p})) - \nabla_{\Theta} \log p(\Theta | \mathbf{p}, \mathbf{W}) - \underbrace{\nabla_{\Theta} \log p(\mathbf{p}, \mathbf{W})}_0 \\
 &= -\nabla_{\Theta} \log p(D | \Theta, \tau(\mathbf{W}, \mathbf{p})) - \nabla_{\Theta} \log p(\Theta)
 \end{aligned}$$

$\square$

#### B.6. Derivation of ELBO

$$\begin{aligned}
 \log p(\mathbf{p}, \Theta, D) &= \log \int p(\mathbf{p}, \Theta, D, \mathbf{W}) d\mathbf{W} \\
 &= \log \int \frac{q_{\phi}(\mathbf{W} | \mathbf{p})}{q_{\phi}(\mathbf{W} | \mathbf{p})} p(\mathbf{p}, \Theta, D, \mathbf{W}) d\mathbf{W} \\
 &\geq \int q_{\phi}(\mathbf{W} | \mathbf{p}) \log p(\mathbf{p}, \Theta, D | \mathbf{W}) d\mathbf{W} + \int q_{\phi}(\mathbf{W} | \mathbf{p}) \log \frac{p(\mathbf{W})}{q_{\phi}(\mathbf{W} | \mathbf{p})} d\mathbf{W} \quad (20) \\
 &= \mathbb{E}_{q_{\phi}(\mathbf{W} | \mathbf{p})} [\log p(\mathbf{p}, \Theta, D | \mathbf{W})] - D_{\text{KL}} [q_{\phi}(\mathbf{W} | \mathbf{p}) \| p(\mathbf{W})]
 \end{aligned}$$

where the Equation (20) is obtained by Jensen's inequality.

#### C. SG-MCMC Update

Assume we want to draw samples  $\mathbf{p} \sim p(\mathbf{p} | D, \mathbf{W}, \Theta) \propto \exp(-U(\mathbf{p}, \mathbf{W}, \Theta))$  with  $U(\mathbf{p}, \mathbf{W}, \Theta) = -\log p(\mathbf{p}, \mathbf{W}, \Theta)$ , we can compute  $U$  by

$$U(\mathbf{p}, \mathbf{W}, \Theta) = - \sum_{n=1}^N \log p(\mathbf{x}_n | \mathbf{G} = \tau(\mathbf{W}, \mathbf{p}), \Theta) - \log p(\mathbf{p}, \mathbf{W}, \Theta) \quad (21)$$

In practice, we typically use mini-batches  $\mathcal{S}$  instead of the entire dataset  $D$ . Therefore, an approximation is

$$\tilde{U}(\mathbf{p}, \mathbf{W}, \Theta) = -\frac{|D|}{|\mathcal{S}|} \sum_{n \in \mathcal{S}} \log p(\mathbf{x}_n | \mathbf{G} = \tau(\mathbf{W}, \mathbf{p}), \Theta) - \log p(\mathbf{p}, \mathbf{W}, \Theta) \quad (22)$$

where  $|\mathcal{S}|$  and  $|D|$  are the minibatch and dataset sizes, respectively.

(Gong et al., 2019) uses the preconditioning technique on *stochastic gradient Hamiltonian Monte Carlo* (SG-HMC), similar to the preconditioning technique in (Li et al., 2016). In particular, they use a moving-average approximation of diagonal Fisher information to adjust the momentum. The transition dynamics at step  $t$  with EM discretization is

$$\begin{aligned} B &= \frac{1}{2}l \\ \mathbf{V}_t &= \beta_2 \mathbf{V}_{t-1} + (1 - \beta_2) \nabla_{\mathbf{p}} \tilde{U}(\mathbf{p}, \mathbf{W}, \Theta) \odot \nabla_{\mathbf{p}} \tilde{U}(\mathbf{p}, \mathbf{W}, \Theta) \\ g_t &= \frac{1}{\sqrt{1 + \sqrt{\mathbf{V}_t}}} \\ \mathbf{r}_t &= \beta_1 \mathbf{r}_{t-1} - l g_t \nabla_{\mathbf{p}} \tilde{U}(\mathbf{p}, \mathbf{W}, \Theta) + l \frac{\partial g_t}{\partial \mathbf{p}_t} + s \sqrt{2l \left( \frac{1 - \beta_1}{l} - B \right) \eta} \\ \mathbf{p}_t &= \mathbf{p}_{t-1} + l g_t \mathbf{r}_t \end{aligned} \quad (23)$$

where  $l^2$  is the learning rate;  $(\beta_1, \beta_2)$  controls the preconditioning decay rate,  $\eta$  is the Gaussian noise with 0 mean and unit variance, and  $s$  is the hyperparameter controlling the level of injected noise to SG-MCMC. Throughout the paper, we use  $(\beta_1, \beta_2) = (0.9, 0.99)$  for all experiments.

The following proposition specifies the gradients required by SG-MCMC:  $\nabla_{\mathbf{p}, \Theta} U(\mathbf{p}, \mathbf{W}, \Theta)$ .

**Proposition C.1.** *We have the following gradient equations for SG-MCMC:*

$$\nabla_{\mathbf{p}} U = -\nabla_{\mathbf{p}} \log p(\mathbf{p}) - \nabla_{\mathbf{p}} \log p(D | \Theta, \tau(\mathbf{W}, \mathbf{p})) \quad (24)$$

and

$$\nabla_{\Theta} U = -\nabla_{\Theta} \log p(\Theta) - \nabla_{\Theta} \log p(D | \Theta, \tau(\mathbf{p}, \mathbf{W})) \quad (25)$$

Refer to [Appendix B.5](#) for details.

## D. Gumbel-Sinkhorn Operator

The Sinkhorn operator  $\mathcal{S}(M)$  on a matrix  $M$  (Adams & Zemel, 2011) is defined as a sequence of row and column normalizations, each is called Sinkhorn iteration. (Mena et al., 2018) showed that the non-differentiable  $\arg \max$  problem

$$\boldsymbol{\sigma} = \arg \max_{\boldsymbol{\sigma}' \in \Sigma_d} \langle \boldsymbol{\sigma}', M \rangle \quad (26)$$

can be relaxed through an entropy regularizer with its solution being expressed by  $\mathcal{S}(\cdot)$ . In particular, they showed that  $\mathcal{S}(M/t) = \arg \max_{\boldsymbol{\sigma}' \in \mathcal{B}_d} \langle \boldsymbol{\sigma}', M \rangle + th(\boldsymbol{\sigma}')$ , where  $h(\cdot)$  is the entropy function. This regularized solution converges to the solution of [Equation \(26\)](#) when  $t \rightarrow 0$ , i.e.  $\lim_{t \rightarrow 0} \mathcal{S}(M/t)$ . Since the Sinkhorn operator is differentiable,  $\mathcal{S}(M/t)$  can be viewed as a differentiable approximation to [Equation \(26\)](#), which can be used to obtain the solution of [Equation \(5\)](#). Specifically, we have

$$\arg \max_{\boldsymbol{\sigma}' \in \Sigma_d} \mathbf{p}^T(\boldsymbol{\sigma}' \mathbf{o}) = \arg \max_{\boldsymbol{\sigma}' \in \Sigma_d} \langle \boldsymbol{\sigma}', \mathbf{p} \mathbf{o}^T \rangle = \lim_{t \rightarrow 0} \mathcal{S}\left(\frac{\mathbf{p} \mathbf{o}^T}{t}\right) \quad (27)$$

In practice, we approximate it with  $t > 0$ , resulting in a doubly stochastic matrix. To get the binary permutation matrix, we apply the Hungarian algorithm (Munkres, 1957). During the backward pass, we use a straight-through estimator for  $\mathbf{p}$ .

Some of the previous works (Charpentier et al., 2022; Cundy et al., 2021) have leveraged the Sinkhorn operator to model variational distributions over permutation matrices. However, they start with a full rank  $M$ , which has been reported to require over **1000** Sinkhorn iterations to converge (Cundy et al., 2021). However, our formulation, based on explicit node potential  $\mathbf{p} \mathbf{o}^T$ , generates a rank-1 matrix, requiring much fewer Sinkhorn steps (around **300**) in practice, saving two-thirds of the computational cost.



## E. Related Work

Bayesian causal discovery literature has primarily focused on inference in linear models with closed-form posteriors or marginalized parameters. Early works considered sampling directed acyclic graphs (DAGs) for discrete (Cooper & Herskovits, 1992; Madigan et al., 1995; Heckerman et al., 2006) and Gaussian random variables (Friedman & Koller, 2003; Tong & Koller, 2001) using Markov chain Monte Carlo (MCMC) in the DAG space. However, these approaches exhibit slow mixing and convergence (Eaton & Murphy, 2012; Grzegorzczuk & Husmeier, 2008), often requiring restrictions on number of parents (Kuipers & Moffa, 2017).

Recent advances in variational inference (Zhang et al., 2018) have facilitated graph inference in DAG space, with gradient-based methods employing the NOTEARS DAG penalty (Zheng et al., 2018). (Annadani et al., 2021) samples DAGs from autoregressive adjacency matrix distributions, while (Lorch et al., 2021) utilizes Stein variational approach (Liu & Wang, 2016) for DAGs and causal model parameters. (Cundy et al., 2021) proposed a variational inference framework on node orderings using the gumbel-sinkhorn gradient estimator (Mena et al., 2018). (Deleu et al., 2022; Nishikawa-Toomey et al., 2022) employ the GFlowNet framework (Bengio et al., 2021) for inferring the DAG posterior. Most methods, except (Lorch et al., 2021) are restricted to linear models, while (Lorch et al., 2021) has high computational costs and lacks DAG generation guarantees compared to our method.

In contrast, *quasi-Bayesian* methods, such as DAG bootstrap (Friedman et al., 2013), demonstrate competitive performance. DAG bootstrap resamples data and estimates a single DAG using PC (Spirtes et al., 2000), GES (Chickering, 2002), or similar algorithms, weighting the obtained DAGs by their unnormalized posterior probabilities. Recent neural network-based works employ variational inference to learn DAG distributions and point estimates for nonlinear model parameters (Charpentier et al., 2022; Geffner et al., 2022).

## F. Experimental Settings

### F.1. Baselines

For all the experimental settings, we compare with the following baselines:

- **Bootstrap GES (BGES)** (Friedman et al., 2013; Chickering, 2002) is a bootstrap based quasi-Bayesian approach for linear Gaussian models which first resamples with replacement data points at random and then estimates a linear SCM using the GES algorithm (Chickering, 2002) for each bootstrap set. GES is a score based approach to learn a point estimate of a linear Gaussian SCM. For all the experimental settings, we use 50 bootstrap sets.
- **Differentiable DAG Sampling (DDS)** is a VI based approach to learn distribution over DAGs and a point estimate over the nonlinear functional parameters. DDS performs inference on the node permutation matrices, thus directly generating DAGs. Gumbel-sinkhorn (Mena et al., 2018) is used for obtaining valid gradients and Hungarian algorithm is used for the straight-through gradient estimator. In the author provided implementation, for evaluation, a single permutation matrix is sampled and the logits of the edge beliefs are directly thresholded. In this work, in order to make the comparison fair to Bayesian learning methods, we directly sample the binary adjacency matrix based on the edge logits.
- **BCD Nets** (Cundy et al., 2021) is a VI based fully Bayesian structure learning approach for linear causal models. BCD performs inference on both the node permutations through the Gumbel-sinkhorn (Mena et al., 2018) operator as well as the model parameters through a VI distribution. Both DDS and BCD nets operate directly on full rank initializations to the Gumbel-sinkhorn operator, unlike our rank-1 initialization, which saves computations in practice.
- **DIBS** (Lorch et al., 2021) uses SVGD (Liu & Wang, 2016) with the DAG regularizer (Zheng et al., 2018) and bilinear embeddings to perform inference over both linear and nonlinear causal models. As our data generation process involves SCM with unequal noise variance, we extend DIBS framework with an inference over noise variance using SVGD, similar to the original paper.

While DIBS and DDS can handle nonlinear parameterization, approaches like BGES and BCD, which are primarily designed for linear models still give competitive results when applied on nonlinear data. Given that there are limited number of baselines in the nonlinear case, and DIBS being the only fully Bayesian nonlinear baseline, we compare with BGES and BCD for all settings despite their model misspecification.

## F.2. Evaluation Metrics

For higher dimensional settings with nonlinear models, the true posterior is intractable. While in general it is hard to evaluate the posterior inference quality in high dimensions, prior work has suggested to evaluate on proxy metrics which we adopt in this work as well (Lorch et al., 2021; Geffner et al., 2022; Annadani et al., 2023). In particular, we evaluate the following metrics:

- **$\mathbb{E}$ -SHD:** Structural Hamming Distance (SHD) measures the hamming distance between graphs. In particular, it is a measure of number of edges that are to be added, removed or reversed to get the ground truth from the estimated graph. Since we have a posterior distribution  $q(\mathbf{G})$  over graphs, we measure the *expected* SHD:

$$\mathbb{E}\text{-SHD} := \mathbb{E}_{\mathbf{G} \sim q(\mathbf{G})}[\text{SHD}(\mathbf{G}, \mathbf{G}^{GT})] \approx \frac{1}{N_e} \sum_{i=1}^{N_e} [\text{SHD}(\mathbf{G}^{(i)}, \mathbf{G}^{GT})] \quad , \text{ with } \mathbf{G}^{(i)} \sim q(\mathbf{G})$$

where  $\mathbf{G}^{GT}$  is the ground-truth causal graph.

- **Edge F1:** It is F1 score of each edge being present or absent in comparison to the true edge set, averaged over all edges.
- **NLL:** We also measure the negative log-likelihood of the held-out data.

The first two metrics measure the goodness of the graph posterior while the NLL measures the goodness of the joint posterior over the entire causal model.

## F.3. Synthetic Data

As knowledge of ground truth graph is not possible in many real world settings, it is standard across causal discovery to benchmark in synthetic data settings. Following prior work, we generate synthetic data by first sampling a DAG at random from either Erdos-Rényi (ER) (Erdős et al., 1960) or Scale-Free (SF) (Barabási & Albert, 1999) family. We ensure that the graphs have  $2d$  edges in expectation. For nonlinear models, the nonlinear functions are defined by randomly initialized Multi-Layer Perceptrons (MLP) with a single hidden layer of 5 nodes and ReLU nonlinearity. The variance of the exogenous Gaussian noise variable is drawn from an Inverse Gamma prior with concentration  $\alpha = 1.5$  and rate  $\beta = 1$ . For higher dimensional settings, we consider  $N = 5000$  random samples for training and  $N = 1000$  samples for held-out evaluation. For all settings, we evaluate on 30 random datasets.

## F.4. Hyperparameter Selection

In this section, we will give the details our how to select the hyperparameters for our method and all the baseline models.

We employ a cross-validation-like procedure for hyperparameter tuning in BayesDAG and DIBS to optimize  $\mathbb{E}$ -SHD value (for nonlinear setting). For each ER and SF dataset with varying dimensions, we initially generate five tuning datasets. After determining the optimal hyperparameters, we fix them and evaluate the models on 30 test datasets. For DDS, we adopt the hyperparameters provided in the original paper (Charpentier et al., 2022). BCD and BGES do not necessitate hyperparameter tuning since BCD already incorporates the correct prior graph for ER and SF datasets. For semi-synthetic Syntren and real world Sachs protein cells datasets, we assume the number of edges in the ground truth graphs are known and we tune our hyperparameters to produce roughly correct number of edges. BCD and DIBS also assume access to the ground truth edge number and use the graph prior to enforce the number of edges.

**Network structure** We use one hidden layer MLP with hidden size of  $\max(4 * d, 64)$  for the nonlinear functional relations, where  $d$  is the dimensionality of dataset. We use **LeakyReLU** as the activation function. We also enable the **LayerNorm** and **residual connections** in the network. In particular, for variational network  $\mu_\phi$  in BayesDAG, we apply the **LayerNorm** on  $\mathbf{p}$  before inputting it to the network. We use 2 hidden layer MLP with size 48, **LayerNorm** and **residual connections** for  $\mu_\phi$ .

**Sparse initialization for BayesDAG** For BayesDAG, we additionally allow sparse initialization by sampling a sparse  $\mathbf{W}$  from the  $\mu_\phi$ . This can be achieved by subtracting a constant 1 from the existing logits (i.e. the output from  $\mu_\phi$ ).

BayesDAG				
	$\lambda_s$	Scale $\mathbf{p}$	Scale $\Theta$	Sparse Init.
nonlinear ER $d = 20$	300	0.01	0.01	False
nonlinear SF $d = 20$	200	0.1	0.1	False
nonlinear ER $d = 30$	500	1	0.01	False
nonlinear SF $d = 30$	300	0.01	0.01	False
nonlinear ER $d = 50$	500	0.01	0.01	True
nonlinear SF $d = 50$	300	0.1	0.01	False
nonlinear ER $d = 70$	700	0.1	0.01	True
nonlinear SF $d = 70$	300	0.01	0.01	False
nonlinear ER $d = 100$	700	0.1	0.01	False
nonlinear SF $d = 100$	700	0.1	0.01	False
SynTren	300	0.1	0.01	False
Sachs Protein Cells	1200	0.1	0.01	False

Table 3: The hyperparameter selection for BayesDAG for each setting.

DIBS				
	$\alpha$	$h$ latent	$h_\theta$	$h_\sigma$
nonlinear ER $d = 20$	0.02	5	1500	10
nonlinear SF $d = 20$	0.2	5	1500	10
nonlinear ER $d = 30$	0.2	5	500	1
nonlinear SF $d = 30$	0.2	5	1000	1
nonlinear ER $d = 50$	0.2	5	500	10
nonlinear SF $d = 50$	0.2	5	1500	1
SynTren	0.2	5	500	10
Sachs Protein Cells	0.2	5	500	10

Table 4: The hyperparameter selection for DIBS for each setting.

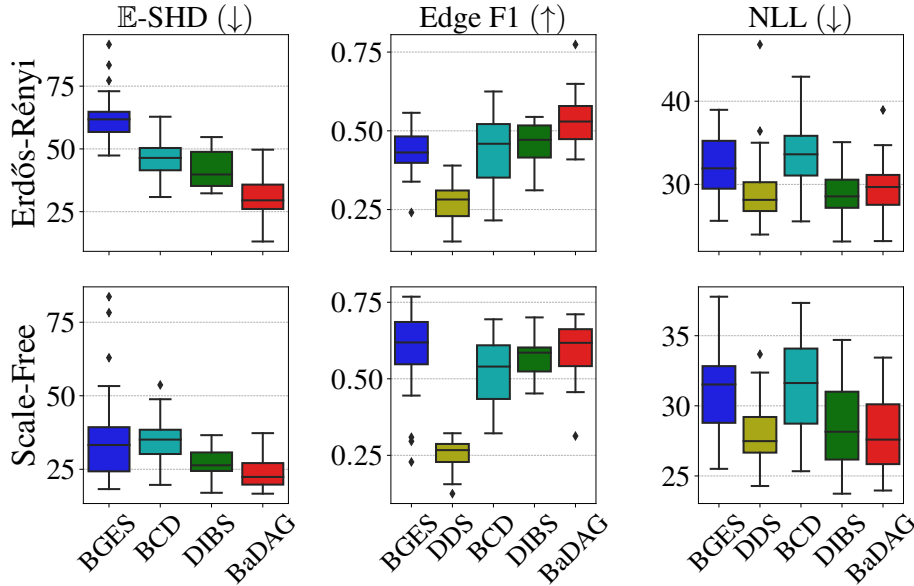


Figure 4: Posterior inference of both graph and functional parameters on synthetic datasets of nonlinear causal models with  $d = 20$  variables. `BayesDAG` gives best results across all metrics.  $\downarrow$  denotes lower is better and  $\uparrow$  denotes higher is better. For the sake of clarity, `DDS` has been omitted for  $\mathbb{E}$ -SHD due to its significantly inferior performance on this metric.

**Other hyperparameters** For `BayesDAG`, we run 10 parallel SG-MCMC chains for  $\mathbf{p}$  and  $\Theta$ . We implement an adaptive sinkhorn iteration where the iteration automatically stops when the sum of rows and columns are closed to 1 within the threshold 0.001 (upto a maximum of 3000 iterations). Typically, we found this to require only around 300 iterations. We set the sinkhorn temperature  $t$  to be 0.2. For the reparametrization of  $\mathbf{W}$  matrix with Gumbel-softmax trick, we use temperature 0.2. During evaluation, we use 100 SG-MCMC particles extracted from the particle buffer. We use 0.0003 for SG-MCMC learning rate  $l$  and batch size 512. We run 700 epochs to make sure the model is fully converged.

For `DIBS`, we can only use 20 SVGD particles for evaluation due to the quadratic scaling with the number of particles. We use 0.1 for Gumbel-softmax temperature. We run 10000 epochs for convergence. The learning rate is selected as 0.01.

[Table 3](#) shows the hyperparameter selection for `BayesDAG`. [Table 4](#) shows the hyperparameter selection for `DIBS`.

## G. Additional Results

### G.1. Performance with higher dimensional datasets

Full results for all the metrics for settings  $d = 20$ ,  $d = 70$  and  $d = 100$  for nonlinear settings are presented in [Figure 4](#), [Figure 5](#) and [Figure 6](#). We find that our method consistently outperforms the baselines with  $d = 70$  and in terms of  $\mathbb{E}$ -SHD with  $d = 100$ . Competitive performance for  $d > 50$  in nonlinear settings further demonstrates the applicability and computational efficiency of the proposed approach. In contrast, the only fully Bayesian nonlinear method, `DIBS`, is not computationally efficient to run for  $d > 50$ .

### G.2. Performance of SG-MCMC with Continuous Relaxation

We compare the performance of SG-MCMC+VI and SG-MCMC with  $\tilde{\mathbf{W}}$  on  $d = 10$  ER and SF graph settings. [Figure 7](#) shows the performance comparison. We can observe that SG-MCMC+VI generally outperforms its counterpart in most of the metrics. We hypothesize that this is because VI network  $\mu_\phi$  couples  $\mathbf{p}$  and  $\mathbf{W}$ . This coupling effect is crucial since the changes in  $\mathbf{p}$  results in the change of permutation matrix, where the  $\mathbf{W}$  can immediately respond to this change through  $\mu_\phi$ . On the other hand,  $\tilde{\mathbf{W}}$  can only respond to this change through running SG-MCMC steps on  $\tilde{\mathbf{W}}$  with fixed  $\mathbf{p}$ . In theory, this is the most flexible approach since this coupling do not requires parametric form like  $\mu_\phi$ . However in practice, we cannot

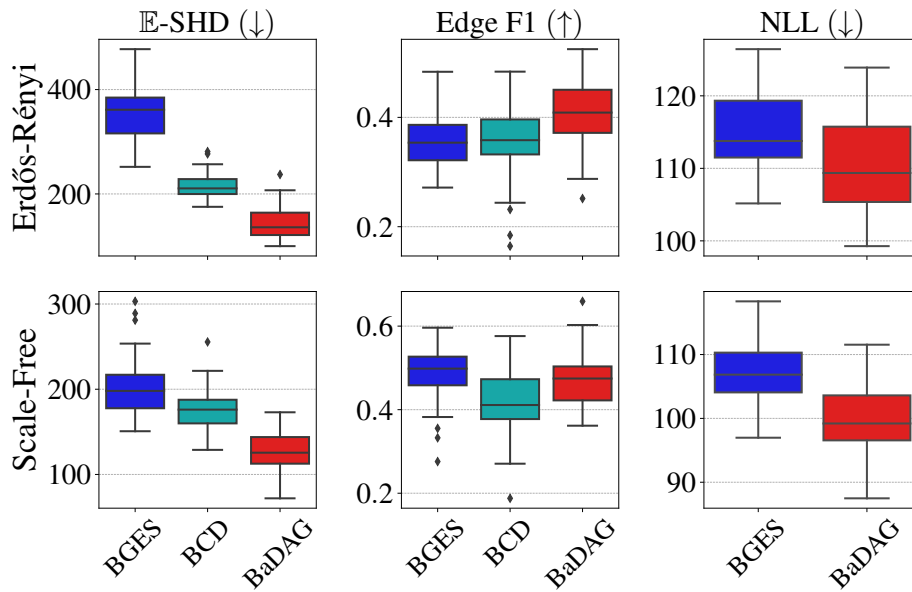


Figure 5: Posterior inference of both graph and functional parameters on synthetic datasets of nonlinear causal models with  $d = 70$  variables. BayesDAG gives best results across most metrics.  $\downarrow$  denotes lower is better and  $\uparrow$  denotes higher is better. As DIBS and DDS are computationally prohibitive to run for this setting, it has been omitted. BCD has been omitted for NLL as we observed that it performs significantly worse.

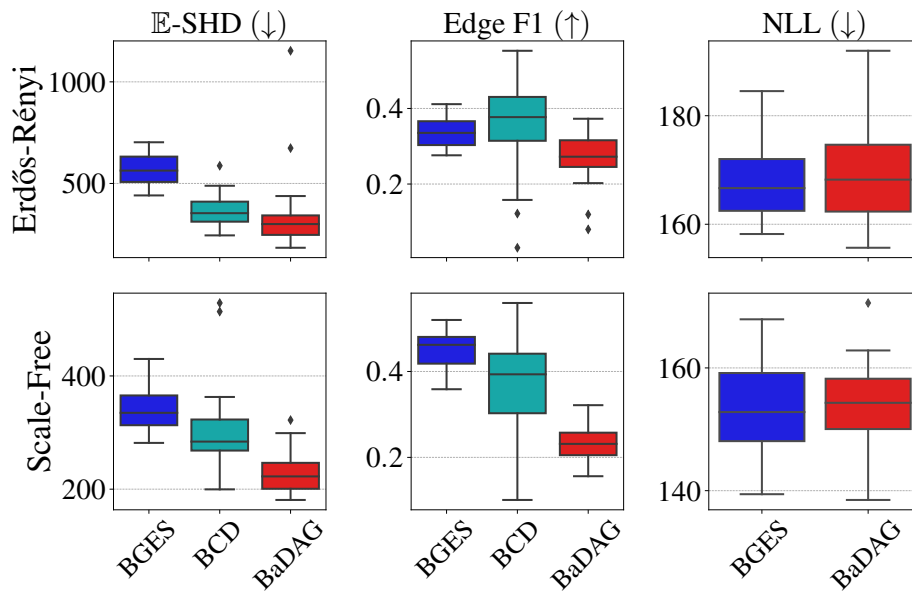


Figure 6: Posterior inference of both graph and functional parameters on synthetic datasets of nonlinear causal models with  $d = 100$  variables. BayesDAG gives best results across  $\mathbb{E}$ -SHD, comparable across NLL but slightly worse for Edge F1.  $\downarrow$  denotes lower is better and  $\uparrow$  denotes higher is better. As DIBS and DDS are computationally prohibitive to run for this setting, it has been omitted. BCD has been omitted for NLL as we observed that it performs significantly worse.



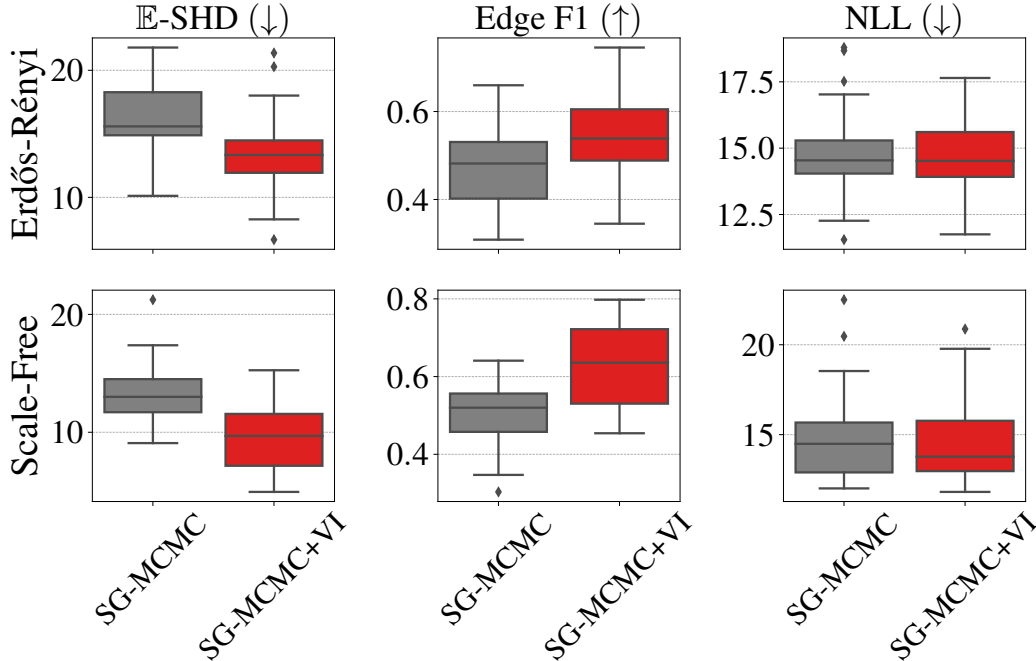


Figure 7: Performance comparison of SG-MCMC+VI v.s. fully SG-MCMC with  $\tilde{W}$  for  $d = 10$  variables.

run many SG-MCMC steps with fixed  $p$  for convergence, which results in the inferior performance.

### G.3. Ablation Study

We conduct ablation studies on our method using the nonlinear ER  $d = 30$  dataset. Key findings are summarized below. All ablation studies are conducted using ER  $d = 30$  datasets with the same hyperparameters reported in Appendix F.4.

**Initialized  $p$  scale** Figure 8 investigates the influence of the initialized scale of  $p$ . We found that the performance is the best with  $\alpha = 0.01$  or  $10^{-5}$ , and deteriorates with increasing scales. This is because with larger initialization scale, the absolute value of the  $p$  is large. Longer SG-MCMC updates are needed to reverse the node potential order, which hinders the exploration of possible permutations, resulting in the convergence to poor local optima.

**Number of SG-MCMC chains** We examine the impact of the number of parallel SG-MCMC chains in Figure 9. We observe that the number of chains does not have a significant impact on the E-SHD and Edge F1 scores.

**Injected noise level for SG-MCMC** In Figures 10 and 11, we study the performance differences arising from various injected noise levels for  $p$  and  $\Theta$  in the SG-MCMC algorithm (i.e.  $s$  of the SG-MCMC formulation in Appendix C). Interestingly, the noise level of  $p$  does not impact the performance as much as the level of  $\Theta$ . Injecting noise helps improve the performance, but a smaller noise level should be chosen for  $\Theta$  to avoid divergence from optima.

## H. Code and License

For the baselines, we use the code from the following repositories:

- BGES: We use the code from (Agrawal et al., 2019) from the repository [https://github.com/agrawalraj/active\\_learning](https://github.com/agrawalraj/active_learning) (No license included).
- DDS: We use the code from the official repository <https://github.com/sharpenb/Differentiable-DAG-Sampling> (No

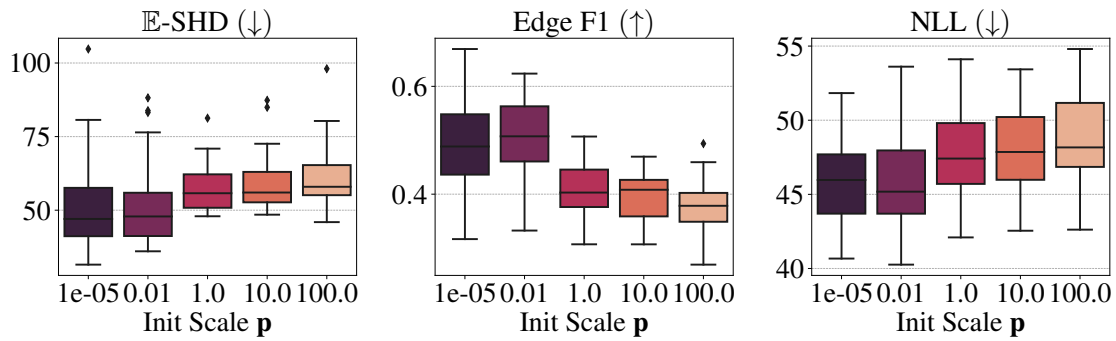


Figure 8: Posterior inference quality for  $d = 30$  ER synthetic datasets with different initialized  $p$  scale. Best performance is obtained for 0.01, with larger scales leading to progressively worse performance.

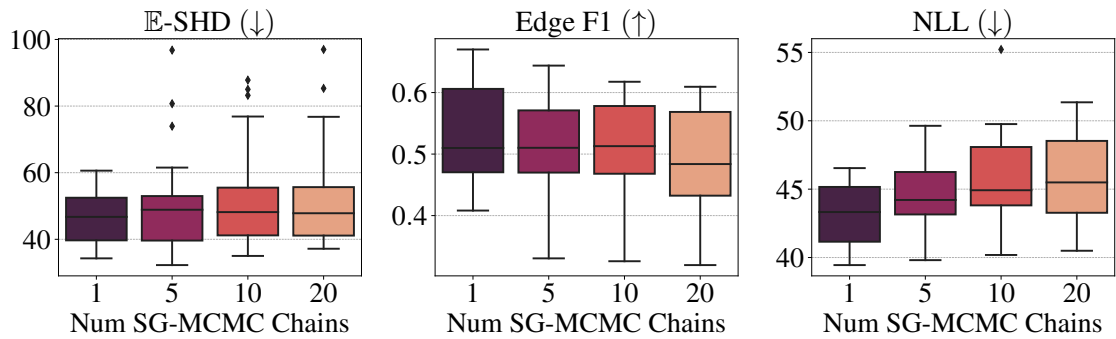


Figure 9: Posterior inference quality for  $d = 30$  ER synthetic datasets with different number of parallel SG-MCMC chains. We find that the method performs the best with fewer chains, possibly due to easier optimization and computational complexity.

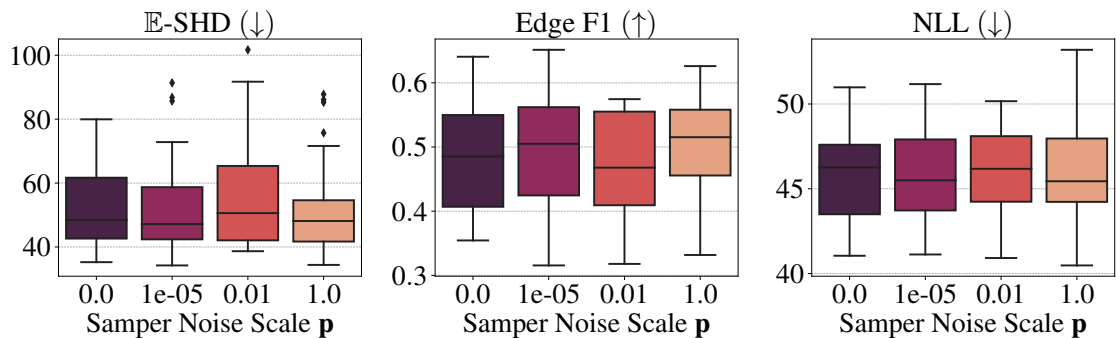


Figure 10: Posterior inference quality for  $d = 30$  ER synthetic datasets with different level of injected noise scale for  $p$ . We notice that the method is fairly insensitive to the level of noise injection.

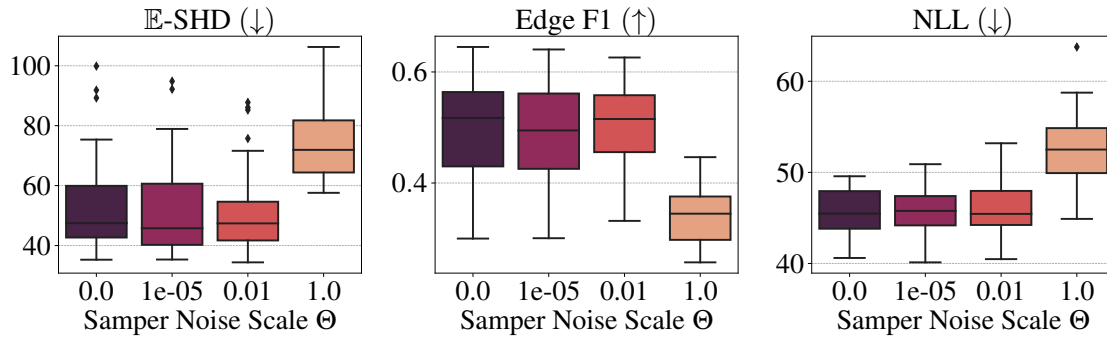


Figure 11: Posterior inference quality for  $d = 30$  ER synthetic datasets with different level of injected noise scale for  $\Theta$ . We notice that the method converges to good solutions with low level of noise injections, leading to superior performance.

license included).

- BCD: We use the code from the official repository <https://github.com/ermongroup/BCD-Nets> (No license included).
- DIBS: We use the code from the official repository <https://github.com/larslorch/dibs> (MIT license).

Additionally for the Syntren (Van den Bulcke et al., 2006) and Sachs Protein Cells (Sachs et al., 2005) datasets, we use the data provided with repository <https://github.com/kurowasan/GraN-DAG> (MIT license).