

# Distributions over DAGs for Causal Discovery: Limitations of Expressiveness

Simon Rittel<sup>1,2,3</sup> (✉) and Sebastian Tschiatschek<sup>4</sup>

<sup>1</sup> Ludwig-Maximilians-Universität München, Germany [simon.rittel@lmu.de](mailto:simon.rittel@lmu.de)

<sup>2</sup> Munich Center for Machine Learning (MCML), Munich, Germany

<sup>3</sup> UniVie Doctoral School Computer Science, Vienna, Austria

<sup>4</sup> Faculty of Computer Science, University of Vienna, Austria  
[sebastian.tschiatschek@univie.ac.at](mailto:sebastian.tschiatschek@univie.ac.at)

**Abstract.** Bayesian approaches for causal discovery can—in principle—quantify uncertainty in the prediction of the underlying causal structure, typically modeled by a directed acyclic graph (DAG). Various semi-implicit models for parametrized distributions over DAGs have been proposed, but their limitations have not been studied thoroughly. In this work, we focus on the expressiveness of parametrized distributions over DAGs in the context of causal discovery. We show several limitations of candidate models in a theoretical analysis and validate them empirically in supervised settings. To overcome these limitations, we propose using mixture models of the considered distributions over DAGs.

*This workshop paper has been updated to an extended version [25].*

## 1 Introduction

Causal discovery, also known as *causal structure learning* (CSL), is the task of uncovering cause and effect relations among modeled variables based on observed data [7, 29]. The inferred structures govern the translation of a causal estimand into a statistical one that can be measured from data and, hence, provide a basis for causal inference [15, 19]. Errors in causal discovery can imply different statistical estimands and bias the analysis of causal queries. This dependence highlights the importance of quantifying the uncertainty of a predicted causal structure. Bayesian CSL goes beyond identifying only a single, potentially incorrect, causal structure and rather models the uncertainty over the true causal structure by a distribution. It incorporates prior knowledge in the form of a prior distribution and in principle enables the computation of a posterior distribution over possible causal structure when new evidence is presented. Typical assumptions for causal discovery are causal sufficiency, the absence of selection bias, and acyclicity of the causal graph, allowing to represent the causal structure by a *directed acyclic graph* (DAG). When modeling the data-generating mechanism for such acyclic causal model, two sources of error can be distinguished. The *approximation error* results from the finite amount of evidence for the prediction of the possibly nonlinear relations among the observed variables and decreases

with increasing data size. Conversely, the model choices for the functional relationships and the distribution over the causal graph introduce a *model error* that cannot be overcome with more data. Novel Bayesian CSL algorithms often differ not only in a single model choice, but multiple ones and seldom identify in ablation studies for all of them which key detail is responsible for the reported increase in performance in contrast to competing algorithms. While a branch of research focuses exclusively on the functional relationships and parametric forms that allow identification of cause and effect pairs [8, 9, 12, 27, 34], distributions over causal graphs, e.g., DAGs, have received only little attention yet.

*Contributions* In our work, we investigate and compare the expressiveness of distributions over DAGs used in recent Bayesian CSL algorithms [1, 3–5, 13, 23, 32]. We highlight their limitations to assign equal probabilities to graphs of the same Markov equivalence class (Example 1) or capture dependencies between edges in the graphs (Example 2) and compare their ability to match synthetic graph distributions. In addition, we provide experimental and theoretical evidence that probabilistic models [5, 23, 32] and mixtures of them are more expressive than particle distributions over a substantially higher number of graphs. We believe our work will be helpful to researchers and practitioners alike by demonstrating shortcomings of recently proposed distributions over DAGs that can limit the applicability of Bayesian CSL and proposing mixture models as an effective countermeasure. In addition, we provide a recipe to evaluate the probability mass function of a given graph, a directed path or a subgraph, that is induced by a generative model by marginalization over an auxiliary random structure using importance samples.

*Structure of the paper* We begin in section 2 with the presentation of the considered distributions over DAGs and discuss their theoretical limitations. In section 3, we outline an efficient evaluation of generative models based on importance sampling and derive analytically the minimal statistical divergences for particle distributions with  $K$  graphs. Finally, we apply both techniques to supervised learning of different target distributions in section 4 and conclude our paper in section 5. Due to page limitation, we explain in Appendix section A the used notation and provide necessary background on causal discovery.

## 2 Distributions over DAGs

The number of DAGs for a given number of variables  $D$  is finite, but super-exponential in  $D$  [18, 30]. Consequently, using categorical distributions that distribute the probability mass freely among all graphs quickly becomes infeasible, even for small  $D$ . However, for many applications, it is not necessary to specify arbitrary probabilities to all DAGs, i.e. a smaller degree of freedom may suffice. For causal discovery, the posterior distribution is expected to concentrate on graphs that are *similar* to the ground truth graph according to which the observed data was generated. This motivates the design of probabilistic models

Table 1: Overview of different candidates for distributions over DAGs alongside the number of learnable parameters which depends of the number of variables  $D$ , particles  $K$ , hidden neurons  $H_N$  as well as embedding and key size  $H_E$  and  $H_K$ . Their default values are reported in section E

Graph model	Figure	Equation	# Learnable parameters
DAG with independent edges	1a	(2)	$D(D-1)$
Graph particles	1b	(3)	$K(D(D-1)+1)$
DPM-DAG [23]	1c	(5)	$(D-1)+D(D-1)$
ARCO-DAG [32]	1c	(6)	$H_N(D^2+1)+(H_N+1)D+D(D-1)$
GFlowNet-DAG [5]	1e	(7)	$2DH_E+7(H_ED^2+16(H_EH_K+H_K^2))+4H_E^2$

for graphs that are flexible enough to model any possible DAG  $\mathbf{G}$  and assign probability mass to some candidate graphs but require substantially fewer parameters than a general categorical distribution.

In the following, we discuss different models for distributions over DAGs proposed in recent works and investigate their expressiveness. To enhance comprehensibility, we introduce them by their generative model and provide the corresponding probability mass functions in section B. As a summary, we visualize their generative models in Figure 1 and provide a list of the number of learnable parameters for each model in Table 1. To avoid limitations by specific functional relationships or subjective prior knowledge, we consider independence relations as a general desideratum. In particular, we begin the investigation of the expressiveness of the candidate graph distributions with the following Markov equivalence class that serves as a running example.

*Example 1.* For a parametric linear model with exogenous Gaussian noise, the true causal graph is identifiable only up to its Markov equivalence class [20]. Consider the simple case with three variables and the chain graph  $X_1 - X_2 - X_3$  as the Markov equivalence class (MEC). Its edges cannot be further oriented even in the asymptotic limit of infinite data. To represent the corresponding uncertainty over the true DAG, all graphs of this MEC, the common cause,  $\mathbf{G}^{(1)}: X_1 \leftarrow X_2 \rightarrow X_3$ , and the two causal chains,  $\mathbf{G}^{(2)}: X_1 \rightarrow X_2 \rightarrow X_3$  and  $\mathbf{G}^{(3)}: X_1 \leftarrow X_2 \leftarrow X_3$ , should be assigned the probability  $\frac{1}{3}$ .

## 2.1 Independent edges

The arguably simplest probabilistic model for a distribution over a directed, not necessarily acyclic graph  $\mathbf{A}$  consists of a product over independent Bernoulli probabilities, each of them modeling a possible directed edge  $X_i \rightarrow X_j$ :

$$\mathbf{A} \sim q_\phi(\mathbf{A}) \quad \text{with} \quad A_{ij} \sim q_{\phi_{ij}}(A_{ij}), \quad (1)$$

where  $\phi$  are the parameters of the Bernoulli distributions over the edges. While self-loops can be directly ruled out by setting  $\forall i \in [D]: A_{ii} = 0$ , the random graphs  $\mathbf{A}$  drawn from this distribution can still have cycles. For continuous optimization, acyclicity of a point predictor can be enforced using a nonnegative,

differentiable constraint function  $h : \mathcal{R}^{D \times D} \rightarrow \mathcal{R}_+$  that evaluates to zero for any acyclic graph and otherwise to some positive value that quantifies the deviation, e.g., number of closed cycles [2, 17, 33, 35]. For probabilistic models,  $h$  can be introduced within an exponential prefactor to the unconstrained probability  $q_\phi(\mathbf{A})$ , i.e.,

$$q_{\phi,\lambda}(\mathbf{G}) \propto \exp(-\lambda h(\mathbf{G})) q_\phi(\mathbf{G}) . \quad (2)$$

For a sufficiently high prefactor  $\lambda$ , the resulting distribution  $\tilde{q}_{\phi,\lambda}(\mathbf{G})$  assigns only negligible probability mass to any cyclic graph. For an independent factorization over the edges of the graph, this comes at the cost of expressivity, since the resulting model is locked to some ordering of the nodes [24]. Note that having a set of interdependent parameters  $\phi_{ij}$  does not prevent this, since the realizations of the edges have to be coupled, e.g., two non-zero probabilities for  $P_{\phi_{ij}}(G_{ij})$  and  $P_{\phi_{ji}}(G_{ji})$  result in a positive probability for the cycle  $X_i \rightarrow X_j \rightarrow X_i$  of length two. In Example 1, each chain implies a total order and the common cause a partial order, all of them being incompatible with each other. Hence, such distribution with independent edge probabilities can only concentrate its probability mass on one of the three graphs resulting in a skewed uncertainty measure.

## 2.2 Particle representations & mixture model

Lorch et al. [13] circumvents this limitation by modeling the graph posterior  $p(\mathbf{G}|\mathcal{D})$  by a set of  $K$  particles, each representing a single DAG  $\mathbf{G}^{(k)}$ . The term particle distribution originates from the idea of approximating a continuous density function by discrete probability masses. In the context of approximating a discrete distribution, it rather refers to a simplified model that constrains the support to fewer possible outcomes, e.g.,  $K$  DAGs. The corresponding generative model can be expressed as

$$\mathbf{k} \sim q_{\mathbf{w}}(k), \quad \mathbf{G} = \mathbf{G}^{(k)}, \quad (3)$$

where  $\mathbf{w}$  are the unnormalized weights defining the probability of each particle. As an alternative to the constraint function  $h(\mathbf{G})$ , a total order induced by the permutation  $\pi$  over  $\mathbf{X}$  naturally constrains a graph to be acyclic, i.e., the adjacency matrix of the permuted graph is upper-triangular. Instead of modeling  $K$  deterministic graphs with some positive probability mass,  $K$  permutations allow to model a random upper-triangular matrix  $\mathbf{U}$ . Consequently, we no longer refer to it as a particle distribution, but a probabilistic mixture model. Denoting  $\Pi$  as the permutation matrix corresponding to  $\pi$ , the permuted random matrix  $\mathbf{U}^{(\Pi^{(k)})}$  generates a DAG  $\mathbf{G}$ :

$$\mathbf{k} \sim q_{\mathbf{w}}(k), \quad \mathbf{U} \sim q_\phi(\mathbf{U}), \quad \mathbf{G} = \mathbf{U}^{(\Pi^{(k)})} := \Pi^{(k)T} \mathbf{U} \Pi^{(k)}. \quad (4)$$

The expressivity of the mixture model can be further increased by also modeling different distributions for the upper-triangular matrix, i.e.,  $\mathbf{U}^{(k)}$ . Although

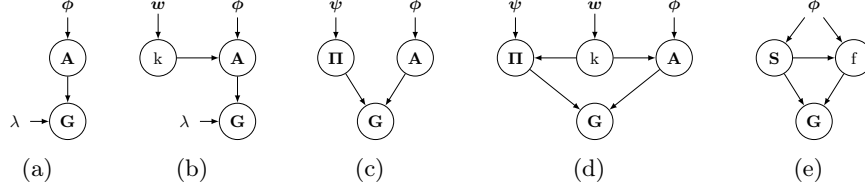


Fig. 1: Generative models of the candidate distributions over DAGs in section 2. (a) Adjacency matrix  $\mathbf{A}$  with independent edges, (b) K graph distribution extending the former by a categorical distribution over  $k$ , (c) order-based model with permutation matrix  $\Pi$  to mask  $\mathbf{A}$ , (d) K order-based model, (e) sequence-based model adding edges autoregressively until the DAG is final ( $f = 1$ ).

smaller than the number of possible DAGs, with  $D!$  the number of possible permutations is still super-exponential. Hence, only a very small fraction of all  $K$  permutations can be represented by the previous model highlighting the limitation of particle distributions. In the running Example 1 there are 25 possible DAGs and the stated MEC consisting only of the 3 graphs, however, they are each admissible under a different total order. This illustrates that here at least  $K = 3$  particles are required. Note that while the number of different MECs is still super-exponential, most of them contain only very few graphs and almost 27 % of all classes consist of only a single graph [6]. The upper bound on the size of a MEC is given by the completely connected chain graph that summarizes  $D!$  graphs each of them following a different permutation. Evidently, splitting the probability mass equally among a huge number of graphs may be analytically correct, albeit it provides a poor trade-off between expressivity and computational efficiency represented by the number of learnable parameters that scales with  $K$ . In case the distribution concentrates its probability mass on graphs that are admissible for a low number of permutations, e.g., sparse graphs with multiple components, such a model may suffice.

### 2.3 Probabilistic models over orders

Several works apply the idea of an order-based search [31] in a probabilistic generative model [1, 3, 4, 23]. Their shared underlying idea is to learn a total order of the variables that imposes an acyclicity constraint on the (sampled) adjacency matrix  $\mathbf{A}$ . Following the description of DPM-DAG [23], we outline its generative model and provide the corresponding graphical model in Figure 1c. The generative process starts by drawing a total order that defines a permutation  $\Pi$  of the variables. A permuted upper-triangular matrix of ones  $\mathbf{M}$  then defines a random acyclicity matrix  $\mathbf{M}^{(\Pi)}$  that is used to mask a sample of an unconstrained adjacency matrix  $\mathbf{A}$  as modeled in Equation (2), i.e.,

$$\Pi \sim q_{\psi}(\Pi), \quad \mathbf{A} \sim q_{\phi}(\mathbf{A}), \quad \mathbf{G} = \mathbf{M}^{(\Pi)} := (\Pi^T \mathbf{M} \Pi). \quad (5)$$

Example 1 shows the limitation of the *Plackett-Luce* (PL) distribution [21] for causal discovery [23, 32] that is used for sampling the permutation. Both chains imply a total order of the three variables, but reversed ones. In the PL model, their permutations can only receive the same probability mass in the case of uniform weights for all three variables, yet the probability for the two causal chains is then upper-bounded by  $\frac{1}{6}$  and cannot take the value of  $\frac{1}{3}$ . Some remaining probability mass is then concentrated on other graphs that do not belong to the MEC. This motivated the ARCO-DAG model [32], an autoregressive model over causal orders that computes the weights for the categorical sampling without replacement conditionally on the previous drawn sequence, i.e.,

$$q_{\psi}(\pi) = \prod_{d=1}^D q_{\psi^{(d)}}(\pi(d)) \quad \text{with} \quad \psi^{(d)} = f_{\psi}(\{\pi(i)\}_{i=1}^d) . \quad (6)$$

The function  $f_{\psi} : \mathcal{R}^{D \times D} \mapsto [D]$  takes a permutation matrix where some rows are still zeros as input and predicts the weights for sampling the next variable in the order. For their experiments, the authors of ARCO applied a multilayer perceptron with a single hidden layer with  $H > D$  neurons. To avoid marginalizing over the set of possible parents under a sampled permutation as in Toth et al. [32], the same distribution over an unconstrained graph may be used yielding Equation (5) but with an autoregressive distribution over the permutation. In case of Example 1 where the probability mass should be split equally on all graphs of the MEC  $X_1 - X_2 - X_3$ , this is consistent with setting  $P_{\mathbf{A}_{13}}(1) = P_{\mathbf{A}_{31}}(1) = 0$  and  $P_{\mathbf{A}_{ij}}(1) = 1$  otherwise.

## 2.4 Fully autoregressive model

In all covered candidate models for distributions over DAGs, the unmasked edges are still sampled from independent Bernoulli distributions as stated in Equation (2). Besides acyclicity, the edges of the causal graph drawn from these distributions are not coupled limiting the models' expressivity.

*Example 2.* Consider the four DAGs depicted in Figure 2 with  $\mathbf{G}^{(1)}$  being the true causal graph. Assume the functional relationship between  $X_1$  and  $X_2$  is identifiable and, hence, the cumulative probability of all graphs that contain this edge should be very high in the posterior. If the structural equation for variable  $X_3$ ,  $f_{X_3}(X_1, X_2, \epsilon_3)$ , depends only on both variables  $X_1$  and  $X_2$ , but not a single one alone, e.g.,  $f_3 = [X_1 > a][X_2 > b] + \epsilon_3$ , then the edges are consequently coupled. This implies that the posterior concentrate its probability mass on  $\mathbf{G}^{(1)}$  and  $\mathbf{G}^{(2)}$  that misses  $X_1$  and  $X_3$  as causes of  $X_2$ . The graphs  $\mathbf{G}^{(3)}$  and  $\mathbf{G}^{(4)}$  with only  $X_1$  or  $X_2$  as causes of  $X_3$  should be assigned small probability in this setting.

Deleu et al. [5] construct a DAG by adding edges sequentially, we refer to their graphical model as the GFlowNet-DAG model. A sampled sequence of distinct edges  $\mathbf{S} = (S_1, \dots, S_E) \in \mathcal{S}$  uniquely defines a DAG  $\mathbf{G}$  over a mapping

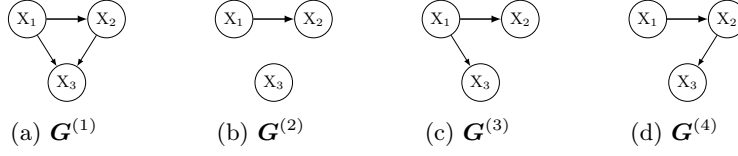


Fig. 2: Graphs of Example 2 that share the highlighted edge between  $X_1$  and  $X_2$ .

$g : \mathcal{S} \mapsto \mathcal{G}$ . Note that  $E!$  different sequences can yield the same graph  $G$  where  $E := |\mathcal{S}|$  equals its number of edges  $|\mathcal{G}|$ . The underlying sampling is based on a transformer architecture  $t_\phi : \mathcal{G} \mapsto \mathcal{R}^{D \times D}$  that autoregressively computes the parameters  $\varphi$  defining the probabilities for potential new edges,  $\varphi := t_\phi(S_{:i-1})$ . At each sampling step  $i$ , acyclicity is enforced by masking of the edges that would create a cycle, i.e., the transitive closure of the adjacency matrix of the graph that decodes ancestral relations. After the addition of a new edge, the ancestral mask is updated and the iterative sampling continues until a (binary) stop signal  $f_i = 1$  is sampled indicating that the graph sample is final, or the DAG is fully connected:

$$\mathbf{S}_i \sim q_\phi(\mathbf{S}_i | \mathbf{S}_{:i-1}), \quad f_i \sim q_\phi(f_i | \mathbf{S}_{:i-1}) \quad (7)$$

The probability of stopping  $q_\phi(1 | \mathbf{S}_{:i-1})$  theoretically guarantees that each DAG can be assigned an arbitrary probability mass provided that the function  $t_\phi$  has the capacity to model the parameter for the transition probabilities between different states exactly.

### 3 Evaluation

In contrast to Bayesian causal discovery algorithms where an observed data set is the basis for unsupervised learning of the causal graph, in this work we focus on a comparison of different models for distributions over DAGs w.r.t. their expressiveness for known target distributions. The supervised setting is motivated by shielding of other sources of error, including any bias due to the choice of functional relationships, prior distributions over the data, or the size of the observed training data. To demonstrate the limitations of the candidate models for distribution over DAGs obtained in section 2, we fit each candidate model  $q_\mathbf{G}$  to a specified target distribution  $p_\mathbf{G}$ . The target distribution is either derived by the MEC class from Example 1, by the coupling of edges from Example 2 or a synthetically generated distribution that arises by concentrating the probability mass around a target graph using the structural Hamming distance (SHD). For training of the parameters  $\phi$  with gradient descent, we take the forward KL divergence between the target distribution  $p_\mathbf{G}$  and the candidate distribution  $q_\mathbf{G}$  as loss function and approximate it using samples from the target distribution. Due to the limited support of particle distributions, we evaluate the fitted candidate distribution with the reverse KL divergence  $D_{\text{KL}}(q_\mathbf{G} \| p_\mathbf{G})$  and the total variation

distance  $D_{\text{TV}}(q_{\mathbf{G}} \| p_{\mathbf{G}})$ . In contrast to its backward formulation, the forward KL divergence estimates the log-likelihood ratio using samples from the target distribution. During supervised training with gradient descent, this ensures finding the region of the support where the probability mass is concentrated. To illustrate this important detail consider the GFlowNet model where the generation of a DAG starts with the empty graph. If the probability mass is concentrated around graphs with a moderate number of edges, almost no signal for learning is provided for graphs with only one or two of its edges.

### 3.1 Particle distributions

**Lemma 1 (Minimal statistical distances for the particle distribution).** *The reverse Kullback-Leibler (KL) divergence as well as the total variation, Hellinger and Bhattacharyya distance between a discrete target distribution  $p_{\mathbf{G}}$  and a particle representation as candidate model  $q_{\mathbf{G}}$  are minimized if the particles of the candidate  $q_{\mathbf{G}}$  represent the items with the highest probability mass in the target distribution  $p_{\mathbf{G}}$  and have their normalized probability mass.*

The minimal reverse KL divergence and total variation distance are given by:

$$\min_{q \in \mathcal{Q}} D_{\text{KL}}(q_{\mathbf{G}} \| p_{\mathbf{G}}) = -\log \sum_{k=1}^K p(\mathbf{G}^{(k)}), \quad \min_{q \in \mathcal{Q}} D_{\text{TV}}(q_{\mathbf{G}} \| p_{\mathbf{G}}) = 1 - \sum_{k=1}^K p(\mathbf{G}^{(k)})$$

$$\text{with } \forall k \in [1, K] : \mathbf{G}^{(k)} = \arg \max_{\mathbf{G} \in \mathcal{G} \setminus \{\mathbf{G}^{(i)}\}_{i < k}} p(\mathbf{G}). \quad (8)$$

We present the derivation of all minima and the proof of Lemma 1 in section C.

### 3.2 Semi-implicit generative models

A latent generative model for  $\mathbf{G}$  defines a sampling procedure using some auxiliary random objects  $\mathbf{Y} \sim q_{\phi}$ . A mechanism  $g$  maps the latent samples  $\mathbf{Y}$  to a DAG, i.e.,  $\mathbf{G} = g(\mathbf{Y})$ , and induces the joint distribution  $q(\mathbf{G}, \mathbf{Y})$ . For DAGs  $\mathbf{Y}$  can be a combination of a permutation  $\pi$  and an unconstrained adjacency matrix  $\mathbf{A}$  or alternatively a sequence of edges  $\mathbf{S}$ . Such implicit distribution cannot be directly evaluated in a graph  $\mathbf{G}^*$  sampled from the target distribution since  $g$  is in general not injective. The joint distribution has to be marginalized either by explicit summation over  $\mathbf{Y}$  or approximated by samples of  $\mathbf{Y}$ , i.e.,

$$q_{\phi}(\mathbf{G}^*) = \sum_{\mathbf{Y}} [\mathbf{G}^* = g(\mathbf{Y})] q_{\phi}(\mathbf{G}, \mathbf{Y}). \quad (9)$$

We leverage importance sampling to efficiently estimate the marginal probability of  $\mathbf{G}^*$  using a low number of samples of  $\mathbf{Y}^*$  that only generates samples of the target graph  $\mathbf{G}^* = g(\mathbf{Y}^*)$ :

$$q_{\phi}(\mathbf{G}^*) = \sum_{\mathbf{Y}^*: f(\mathbf{Y}^*) = \mathbf{G}^*} q_{\phi}(\mathbf{G}^*, \mathbf{Y}^*) = \mathbb{E}_{(\mathbf{G}^*, \mathbf{Y}^*) \sim q_{\mathbf{G}}^*} \left[ \frac{q_{\phi}(\mathbf{G}^*, \mathbf{Y}^*)}{q_{\phi}^*(\mathbf{G}^*, \mathbf{Y}^*)} \right]. \quad (10)$$



The proposal distribution  $q_\phi^*$  is derived from the distribution  $q_\phi$  by constraining the discrete samples of  $\mathbf{Y}$  to  $\mathbf{Y}^*$ . Note that due to the parameter sharing of  $\phi$ , the proposal distribution is not constant during training and is updated jointly with the candidate distribution at every iteration. We provide further details on the implementation in section D.

## 4 Experiments

In our experiments, we compare the following models from section 2 and Table 1:

- distributions over a fixed number of DAGs ( $K$  graph particle distribution),
- *reinforced probabilistically masked DAG* (RPM-DAG), a discretized version of DPM-DAG [23] that is trained with the score-function gradient estimator,
- ARCO-DAG, the autoregressive model over the causal order [32],
- Mixture models of RPM- and ARCO-DAG, as depicted in Figure 1d,
- GFlowNet-DAG[5], an autoregressive model adding edges sequentially.

We minimize the forward KL divergence with the AdamW Optimizer [14] over 1000 optimization steps. Further details and the used hyperparameters for the supervised training of each model are provided in section E.

### 4.1 Markov equivalent graphs

The first experiment consists of the MEC represented by the chain graph  $A-B-C$  motivated in Example 1. The target distribution assigns all three graphs of the equivalence class the same probability of a third. In Table 2a we report the mean of the reverse KL divergence and the total variation distance together with their standard deviations error over 20 independent runs for the RPM-DAG, ARCO-DAG, and GFlowNet-DAG model. For the particle distributions we compute the optimal empirical values according to Equation (8). The results align with the theoretical analysis that RPM-DAG fails to distribute the probability equally among all graphs of the MEC. The observed preference for the common cause  $\mathbf{G}^{(1)}$  is a model bias arising from the fact that  $\mathbf{G}^{(1)}$  induces only a partial order that is compatible with two total ones. For  $K = 3$ , a particle distribution can theoretically match the target distribution, whereas a reduced number of modeled graphs is limited to an overconfidence estimation that is outperformed by ARCO-DAG and GFlowNet-DAG. Our experimental result for a mixture of 3 RPM-DAG models further indicates that these optimal theoretical values are not reached during model training. With increasing mixture models, the statistical divergence only decreases slowly, while their standard deviations increases initially, before they decrease again. The standard errors underline the validity of the reported mean values, since they are by a factor of  $\sqrt{20}$  smaller than the reported standard deviations. The high variance is due to some outlier, illustrating that there is a regime where training is less stable. By contrast, the results for ARCO-DAG demonstrate that it can overcome the limitations of RPM-DAG by a simple autoregressive model for the permutation weights. In comparison to the

Table 2: Statistical divergences and graph probabilities for different candidate models  $q_G$  and underlying target distribution  $p_G$ . Empirical means and standard deviations are reported for 20 runs, except for the particle representations (analytical values using Equation (8)).

(a) Example 1: MEC class  $X_1 - X_2 - X_3$ .

$q_G$	$D_{KL}(q_G \  p_G)$	$D_{TV}(q_G \  p_G)$	$q(G_1)$	$q(G_2)$	$q(G_3)$
$p_G$	0	0	0.3	0.3	0.3
1 graph particle	1.0986	0.6	1	0	0
2 graph particles	0.4055	0.3	0.5	0.5	0
3 graph particles	0	0	0.3	0.3	0.3
RPM-DAG	0.32673 ± 0.00004	0.34377 ± 0.00003	0.53318 ± 0.00001	0.16145 ± 0.00002	0.16145 ± 0.00002
5-RPM-DAG	0.05 ± 0.12	0.05 ± 0.11	0.36 ± 0.07	0.31 ± 0.06	0.07 ± 0.06
10-RPM-DAG	0.0013 ± 0.0009	0.003 ± 0.008	0.335 ± 0.007	0.333 ± 0.001	0.331 ± 0.009
ARCO-DAG	0.00370 ± 0.00005	0.0038 ± 0.0003	0.3322 ± 0.0007	0.3320 ± 0.0007	0.3321 ± 0.0006
5-ARCO-DAG	0.0021 ± 0.0006	0.003 ± 0.004	0.333 ± 0.004	0.332 ± 0.005	0.333 ± 0.002
10-ARCO-DAG	0.0015 ± 0.0006	0.022 ± 0.0021	0.333 ± 0.002	0.3327 ± 0.0017	0.3324 ± 0.0012
GFlowNet-DAG	0.004 ± 0.015	0.02 ± 0.04	0.328 ± 0.026	0.34 ± 0.04	0.328 ± 0.023

(b) Example 2: Coupled edges  $X_1 \rightarrow X_3 \leftarrow X_2$  and identifiable causal effect  $X_1 \rightarrow X_2$ .

$q_G$	$D_{KL}(q_G \  p_G)$	$D_{TV}(q_G \  p_G)$	$q(G_1)$	$q(G_2)$	$q(G_3) + q(G_4)$
$p_G$	0	0	0.3	0.6	0.1
1 graph particle	0.5108	0.4	0	1	0
2 graph particles	0.1054	0.1	0.3	0.6	0
3 graph particles	0.0513	0.05	0.3158	0.6316	0.0526
4 graph particles	0	0	0.3	0.6	0.1
RPM-DAG	0.333 ± 0.009	0.36 ± 0.01	0.126 ± 0.012	0.414 ± 0.023	0.44 ± 0.02
2-RPM-DAG	0.0024 ± 0.0008	0.0046 ± 0.0047	0.2986 ± 0.0006	0.5970 ± 0.0046	0.1004 ± 0.0011
ARCO-DAG	0.328846 ± 0.000002	0.355704 ± 0.000012	0.122402 ± 0.000015	0.42189 ± 0.00003	0.45447 ± 0.00002
2-ARCO-DAG	0.00126 ± 0.00005	0.0021 ± 0.0017	0.2999 ± 0.0020	0.5993 ± 0.0020	0.1004 ± 0.0002
GFlowNet-DAG	0.0001 ± 0.0003	0.004 ± 0.004	0.30 ± 0.05	0.599 ± 0.004	0.099 ± 0.004

GFlowNet-DAG model, it yields a considerable lower total variation distance, is more stable and much faster in training as well as evaluation, due to the reduced number of model parameters (see Table 1). For mixtures of 10 models, the situation slightly reverses with 10-RPM-DAG yielding a smaller statistical distance than 10-ARCO-DAG. Note that the structure consisting of only three variables also appears in graphs with more variables and comes w.l.o.g.. In the case of a graph with  $C$  unconnected graph components, each consisting of a graph of the simple MEC class, its MEC contains  $3^C$  graphs. Under assumed ideal conditions and a supervised setting, a particle method picks  $K$  graphs of them with equal probability. Then Equation (8) implies high statistical divergences that are not competitive with parametrized distributions over orders, e.g., RPM-DAG, that scale well for unconnected components.

## 4.2 Coupled edges

Dependent edges in the posterior as in Example 2 motivate the second experiment in which we consider the corresponding graphs of Figure 2 with the following probabilities. The target distribution concentrates 60% of the probability mass on the graph  $G^{(2)}$  that misses  $X_1$  and  $X_2$  as parents of  $X_3$ , further 30% are assigned to the true causal graph  $G^{(1)}$ . To represent a coupling of the two causes, the graphs  $G^{(3)}$  and  $G^{(4)}$  that miss either of these parents get only a

Table 3: Target distribution around the MAP graph  $\mathbf{G}_0$ . Probabilities degrees with decreasing SHD.

SHD	$p(\mathcal{G}_{\text{SHD}})$	$ \mathcal{G}_{\text{SHD}} $	$p(\mathbf{G} \in \mathcal{G}_{\text{SHD}})$
0	0.15000	1	0.15000
1	0.42500	8	0.05313
2	0.22500	28	0.00804
3	0.10000	61	0.00164
4	0.05000	94	0.00053
5	0.01581	111	0.00014
6	0.01439	101	0.00014
7	0.01068	75	0.00014
8	0.00627	44	0.00014
9	0.00242	17	0.00014
10	0.00043	3	0.00014

Table 4: Statistical divergences of the candidate models for the posterior distribution in Table 3.

$q_{\mathbf{G}}$	$D_{\text{KL}}(q_{\mathbf{G}} \  p_{\mathbf{G}}) \downarrow$	$D_{\text{TV}}(q_{\mathbf{G}} \  p_{\mathbf{G}}) \downarrow$
10 graph particles	0.5395	0.4170
50 graph particles	0.1969	0.1787
100 graph particles	0.1042	0.0989
250 graph particles	0.0426	0.0417
RPM-DAG	$0.133 \pm 0.008$	$0.166 \pm 0.007$
2-RPM-DAG	$0.107 \pm 0.009$	$0.149 \pm 0.009$
10-RPM-DAG	$0.041 \pm 0.005$	$0.09 \pm 0.01$
ARCO-DAG	$0.087 \pm 0.009$	$0.146 \pm 0.007$
2-ARCO-DAG	$0.079 \pm 0.010$	$0.137 \pm 0.008$
10-ARCO-DAG	$0.058 \pm 0.007$	$0.114 \pm 0.007$
GFlowNet-DAG	$0.25 \pm 0.17$	$0.25 \pm 0.08$

probability of 0.05 each. In Table 2b the corresponding metrics for 20 runs are listed, following the same reporting as in subsection 4.1. Both RPM-DAG and ARCO-DAG, assign the graphs  $\mathbf{G}^{(3)}$  and  $\mathbf{G}^{(4)}$  too high probability as they fail to account for the coupling of the edges  $X_1 \rightarrow X_3$  and  $X_2 \rightarrow X_3$ . Due to its autoregressive model over edges, GFlowNet-DAG can approximate the target posterior with very high accuracy in terms of individual graph probabilities as well as the two evaluated statistical divergences. By design, a particle distribution is not constrained by dependent edges. Since 90% of the probability is concentrated on two graphs, a reasonable approximation can be obtained by using a mixture of only two graphs in this particular example. However, our results show that both probabilistic mixture model, with two components of RPM- or ARCO-DAG models each, outperform the simple mixture of graphs substantially.

### 4.3 Concentration of posterior mass

For an identifiable FCM and a high number of samples, the posterior should be concentrated on graphs that show a high likelihood. In the idealized setting of perfect regression of a child on its parents, the likelihood can be expected to peak for the true underlying graph  $\mathbf{G}_0$  provided that parameter uncertainty or regularization prevents superfluous edges that are not in  $\mathbf{G}_0$ . This implies that high-scoring graphs are 'similar' to  $\mathbf{G}_0$  implied by the FCM. In the absence of an analytic posterior that motivates such similarity, we generate a synthetic target distribution around the assumed *maximum-a-posteriori* (MAP) graph  $\mathbf{G}_0$  depicted in Figure 5 that has positive support for all 543 possible DAGs with 4 nodes. We start by assigning probability masses to the groups of graphs that have a *structural Hamming distance* (SHD) up to 4, i.e., 4 different entries in the adjacency matrix, and split the remaining probability mass equally among all graphs with greater SHD. For the five groups with SHD up to 4, we distribute their cumulated probability masses equally among all graphs within the respective group. The target distribution is summarized in Table 3 which lists the probability masses for the groups of graphs with the same SHD along side individual ones. We report in Table 4 the mean values and standard deviations for the statistical divergences over 20 independent runs for all candidate

models, except the particle representations for which the optimal values are computed analytically using Equation (8). The results highlight the expressivity of the RPM-DAG and its extension ARCO-DAG that both outperform a particle distribution consisting of the 50 graphs with the highest probability in the target distribution that account for a cumulative probability of 86.3%. A single ARCO-DAG yields lower statistical divergences than a mixture of 2 RPM-DAG models in this very setting. When comparing mixture models with  $K \geq 3$ , it can be observed that mixtures of RPM-DAG yields better results than ARCO-DAG. Although we trained in an idealized very low-dimensional setting with samples drawn from the target distribution, the performance of the GFlowNet-DAG model does not match its performance from subsection 4.1 and 4.2. We conjecture that the gradients from a changing variety of different graphs do not yield a sufficiently stable training signal that is necessary to tune the transformer network.

## 5 Conclusion

Bayesian causal discovery promises uncertainty quantification in the prediction of the causal graph. We reviewed several candidate models for distribution over DAGs and investigated the limitations of their expressivity. To minimize confounding factors such as the influence of functional relationships or the size of the training set, we considered an idealized supervised setting that can also be used to specify prior distributions. We showed in our theoretical analysis and experimental results that all considered candidate models, besides the autoregressive GFlowNet-DAG model, are theoretically not sufficiently expressive to match simple target distributions in which edges are coupled. While GFlowNet-DAG satisfies the theoretical requirements, we cannot verify its expressiveness in a low-dimensional experiment where graph samples are drawn from the target distribution with support over all possible DAGs. Since causal structure learning is typically an unsupervised problem, this poses a major limitation to this model.

Coupled edges due to interaction effects pose a general challenge in causal structure learning and affect most algorithms—not only Bayesian approaches. Our results further suggest that a mixture of a small number of simple probabilistic models as RPM-DAG models may approximate distributions with coupled edges sufficiently well in practical applications and outperform particle distributions with a moderate number of modeled graphs. Its extension, ARCO-DAG is more expressive w.r.t. probabilistic causal orders and shows competitive performance to GFlowNet-DAG, while having a substantially lower number of parameters and is more stable in training. When comparing mixture models of RPM- and ARCO-DAG, the simpler model outperforms the latter more expressive model for 10 mixture components consistently in all experiments. Therefore, we conjecture that a mixture of RPM-DAG models is particularly suited to scale Bayesian causal discovery to higher numbers of variables.

## Bibliography

- [1] Annadani, Y., Pawlowski, N., Jennings, J., Bauer, S., Zhang, C., Gong, W.: BayesDAG: Gradient-based posterior inference for causal discovery. In: Advances in Neural Information Processing System, vol. 36 (2023)
- [2] Bello, K., Aragam, B., Ravikumar, P.: DAGMA: learning DAGs via M-matrices and a log-determinant acyclicity characterization. In: Advances in Neural Information Processing Systems, vol. 35 (2022)
- [3] Charpentier, B., Kibler, S., Günnemann, S.: Differentiable DAG sampling. In: International Conference on Learning Representations, 10, Openreview (2022)
- [4] Cundy, C., Grover, A., Ermon, S.: BCD nets: Scalable variational approaches for Bayesian causal discovery. In: Advances in Neural Information Processing Systems, vol. 34, pp. 7095–7110 (2021)
- [5] Deleu, T., Góis, A., Emezue, C.C., Rankawat, M., Lacoste-Julien, S., Bauer, S., Bengio, Y.: Bayesian structure learning with generative flow networks. In: Conference on Uncertainty in Artificial Intelligence, Proceedings of Machine Learning Research, vol. 180, pp. 518–528, PMLR (2022)
- [6] Gillispie, S.B., Perlman, M.D.: Enumerating markov equivalence classes of acyclic digraph models. In: Proceedings of the 17th Conference in Uncertainty in Artificial Intelligence, pp. 171–177, Morgan Kaufmann (2001)
- [7] Glymour, C., Zhang, K., Spirtes, P.: Review of causal discovery methods based on graphical models. *Frontiers in Genetics* **10**, 524 (2019)
- [8] Hoyer, P.O., Janzing, D., Mooij, J.M., Peters, J., Schölkopf, B.: Nonlinear causal discovery with additive noise models. In: Advances in Neural Information Processing Systems, vol. 21, pp. 689–696 (2008)
- [9] Immer, A., Schultheiss, C., Vogt, J.E., Schölkopf, B., Bühlmann, P., Marx, A.: On the identifiability and estimation of causal location-scale noise models. In: International Conference on Machine Learning, Proceedings of Machine Learning Research, vol. 202, pp. 14316–14332, PMLR (2023)
- [10] Jang, E., Gu, S., Poole, B.: Categorical reparameterization with gumbel-softmax. In: International Conference on Learning Representations, 5, Openreview (2017)
- [11] Koller, D., Friedman, N.: Probabilistic Graphical Models - Principles and Techniques. MIT Press (2009)
- [12] Loh, P., Bühlmann, P.: High-dimensional learning of linear causal networks via inverse covariance estimation. *Journal of Machine Learning Research* **15**(1), 3065–3105 (2014)
- [13] Lorch, L., Rothfuss, J., Schölkopf, B., Krause, A.: DiBS: Differentiable Bayesian structure learning. In: Advances in Neural Information Processing Systems, vol. 34, pp. 24111–24123 (2021)
- [14] Loshchilov, I., Hutter, F.: Decoupled weight decay regularization. In: International Conference on Learning Representations, 7, OpenReview (2019)

- [15] Lundberg, I., Johnson, R., Stewart, B.M.: What is your estimand? Defining the target quantity connects statistical evidence to theory. *American Sociological Review* **86**(3), 532–565 (2021)
- [16] Maddison, C.J., Mnih, A., Teh, Y.W.: The concrete distribution: A continuous relaxation of discrete random variables. In: *International Conference on Learning Representations*, 5, Openreview (2017)
- [17] Nazaret, A., Hong, J., Azizi, E., Blei, D.: Stable differentiable causal discovery. In: *International Conference on Machine Learning*, OpenReview (2024)
- [18] OEIS Foundation Inc.: The On-Line Encyclopedia of Integer Sequences (2023), published electronically at <http://oeis.org/A003024>
- [19] Pearl, J.: *Causality*. Cambridge University Press, 2nd edn. (2009)
- [20] Peters, J., Janzing, D., Schölkopf, B.: *Elements of causal inference: Foundations and learning algorithms*. The MIT Press (2017)
- [21] Plackett, R.L.: The analysis of permutations. *Journal of the Royal Statistical Society Series C: Applied Statistics* **24**(2), 193–202 (1975)
- [22] Prillo, S., Eisenschlos, J.: Softsort: A continuous relaxation for the argsort operator. In: *International Conference on Machine Learning, Proceedings of Machine Learning Research*, vol. 119, pp. 7793–7802, PMLR (2020)
- [23] Rittel, S., Tschitschek, S.: Specifying prior beliefs over DAGs in deep Bayesian causal structure learning. In: *European Conference on Artificial Intelligence, Frontiers in Artificial Intelligence and Applications*, vol. 372, pp. 1962–1969, IOS Press (2023)
- [24] Rittel, S., Tschitschek, S.: On differentiable Bayesian causal structure learning. In: *UAI Workshop on Causal Inference* (2024)
- [25] Rittel, S., Tschitschek, S.: Expressiveness of parametrized distributions over DAGs for causal discovery. *Transactions on Machine Learning Research* (2025)
- [26] Shah, R.D., Peters, J.: The hardness of conditional independence testing and the generalised covariance measure. *The Annals of Statistics* **48**(3), 1514–1538 (2020)
- [27] Shimizu, S., Hoyer, P.O., Hyvärinen, A., Kerminen, A.J.: A linear non-Gaussian acyclic model for causal discovery. In: *Journal of Machine Learning Research*, vol. 7, pp. 2003–2030, JMLR (2006)
- [28] Spirtes, P., Glymour, C., Scheines, R.: *Causation, Prediction and Search. Adaptive Computation and Machine learning*, MIT Press, 2nd edn. (2000)
- [29] Squires, C., Uhler, C.: Causal structure learning: A combinatorial perspective. *Foundations of Computational Mathematics* pp. 1–35 (2022)
- [30] Stanley, R.P.: Acyclic orientations of graphs. *Discrete Mathematics* **5**(2), 171–178 (1973)
- [31] Teyssier, M., Koller, D.: Ordering-based search: A simple and effective algorithm for learning Bayesian networks. In: *Proceedings of the 21st Conference in Uncertainty in Artificial Intelligence*, pp. 548–549, AUAI Press (2005)
- [32] Toth, C., Knoll, C., Pernkopf, F., Peharz, R.: Effective Bayesian causal inference via structural marginalisation and autoregressive orders. In: *ICML Workshop on Structured Probabilistic Inference & Generative Modeling* (2024)

- [33] Yu, Y., Chen, J., Gao, T., Yu, M.: DAG-GNN: DAG structure learning with graph neural networks. In: International Conference on Machine Learning, Proceedings of Machine Learning Research, vol. 97, pp. 7154–7163, PMLR (2019)
- [34] Zhang, K., Hyvärinen, A.: On the identifiability of the post-nonlinear causal model. In: Proceedings of the 25th Conference on Uncertainty in Artificial Intelligence, pp. 647–655, AUAI Press (2009)
- [35] Zheng, X., Aragam, B., Ravikumar, P., Xing, E.P.: DAGs with NO TEARS: continuous optimization for structure learning. In: Advances in Neural Information Processing Systems, vol. 31, pp. 9492–9503 (2018)

## A Preliminaries

### A.1 Probability and random variables

We denote scalar random variables as  $y$ , random vectors as  $\mathbf{y}$ , and random matrices as  $\mathbf{Y}$ . Single elements of a random vector or matrix are scalars and written as  $y$  or  $Y$  respectively. The probability of a discrete random variable  $y$  taking the value  $y$  is expressed by the probability mass function  $p_y(y) := \mathbb{P}(y = y)$ . When clear from the context, we omit the random variable in the subscript of a probability mass function in favor of the parameters  $\boldsymbol{\theta}$  of the distribution to increase readability, e.g.,  $p_{\boldsymbol{\theta}}(y)$ . We mainly apply this concise notation in the main text to introduce the generative models. With a slight abuse of terminology, we refer to a distribution  $P_y$  by its induced probability mass function  $p_y(y) = P_y(\{y\})$ . To distinguish approximations of the true target distribution, we denote model distributions by  $q_y(y) = Q_y(y)$  or simply  $q_{\boldsymbol{\theta}}(y)$ .

### A.2 Functional causal models

A *functional causal model* (FCM) is defined as the triple  $\mathcal{M}_{\mathbf{X}} := \{\mathbf{X}, (\boldsymbol{\epsilon}, P_{\boldsymbol{\epsilon}}), \mathbf{f}\}$  consisting of a set of endogenous random variables  $\mathbf{X}$ , a set of exogenous noise variables  $\boldsymbol{\epsilon}$  with joint probability distribution  $P_{\boldsymbol{\epsilon}}$  and a set of deterministic functions  $\mathbf{f}$ , all three indexed by  $[D] := \{1, \dots, D\}$ . Each endogenous variable  $X_d$  of the model is generated by a function of a subset of the endogenous variables  $\mathbf{X}$  and its exogenous noise  $\varepsilon_d$ , i.e.,

$$X_d := f_d(\mathbf{X}, \varepsilon_d) = f_d(\text{Pa}_{\mathbf{G}}(X_d), \varepsilon_d) \quad \forall d \in [D]. \quad (11)$$

The distribution of  $X_d$  is implicitly defined as the pushforward measure of  $P_{\boldsymbol{\epsilon}}$  through the *causal mechanism*  $f_d$ . The structure induced by the direct functional dependencies is often restricted to be acyclic such that it can be represented by a *directed acyclic graph* DAG or—equivalently—its adjacency matrix  $\mathbf{G} \in \mathcal{G} \subset \{0, 1\}^{D \times D}$  with a one-to-one correspondence between random variables and nodes. Let  $\sim d$  denote the index set  $[D] \setminus \{d\}$ , then the  $d$ -th column of  $\mathbf{G}$  encodes the parents  $\text{Pa}_{\mathbf{G}}(X_d) \subseteq \mathbf{X}_{\sim d}$  of a node/random variable  $X_d$ , i.e., the subset of  $\mathbf{X}_{\sim d}$  that has a direct influence on  $X_d$  via  $f_d$ . Ancestors of a random variable,  $\text{An}_{\mathbf{G}}(X_d)$ , have a directed path to  $X_d$  in the causal graph  $\mathbf{G}$  and indirectly influence  $X_d$ . Their causal effect is mediated by at least one parent of  $X_d$  (possibly themselves), it can be blocked by conditioning on all its parents  $\text{Pa}_{\mathbf{G}}(X_d)$ . Children  $\text{Ch}_{\mathbf{G}}(X_d)$  and descendants  $\text{De}_{\mathbf{G}}(X_d)$  are affected by changes of the corresponding node  $X_d$  and are defined complementary to parents and ancestors.



### A.3 Causal structure learning

The objective of causal discovery is to learn the underlying causal graph  $\mathbf{G}$  from observed random variables  $\mathbf{X}$  that encodes the causal effects implied by the FCM  $\mathcal{M}_{\mathbf{X}}$ . Throughout this work, we assume causal sufficiency, i.e., all endogenous variables  $\mathbf{X}$  are observable and the exogenous noise variables  $\epsilon$  are mutually independent. This implies that all dependencies and independencies between the observed values of the random variables  $\mathbf{X}$  result from their causal effects over the functions  $\mathbf{f}$  and not from some latent common causes, i.e., an unobserved shared ancestors of them. In addition, we assume that all samples in the data set  $\mathcal{D} := \{\mathbf{X}^{(n)}\}_{n=1}^N$  are generated i.i.d. from the FCM  $\mathcal{M}_{\mathbf{X}}$  without any selection bias for the generated samples, e.g., there is no conditioning on unobserved confounders.

Without any assumptions on the functions  $\mathbf{f}$ , the *Markov equivalence class* (MEC) of a causal graph  $\mathbf{G}$  can be consistently identified from  $\mathcal{D}$  using *conditional independence* (CI) tests [28]. This equivalence class contains all DAGs that entail the same observed independence relations and can be compactly represented by a *completely partially directed acyclic graph* (CPDAG), an acyclic mixed graph that has the same adjacencies,  $X_i - X_j$ , and unshielded colliders,  $X_i \rightarrow X_k \leftarrow X_j$ , as the underlying true graph  $\mathbf{G}$ , but with some of its edges remaining undirected. While CI tests can be easily parallelized and their required overall number can be sequentially restricted by the individual test results as in the PC algorithm [28], the combination of the uncertainty attached to each test is non-trivial. Moreover, CI testing for continuous random variables lacks statistical power against alternatives and suffers from the curse of dimensionality [26].

### A.4 Bayesian causal discovery

Alternatively, score-based algorithms for causal discovery optimize a scalar quantity that is typically derived from the likelihood of the observed data  $p_{\boldsymbol{\theta}}(\mathcal{D}|\mathbf{G})$ . Since the maximum likelihood  $p_{\boldsymbol{\theta}^*}(\mathcal{D}|\mathbf{G})$  can only improve with increasing number of allowed parents as covariates, regularization of the number of edges is required to prevent estimating a fully connected DAG. Bayesian causal structure learning is also based on the likelihood, but models the full data generating process for the so called evidence  $p(\mathcal{D})$  that includes the uncertainty over model parameters  $\boldsymbol{\theta}$  and a prior distribution over the causal graph  $p(\mathbf{G})$ :

$$p(\mathcal{D}) = \int \int \prod_{n=1}^N p(\mathbf{X}^{(n)}|\boldsymbol{\theta})p(\boldsymbol{\theta}|\mathbf{G})p(\mathbf{G}) \, d\boldsymbol{\theta} \, d\mathbf{G} = \int p(\mathcal{D}|\mathbf{G})p(\mathbf{G}) \, d\mathbf{G} \quad (12)$$

The marginal likelihood  $p(\mathcal{D}|\mathbf{G})$  arises by averaging over the (conditional) prior distribution of model parameters  $p(\boldsymbol{\theta}|\mathbf{G})$ . In contrast to the maximum likelihood estimate  $p_{\boldsymbol{\theta}^*}(\mathcal{D}|\mathbf{G})$ , it avoids overfitting to the noise of the data  $\mathcal{D}$  [11]. The posterior distribution over the causal  $p(\mathbf{G}|\mathcal{D})$  that quantifies the uncertainty over the true causal graph is obtained by Bayes' formula:

$$p(\mathbf{G}|\mathcal{D}) = \frac{p(\mathcal{D}|\mathbf{G})p(\mathbf{G})}{p(\mathcal{D})}. \quad (13)$$

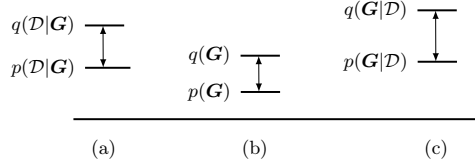


Fig. 3: Modeling error in Bayesian causal discovery due to approximations  $q$  of (a) the marginal likelihood, (b) the model for the prior and (c) the posterior over the DAG  $\mathbf{G}$ .

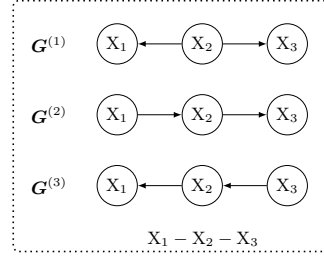


Fig. 4: Graphs of the MEC from Example 1 along its CPDAG.

Computing the evidence involves marginalization over all DAGs and parameters of the likelihood and is intractable. Therefore, approximate methods for Bayesian inference are needed. In the following we focus on variational inference where the joint distribution  $p(\mathbf{G}, \mathcal{D})$  is approximated by a generative model  $q(\mathcal{D}|\mathbf{G})p(\mathbf{G})$ , and a variational family of distributions  $\mathcal{Q} := \{q_\phi(\mathbf{G}|\mathcal{D}) \mid \phi \in \mathbb{R}\}$  is specified for the posterior  $q(\mathbf{G}|\mathcal{D})$ . Both steps involve parametric distributions over DAGs.

### A.5 Bayesian model error

Given a fixed data set  $\mathcal{D}$ , deviations of a modeled posterior  $q(\mathbf{G}|\mathcal{D})$  to the true posterior  $p(\mathbf{G}|\mathcal{D})$  can be backtracked to three different model errors as depicted in Figure 3. The modeled marginal likelihood  $q(\mathcal{D}|\mathbf{G})$  may induce a bias, with the particular case of the maximum likelihood estimate leading to an overconfident prediction. Throughout this work, we assume a flexible model for the marginal likelihood and focus on the role of the distribution over DAGs instead. The model for the prior distribution  $q(\mathbf{G})$  directly constrains the belief that domain experts can express over the causal structure, but the subjective bias vanishes in the asymptotic limit as long as it assigns each DAG some positive probability mass. Lastly, the true posterior may not be in the variational family  $\mathcal{Q}$ , i.e. cannot be expressed by  $q(\mathbf{G}|\mathcal{D})$  which poses a direct constrain for the outcome of the Bayesian analysis.

## B Probability mass functions over DAGs

For conciseness, we present in section 2 only the generative models for DAGs that implicitly induce a probability distribution. In the following, we provide their analytic probability mass functions.

### B.1 Independent Edges

$$q_{\mathbf{A}}(\mathbf{A}) = \prod_{i \neq j}^D q_{A_{ij}}(A_{ij}) \quad (14)$$

$$q_{\mathbf{G}}(\mathbf{G}) = \frac{1}{Z} \exp(-\lambda h(\mathbf{G})) q_{\mathbf{A}}(\mathbf{G})$$

with  $Z := \sum_{\mathbf{G} \in \mathcal{G}} \exp(-\lambda h(\mathbf{G})) q_{\mathbf{A}}(\mathbf{G})$ . (15)

### B.2 Particle distribution

$$q_{\mathbf{G}}(\mathbf{G}) = \sum_{k=1}^K \left[ \mathbf{G}^{(k)} = \mathbf{G} \right] q_k(k), \quad (16)$$

### B.3 Probabilistically masked DAG

$$q_{\mathbf{G}}(\mathbf{G}) = \sum_{\mathbf{\Pi}, \mathbf{A}} \left[ \mathbf{M}^{(\mathbf{\Pi})} \circ \mathbf{A} = \mathbf{G} \right] q_{\mathbf{\Pi}}(\mathbf{\Pi}) q_{\mathbf{A}}(\mathbf{A})$$

with  $\mathbf{M}^{(\mathbf{\Pi})} := \mathbf{\Pi}^T \mathbf{M} \mathbf{\Pi}$ . (17)

Sampling without replacement from a categorical distribution with some fixed weights  $\{w_d\}_{d=1}^D$  is equivalent to drawing a permutation over  $[D]$  and known as the Plackett-Luce distribution [21]. At each sampling stage  $d$  the categorical weights for the remaining variables that were not yet sampled are normalized :

$$q_{\pi}(\pi) = \prod_{d=1}^{D-1} q_{\pi(d)}(\pi(d)) = \prod_{d=1}^{D-1} \frac{\psi_{\pi(d)}}{\sum_{i=1}^D \psi_i - \sum_{j < d} \psi_{\pi(j)}}. \quad (18)$$

A differentiable, continuous relaxation of the discrete sampling of permutation matrices  $\mathbf{\Pi}$  can be obtained by pairing the Gumbel-Softmax trick [10, 16] with Softsort [22].

#### B.4 Mixture of DAGs with probabilistic entries

$$q_{\mathbf{G}}(\mathbf{G}) = \sum_{k=1}^K \sum_{\mathbf{U}} \left[ \mathbf{U}^{(\boldsymbol{\Pi}^{(k)})} = \mathbf{G} \right] q_{\mathbf{U}}(\mathbf{U}) q_k(k)$$

with  $\mathbf{U}^{(\boldsymbol{\Pi}^{(k)})} := \boldsymbol{\Pi}^{(k)T} \mathbf{U} \boldsymbol{\Pi}^{(k)}$ .

(19)

#### B.5 Autoregressive model over all potential edges

Denoting the mapping between the sequence of edges  $\mathbf{S} \in \mathcal{S}$  to a graph  $\mathbf{G}$  by  $g : \mathcal{S} \mapsto \{0, 1\}^{D \times D}$ , the probability of sampling a graph  $\mathbf{G}$  equals:

$$q_{\mathbf{G}}(\mathbf{G}) = \sum_{\mathbf{S} \in \mathcal{S}} q_{\mathbf{G}, \mathbf{S}}(\mathbf{G}, \mathbf{S}) = \sum_{\mathbf{S} \in \mathcal{S}} [\mathbf{G} = g(\mathbf{S})] q_{\mathbf{S}}(\mathbf{S}). \quad (20)$$

$$q_{\mathbf{S}}(\mathbf{S}) = \left( \prod_{i=1}^E q_{S_i | \mathbf{S}_{:i-1}}(S_i | \mathbf{S}_{:i-1}) \right) \left( \prod_{j=1}^{E-1} q_{f_j | \mathbf{S}_{:j-1}}(0 | \mathbf{S}_{:j-1}) \right) q_{f_E | \mathbf{S}_{:E-1}}(1 | \mathbf{S}_{:E-1}). \quad (21)$$

The sequential process of a GFlowNet (GFN) itself defines a DAG where all non-leaf nodes represent the states of an unfinished graph  $\mathbf{G}^{(i)} = g(\mathbf{S}_{:i})$ . Starting with the empty graph state as the root node and probability 1, it splits the probability of a state representing a preliminary graph among its children who are either other preliminary graphs with an additional edge or finalized graphs.

## C Particle distributions

In the following we present a proof of Lemma 1 and derive the minimal values for the four statistical distances between a discrete target distribution  $p$  and a reduced distribution  $q$  with only  $K$  particles.

#### C.1 Minimal total variation distance

The total variation distance  $D_{\text{TV}}$  between the candidate and target distribution,  $q_{\mathbf{G}}$  and  $p_{\mathbf{G}}$ , can be rewritten as:

$$\begin{aligned} D_{\text{TV}}(q_{\mathbf{G}} \| p_{\mathbf{G}}) &= \frac{1}{2} \sum_{\mathbf{G} \in \mathcal{G}} |q(\mathbf{G}) - p(\mathbf{G})| \\ &= \frac{1}{2} \left( \underbrace{\sum_{\mathbf{G} \in \mathcal{G} \setminus \{\mathbf{G}^{(k)}\}_{k=1}^K} p(\mathbf{G})}_{\sum_{k=1}^K p(\mathbf{G}^{(k)})} + \sum_{k=1}^K \left| q(\mathbf{G}^{(k)}) - p(\mathbf{G}^{(k)}) \right| \right). \quad (22) \\ &= 1 - \sum_{k=1}^K p(\mathbf{G}^{(k)}) \end{aligned}$$

In the RHS of Equation (22) the sum over all graphs is split into two terms. The graphs that are not in the support of the particle distribution  $q_{\mathbf{G}}$  lead to a distance that is independent of the weights for the particles. The second term measures the absolute deviations for the probabilities of the  $K$  graphs that are assigned non-zero probabilities in  $q_{\mathbf{G}}$ .

For equality of the probability masses  $p(\mathbf{G}^{(k)})$  and  $q(\mathbf{G}^{(k)}) \forall k \in [K]$ , it would be minimized. Due to the normalization constraint, the accumulated probability mass of the not modeled graphs has to be distributed among the  $K$  graphs in  $q_{\mathbf{G}}$ , leading to  $q(\mathbf{G}^{(k)}) \geq p(\mathbf{G}^{(k)}) \forall k \in [K]$ :

$$\begin{aligned} \min_{q \in \mathcal{Q}} D_{\text{TV}}(q_{\mathbf{G}} \| p_{\mathbf{G}}) &= \frac{1}{2} \left( 1 - \sum_{k=1}^K p(\mathbf{G}^{(k)}) + \left( \sum_{k=1}^K q(\mathbf{G}^{(k)}) - p(\mathbf{G}^{(k)}) \right) \right) \\ &= 1 - \sum_{k=1}^K p(\mathbf{G}^{(k)}). \end{aligned} \quad (23)$$

Consequently, both sums in the previous equation amount precisely to the loss of the accumulated probability mass for the graphs not modeled by  $q_{\mathbf{G}}$  and the total variation distance is minimized when  $q_{\mathbf{G}}$  assigns non-zero probability to the  $K$  graphs with the highest probability in  $p_{\mathbf{G}}$ :

$$\forall k \in [K] : \quad \mathbf{G}^{(k)} = \arg \max_{\mathbf{G} \in \mathcal{G} \setminus \{\mathbf{G}^{(i)}\}_{i < k}} p_{\mathbf{G}}(\mathbf{G}). \quad (24)$$

In general their corresponding probabilities  $q(\mathbf{G}^{(k)})$  are not unique. The approximated probability distribution  $q_{\mathbf{G}}$  with the minimal  $D_{\text{TV}}$  to  $p_{\mathbf{G}}$  does not constrain how the additional probability mass is distributed.

In principle the probability of a single graph could account for the missing accumulated probability mass yielding only a single biased value, but distorting the relative probabilities. This limits the utility of such an approximated distribution heavily and motivates the following constraint:

$$\forall i \neq j \in [K] : \quad \frac{q(\mathbf{G}^{(i)})}{q(\mathbf{G}^{(j)})} = \frac{p(\mathbf{G}^{(i)})}{p(\mathbf{G}^{(j)})}. \quad (25)$$

Preserving the relative probabilities of  $p_{\mathbf{G}}$  in  $q_{\mathbf{G}}$  uniquely defines  $q_{\mathbf{G}}$  by normalized probabilities

$$q(\mathbf{G}^{(k)}) = \frac{p(\mathbf{G}^{(k)})}{\sum_{j=1}^K p(\mathbf{G}^{(j)})}, \quad (26)$$

and yields a approximation that is for all  $K$  graphs overconfident w.r.t to the approximated distributed  $p_{\mathbf{G}}$ , i.e.,  $q(\mathbf{G}^{(k)}) > p(\mathbf{G}^{(k)}) \forall k \in [K] \subset |\mathcal{G}|$ .  $\square$

## C.2 Minimal Kullback-Leibler divergence

In contrast to the total variation distance, the KL divergence is not symmetric. Due to the smaller support of the candidate distribution  $q_{\mathbf{G}}$  compared to the

target distribution over all DAGs  $p_{\mathbf{G}}$  only the reverse KL divergence,  $D_{\text{KL}}(q|p)$ , is properly defined. Imposing normalized probabilities as in Equation (26) to preserve the relative probabilities of the graphs in the target distribution  $p_{\mathbf{G}}$  simplifies the KL divergence to:

$$\begin{aligned} D_{\text{KL}}(q_{\mathbf{G}} \| p_{\mathbf{G}}) &= \sum_{\mathbf{G} \in \mathcal{G}} q(\mathbf{G}) \log \frac{q(\mathbf{G})}{p(\mathbf{G})} = \sum_{k=1}^K q(\mathbf{G}^{(k)}) \log \frac{q(\mathbf{G}^{(k)})}{p(\mathbf{G}^{(k)})} \\ &= - \sum_{k=1}^K q(\mathbf{G}^{(k)}) \log \sum_{j=1}^K p(\mathbf{G}^{(j)}) = - \log \sum_{j=1}^K p(\mathbf{G}^{(j)}) . \end{aligned} \quad (27)$$

Hence, the constrained KL divergence is minimized, when  $q_{\mathbf{G}}$  assigns non-zero probabilities to the  $K$  graphs with the highest probabilities in  $p_{\mathbf{G}}$  as in Equation (3).  $\square$

### C.3 Minimal Hellinger & Bhattacharyya distance

The derivation of the minimal Hellinger distance  $D_{\text{H}}$  follows analog to the one in subsection C.1 with the same constraint of preserved relative probabilities from Equation (26):

$$\begin{aligned} D_{\text{H}}^2(q_{\mathbf{G}} \| p_{\mathbf{G}}) &= \frac{1}{2} \sum_{\mathbf{G}} \left( \sqrt{q(\mathbf{G})} - \sqrt{p(\mathbf{G})} \right)^2 \\ &= \frac{1}{2} \left( \sum_{\mathbf{G} \in \mathcal{G} \setminus \{\mathbf{G}^{(k)}\}_{k=1}^K} p(\mathbf{G}) + \sum_{k=1}^K \left( \sqrt{q(\mathbf{G}^{(k)})} - \sqrt{p(\mathbf{G}^{(k)})} \right)^2 \right) \\ &= \frac{1}{2} \left( 1 + \sum_{k=1}^K q(\mathbf{G}^{(k)}) - 2 \sqrt{q(\mathbf{G}^{(k)}) p(\mathbf{G}^{(k)})} \right) \\ &= \frac{1}{2} \left( 2 - \frac{2}{\sqrt{\sum_{j=1}^K p(\mathbf{G}^{(j)})}} \sum_{k=1}^K p(\mathbf{G}^{(k)}) \right) \\ &= 1 - \sqrt{\sum_{k=1}^K p(\mathbf{G}^{(k)})} . \end{aligned} \quad (28)$$

The Bhattacharyya distance  $D_{\text{B}}$  can be directly derived from the Hellinger distance  $D_{\text{H}}$  by:

$$\begin{aligned} D_{\text{B}}(q_{\mathbf{G}} \| p_{\mathbf{G}}) &= - \ln \left( \sum_{\mathbf{G} \in \mathcal{G}} \sqrt{q(\mathbf{G}) p(\mathbf{G})} \right) = - \ln (1 - D_{\text{H}}^2(q_{\mathbf{G}} \| p_{\mathbf{G}})) \\ &= - \frac{1}{2} \ln \sum_{k=1}^K p(\mathbf{G}^{(k)}) . \end{aligned} \quad (29)$$

Both are minimized by selecting the graphs with the highest probability in  $p_{\mathbf{G}}$  as in Equation (24).  $\square$

## D Importance sampling for parametrized distributions with an auxiliary discrete structure

Implicit generative models for discrete structures such as DAGs samples substructures from categorical distributions. The weights of discrete objects, e.g., edges, that do not appear in some target structure can be masked to 0. Since the relative probabilities of admissible substructures are preserved then, it is equal in probability to sampling from the unconstrained distribution and rejecting any sample that is not admissible under the target structure. The resulting proposal distribution  $q_\phi^*$  restricts the sample space to this target structure and is consequently optimal to the candidate distribution  $q_\phi$  at any optimization step due to the shared parameters. In comparison direct sampling has higher variance, in particular for structures with very low probability in the target distribution. While we present the proposed importance sampling to evaluate the probability of a full graph, its application to directed paths or subgraphs is straight-forward, following the same rational.

*Order-based model* For order-based models as described in subsection 2.3 and depicted in Figure 1c, e.g., RPM-DAG and ARCO-DAG, the weights of the Bernoulli distribution for any edges that do not appear in the target graph can be set to zero, i.e.,  $\forall i, j$  with  $\mathbf{G}_{ij}^* = 0 : \phi_{ij} = 0$ . In addition, the sampling of a total order in Equation (18) can be constrained to the (potentially partial) order induced by the target graph by setting the weights of all variables to 0 as long as their parents have not been sampled as a predecessor in the order.

*Sequence-based model* Sampling of the next edge in a sequence as described in subsection 2.4 and depicted in Figure 1e can be directly constrained to the edges that appear in the target graph. The ancestral masking that ensures acyclicity remains unaffected.

## E Experimental Parameters

In the experiments in subsection 4.1 and 4.2, we evaluate the forward KL divergence in all three, respectfully four, graphs with positive probability mass. The probability of a graph under the candidate model is approximated by Equation (10) using *importance samples* (IS). To speed up the calculations during training, the number of IS is chosen to be smaller than for the final evaluation of the trained model. The learning rate was chosen by a grid search over  $\{1, 5 \times 10^{-1}, 1 \times 10^{-1}, 5 \times 10^{-2}, 1 \times 10^{-2}, 5 \times 10^{-3}, 1 \times 10^{-3}, 5 \times 10^{-4}, 1 \times 10^{-4}\}$ . Training was performed over 1000 optimization steps to account for the instability of GFlowNet-DAG. We document the used hyperparameters in Table 5. Due to the low number of DAGs, our first experiment (Section 5.1) allows to evaluate all 25 graphs within a single forward pass during training (hence, 1 optimization step equals 1 epoch). To account for the higher number of possible sequences, we average the probability of sampling a target graph from the model

distribution using 10 importance samples for GFlowNet-DAG in contrast to a single one. For evaluation of each graph probability we drew 100 importance samples for all methods. In the second experiment (Section 5.2) we had to limit the number of samples from the target distribution for a single optimization step. Due to the high number of trainable parameters we decided to go with 25 instead of 100 samples from the target distribution and 5 instead of 10 for the GFlowNet-DAG model. By contrast, we doubled the number of importance samples to 20 during evaluation of each of the 543 DAGs.

*GFlowNet-Dag* The transformer architecture for the GFlowNet-Dag model follows the official implementation provided on <https://github.com/tristandeleu/jax-dag-gflownet>, as default parameters it uses an embedding size of  $H_E = 128$ , a key size of  $H_K = 64$  and  $H_L = 7$  layers of Transformer blocks.

*ARCO-DAG* The autoregressive neural network of ARCO-DAG consists of a simple two layer perceptron with  $H_N = 30$  hidden neurons and ReLU-activations. It follows the official implementation provided on <https://github.com/chritoth/bci-arco-gp/>.

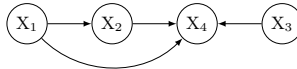
Table 5: Hyperparameters.

(a) For the experiments in subsection 4.1 and 4.2.

Model $q_G$	Learning rate	# forward KL samples	# IS for training	# IS for evaluation
RPM-DAG	$5 \times 10^{-1}$	"all"	1	100
ARCO-DAG	$5 \times 10^{-1}$	"all"	1	100
GFlowNet-DAG	$1 \times 10^{-3}$	"all"	10	100

(b) For the experiments in subsection 4.3.

Model $q_G$	Learning rate	# forward KL samples	# IS for training	# IS for evaluation
RPM-DAG	$5 \times 10^{-2}$	100	10	10
ARCO-DAG	$5 \times 10^{-2}$	100	10	10
GFlowNet-DAG	$1 \times 10^{-3}$	25	5	20

Fig. 5: True graph  $G_0$  as the MAP graph of the posterior in Table 3.