CAUSAL GRAPH LEARNING VIA DISTRIBUTIONAL INVARIANCE OF CAUSE-EFFECT RELATIONSHIP

Anonymous authors

Paper under double-blind review

Abstract

This paper introduces a new framework for recovering causal graphs from observational data, leveraging the fact that the distribution of an effect, conditioned on its causes, remains invariant to changes in the prior distribution of those causes. This insight enables a direct test for potential causal relationships by checking the variance of their corresponding effect-cause conditional distributions across multiple downsampled subsets of the data. These subsets are selected to reflect different prior cause distributions, while preserving the effect-cause conditional relationships. Using this invariance test and exploiting an (empirical) sparsity of most causal graphs, we develop an algorithm that efficiently uncovers causal relationships with quadratic complexity in the number of observational features/variables, reducing the processing time by up to 25x compared to state-of-the-art methods. Our empirical studies on a diverse benchmark of large-scale datasets demonstrate that the developed algorithm consistently performs better or comparable to existing works while generally achieving better scalability. Our code is publicly accessible at https://anonymous.4open.science/r/GLIDE-DC57.

1 INTRODUCTION

A core challenge in causal learning is finding a directed acyclic graph (DAG) that captures the causeeffect relationships between variables in an observational dataset (Zanga et al., 2022). A direct approach to uncover these causal links is through intervention which conducts more experiments to confirm whether changes in one set of variables will consequently affect the outcome distribution of another variable (Peters et al., 2016; Guo et al., 2024). However, such interventional experiments can be prohibitively expensive. Furthermore, recovering the entire causal graph requires running intervention on an exponential number of candidate cause-effect relationships among subsets of variables, thus rendering this approach both costly and impractical (Pearl, 2009).

To sidestep such expensive interventions, numerous approximation approaches have been developed to instead find an equivalence class of DAGs (Pearl et al., 2016; Peters et al., 2017), all of which are compatible with a set of statistical evidence or constraints derived from the observational data¹. Most of these approaches are formulated as graph searches that proceed either by finding statistical evidence to eliminate incompatible graph candidates (Spirtes & Glymour, 1991; Spirtes, 2001) or by optimizing for a heuristic score defined on graphs (Hauser & Bühlmann, 2012; Rolland et al., 2022; Montagna et al., 2023). Such approaches however are either (a) not scalable due to the expensive computation cost of running statistical tests on an exponentially large number of cause-effect relationship candidates (i.e., subsets of variables); or (b) less accurate due to the heuristic nature of the score function defined over the graph space.

For examples, Spirtes & Glymour (1991) and Spirtes (2001) required performing local conditional independence tests between effect variables and their corresponding sets of causes across all candidate causal relationships, incurring a process that scales exponentially with the number of features/variables. In contrast, Hauser & Bühlmann (2012), Rolland et al. (2022), and Montagna et al. (2023) proposed direct optimization of a global heuristic score on graphs, avoiding the need to solve such constraint satisfaction problems involving an exponentially large set of local constraints. There are also other approaches in this direction which further impose simplified assumptions on the

¹Finding the true causal graph is not possible without running intervention since the empirical distribution over the observational data might admit different sets of statistical constraints entailed by different graphs.

functional structure of the causal relationship (e.g., linear relationship between effects and causes perturbed with Gaussian noises). Base on those assumptions, the problem of causal learning can be formulated as a continuous optimization task and can be solved by more effective solution techniques (Kalainathan et al., 2022; Lachapelle et al., 2019; Ng et al., 2022; 2019). However, these approaches often perform less robustly when the heuristic scores or modeling assumptions do not fit well with the nature of the observational data (see Section 5).

To mitigate the aforementioned limitations of existing causal learning methods, we propose a new solution perspective based on a new causal test. Such a method that avoids imposing additional assumptions on the effect-cause data generation process while also achieving better performance with affordable computational costs. Our approach leverages the invariance of the effect-cause conditional distribution $P(\text{effect} \mid \text{cause})$ to changes in the prior cause distribution P(cause). This inspires a principled test for causal relationships via estimating the variance of $P(X \mid Z)$ across synthetic data augmentations that reflect different cause distributions P(Z). That is, $(Z \to X)$ represents a causal relationship if the data-induced $P(X \mid Z)$ does not change much² across different choices of P(Z). This test can then be integrated with a systematic search that identifies candidates for the causal parents of all variables with quadratic complexity, achieving improved scalability and performance over previous methods. In particular, our technical contributions that substantiate the above include:

1. An invariance test that reliably determines if Z = Pa[X] for each effect variable X and a candidate set of causes Z via checking the variance of $P(X \mid Z)$ against changes in P(Z). In particular, the invariance test identifies and constructs the most informative data augmentations to reliably approximate the variance of $P(X \mid Z)$ across potential changes to P(Z). This is achieved via a downsampling scheme of the observational data that modifies P(cause) without changing the conditional $P(\text{effect} \mid \text{cause})$ for true (cause, effect) tuples (Section 4.2).

2. A practical parent-finding algorithm that (i) adopts a previous approach (Edera et al., 2014) to find the Markov blankets of all variables using observational data; and (ii) uses this information to construct an augmented bidirectional graph for each variable whose maximal cliques correspond to its most plausible candidate parent sets. Due to the sparsity of such augmented graphs, the number of maximal cliques is quadratic in the number of variables and an effective depth-first search (DFS) algorithm can be devised to enumerate through all these cliques and hence, the corresponding parent candidates (Section 4.3). The developed invariance test (Section 4.2) can then be used to find the true parent set among the plausible candidates for each variable, thus recovering the true causal graph.

3. An extensive empirical evaluation of our proposed causal discovery framework on a variety of diverse benchmark datasets including both synthetic and large-scale real-world datasets. The reported results consistently show that our framework performs better or comparable to prior work in terms of causal discovery performance while achieving generally better scalability, with an average of up to 73.3% reduction in spurious rate and (up to) $25 \times$ reduction in processing time (Section 5).

2 RELATED WORKS

Existing causal discovery methods that aim to recover causal graphs from observational data without using interventional data can be categorized in three main groups (Glymour et al., 2019):

First, constraint-based methods aim to recover an equivalence class of causal graphs via deriving statistical evidence from the observational data to eliminate incompatible candidates as much as possible. Their main goal is to minimize the chance of mistaking correlation for causation and hence, maximizing the reliability of the output graph. For example, Peter-Clark (PC) (Spirtes & Glymour, 1991) and Fast Causal Inference (FCI) (Spirtes, 2001; Spirtes et al., 2013) run all possible conditional independence tests between all potential (effect, cause) tuples to find the most reliable relationship candidates that pass all tests. In practice, while such approaches often produce reliable results, they do not scale well to high-dimensional datasets since the number of (effect, cause) candidate tuples often grows exponentially in the number of variables/features (Spirtes et al., 2001).

²The variance of the true effect-cause conditional distribution P(effect | cause) with respect to changes in P(cause) is theoretically zero but in practice, P(effect | cause) has to be estimated using observational data which causes (small) additional variance due to (slight) variations in its estimation across augmented datasets.



Figure 1: Overall workflow of our proposed **GLIDE** framework which comprises two main steps: (a) an algorithmic configuration –the key effect-cause distributional invariance test that helps test potential parent-child relationships (Section 4.2); and (b) a graph search algorithm exploiting prior knowledge of each node's Markov blanket and an (empirically verified) sparsity of causal graphs to provably reduces the number of tests to recover the true causal graph (Section 4.3).

Alternatively, score-based methods instead use heuristic scores defined on graphs to reformulate causal learning as an optimization task which associate true causal graphs with those that maximize the score (Heckerman et al., 1995; Chickering, 2002; Teyssier & Koller, 2012; Solus et al., 2021). Thus, unlike constraint-based methods which cast causal learning as a constraint satisfaction task that involves an exponentially large set of local constraints, scored-based methods recast it as a global optimization task, which often admits more scalable solutions. As a result, existing score-based methods (Hauser & Bühlmann, 2012; Rolland et al., 2022; Montagna et al., 2023) often scale better to larger datasets. However, the heuristic design of the score function imposes implicit assumptions on the causal structure of data which are often violated in practice. Consequently, the reliability (i.e., not mistaking correlation for causation) of methods in this group is relatively lower than those of constraint-based methods. This is also verified in Section 5.1 (see Figures 2 and 3).

To reconcile the above conflicting goals of reliability and scalability, model-based methods adopt additional assumptions on the effect-cause generation model of the observed data (e.g., linear relationship perturbed with Gaussian noise (Zheng et al., 2018)). This often allows a provable reformulation of the true causal structure as an optimal solution to a continuous optimization task that can be solved efficiently with numerous modern and scalable machine learning algorithms (Scanagatta et al., 2015; Zheng et al., 2018; 2020). However, their performance is often not stable in scenarios where such assumptions on the data generation process do not hold, e.g., non-linear effect-cause relationships perturbed with non-Gaussian noises.

Existing Limitations. Overall, existing methods are either limited by expensive processing costs (as seen in constraint-based approaches), unreliable performance due to the use of a heuristic score function and greedy nature of the optimization algorithms (as with score-based methods), or inapplicable to scenarios involving unknown data-generation models (as is the case with model-based approaches). To mitigate these limitations, we investigate an alternative approach to causal learning by exploiting the invariance of the effect-cause conditional distribution across different data augmentations that induce changes to the prior distribution of causes. This, in turn, inspires a highly scalable graph search for effective causal discovery, as detailed in Section 4.

3 PROBLEM FORMULATION AND BACKGROUND

Let *D* denote a dataset comprising *n* observations $\mathbf{X}^{(1)}, \ldots, \mathbf{X}^{(n)}$ of a set $\mathbf{X} = (X_1, X_2, \ldots, X_d)$ of *d* random variables. Suppose these observations are drawn independently from an unknown distribution $P(\mathbf{X}) = P(X_1, X_2, \ldots, X_d)$, we want to learn from *D* a DAG $G = (\mathbf{X}, \mathbf{E})$ over (X_1, X_2, \ldots, X_d) which is a causal model of $P(\mathbf{X})$ following Definition 1.

Definition 1 (Causal Model). A direct acylic graph G = (X, E) is a causal model of P(X) if every conditional independence derived from P(X) can be derived from G, and vice versa.

A conditional independence $X \perp_P Y \mid Z$ derived from P(X) means $P(X \mid Y, Z) = P(X \mid Z)$ where $P(X \mid Y, Z)$ and $P(X \mid Z)$ are marginal likelihoods derived from P(X). On the other hand, a conditional independence $X \perp_G Y \mid Z$ derived from G means given Z, X and Y are d-separated following the below definition of d-separation.

Definition 2 (*D*-separation). Given two variables $X, Y \in \mathbf{X}$ and a set of variables $\mathbf{Z} \subseteq \mathbf{X} \setminus \{X, Y\}$, X and Y are d-separated given \mathbf{Z} (i.e., $X \perp _G Y \mid \mathbf{Z}$) if any path between them contains either: (a) a fork $A \leftarrow B \rightarrow C$ with $B \in \mathbf{Z}$, (b) a chain $A \rightarrow B \rightarrow C$ such that $B \in \mathbf{Z}$; and (c) a collider $A \rightarrow B \leftarrow C$ such that B or any of its descendants Desc[B] is not in \mathbf{Z} .

Definition 1 implies G is a causal model of $P(\mathbf{X})$ when $X \perp _P Y \mid \mathbf{Z} \Leftrightarrow X \perp _G Y \mid \mathbf{Z}$. Thus, suppose G is a causal model of $P(\mathbf{X})$, $X \perp _P (\mathbf{X} \setminus \text{Desc}[X]) \mid \text{Pa}[X]$ since $X \perp _G (\mathbf{X} \setminus \text{Desc}[X]) \mid \text{Pa}[X]$ due to d-separation – see Definition 2 – where Pa[X] denote the set of parents of X in G.

Local Markov Conditions. The above means each node X is independent of its non-descendants $X \setminus \text{Desc}[X]$ given its parents Pa[X] according to P(X), resulting in the following factorization:

(I)
$$P(\mathbf{X}) = \prod_{i=1}^{n} P(X_i | \operatorname{Pa}[X_i])$$
 which implies (II) $X \perp P(\mathbf{X} \setminus \mathbb{M}(X)) | \mathbb{M}(X)$,

where $\mathbb{M}(X)$ denotes the Markov blanket of a variable X that consists of its immediate parents, its children, and its children's other parents. (II) can also be verified via checking *d*-separation on G.

Core Idea. Suppose the Markov blanket in condition (II) is known, an augmented bidirectional graph for each variable X can be constructed such that each of its maximal cliques corresponds to a plausible candidate Z of the true causal parents Pa[X]. We can leverage this representation to develop an effective DFS procedure that systematically enumerates through each candidate with a quadratic time complexity in the number of features/variables (Section 4.3). Each candidate can be tested using the developed causal test via synthetic data augmentation (Section 4.2). The feasibility of this approach is enabled by adopting an existing algorithm that can recover the Markov blankets for all $X \in G$ from D with $O(d^2)$ complexity (Edera et al., 2014). This is guaranteed under the causal sufficient assumption (Pearl et al., 2016) that there is no unobserved variables (i.e., confounders) that affect the causal mechanism that generate the observational data.

4 INVARIANCE OF EFFECT-CAUSE CONDITIONAL DISTRIBUTION

Here, we first present a bird's eye view on our proposal in Section 4.1, and then navigate into the details of the method in Sections 4.2 and 4.3, as depicted in Figure 1.

4.1 AN OVERVIEW

The core idea of our causal discovery framework is based on the invariance of the (marginal) effectcause conditional distribution $P(X | \operatorname{Pa}[X])$ induced from P(X) across different choices of the prior P(B) over the source variables $B \triangleq \{B | B \in X, \operatorname{Pa}[B] = \emptyset\} \subseteq X$ of the true causal graph G = (X, E) – see Theorem 1 below (a detailed proof is provided in Appendix A.1).

Theorem 1 (Effect-Cause Distributional Invariance). Let $P_1(\mathbf{B}), P_2(\mathbf{B}), \ldots, P_m(\mathbf{B})$ denote a set of *m* different priors over **B**. Let $P_i(\mathbf{X}) = P_i(\mathbf{B}) \cdot P(\mathbf{X} \setminus \mathbf{B} \mid \mathbf{B})$ denote the corresponding augmentation of the true data distribution $P(\mathbf{X})$ when we replace its marginal prior $P(\mathbf{B})$ with $P_i(\mathbf{B})$. For each variable $X \in \mathbf{X} \setminus \mathbf{B}$, its true causal parents Pa[X] and a candidate \mathbf{Z} , we have

$$\mathbb{V}_{P_{+}(\boldsymbol{X}) \sim \mathcal{P}}\left[P_{+}\left(\boldsymbol{X} \mid \boldsymbol{Z}\right)\right] > 0 \quad \Rightarrow \quad \boldsymbol{Z} \neq \operatorname{Pa}[\boldsymbol{X}], \tag{1}$$

regardless of how $P_+(\mathbf{X})$ is drawn from $\mathcal{P} \triangleq (P_1(\mathbf{X}), P_2(\mathbf{X}), \dots, P_m(\mathbf{X}))$. Here, $P_+(X \mid \mathbf{Z})$ is induced from $P_+(\mathbf{X})$.

This result reveals a principled test of whether a subset $Z \subseteq X$ is the parent set of $X \in X$ according to the true causal graph G. The intuition is if we can re-sample D_i from $D \sim P(X)$ such that $D_i \sim P_i(X)$, we can test whether Z = Pa[X] via checking whether the sample variance of the empirical conditional $P_i(X \mid Z)$ is distinguishably small. This is formalized below:

A. Effect-Cause Invariance Test. Given m augmented datasets D_1, D_2, \ldots, D_m which are resampled from $D \sim P(\mathbf{X})$ such that $D_i \sim P_i(\mathbf{X})$, then $\mathbf{Z} = Pa[X]$ when

$$\frac{1}{m}\sum_{i=1}^{m} \left\| P_i \left(X \mid \mathbf{Z} \right) - \overline{P} \left(X \mid \mathbf{Z} \right) \right\|^2 \simeq 0 \quad \text{where} \quad \overline{P} \left(X \mid \mathbf{Z} \right) = \frac{1}{m}\sum_{i=1}^{m} P_i \left(X \mid \mathbf{Z} \right). \tag{2}$$

This test will become more accurate with more source priors. When m is infinitely large, the implication in Eq. (1) becomes bi-directional and $\mathbb{V}[P_+(X \mid Z)] = 0$ definitively implies $Z = \operatorname{Pa}[X]$. Otherwise, when m is sufficiently large and $\{P_i(B)\}_{i=1}^m$ are sufficiently representative of the entire space of source priors, the test might not be perfect but remains highly accurate (see Section 5).

For more convenience, we note that the above test can also be rephrased as a selection criteria to find the true causal parents Pa[X] for each variable X as detailed next.

B. Parent-Finding via Effect-Cause Invariance. Given *m* augmented datasets D_1, D_2, \ldots, D_m which are resampled from $D \sim P(\mathbf{X})$ such that $D_i \sim P_i(\mathbf{X})$. Let $\mathbf{Z}_1, \mathbf{Z}_2, \ldots, \mathbf{Z}_p$ denote the set of plausible parent candidates containing the true parent set, then Pa[X] can be determined via

$$\operatorname{Pa}[X] = \mathbf{Z}_{r} \quad \text{where} \quad r \triangleq \operatorname{arg\,min}_{t} \left(\frac{1}{m} \sum_{i=1}^{m} \left\| P_{i} \left(X \mid \mathbf{Z}_{t} \right) - \overline{P} \left(X \mid \mathbf{Z}_{t} \right) \right\|^{2} \right), \tag{3}$$

when the number m of augmented datasets is sufficiently large and diversified (see Section 4.2).

C. High-Level Framework. Suppose we know how to generate D_1, D_2, \ldots, D_m (see Section 4.2) for which $D_i \sim P_i(\mathbf{X})$ as required in Theorem 1, optimizing Eq. (3) can be achieved via exhaustively checking all subsets \mathbf{Z} as candidates for Pa[X]. Repeating this for all X allows us to recover the causal graph. This process is however impractical since its complexity is exponential in d.

Fortunately, this can be avoided using the local Markov condition (II) (see Section 3) which implies that (i) $Z \subseteq \mathbb{M}(X)$ if $Z = \operatorname{Pa}[X]$ and (ii) there are at most p = O(d) candidates for $\operatorname{Pa}[X]$ which can be provably identified via finding maximum cliques on augmented bidirectional graphs. Under an empirically verified assumption on causal graph sparsity, this can be achieved with a customized DFS procedure with $O(d^2)$ complexity (see Theorem 7, Section 4.3). This results in an effective $O(d^2)$ total complexity for our causal discovery framework as detailed below.

D. Time Complexity. As there are d observational variables, our framework needs to make d calls to the parent-finding routine in part (**B**). As each routine will consider at most p = O(d) plausible parent candidates, the per-call complexity is O(md). This amounts to a total cost of $O(md^2)$ to recover the full causal graph. This can be enabled with (a) an $O(d^2)$ overhead to find the Markov blankets for all variables following the prior work of Edera et al. (2014); and (b) another $O(d^2 + m|D| + m|B|)$ overhead to construct D_1, \ldots, D_m for Eq. (3) (Section 4.2.4). The overall complexity (including overhead) is therefore $O(md^2 + m|D| + m|B|)$.

To substantiate the above high-level framework, we need to (i) choose representative variants of the source priors to improve the test reliability (Section 4.2); and (ii) find all plausible sets of candidate for Pa[X] for all $X \in \mathbf{X}$ which are guaranteed to have sizes at most O(d) (Section 4.3).

4.2 AUGMENTING SOURCE PRIOR

To enable practical use of the parent-finding routine in Eq. (3), we need to (i) find the set B of sources, (ii) choose a representative set of source priors $P_1(B), \ldots, P_m(B)$ so that the invariance test is reliable; and (iii) re-sample D_i from D so that $D_i \sim P_i(X) = P_i(B) \cdot P(X \setminus B \mid B)$.

The above (iii) is essential because we do not have direct access to $P(\mathbf{X} \setminus \mathbf{B} \mid \mathbf{B})$. This means we cannot compute $P_i(\mathbf{X})$ and induce $P_i(X \mid \mathbf{Z})$ to assess its variance directly. However, as $P(\mathbf{X} \setminus \mathbf{B} \mid \mathbf{B})$ can be simulated using the original data D, we can re-sample D_i from D so that $D_i \sim P_i(\mathbf{X})$, which then allows us to use D_i to estimate $P_i(X \mid \mathbf{Z})$ and consequently, its variance.

To achieve the above, we will develop an algorithm to find the source variables B (Section 4.2.1). We will then derive a re-sampling procedure to obtain a downsampled dataset D_i from D such that $D_i \sim P_i(X)$ (Section 4.2.2). Last, we will detail a criterion to choose informative $P_i(B)$ (for the invariance test) based on the above re-sampling procedure and how to determine a set of representative source priors (Section 4.2.3).

4.2.1 FINDING SOURCE VARIABLES

To find the set of source variables (i.e., those with no parent in the true causal graph), we introduce below the concept of a basis of a DAG.

Definition 3. A basis $B \subseteq X$ of a DAG G = (X, E) is a set of mutually *d*-separated (independent) variables (see Definition 2) such that for each $X \notin B$, there exists $X' \in B$ and $X' \not \perp X$.

As finding the true sources is not possible without intervention, we will use a basis set as a surrogate. This is because the basis has similar properties to the set of sources. First, similar to source variables, basis variables are mutually independent. Second, each source variable X is either in the basis B or shares the same dependence set $\Phi(X) \equiv \Phi(X')$ – as defined in Theorem 3 – with another basis variable $X' \in B$ (see Appendix B.1). Furthermore, Theorem 2 confirms that the maximum size of a basis set is equal to the number of sources in the true causal graph. This means changing the prior over basis variables will have similar effect to changing prior over source variables. Hence, we can use the priors over basis variables instead of priors over sources to test the effect-cause invariance.

Theorem 2. The maximum size of a basis set of a DAG is equal to the number of its sources.

The proof of Theorem 2 is deferred to Appendix A.2. Another advantage of the maximum basis set is that it can be provably identified with $O(d^2)$ complexity following the procedure detailed in Theorem 3 below. See its detailed proof in Appendix A.3.

Theorem 3. Let V = X and $\Phi(X) \triangleq \{Y \mid Y \in V \setminus \{X\} : Y \not\perp X\}$. The maximum basis can be constructed via continually (i) selecting X in V with lowest $|\Phi(X)|$; (ii) setting $V \leftarrow V \setminus (\Phi(X) \cup \{X\})$; and (iii) stopping when $V = \emptyset$. The selected X forms the maximum basis.

Computing $\Phi(X)$ requires checking $Y \not\perp X$ which can be practically achieved using existing reliable pairwise independence tests provided by the causal-learn open library (Zheng et al., 2024).

4.2.2 RE-SAMPLING OBSERVATIONAL DATA

Given a target $P_i(B)$, we want to find a resampled dataset D_i from the original observation data D such that $D_i \sim P_i(X)$. To guarantee this, D_i must be a downsampled version of D via sampling with no replacement to avoid introducing duplicates and hence, false causal biases into D_i .

On the other hand, among valid downsamples $D_i \sim P_i(\mathbf{X})$, we want to choose the one that has the minimal downsampled rate $|D|/|D_i|$ to preserve (as much as possible) observations of the effectcause conditionals in $P(\mathbf{X} \setminus \mathbf{B} \mid \mathbf{B})$. Interestingly, it can be shown that given the target $P_i(\mathbf{B})$, the minimum downsampling rate can be computed and achieved via the following results.

Theorem 4. Suppose D_i is the downsampled dataset with minimum downsampling rate that satisfies the condition $D_i \sim P_i(\mathbf{X})$. Then, it follows that

$$|D|/|D_i| = \gamma_i^{-1}, \quad \text{where} \quad \gamma_i = \min_{\boldsymbol{b}} \left(P(\boldsymbol{B} = \boldsymbol{b}) / P_i(\boldsymbol{B} = \boldsymbol{b}) \right). \tag{4}$$

Given the (computable) optimal downsampling rate in Theorem 4, the corresponding downsampling procedure that achieves it can be derived via Theorem 5 below.

Theorem 5. Let γ_i be defined in Theorem 4. Suppose D_i is created via sampling without replacement $P_i(\mathbf{B} = \mathbf{b}) \cdot |D| \cdot \gamma_i$ points from D where $\mathbf{B} = \mathbf{b}$. Then, $|D|/|D|_i = 1/\gamma_i$ and $D_i \sim P_i(\mathbf{X})$.

The derivations of Theorems 4 and 5 are deferred to Appendices A.4 and A.5. For a practical implementation, we associate P(B = b) with its empirical estimate P(B = b) = |D[B = b]|/|D|. To accommodate for continuous data, we use binning with fixed bin width to make the data categorical.

4.2.3 CHOOSING SOURCE PRIORS

We can now leverage the insight of Theorem 4 to choose the most informative source priors to enhance the reliability of the invariance test in Eq. (3). Intuitively, we want $D_i \neq D$ so γ_i should not be too large. Otherwise, as $\gamma_i \rightarrow 1$, $D_i \rightarrow D$ and $P_i(\mathbf{B}) \rightarrow P(\mathbf{B})$ which cannot be used to test the effect-cause invariance against changes in $P(\mathbf{B})$. On the other hand, if γ_i is too small, D_i might drop too much information from D which might obscure some effect-cause relationship.

As such, we want to choose $P_i(B)$ such that its inverse downsampling rate γ_i (see Theorem 4) is above a certain threshold γ_o where $\gamma_o \in (0, 1)$ is an adjustable parameter that we can experiment with. Our ablation studies in Section 5 shows the impact of γ_o on the overall causal discovery performance. The question is now how to sample representative $P_i(B)$ from the subspace of source priors whose (inverse) optimal downsampling rate $\gamma_i \geq \gamma_o$. To answer this question, we establish Theorem 6 below which characterizes its convex hull (see a detailed derivation in Appendix A.6). **Theorem 6.** Let $r \triangleq |\text{Dom}(B)|$ denote the number of (categorical) candidate values of B. The subspace of $P_i(B)$ that satisfies $\gamma_i \ge \gamma_o$ is a convex subspace $C_r(\gamma_o)$ of the r-dimensional simplex Δ_r over Dom(B) which cuts Δ_r at r points $P^{(1)}(B), \ldots, P^{(r)}(B)$ representing its convex hull:

$$P^{(k)}(\boldsymbol{B}) = \alpha_k \cdot P(\boldsymbol{B}) + (1 - \alpha_k) \cdot \delta_k(\boldsymbol{B}), \text{ where}$$

$$\alpha_k = \left(1 - P\left(\boldsymbol{B} = \boldsymbol{b}^{(k)}\right) \gamma_o^{-1}\right) / \left(1 - P\left(\boldsymbol{B} = \boldsymbol{b}^{(k)}\right)\right), \tag{5}$$

 $\delta_k(\mathbf{B})$ is a point mass function that assigns 1 when $\mathbf{B} = \mathbf{b}^{(k)}$ and 0 otherwise. Here, $\mathbf{b}^{(k)}$ is the *k*-th candidate value in Dom(\mathbf{B}).

Each source prior $P_i(B)$ with $\gamma_i \ge \gamma_o$ belongs to this convex set and can be represented as a linear combination of the above points:

$$P_i(\mathbf{B}) = \sum_{k=1}^r a_k \cdot P^{(k)}(\mathbf{B}), \text{ where } \sum_{k=1}^r a_k = 1 \text{ and } a_k \ge 0.$$
 (6)

Hence, $P_i(B)$ can be sampled via drawing $a = (a_1, a_2, \dots, a_r)$ from Δ_r and using Eq. (6).

4.2.4 PRACTICAL IMPLEMENTATION

The parent-finding procedure involves three main steps: finding basis variables, sampling source priors, and re-sampling observational data. First, finding the basis variables has a time complexity of $O(d^2)$. Next, we sample a from a Dirichlet distribution to compute $P_i(B)$, with 10^4 samples clustered using K-means. The K = m centroids are selected as source priors, with a complexity of $O(mr) = O(mc^{|B|})$ where c is the maximum number of candidate values of a single variable. To avoid exponential costs in |B|, we resample for each basis variable individually, reducing the complexity to O(mc|B|). Finally, generating the augmented dataset for each $P_i(B)$ incurs a cost of O(m|D|) via Theorem 5. This amounts to $O(d^2 + mc|B| + m|D|)$ total complexity.

4.3 FINDING PLAUSIBLE PARENT SETS

To ensure the parent-finding routine is scalable, we customize a DFS algorithm which leverages prior work in Markov blanket identification to provably find O(d) candidate parent sets for each effect variable X. The true causal graph can then be recovered via solving Eq. (3) for each X with respect to the discrete set of O(d) plausible parents found above.

As previous work in Markov blanket identification incurs an $O(d^2) \operatorname{cost} (\operatorname{Edera} \operatorname{et} \operatorname{al.}, 2014)$, the total cost of our parent finding phase is also $O(d^2)$. This helps avoid the brute-force search through all subsets of variables as candidates for $\operatorname{Pa}[X]$ whose complexity is otherwise exponential in d (Peters et al., 2016). To achieve this, we will leverage the following result which relates candidate sets of causal parents to maximal cliques on an augmented bidirectional graph.

Theorem 7 (Plausible Parent Sets). For each variable X, let G'(X) = (V, E') denote a bidirectional graph where $V = \mathbb{M}(X)$ and $(U, V) \in E$ iff $V \in \mathbb{M}(U)$ and $U \in \mathbb{M}(V)$. Then, $\operatorname{Pa}[X] \subseteq \mathbb{M}(X)$ corresponds to a clique in G'(X).

See Appendix A.7 for a detailed proof. Leveraging Theorem 7, we can find all plausible parent sets of X as subsets in $\mathbb{M}(X)$, each corresponding to a clique in G'(X). Restricting the set of plausible parent sets to that of maximal cliques in G', we can reduce the task of finding plausible parent sets to finding maximal cliques in a bidirectional graph. For a sparse graph G'(X) with a low degeneracy constant p (see Definition 4), this can be achieved effectively via a customized version of the DFS-based Bron-Kerbosch (Bron & Kerbosch, 1973) algorithm (see Appendix B.2) which finds all maximal cliques with $O(dp3^{p/3})$ time complexity.

Definition 4 (Degeneracy). A bidirectional graph is *p*-degenerate if every subgraph has at least one node with degree $\leq p$. The graph's degeneracy is the smallest *p* for which it is *p*-degenerate.

Consider p as a small constant compared to d, the cost of finding all plausible parent sets for each variable X is effectively O(d). Hence, the total cost of finding all plausible parent sets for all variables is $O(d^2)$. Furthermore, it is also established in Bron & Kerbosch (1973) that the (worst-case) number of maximum cliques is $O((d - p) \cdot 3^{p/3})$. Our empirical studies in fact show that

Data Models	PC	GIES	FCI	Notears	MLP-Notears	DAS	SCORE	GLIDE(Ours)
Linear, Gaussian (L-G)		\checkmark	\checkmark	\checkmark	 	\checkmark	\checkmark	<hr/>
Non-linear, non-Gaussian (nL-nG)		\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	 ✓
Categorical, Synthetic	🗸	\checkmark						 ✓
Categorical, Real-world	🗸	\checkmark						 ✓

Table 1: **Baselines and Data Models.** The table indicates the applicability of each baseline to each data models. Applicability is determined based on the baseline's empirical effectiveness on recovering the causal graph from observational data under the corresponding data model.

 $p \le 13$ across all benchmark datasets even for the largest graphs, confirming that the set of maximal cliques in G'(X) is practically O(d) with a constant factor smaller than $3^{13/3} < 117$. Our additional ablation studies in Appendix C.2.3 also show that for a variety of graph topologies, random DAGs with an increasing number of causal edges, the sparsity constant p of the corresponding augmented (bidirectional) graph remains insignificant compared to d. Exploiting this, our proposed method is able to scale efficiently on these large graphs where previous methods fail to do so.

5 EXPERIMENT

This section evaluates and compares the performance of our proposed method, Causal <u>G</u>raph <u>L</u>earning via <u>D</u>istributional Invariance of Cause-<u>E</u>ffect Relationship (**GLIDE**), against existing state-of-the-art baselines in causal graph learning.

Baselines. Our empirical evaluations are conducted with respect to a variety of data models as summarized in Table 1³ and a diverse suite of both classical and recent baselines, including PC (Spirtes & Glymour, 1991), GIES (Hauser & Bühlmann, 2012), FCI (Spirtes, 2001), NOTEARS (Zheng et al., 2018), MLP-NOTEARS (Zheng et al., 2020), DAS (Montagna et al., 2023), and SCORE (Rolland et al., 2022). Each baseline is configured with its best hyper-parameters (see Appendix C.1).

Datasets. Our experiments are based on both synthetic and real-world datasets. For synthetic experiments (Sections 5.1 and 5.2), we follow the commonly used protocols in previous work to generate observationals data based on the Erdos-Renyi (Zheng et al., 2018; Heinze-Deml et al., 2018), bipartite and scale-free (Zheng et al., 2020) classes of causal graphs. Due to limited space, we only present the empirical results on synthetic data generated from the Erdos-Renyi class of graphs. Our other experiments with the bipartite and scale-free classes of graphs are deferred to Appendix C.2. For real-world experiments (Section 5.3), we use the datasets provided in the bnlearn package (Scutari, 2009), which includes Sachs, Insurance, Water, Alarm, Barley, Pathfinder, and Munin. Notably, the Munin dataset features observational data from a large-scale graph comprising 1041 variables.

Evaluation Metrics. We use the structural Hamming distance (SHD), spurious rate (percentage), and running time (minutes) to measure both the effectiveness and scalability of our proposed method. The SHD and running time metrics are commonly used in the evaluation of existing causal discovery methods (Zheng et al., 2018; Montagna et al., 2023). The spurious rate measures the ratio of the number of false causal relationships among all causal relationships discovered by each algorithm. This helps compare the reliability of causal prediction across baselines. All performance metrics are reported with their mean and confidence interval 95% averaged over 10 independent runs.

5.1 SYNTHETIC CONTINUOUS DATA

We consider two experiment settings. First, in the normal setting, the graph's number of edges is the same as its number of nodes, which varies from 100 to 500. Second, in the extreme setting, we fix the number of nodes at 500 and increase the number of edges from 600 to 1000. In each setting, we further consider two data-generation models for the effect-cause relationship: (1) linear relationship perturbed with Gaussian noise (annotated as **L-G**); and (2) non-linear and non-Gaussian relationship (annotated as **nL-nG**). Such data generation mechanism is implemented using the public code in (Zheng et al., 2018; 2020). In each experiment, we first generate a random Erdos-Renyi DAGs and then simulate (synthetic) observational data from it using the above mechanisms. Our reported results and observations in each setting are detailed below.

³Note that except for our versatile method, most other baselines are not applicable to all data models (i.e., producing very poor performance that is not meaningful for comparison).



Figure 2: Baseline performance in normal setting (continuous data). Lower metrics are better.

Figure 3: Baseline performance in extreme setting (continuous data). Lower metrics are better.

Normal Setting. Figure 2 reports the performance of our method GLIDE in comparison to those of other baselines. In the L-G cases, GLIDE, NOTEARS, and MLP-NOTEARS outperform the rest of the baselines significantly in terms of both SHD and spurious rate. Notably, **GLIDE** achieves a remarkable 64.17% reduction rate on SHD over FCI. Furthermore, GLIDE incurs much less computational cost than both NOTEARS and MLP-NOTEARS (see the runtime plots) while achieving comparable or better performance. For example, **GLIDE** runs $96.54 \times$ and $15.66 \times$ faster than both NOTEARS and MLP-NOTEARS in in 100- and 500-node graphs while being second to MLP-NOTEARS in terms of SHD (with a small gap of 16.2%) and achieving best spurious rate in the largest graph setting with 500 nodes (4.2% vs the second best of 11.75% of NOTEARS and third best of 13.19% of MLP-NOTEARS). In the **nL-nG** cases, we also have a similar observations where GLIDE is again the second best in SHD and best in spurious rate (in the largest graph settings) while being much more scalable than the best and second best baselines in SHD and spurious rate, respectively. Our reported results also confirm our intuition earlier that score-based methods, despite being more scalable than constraint-based approaches, tend to perform much less robustly in large graph settings where the score function becomes less accurate. For example, both SCORE and DAS reach over 30% of false causal detection rate (i.e., spurious rate) in 500-node graphs.

Extreme Setting. Figure 3 reports the performance of **GLIDE** in comparison to other baselines in this setting. In the **L-G** cases, it is observed that **GLIDE** achieves the best performance in both SHD and spurious rate while also achieving the fastest running time in most graph settings. In terms of SHD, the performance gaps between **GLIDE** and the second best (NOTEARS) and worst baselines (SCORE) are 11.74% and 108.5%, respectively. In terms of spurious rate, **GLIDE** also improves over the baselines with substantial performance gaps ranging from 5.46% and 48.61% against the second best and worst baselines, respectively. In the **nL-nG** cases, **GLIDE** and **MLP-NOTEARS** perform comparably as best baselines in SHD but **GLIDE** outperforms it by a gap of 8.23% in term of spurious rate. Furthermore, against the most high-performing baselines in this case, NOTEARS and MLP-NOTEARS, **GLIDE** achieve $9.6 \times$ and $25.52 \times$ faster processing time, respectively.

5.2 SYNTHETIC CATEGORICAL DATA

Categorical data are generated in the same manner as continuous data in Section 5.1. The only difference is that instead of using Gaussian models, we first randomize conditional probabilities

Metrics	Method	Sachs (11)	Insurance (27)	Water (32)	Alarm (37)	Barley (48)	Pathfinder (186)	Munin (1041)
SHD	GIES PC GLIDE	$ \begin{vmatrix} 14.0 \pm 0.0 \\ 10.3 \pm 0.7 \\ \textbf{5.2} \pm \textbf{0.4} \end{vmatrix} $	$\begin{array}{c} 36.0\pm0.0\\ 33.7\pm0.7\\ \textbf{18.0}\pm\textbf{2.8} \end{array}$	$\begin{array}{c} 49.0\pm 0.0\\ 54.3\pm 0.7\\ \textbf{41.6}\pm \textbf{1.8} \end{array}$	$\begin{array}{c} 44.0\pm 0.0\\ 35.7\pm 2.4\\ \textbf{27.8}\pm \textbf{2.0} \end{array}$	$\begin{array}{c} 65.0 \pm 0.0 \\ 58.3 \pm 1.7 \\ \textbf{45.8} \pm \textbf{2.4} \end{array}$	$ \begin{vmatrix} 1156.0 \pm 0.0 \\ \text{N/A} \\ \textbf{59.1} \pm \textbf{1.9} \end{vmatrix} $	$ \begin{vmatrix} 1235.0 \pm 0.0 \\ \text{N/A} \\ \textbf{883.2} \pm \textbf{21.8} \end{vmatrix} $
Spurious rate (%)	GIES PC GLIDE	$ \begin{vmatrix} 34.8 \pm 0.0 \\ 4.2 \pm 4.1 \\ \textbf{0.0} \pm \textbf{0.0} \end{vmatrix} $	$\begin{array}{c} 31.9 \pm 0.0 \\ \textbf{0.0} \pm \textbf{0.0} \\ 3.6 \pm 2.7 \end{array}$	$\begin{array}{c} 35.8 \pm 0.0 \\ \textbf{17.1} \pm \textbf{2.1} \\ 23.0 \pm 0.9 \end{array}$	$\begin{array}{c} 44.3 \pm 0.0 \\ 27.7 \pm 3.6 \\ \textbf{13.1} \pm \textbf{0.9} \end{array}$	$\begin{array}{c} 35.8 \pm 0.0 \\ 17.8 \pm 1.9 \\ \textbf{8.2} \pm \textbf{1.1} \end{array}$	$\begin{array}{c} 87.0\pm0.0\\ \text{N/A}\\ \textbf{1.9}\pm\textbf{0.5} \end{array}$	$\begin{vmatrix} 42.4 \pm 0.0 \\ \text{N/A} \\ \textbf{1.8} \pm \textbf{0.2} \end{vmatrix}$
Runtime (seconds)	GIES PC GLIDE	$\begin{array}{c c} \textbf{1.5} \pm \textbf{0.4} \\ 52.0 \pm 1.8 \\ 23.0 \pm 1.1 \end{array}$	$\begin{array}{c} \textbf{2.7} \pm \textbf{0.6} \\ 450.3 \pm 19.4 \\ 34.2 \pm 0.9 \end{array}$	$\begin{array}{c} \textbf{3.1} \pm \textbf{0.6} \\ 573.8 \pm 16.1 \\ 49.1 \pm 3.3 \end{array}$	$\begin{array}{c} \textbf{3.3} \pm \textbf{0.4} \\ 503.2 \pm 21.2 \\ 43.3 \pm 0.4 \end{array}$	$\begin{array}{c} \textbf{3.4} \pm \textbf{0.8} \\ \textbf{761.1} \pm \textbf{9.5} \\ \textbf{61.7} \pm \textbf{2.0} \end{array}$	$\begin{array}{ } \textbf{19.3} \pm \textbf{0.4} \\ \text{N/A} \\ 197.0 \pm 21.2 \end{array}$	$\begin{vmatrix} \textbf{61.5} \pm \textbf{15.2} \\ \text{N/A} \\ \textbf{6200.3} \pm \textbf{88.7} \end{vmatrix}$
$ \begin{array}{c} \hline \hline$								

Table 2: Performance on real-world datasets. N/A denotes no results within 48 hours.



SHD (×1

Spurious

for each node and then use Gibbs sampling to simulate the data. The number of categories are randomly chosen from 2 to 5 for each variable. We compare our proposal with two competitive baselines, GIES and PC, while excluding the rest since they fail to produce meaningful results for comparison (i.e., very poor performance). To focus on settings where the PC baseline can produce results within the time limit of 48 hours, we restrict our evaluation to two classes of graphs where the number of edges is (i) $3\times$, and (ii) $4\times$ the number of nodes, which ranges between 100 and 500.

The results are reported in Figure 4 which shows that **GLIDE** consistently outperforms the baselines in terms of SHD with substantial gaps of 10.49% and 27.95% in the ($\mathbf{E} = 3\mathbf{V}$) case; and 11.56% and 13.3% in the ($\mathbf{E} = 4\mathbf{V}$) case. In both cases, **GLIDE** is the second best in terms of spurious rate, increasing the false causal detection (spurious) rate of the best baseline (PC) by a mere margin of 4%. In exchange, **GLIDE** achieves a 30× faster processing time than PC. In contrast, GIES achieves the fastest running time but incurs 40% spurious rate (i.e., low reliability). Overall, **GLIDE** has the best trade-off between scalability (processing time) and performance (SHD, spurious rate).

5.3 REAL-WORLD CATEGORICAL DATA

20

0 Durio

Ē

Table 2 reported the performance of PC, GIES, and **GLIDE** on 7 real-world datasets. Note that GIES is a deterministic method and has zero deviation across different seeds on the same test case. It is observed that **GLIDE** achieves the best SHD performance across 7/7 datasets and also achieves the best spurious rate in 5/7 datasets, especially on large datasets such as Barley, Pathfinder, and Munin. On the largest dataset (Munin), **GLIDE** achieves a spurious rate of 1.8% which is remarkably better than GIES' (42.36%). Again, **GLIDE** has the best balance between performance and scalability.

6 CONCLUSION

This paper presents a new perspective of causal learning via a new invariance test for causality that inspires a reliable and scalable algorithm for recovering causal graphs from observational data. Our approach explores a key insight that the effect-cause conditional distribution remain invariant under changes in the prior cause distribution, leading to a parent-finding procedure for each variable via synthetic data-augmentation. This procedure is further coupled with an effective search algorithm that exploits prior knowledge of each effect variable's Markov blanket and an (empirically verified) sparsity of the causal graphs to significantly reduces the overall complexity. The reduction in complexity and the marked improvements in both speed and accuracy, as demonstrated on large-scale benchmark datasets, highlight the potential for this approach to outperform existing methods. These findings suggest that our framework offers a promising direction for further research and practical applications in causal inference, especially in scenarios involving large and complex datasets.

REFERENCES

- Coen Bron and Joep Kerbosch. Algorithm 457: finding all cliques of an undirected graph. *Communications of the ACM*, 16(9):575–577, 1973.
- David Maxwell Chickering. Optimal structure identification with greedy search. *Journal of machine learning research*, 3(Nov):507–554, 2002.
- Alejandro Edera, Yanela Strappa, and Facundo Bromberg. The grow-shrink strategy for learning markov network structures constrained by context-specific independences. In Advances in Artificial Intelligence–IBERAMIA 2014: 14th Ibero-American Conference on AI, Santiago de Chile, Chile, November 24-27, 2014, Proceedings 14, pp. 283–294. Springer, 2014.
- Clark Glymour, Kun Zhang, and Peter Spirtes. Review of causal discovery methods based on graphical models. *Frontiers in genetics*, 10:524, 2019.
- Siyuan Guo, Viktor Tóth, Bernhard Schölkopf, and Ferenc Huszár. Causal de finetti: On the identification of invariant causal structure in exchangeable data. Advances in Neural Information Processing Systems, 36, 2024.
- Alain Hauser and Peter Bühlmann. Characterization and greedy learning of interventional markov equivalence classes of directed acyclic graphs. *The Journal of Machine Learning Research*, 13 (1):2409–2464, 2012.
- David Heckerman, Dan Geiger, and David M Chickering. Learning bayesian networks: The combination of knowledge and statistical data. *Machine learning*, 20:197–243, 1995.
- Christina Heinze-Deml, Jonas Peters, and Nicolai Meinshausen. Invariant causal prediction for nonlinear models. *Journal of Causal Inference*, 6(2):20170016, 2018.
- Diviyan Kalainathan, Olivier Goudet, and Ritik Dutta. Causal discovery toolbox: Uncovering causal relationships in python. *Journal of Machine Learning Research*, 21(37):1–5, 2020.
- Diviyan Kalainathan, Olivier Goudet, Isabelle Guyon, David Lopez-Paz, and Michèle Sebag. Structural agnostic modeling: Adversarial learning of causal graphs. *Journal of Machine Learning Research*, 23(219):1–62, 2022.
- Sébastien Lachapelle, Philippe Brouillard, Tristan Deleu, and Simon Lacoste-Julien. Gradient-based neural dag learning. arXiv preprint arXiv:1906.02226, 2019.
- Francesco Montagna, Nicoletta Noceti, Lorenzo Rosasco, Kun Zhang, and Francesco Locatello. Scalable causal discovery with score matching. In *Conference on Causal Learning and Reasoning*, pp. 752–771. The Proceedings of Machine Learning Research, 2023.
- Ignavier Ng, Shengyu Zhu, Zhitang Chen, and Zhuangyan Fang. A graph autoencoder approach to causal structure learning. *arXiv preprint arXiv:1911.07420*, 2019.
- Ignavier Ng, Shengyu Zhu, Zhuangyan Fang, Haoyang Li, Zhitang Chen, and Jun Wang. Masked gradient-based causal structure learning. In *Proceedings of the 2022 SIAM International Conference on Data Mining (SDM)*, pp. 424–432. SIAM, 2022.
- Judea Pearl. Causality. Cambridge university press, 2009.
- Judea Pearl, Madelyn Glymour, and Nicholas P Jewell. *Causal inference in statistics: A primer*. John Wiley & Sons, 2016.
- Jonas Peters, Peter Bühlmann, and Nicolai Meinshausen. Causal inference by using invariant prediction: identification and confidence intervals. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 78(5):947–1012, 2016.
- Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. *Elements of causal inference: foundations and learning algorithms*. The MIT Press, 2017.

- Paul Rolland, Volkan Cevher, Matthäus Kleindessner, Chris Russell, Dominik Janzing, Bernhard Schölkopf, and Francesco Locatello. Score matching enables causal discovery of nonlinear additive noise models. In *International Conference on Machine Learning*, pp. 18741–18753. The Proceedings of Machine Learning Research, 2022.
- Mauro Scanagatta, Cassio P de Campos, Giorgio Corani, and Marco Zaffalon. Learning bayesian networks with thousands of variables. *Advances in neural information processing systems*, 28, 2015.
- Marco Scutari. Learning bayesian networks with the bnlearn r package. *arXiv preprint* arXiv:0908.3817, 2009.
- Liam Solus, Yuhao Wang, and Caroline Uhler. Consistency guarantees for greedy permutation-based causal inference algorithms. *Biometrika*, 108(4):795–814, 2021.
- Peter Spirtes. An anytime algorithm for causal inference. In *International Workshop on Artificial Intelligence and Statistics*, pp. 278–285. The Proceedings of Machine Learning Research, 2001.
- Peter Spirtes and Clark Glymour. An algorithm for fast recovery of sparse causal graphs. *Social science computer review*, 9(1):62–72, 1991.
- Peter Spirtes, Clark Glymour, and Richard Scheines. *Causation, prediction, and search*. MIT press, 2001.
- Peter L Spirtes, Christopher Meek, and Thomas S Richardson. Causal inference in the presence of latent variables and selection bias. *arXiv preprint arXiv:1302.4983*, 2013.
- Marc Teyssier and Daphne Koller. Ordering-based search: A simple and effective algorithm for learning bayesian networks. *arXiv preprint arXiv:1207.1429*, 2012.
- Alessio Zanga, Elif Ozkirimli, and Fabio Stella. A survey on causal discovery: Theory and practice. International Journal of Approximate Reasoning, 151:101–129, 2022.
- Kun Zhang, Biwei Huang, Jiji Zhang, Clark Glymour, and Bernhard Schölkopf. Causal discovery from nonstationary/heterogeneous data: Skeleton estimation and orientation determination. In *IJCAI: Proceedings of the Conference*, volume 2017, pp. 1347. NIH Public Access, 2017.
- Xun Zheng, Bryon Aragam, Pradeep K Ravikumar, and Eric P Xing. Dags with no tears: Continuous optimization for structure learning. *Advances in neural information processing systems*, 31, 2018.
- Xun Zheng, Chen Dan, Bryon Aragam, Pradeep Ravikumar, and Eric Xing. Learning sparse nonparametric dags. In *International Conference on Artificial Intelligence and Statistics*, pp. 3414– 3425. The Proceedings of Machine Learning Research, 2020.
- Yujia Zheng, Biwei Huang, Wei Chen, Joseph Ramsey, Mingming Gong, Ruichu Cai, Shohei Shimizu, Peter Spirtes, and Kun Zhang. Causal-learn: Causal discovery in python. *Journal* of Machine Learning Research, 25(60):1–8, 2024.

A PROOFS

A.1 PROOF OF THEOREM 1

To prove Theorem 1, it suffices to prove its (equivalent) contrapositive statement that for $X \in \mathbf{X} \setminus \mathbf{B}$,

$$\mathbb{V}_{P_{+}(\boldsymbol{X})}\Big[P_{+}\Big(X \mid \operatorname{Pa}[X]\Big)\Big] = 0, \qquad (7)$$

which in turn can be established via showing that $P_1(X | \operatorname{Pa}[X]) = P_2(X | \operatorname{Pa}[X]) = \ldots = P_m(X | \operatorname{Pa}[X])$ or equivalently, $P_i(X | \operatorname{Pa}[X])$ does not depend on the choice of $P_i(B)$ in $P_i(X) = P_i(B) \cdot P(X \setminus B | B)$. To see this, we first establish the below lemma.

Lemma 1. Suppose $Z \in X \setminus B$, $P_i(Z \mid B)$ remains constant across different choices of $P_i(B)$.

Proof. Since $Z \in X \setminus B$, we have

$$P_i(\boldsymbol{Z} \mid \boldsymbol{B}) = \int_{\boldsymbol{X} \setminus (\boldsymbol{B} \cup \boldsymbol{Z})} P_i(\boldsymbol{X} \setminus \boldsymbol{B} \mid \boldsymbol{B}) d(\boldsymbol{X} \setminus (\boldsymbol{B} \cup \boldsymbol{Z}))$$
(8)

$$= \int_{\boldsymbol{X}\setminus \left(\boldsymbol{B}\cup\boldsymbol{Z}\right)} P\left(\boldsymbol{X}\setminus\boldsymbol{B}\mid\boldsymbol{B}\right) \mathrm{d}\left(\boldsymbol{X}\setminus\left(\boldsymbol{B}\cup\boldsymbol{Z}\right)\right), \tag{9}$$

where the first equality is a direct application of marginalization. The second equality is true due to the definition of $P_i(\mathbf{X}) = P_i(\mathbf{B})P(\mathbf{X} \setminus \mathbf{B} \mid \mathbf{B})$ which implies $P_i(\mathbf{X} \setminus \mathbf{B} \mid \mathbf{B}) = P(\mathbf{X} \setminus \mathbf{B} \mid \mathbf{B})$. Hence, Eq. (7) implies $P_i(\mathbf{Z} \mid \mathbf{B})$ remains constant across different choices of $P_i(\mathbf{B})$.

Now, suppose $Pa[X] = Z_1 \cup Z_2$ where $Z_1 \subseteq B$ and $Z_2 \nsubseteq B$, it follows that

$$P_i\left(X \mid \operatorname{Pa}[X]\right) = P_i\left(X \mid \boldsymbol{Z}_1, \boldsymbol{Z}_2\right) = P_i\left(X \mid \boldsymbol{Z}_1, \boldsymbol{Z}_2, \boldsymbol{B} \setminus \boldsymbol{Z}_1\right) = P_i\left(X \mid \boldsymbol{B}, \boldsymbol{Z}_2\right)$$
(10)
$$= P_i\left(X \mid \boldsymbol{Z}_2 \mid \boldsymbol{B}\right) / P_i\left(\boldsymbol{Z}_2 \mid \boldsymbol{B}\right)$$
(11)

$$= I_i(\boldsymbol{\Lambda}, \boldsymbol{Z}_2 \mid \boldsymbol{D}) / I_i(\boldsymbol{Z}_2 \mid \boldsymbol{D}), \qquad (11)$$

where the first equality is due to the first local Markov condition (I) in Section 3 which stipulates that given the parent $Pa[X] = Z_1 \cup Z_2$, X is independent of its non-descendant $B \setminus Z_1$. We know $B \setminus Z_1$ is X's non-descendant because $Pa[B] = \emptyset$ due to the choice of B as source variables. Then, the second equality follows from a direct application of the Bayes theorem.

Finally, note that by design, both $Z_2 \cup \{X\}$ and Z_2 are subsets of $X \setminus B$. As a result, Lemma 1 can be applied to both the numerator and denominator of the RHS of Eq. (11), asserting that they both remain constant across different choices of $P_i(B)$. This means $P_i(X \mid \text{Pa}[X])$ also remains constant against changes in $P_i(B)$ and hence, its variance in Eq. (7) is indeed zero as expected. \Box

A.2 PROOF OF THEOREM 2

Let $S(G) \subseteq X$ and $B(G) \subseteq X$ denote respectively the set of sources and an arbitrary basis set of the DAG G = (X, E). We will prove the following inductive statement:

$$P(n)$$
: "Given a DAG G of $n \ (n \ge 1)$ nodes, we have $|\mathbf{B}(G)| \le |\mathbf{S}(G)|$ "

To see this, note that in the base case n = 1, G has single node which is obviously both a source and basis node. Hence, |B(G)| = |S(G)| and P(1) is true.

Now, suppose P(n) is true, we will complete induction by showing that P(n + 1) is also true. To show this, let X be a terminal node in G that has no children. Removing X from G thus results in another DAG G' with n nodes.

Let B(G') denote an arbitrary basis set of G'. By definition, nodes in B(G') are mutually independent and for any nodes $X' \notin X \setminus (B(G') \cup \{X\})$, there exists a node Z in B(G') such that $X' \not\perp Z$ following the definition of d-separation.

Choosing $X' \in Pa[X]$, there must exist $Z \in B(G')$ such that Z is connected to X' via d-separation. Since X' is a parent of X, Z is also connected to X via d-separation. This means each node in G is either a part of B(G') or connected to another node in B(G') via d-separation. Hence, B(G') is also a basis set of G. In this case, $|B(G')| = |S(G')| \le |S(G)|$ since P(n) is true. Otherwise, suppose B(G) is a basis set of G that is not in G'. In this case, B(G) must contain X and by definition of basis (see Definition 3) and its choice, X must be an isolated node with no parents and children. This also means $B(G) \setminus \{X\}$ is a basis set of G'. Otherwise, there exists a node in G that is not connected to any nodes in the basis set B(G), resulting in a contradiction. As such, we have $|B(G)| = |B(G) \setminus \{X\}| + 1 \le |S(G')| + 1 = |S(G)|$ where the first inequality is true due to P(n) and the last equality is true since X is an isolated node which is also a source node.

As such, $|B(G)| \leq |S(G)|$ meaning P(n+1) is true if P(n) is true. Induction completes.

A.3 PROOF OF THEOREM 3

To prove that the procedure in Theorem 3 produces a basis set with maximum size, it suffices to prove that each step in this procedure removes exactly one source from the graph (see Lemma 2). In this case, the returned basis set has the same size as the source set, which is also the maximum basis set according to Theorem 2.

Lemma 2. If X has the smallest number of dependent nodes, i.e., $X \triangleq \arg \min_{Y \in V} |\Phi(Y)|$, then X is a source or X is dependent on exactly one source in G = (V, E).

Proof. To prove this lemma, we will use the following definition that distinguishes between two types of collider stemming from the V-structure $A \rightarrow B \leftarrow C$.

Definition 5. A node *B* is a Type-1 collider if and only if:

$$\exists A, C \in \operatorname{Pa}[B] : A \perp \!\!\!\perp C.$$

If such A, C do not exist, B is a Type-2 collider.

Suppose X has the smallest number of dependent nodes in G and X is dependent on more than 1 source. Since sources are mutually independent, this implies that either X is a Type-1 collider or X has an ancestor Y who is a Type-1 collider.

That is, there must exist a V-structure $S_1 \to \ldots \to Y \leftarrow \ldots \leftarrow S_2$ and a path $Y \to \ldots \to X$ where S_1 and S_2 are two sources in G. However, using the result of Lemma 3,

- 1. There must exists such a node Z on the path from S_1 to Y (Z can be S_1 , but not Y) that satisfies: $|\Phi(Z)| < |\Phi(Y)|$. This is because Y is a Type-1 collider.
- 2. Since $Y \to \ldots \to X$, $|\Phi(Y)| \le |\Phi(X)|$.

The above implies $\exists Z : |\Phi(Z)| < |\Phi(X)|$ which contradicts the assumption that X has the smallest number of dependents. Hence, X must be either dependent on at most 1 source, or X is a source variable itself.

Lemma 3. Let $\Phi(X) \triangleq \{Y \in \mathbf{V} : Y \not\perp X\}$. If $\exists A, B \in \mathbf{V} : A \to \ldots \to B$ then $|\Phi(A)| \le |\Phi(B)|$ where the equality "=" occurs when B is not a Type-1 collider.

Proof. For any nodes X in the graph that is dependent on A, there must exist a path that connects X and A without a collider in between. Furthermore, since $A \to \ldots \to B$, the path from X to A and then, to B also does not contain any colliders. Therefore, all nodes dependent on A are also dependent on B. Hence, $\Phi(A) \subseteq \Phi(B)$ or $|\Phi(A)| \leq |\Phi(B)|$.

Following the above logic, suppose B is not a collider, any dependent of B must also be a dependent of A and hence, $|\Phi(B)| = |\Phi(A)|$. Otherwise, if B is a collider but there does not exist $A, C \in Pa[B]$ such that $A \perp C$, then any dependent of B can be reached by A via d-separation and therefore, we have $|\Phi(B)| = |\Phi(A)|$.

Now, suppose *B* is a collider such that there exists $A, C \in Pa[B]$ where $A \perp C$ (Type-1 collider). In this case, $|\Phi(B)| = |\Phi(A) \cup \Phi(C)| > |\Phi(A)|$. Thus, in all cases, $|\Phi(B)| \ge |\Phi(A)|$ if $A \to \ldots \to B$ and the equality occurs only when *B* is not a Type-1 collider.

A.4 PROOF OF THEOREM 4

Let D_i be any downsampled dataset of D. Consider a candidate value b of the source variable B that occurs D. Let n(b) and $n_i(b)$ denote the corresponding numbers of data points in D and D_i that

have B = b. Since D_i is downsampled from D, $n_i(b) \le n(b)$. Hence, let $n \triangleq |D|$ and $n_i \triangleq |D_i|$,

$$P(\boldsymbol{B} = \boldsymbol{b})/P_i(\boldsymbol{B} = \boldsymbol{b}) \geq (n(\boldsymbol{b})/n)/(n_i(\boldsymbol{b})/n_i)$$

= $(n(\boldsymbol{b})/n_i(\boldsymbol{b})) \cdot (n_i/n) \geq n_i/n = |D_i|/|D|$ (12)

which follows from the facts that P(B = b) = n(b)/n, $P_i(B = b) = n_i(b)/n_i$, and $n(b) \ge n_i(b)$ due to the downsampled nature of D_i . Taking the minimum over all candidate values b of B,

$$\gamma_i \triangleq \min_{\mathbf{b}} \left(P(\mathbf{B} = \mathbf{b}) / P_i(\mathbf{B} = \mathbf{b}) \right) \ge |D_i| / |D|$$
 (13)

As a result, for all downsampled dataset D_i of D, the downsampling rate $|D|/|D_i| \ge \gamma_i$ where γ_i is defined in terms of the original and target marginal over the source variable B as stated in Eq. (13). Hence, if D_i is the downsampled data with minimum downsampling rate $|D|/|D_i| \ge \gamma_i^{-1}$. Such a downsampled dataset can indeed be found using the sampling procedure in Theorem 5 below.

A.5 PROOF OF THEOREM 5

Since D_i is created via sampling with no replacement $P_i(\boldsymbol{B} = \boldsymbol{b}) \cdot |D| \cdot \gamma_i$ data points from D where $\boldsymbol{B} = \boldsymbol{b}$, we know that $n_i(\boldsymbol{b}) = P_i(\boldsymbol{B} = \boldsymbol{b}) \cdot |D| \cdot \gamma_i = P_i(\boldsymbol{B} = \boldsymbol{b}) \cdot n \cdot \gamma_i$. Thus,

$$n_{i} \triangleq \left(\sum_{\mathbf{b}} n_{i}(\mathbf{b})\right) = \left(\sum_{\mathbf{b}} P_{i}(\mathbf{B} = \mathbf{b}) \cdot |D| \cdot \gamma_{i}\right)$$
$$= P_{i}(\mathbf{B} = \mathbf{b}) \cdot n \cdot \gamma_{i} = n \cdot \gamma_{i} \cdot \left(\sum_{\mathbf{b}} P_{i}(\mathbf{B} = \mathbf{b})\right) = n \cdot \gamma_{i} \quad (14)$$

As a result, $|D_i|/|D| = n_i/n = \gamma_i$ and $n_i(\mathbf{b})/n_i = (P_i(\mathbf{B} = \mathbf{b}) \cdot n \cdot \gamma_i)/(n \cdot \gamma_i) = P_i(\mathbf{b})$. Thus, $D_i \sim P_i(\mathbf{X})$ and the downsampling rate $|D|/|D_i|$ achieves minimum as expected.

A.6 PROOF OF THEOREM 6

A. Convexity. We will first prove that the subspace of $P_i(B)$ that satisfies $\gamma_i \ge \gamma_o$ is convex. To see this, consider $P_i^1(B)$ and $P_i^2(B)$ whose (inverse) downsampling rates (see Theorem 4) γ_i^1 and γ_i^2 are both larger than γ_o . This means both $P_i^1(B)$ and $P_i^2(B)$ belong to the aforementioned subspace.

Now, let $\alpha \in (0,1)$ and $P_i^{\alpha}(\mathbf{B}) = \alpha \cdot P_i^1(\mathbf{B}) + (1-\alpha) \cdot P_i^2(\mathbf{B})$. The (inverse) downsampling rate of $P_i^{\alpha}(\mathbf{B})$ is defined as

$$\gamma_i^{\alpha} \triangleq \min_{\boldsymbol{b}} \left(\frac{P(\boldsymbol{B} = \boldsymbol{b})}{P_i(\boldsymbol{B} = \boldsymbol{b})} \right) = \min_{\boldsymbol{b}} \left(\frac{P(\boldsymbol{B} = \boldsymbol{b})}{\alpha \cdot P_i^1(\boldsymbol{B} = \boldsymbol{b}) + (1 - \alpha) \cdot P_i^2(\boldsymbol{B} = \boldsymbol{b})} \right)$$
(15)

$$\geq \min_{\boldsymbol{b}} \left(\frac{P(\boldsymbol{B} = \boldsymbol{b})}{\alpha \cdot \gamma_o^{-1} \cdot P(\boldsymbol{B} = \boldsymbol{b}) + (1 - \alpha) \cdot \gamma_o^{-1} \cdot P(\boldsymbol{B} = \boldsymbol{b})} \right) = \frac{1}{\gamma_o^{-1}} = \gamma_o , \quad (16)$$

where the inequality follows from the facts that (1) $P(\mathbf{B} = \mathbf{b})/P_i^1(\mathbf{B} = \mathbf{b}) \ge \gamma_i^1 \ge \gamma_o$; and $P(\mathbf{B} = \mathbf{b})/P_i^1(\mathbf{B} = \mathbf{b}) \ge \gamma_i^2 \ge \gamma_o$ which follows from the (inverse) downsampling rate's definition in Theorem 4. This means $P_i^1(\mathbf{B} = \mathbf{b}) \le \gamma_o^{-1} \cdot P(\mathbf{B} = \mathbf{b})$ and $P_i^2(\mathbf{B} = \mathbf{b}) \le \gamma_o^{-1} \cdot P(\mathbf{B} = \mathbf{b})$, which can plugged into Eq. (15) to arrive at Eq. (16). This in turn implies $P_i^{\alpha}(\mathbf{B})$ belongs to the subspace of source priors with the induced (inverse) downsampling rate above γ_o . Thus, following the definition of convexity, we know that this subspace is convex.

B. Convex Hull. Note that any source prior $P_i(B)$ is a point in an *r*-dimensional simplex where *r* is the number of candidate values of the source variables *B*. The convex hull of the aforementioned convex set can then be determined by finding its intersection with the edges of the simplex. This comprises *r* points of the following form:

$$P^{(k)}(\boldsymbol{B}) = \alpha_k \cdot P(\boldsymbol{B}) + (1 - \alpha_k) \cdot \delta_k(\boldsymbol{B}), \qquad (17)$$

where α_k is selected such that $\gamma_i^{(k)} \triangleq \min_{\boldsymbol{b}} (P(\boldsymbol{B} = \boldsymbol{b})/P^{(k)}(\boldsymbol{B} = \boldsymbol{b})) = \gamma_o$. Solving for α_k results in the following equation:

$$\gamma_o = \min\left(\frac{1}{\alpha_k}, \frac{P(\boldsymbol{B} = \boldsymbol{b}^{(k)})}{\alpha_k \cdot P(\boldsymbol{B} = \boldsymbol{b}^{(k)}) + (1 - \alpha_k)}\right) = \frac{P(\boldsymbol{B} = \boldsymbol{b}^{(k)})}{\alpha_k \cdot P(\boldsymbol{B} = \boldsymbol{b}^{(k)}) + (1 - \alpha_k)}, \quad (18)$$

where $\mathbf{b}^{(k)}$ denote the k-th candidate value of \mathbf{B} . The second equality in the above holds since $\alpha_k \in (0, 1)$ which means $\alpha_k \cdot P(\mathbf{B} = \mathbf{b}) \leq \alpha_k \cdot P(\mathbf{B} = \mathbf{b}) + 1 - \alpha_k$ or equivalently, $1/\alpha_k \geq P(\mathbf{B} = \mathbf{b})/(\alpha_k \cdot P(\mathbf{B} = \mathbf{b}) + 1 - \alpha_k)$. This allows us to get rid of the min operator in Eq. (18) and consequently, compute α_k in closed-form:

$$\alpha_k = \left(1 - \gamma_o^{-1} \cdot P\left(\boldsymbol{B} = \boldsymbol{b}^{(k)}\right)\right) / \left(1 - P\left(\boldsymbol{B} = \boldsymbol{b}^{(k)}\right)\right) .$$
(19)

This completes our derivation for the convex hull in Theorem 6.

A.7 PROOF OF THEOREM 7

Suppose $U, V \in Pa[X]$, it must follow that $U \in \mathbb{M}(V)$ and $V \in \mathbb{M}(U)$ since U and V are spouses. This means $(U, V) \in E$. Hence, there is a bidirectional edge in G'(X) between any two parents of X which means Pa[X] corresponds to a clique in G'(X).

B ALGORITHMS & ANALYSIS

B.1 FIND THE BASIS OF A GRAPH

This section provides visualization and pseudocode of the basis finding procedure in Theorem 3.

A. Pseudocode. First, the pseudocode of the basis finding algorithm is detailed below.

Algorithm 1 Maximum-Sized Basis Search

Input: Data *D* and the set of all variable indices *V*. Output: Maximum-Sized Basis *B*. 1: Compute $[\Phi]_{ij} = \mathbb{I}(X_i \not \perp X_j)$ using existing statistical independence test. 2: while |V| > 0 do 3: Choose $t = \arg \min_{i \in V} \sum_{j} [\Phi]_{ij}$ 4: Update $B = B \cup \{X_t\}$ 5: Update $V \leftarrow V \setminus (\{i : [\Phi]_{ti} = 1\} \cup \{X_t\})$ 6: end while 7: return *B*



Figure 5: A step-by-step visualization of Algorithm 1. From left to right: each step finds the node with minimum number of dependents and remove it along with its dependent set from the graph. Each node is accompanied by a number in red that indicates the number of variables being dependent on that node. The removed nodes are colored gray while active nodes are blue. The example graph is taken from the ASIA dataset (Scutari, 2009).

This algorithm is guaranteed to find the maximum-sized of basis variables since each of its iteration will remove exactly one source node from V as proved in Appendix A.3. Thus, the number of nodes added to the basis is also the number of source nodes. Such basis is guaranteed to has maximum size following Theorem 2. A visualization of this algorithm is provided in Figure 5.

B. Time Complexity. Algorithm 1 consists of (i) computing the dependence network matrix Φ , and (ii) selecting the node with minimum number of dependents from V per iteration. The complexity for step (i) is $O(d^2)$ while the complexity for step (ii) is O(d) per iteration. As there are at most d iterations, the total cost of Algorithm 1 is therefore $O(d^2)$.

B.2 FIND PLAUSIBLE PARENT SETS

Algorithm 2 Finding Plausible Parent Sets

Inputs: root node r_o containing the virtual variable X_o , Markov blanket $\mathbb{M}(X)$ of all variable X. **Outputs:** a plausible-parent-set tree, the path from each intermediate search node of variable X to each leaf node represents a plausible candidate for $\operatorname{Pa}[X]$

```
1: Initialize an empty list of leaves L \leftarrow \emptyset;
 2: search(r_o) \leftarrow the set of all variables X;
 3: procedure RECURSIVE-BUILD(r)
 4:
           for \ell in L do
 5:
                if path(r_o \to r) \cup search(r) \in path(r_o \to \ell) then return
 6:
                end if
 7:
           end for
 8:
           S \leftarrow \text{sort } X_i \in \text{search}(r) \text{ in the decreasing order of } |\text{search}(r) \cap \mathbb{M}(X_i)|
 9:
           oldsymbol{V} \leftarrow \emptyset
                             #initialize the set of visited nodes
10:
           if |\boldsymbol{S}| > 0 then
11:
                while |\boldsymbol{S}| > 0 do
12:
                      X \leftarrow S.pop()
13:
                      q \leftarrow \operatorname{node}(X)
                      \operatorname{search}(q) \leftarrow (\mathbb{M}(X) \cap \operatorname{search}(r)) \setminus V
14:
                      path(q) \leftarrow path(r)) \cup \{X\}
15:
                      Recursive-Build(q)
16:
                      V \leftarrow V \cup \{X\}
17:
18:
                end while
19:
           else
                L \leftarrow L \cup \{r\}
20:
           end if
21:
22: end procedure
```



Figure 6: Step-by-step visualization of Algorithm 2 showing (a) a Markov blanket of a search node, (b) a plausible parent tree resulting from a search at a particular node, and (c) a search tree starting from a virtual variable which are connected to all other (real) variables. The grey-colored nodes correspond to branches which are terminated following lines 4-7 in Algorithm 2. For example, in (b) we see that continuing to explore the branch $X_6 - X_4$ will re-produce the content which is already produced in a previous $X_6 - X_2 - X_4$ branch. Hence, it was terminated.

Theorem 7 implies that if two nodes in the Markov blanket of X do not belong to each other's Markov blanket, then they can not both be parents of X. According to this logic, we wish to find all sets of nodes whose Markov blanket contains one another. To do this, we develop a tree-based recursive algorithm that works as follows. Given the Markov blanket of all nodes, we perform

Algorithm 2, which is a modification of the Bron-Kerbosch algorithm (Bron & Kerbosch, 1973). Instead of running Bron-Kerbosch for each node X and its Markov blanket $\mathbb{M}(X)$, we add a virtual node that directly connects to all other nodes in the graph and start building the set of plausible parent sets for all nodes recursively from this virtual node.

Note that Algorithm 2 does not create a sub-tree (branch) if the search space of the root of that sub-tree is guaranteed to re-produce the content of a previously visited leaf and generate no new information – see lines 4-7. Otherwise, a plausible parent set of X is the set of variables in the path from the search node containing X to a leaf node, which is computed following the recursion of Algorithm 2. A step-by-step visualization of this algorithm is illustrated in Figure 6.

C ADDITIONAL EXPERIMENTS

C.1 EXPERIMENT SETTINGS

Baselines & Hyper-parameters for baselines. In this section, we present the experimental settings for each baseline used in the empirical report in the Section V of the main manuscript. Firstly, the code for PC is from the bnlearn package(Scutari, 2009), provided in Python. We reported results of the PC-stable, which is an upgraded version of PC. On the other hand, the implementation of GIES and FCI are provided from the cdt package (Kalainathan et al., 2020) and the causal-learn package (Zheng et al., 2024), respectively. These algorithms require no tuned parameters as input and therefore can be used as recommended from the package from which they are provided. Notears and MLP-Notears are open-source on Github, and are provided by Zheng et al. (2018). For our experiments, we reuse the parameters recommendation in the original papers: Notears ($\lambda_1 = 0.1, w_{th} = 0.3$) and MLP-Notears ($\lambda_1 = \lambda_2 = 0.01, w_{th} = 0.3$). It is worth noting that we do not modify the original architecture of the network used in MLP-Notears. Lastly, DAS and SCORE are also open-source and are available from (Montagna et al., 2023). Both of these baselines share the same operational parameters as follows: $\eta_G = \eta_H = 0.001, K = 10, \text{pns} = 10$, threshold = 0.05 and camcutoff = 0.001.

Hyper-parameters for our proposed method GLIDE. Our proposal has the following hyperparameters:

- (1) The number of prior distributions m used for the invariance test.
- (2) γ factor which controls the ratio of the data generated from weak interventions.

Both of these parameters are concerned with the data sampling procedure. Other than these two parameters, there are two more thresholds which control the level of tolerable variance we consider as invariant and the confidence interval for the Conditional Independence Tests (CITs). These thresholds are set at 10^{-3} and 0.05, respectively, across all experiments reported in this study. We recommend tuning the former threshold depends on the nature of the data whereas the latter should remain unchanged.

C.2 Ablation studies

C.2.1 Studies on the influence of m and γ

The number of prior distributions m corresponds to the number of environments/sub-datasets that our proposal generates. The higher this figure, the better the performance of invariance test. However, this comes with a trade-off on time consumption: the invariance test loops through all generated sub-datasets; therefore, it consumes increasing time linearly with the increase of K. On the other hand, γ factor virtually controls the error in estimations produced by sub-datasets because γ bounds the volume of the downsampled datasets.

In this study, we conduct experiments on the 100-variable Erdos-Renyi categorical data graph and continuous data graph. In each data model, we examine different values of $\gamma \in \{0.2, 0.4, 0.6, 0.8\}$, along with increasing values of $m \in \{10, 20, 30, 40, 50\}$. The input data volume is 10000 and we examine 10 runs to report mean values and 95% confidence intervals as in Table3. The results show a significant decrease (averagely, over 18%) in SHD when we increase the number of environments m from 10 to 50. However, such performance comes at the expense of Runtime. As we discussed,

m	γ	SHD	Spurious rate	Runtime	
	'		(%)	(minutes)	
	0.2	122.2 ± 9.788	5.91 ± 1.678	1.95 ± 0.037	
10	0.4	127.4 ± 10.31	5.41 ± 0.951	2.02 ± 0.127	
	0.6	113.6 ± 2.525	4.64 ± 1.01	2.39 ± 0.113	
	0.8	121.0 ± 7.717	4.19 ± 1.552	2.45 ± 0.114	
	0.2	116.8 ± 9.139	5.53 ± 0.841	4.54 ± 0.221	
20	0.4	109.8 ± 8.796	4.30 ± 0.792	4.18 ± 0.086	
20	0.6	123.0 ± 2.772	5.96 ± 0.957	3.88 ± 0.098	
	0.8	112.8 ± 5.958	3.27 ± 1.448	3.93 ± 0.095	
	0.2	117.0 ± 4.679	5.19 ± 0.505	5.67 ± 0.068	
30	0.4	112.8 ± 7.874	4.37 ± 0.613	5.66 ± 0.012	
50	0.6	102.0 ± 6.263	3.85 ± 1.603	5.58 ± 0.218	
	0.8	106.0 ± 9.461	4.15 ± 0.981	5.75 ± 0.194	
	0.2	114.0 ± 5.714	4.79 ± 0.753	7.44 ± 0.032	
40	0.4	112.8 ± 5.088	3.94 ± 0.529	7.50 ± 0.028	
40	0.6	108.6 ± 8.072	3.97 ± 1.536	7.30 ± 0.051	
	0.8	108.8 ± 4.444	4.03 ± 1.279	7.36 ± 0.048	
	0.2	100.8 ± 2.432	5.03 ± 0.518	8.97 ± 0.024	
50	0.4	106.6 ± 5.834	4.64 ± 0.958	9.35 ± 0.061	
50	0.6	99.01 ± 6.261	3.68 ± 1.589	9.21 ± 0.021	
	0.8	107.6 ± 7.606	4.61 ± 0.955	9.22 ± 0.039	

Table 3: Influence of hyper-parameters. Mean value and confidence interval 95% are reported.

the time complexity is linear with m, thus, we see that the runtime is about 5 times higher when m is 5 times higher. As for γ . The increase of γ does not guarantee better performance, as can be seen in settings m = 10, m = 30, and m = 50 where the best γ is neither the smallest nor the biggest. Despite the effect of γ is secondary to that of m, a good γ can boost the performance roughly 6-8% in term of SHD without incurring additional time complexity.

C.2.2 STUDIES ON DIFFERENT TOPOLOGIES

In this section, we test the performance of baselines and our proposal on Scale-free graphs and Bipartite graphs (Zheng et al., 2018). The code used to generate the datasets for experiments is introduced by Zheng et al. (2018). We also evaluate our proposal and baselines in normal and extreme cases - the former has the number nodes increasing from 100 to 500 and the number of edges equals the number of nodes, the latter has the number of nodes fixed at 500 while the number of edges increases from 600 to 1000. In each case, we also generate linear Gaussian data and non-linear non-Gaussian data scenarios. The results are depicted in Figures 7 and 8 for Bipartite graphs and Scale-free graphs, respectively.

Bipartite graphs. Regarding Bipartite graphs, our proposal shows superior performance in almost all cases and data generation scenarios. As can be seen in Figure 7(a)-upper, our proposal achieves relatively similar to the strongest baseline (MLP-Notears) in terms of SHD and spurious rate, but costs over an order of magnitude smaller in runtime, about 15.68 times.

The performance gap is noticeable when it comes to the non-linear non-Gaussian data scenario – Figure 7(a)-lower. Our proposal now has a clear advantage in terms of SHD (roughly 10% performance gap) and spurious rate (17.21% performance gap) against MLP-Notears, while the runtime difference is as significant as the previous scenario. It is worth noticing that, none of the baselines can achieve both fast runtime and relatively high performance. For example, GIES and DAS have the same runtime as our proposal, but their performance in terms of SHD and spurious rate are significantly worse (at an aaverage of 61.55% and 46.22%, respectively).

As for the extreme cases in Bipartite graphs, we isolate best baselines, Notears and MLP-Notears, to compare with our proposal. The overall result in Figure 7(b) is that our proposal outperforms both baselines in both scenarios and in all metrics, albeit with an exception. Figure 7(b)-upper shows the result in linear Gaussian data scenario. We can see that our proposal has a noticeable gap to the other



(a) Evaluation on linear Gaussian (upper) and non-linear non-Gaussian (lower) data models. The number of edges equals the number of nodes.



(b) Evaluation on linear Gaussian (upper) and non-linear non-Gaussian (lower) data models at extremes. The number of nodes is fixed at 500. Best baselines are selected to compare with our proposal.

Figure 7: **Performance on continuous data models in normal and extreme cases on Bipartite graphs.** Apply for all metrics: Lower is better.

2 baselines in both SHD (12.81% and 27.18% lower than MLP-Notears and Notears, respectively) and runtime (16.62 and 14.18 times less than MLP-Notears and Notears, respectively). However, Notears marginally outperforms our proposal in term of spurious rate, about 1%.

Figure 7(b)-lower shows that our proposal returns relatively stable performance even in the nonlinear non-Gaussian data scenario with an average of 0.69% spurious rate. Furthermore, our proposal consistently outperforms both Notears and MLP-Notears in all cases in term of SHD. In contrast, Notears and MLP-Notears seem to be heavily affected by non-linearity in term of spurious rate. MLP-Notears has an average of 10.06% spurious rate – about 5 times higher than linear Gaussian data scenario. Notears also suffers an average 3% higher spurious rate. Interestingly, both baselines seem to benefit from the increasing number of edges in the graph as the results show an downward trend in spurious rate. However, it is worth noticing that Notears incurs a linearly increasing in SHD as the number of edges increases.

Scale-free graphs. Regarding Scale-free graphs, Figure 8 shows the results of algorithms in normal cases. Overall, despite not being the best performer in all metrics, GLIDE maintains a good



Figure 8: Evaluation on linear Gaussian (upper) and non-linear non-Gaussian (lower) data models on Scale-free graphs. The number of edges equals the number of nodes.



Figure 9: Degeneracy measure on different topologies. The number of nodes is fixed at 500.

trade-off between all metrics. Regarding the **L-G** case (Figure 8-upper) shows a different trends in performance of baselines compared to **GLIDE** in term of SHD. In details, the SHD of our proposal gradually increases until plateaus out in the 400 and 500-node graphs contrasting to other baselines (e.g., SCORE, and MLP-NOTEARS) that show an accelerating rate in SHD towards 500-node graphs. Furthermore, these baselines are noticeably worse than **GLIDE** when it comes to spurious rate (**GLIDE** produces an average of 5.52% less spurious relationships than SCORE, 28.57% than MLP-NOTEARS while being comparable to NOTEARS and DAS). As for the **nL-nG** case (Figure 8-lower), we have the same observation where our proposal consistently being comparable with other baselines in term of SHD, having a low spurious rate and runtime, simultaneously.

As regarding the extreme cases on the Scale-free graphs, when the number of edges exceeds 700, we encounter the following problem: the average number of parents for each node becomes exponentially large due to the nature of the graph. Such that the performance of MLP-Notears is heavily impacted: MLP-Notears cannot produce meaningful results within 48 hours. Our proposal - **GLIDE** is also impaired by this setting. The reason is that the data set is not sufficiently large to perform invariance test with adequate accuracy due to the exponentially large number of parents. Therefore, we do not include the report on performance of **GLIDE** as well as other baselines on this particular setting.

Datasets	Sachs	Insurance	Water	Alarm	Barley	Pathfinder	Munin
Number of Nodes	11	27	32	37	48	186	1041
Degeneracy p	3	4	6	4	5	5	4

Table 4: Degeneracy measures on real-world causal graphs.

C.2.3 STUDIES ON THE DEGENERACY MEASURE

In this ablation study, we investigate the degeneracy measure p on different topologies and connect them to the performance of **GLIDE**. As mentioned in Section 4.3, the time complexity that **GLIDE** requires to find the plausible parent sets of all d variables is $O(pd \cdot 3^{p/3})$. As such, we want to assert that, the degeneracy p of most graphs is indeed insignificant compared to the number of variables d of those graphs. To show this, we design the following ablation study: we generate DAG with a fixed number of nodes (500) and an increasing number of edges (500 to 2000). For each setting of the number of edges, we randomly generate 10 graphs and record their degeneracy measure. Notice that most baselines presented in this research are incapable of running on graphs with 500 nodes and more than 1000 edges – as we presented in Main text's Section 5 and Appendix C.2. However, for the sake of the argument of this section, we increases the number of edges to 2000 to investigate the range value of the degeneracy measure p. Figure 9 shows the mean and confidence interval 95% of the degeneracy measure on different topologies (Erdos-Renyi, Bipartite, and Scale-free graphs) as the number of edges increases.

Clearly, Bipartite graphs return the most stable degeneracy measure out of the three topologies. With very small variations, the degeneracy measures on increasingly dense Bipartite graphs demonstrate an almost linear growth. We speculate that this stable behavior of the degeneracy measure of the Bipartite graph benefits the performance of **GLIDE**, as we can see in Figure 7(b) where **GLIDE** consistently and noticeably outperforms the other two prominent baselines, especially in term of runtime. It is worth noting that at the (500-node, 1000-edge) setting, the degeneracy measure is roughly 15, which is much less than the number of nodes. This shows that, at cases where most baselines take hours to solve, **GLIDE** can still achieve almost quadratic time complexity (Section 4.3) and thus only takes minutes to recover the causal graph.

Regarding the Erdos-Renyi graphs, we can see it random nature manifest in a noticeable – yet not too significant – variation of the degeneracy measure. Nonetheless, these values (even at maximum, e.g., p = 13 at 1000-edge or p = 22 at 2000-edge graphs) are insignificant compared to the number of nodes. This partially explains the scalability of **GLIDE** on the Erdos-Renyi graphs as we reported in Main text's Section 5. Contrasting to the previous two, the degeneracy measure on Scale-free graphs shows an unstable behavior and a wide range of fluctuation. The degeneracy measure on 500-edge graphs upto 700-edge graphs has an average of 52.4 ± 14.27 , which is highly unstable and generally much higher than that of the same setting but on other topologies. As the number of edges gradually reaches 700, the degeneracy has a sudden leap upto an average of 112 ± 27.43 . This unusual (and perhaps, non-linear) behavior may stem from the implementation (Zheng et al., 2018) or may need further studies to address. Regardless, we can see that the range for degeneracy in Scale-free graphs are significantly higher than that in Erdos-Renyi or Bipartite graphs, which potentially affect the runtime of Algorithm 2. However, as we will see in Appendix C.2.4, the number of plausible sets in Scale-free graphs is still negligible compared to the number of nodes.

Finally, we also investigate the degeneracy of real-world causal graphs. These graphs are the same as the one presented in the Main text's Section 5.3. As it turns out, real-world causal graphs are indeed sparse (see Table 4). In details, the degeneracy of the causal graph of these dataset are: (Sachs: 3, Insurance: 4, Water: 6, Alarm: 4, Barley: 5, Pathfinder: 5, Munin: 4). In which the Munin dataset has over 1000 variables. This shows that, in practice, we indeed often encounter large causal graphs that are sparse.

C.2.4 Studies on the number of plausible parent sets

The number of plausible parent sets plays an essential role in the time complexity of our proposal. As per Algorithm 2 – an enhanced Bron-Kerbosch algorithm (Bron & Kerbosch, 1973), each plausible set corresponds to a maximal clique. Theoretical results from Bron & Kerbosch (1973) give us the



Figure 10: **Investigation on the number of plausible parent sets in different topologies.** The number of edges equals the number of nodes.



Figure 11: Investigation on the number of plausible parent sets in different topologies. The number of edges doubles the number of nodes.



Figure 12: Investigation on the number of plausible parent sets in different topologies. The number of nodes is fixed at 500.

bound for the number of maximal clique, which is $O((d-p) \cdot 3^{p/3})$ where p is the degeneracy of the graph. As we have shown in Appendix Section C.2.3, p is indeed insignificant compared to d in the common Erdos-Renyi and Bipartite class of graphs. In this section, we empirically investigate the number of plausible parent sets in different graph scenarios. For each scenarios, we randomly generate 10 graphs using the source code provided by (Zheng et al., 2018) and perform Algorithm 2, then record the number of plausible parent sets of each node. We use the violin plot to show the distribution of the number of plausible sets at each setting. The maximum, minimum, mean, and median are reported.



Figure 13: **Preliminary results in Distributed Causal Discovery settings.** The performance is a straightforward application of our proposal to the setting of Distributed Causal Discovery, which might not be the best possible performance.

Figure 10 shows the number of plausible sets on graphs whose number of nodes equals that of edges, which gradually increases from 100 to 1000. Notice that this setting mimics that of experiments that we reported in normal cases in Main text's Section 5.1 but the number of nodes grows from 100 to 500. As we can see in Figure 10(a), the number of plausible sets on Erdos-Renyi graphs peaks at 10, with mean being roughly 2. This supports the excellent runtime of **GLIDE** as we reported. When it comes to Bipartite graphs, the number of plausible sets is approximately double that in the same settings on the Erdos-Renyi graphs. However, in all cases, this figure is much lower than the number nodes. Noticably, in the case of Scale-free graph, we see an almost similar distribution of the number of plausible sets across graphs with different nodes. In details, most nodes have only 1 plausible parent set, and other have 2. This observation combined with nature of the Scale-free graphs, suggest that each plausible set contains a large number of nodes. Consequently, given that the observational data volume is limited at 10000, the accuracy of the invariance test of **GLIDE** reduces significantly, which explains the performance that we reported in Appendix Section C.2.

We also conduct experiments on graphs that are twice as dense (the number of edges is doubled that of nodes) as in the previous setting, and the results are reported in Figure 11. Interestingly, while the number of plausible sets on Erdos-Renyi roughly doubles (both maximum and medium) compared to the previous setting, the that on Bipartite graphs increases significantly. In details, in most setting, the maximum number of plausible ranges from under 100 (on 100-node graphs) to over 300 (on 400-node graphs), which is innegligible compared to the number of nodes. Nontheless, these figures are still strictly less than the number of nodes. It is worth noting that most nodes in the graphs have about 10 to 15 plausible parent sets. As for the Scale-free graphs (Figure 11(c)), despite the median number of plausible sets is 10, the maximum increases linearly with the number of nodes and then plateaus out at about 150 when the number of nodes reaches 900.

Lastly, we study the number of plausible sets at extreme scenarios (as in Main text's Section 5.1): we fix the number of nodes at 500, and increase the number edges from 500 to 1000. The results are depicted by box plot in Figure 12. The number of plausible sets on the Erdos-Renyi graphs is strictly less than 25 even in the hardest setting (1000-edge graphs) while the mean value is roughly 3.8 and 75%-percentile is at 5. In contrast, the number of plausible sets on the Bipartite graphs grows non-linearly, reaching a maximum (outlier) 180 on 1000-edge graphs. However, since most nodes do not have more than 5 plausible parent sets, we see that the outliers have little effect on the overall performance of **GLIDE**, as can be seen in Figure 7(b). On the other hand, the number of plausible sets on Scale-free graphs is very unstable, and (interestingly) has the same trend as the degeneracy (reported in Figure 9(c)). In details, when the number of edges is below 700, we have the number of plausible sets ranges from 1 to 3 with the mean 1.4. But when the number of edges exceeds 700, despite the mean of 5 and 75%-percentile of 10, the maximum number of plausible sets may reach to over 100.

C.2.5 PRELIMINARY RESULTS ON DISTRIBUTED CONFIGURATIONS

Out method benefits from distributed data because the data is already segmented in advance. However, this also poses a new challenge in general: uncontrolled skewness. In this study, we simulate 10 distributed devices with 10 separated datasets, each having 1000 samples of d variables. These datasets share the same causal structure but different data-generating seeds (Zheng et al., 2018). In the realm of distributed causal discovery, CDNOD (Zhang et al., 2017) stands as a seminal representative. We compare the performance of our proposal to that of CDNOD on the same criteria as in the Main text.

Figure 13 shows the results of **GLIDE** as we increase the number of variables d from 10 to 500. Note that the reported performance is a straightforward application of our proposal to the setting of Distributed Causal Discovery, which lacks considerations on the nature of distributed data. Therefore, these results are not the peak performance of **GLIDE**. Nonetheless, it can be seen that our proposal significantly outperforms CDNOD in all criteria, up to over 84.09% in SHD when the number of nodes reaches 500, albeit with relatively poorer performance at small-sized graphs. The difference in term of SHD is staggering - the error of GLIDE grows almost linearly whereas that of CDNOD displays an exponential trend instead. Lastly, the spurious rate of our proposal is stable around 20% - 25% where as that of CDNOD sky-rises up to over 80% when the graph gets large. This study once again demonstrates our capability to scale well with size of the graph.

SUPPORT MATERIALS D



D.1 ON THE GENERATION OF WEAKLY INTERVENTIONAL SAMPLING DISTRIBUTIONS

 $\delta(\boldsymbol{B} = \boldsymbol{b}^3)$

 $D^{(2)}(B)$

 $P_i(\boldsymbol{B})$ (c) Step 2: Performing weighted average of the boundary points to generate midpoints $P_i(\boldsymbol{B})$.

Figure 14: The procedure of generating distributions for downsampling.

The fact that we can sample an infinite number of prior distributions $P_i(B)$ in $C_r(\gamma_0)$ is not helpful from the computational point of view. Because a larger number of prior distributions incurs a corresponding increase in the processing time of the invariance test. Furthermore, when the number of sampled prior distributions grows too large, it is inevitable that there exist similar prior distributions, which in fact does not improve the quality of the invariance test. On the other hand, an insufficient number of prior distributions reduces the reliability of the invariance test. Therefore, we wish to select only an adequate number of representative prior distributions within $C_r(\gamma_0)$.

Inspired by the aforementioned logic, we perform the following:

1. Sampling $P_i(B)$ where i goes up to 10^4 (Figure 15(a)) to overwhelmingly fill in the convex hull $C_r(\gamma_0)$. We use the Dirichlet distribution to sample weighted vectors $a^{(i)}$.



Figure 15: A mock-up experiment: Representative distributions selection via K-means.

2. Using the K-means algorithm with parameter K = m to cluster them into m partitions, each corresponds to a representative centroid. These centroids are the output prior distributions of the procedure (Figure 15(b)).

To this end, we have sampled m prior distributions that act as variants of the original P(B) and can be used to fuel the invariance test by applying Theorem 6.

D.2 FURTHER DISCUSSION

Parallelism Prospect. Our proposed framework **GLIDE** is readily applicable to scenarios where the data is distributed across multiple, private local devices such as the popular federated learning setting. In such scenarios, **GLIDE** will benefit directly from the fact that the distributed datasets, which are presumed to be governed by the same causal model, are already admitting identical underlying conditional distributions. As such, **GLIDE** can use those local datasets as synthetic augmented datasets which are essential to the effect-cause distributional invariance test. This can save time and help avoid unnecessary sampling errors.

Limitations. The proposed framework GLIDE has the following limitations:

1. Our performance partially relies on how well Algorithm 1 can find the source nodes for the basis. As previously mentioned, it is not guaranteed that Algorithm 1 will find all the basis variables. However, our extensive experiments have shown that when the number of data augmentations increases, the chance that non-source nodes are included in the basis set is lessen.

2. GLIDE's performance might be hampered on scenarios where the observational data is insufficient to discern the true complexity of the underlying causal graph – as can be seen with scale-free graphs.

E REPRODUCIBILITY

Software. Our implementation is in Python. The requirements include the installation of the causallearn package (Zheng et al., 2024) for the use of CITs, scikit-learn Python package, along with pandas and numpy, which are standard libraries commonly used in Python.

Support Module for Continuous Data. For the use of our proposal for continuous data models, we use a Discretizer module provided by the Python standard scikit-learn library. The number of discretizing bins is fixed at 4 and the width of bins is equal to capture the marginal distribution of each variable in the observational data. As for the categorical data models, this module is deactivated.

Hardware. All experiments are run on a machine with a 64-core Intel(R) Xeon(R) Gold 6242 CPU @ 2.80GHz. The running of GPU-based baselines is conducted on an NVIDIA GeForce RTX 4090.