
Less Greedy Equivalence Search

Adiba Ejaz Elias Bareinboim
Causal Artificial Intelligence Lab
Columbia University
{adiba.ejaz, eb}@cs.columbia.edu

Abstract

Greedy Equivalence Search (GES) is a classic score-based algorithm for causal discovery from observational data. In the sample limit, it recovers the Markov equivalence class of graphs that describe the data. Still, it faces two challenges in practice: computational cost and finite-sample accuracy. In this paper, we develop Less Greedy Equivalence Search (LGES), a variant of GES that retains its theoretical guarantees while partially addressing these limitations. LGES modifies the greedy step; rather than always applying the highest-scoring insertion, it avoids edge insertions between variables for which the score implies some conditional independence. This more targeted search yields up to a 10-fold speed-up and a substantial reduction in structural error relative to GES. Moreover, LGES can guide the search using prior knowledge, and can correct this knowledge when contradicted by data. Finally, LGES can use interventional data to refine the learned observational equivalence class. We prove that LGES recovers the true equivalence class in the sample limit, even with misspecified knowledge. Experiments demonstrate that LGES outperforms GES and other baselines in speed, accuracy, and robustness to misspecified knowledge. Our code is available at <https://github.com/CausalAILab/lges>.

1 Introduction

Causal discovery, the task of learning causal structure from data, is a core problem in the field of causality [54]. The causal structure may be an end in itself to the scientist, or a prerequisite for downstream tasks such as inference, decision-making, and generalization [2, 42]. Causal discovery algorithms have been used in a range of disciplines that span biology, medicine, climate science, and neuroscience, among others [19, 43, 47, 48].

A hallmark algorithm in the field is Greedy Equivalence Search (GES) [9, 35], which takes as input observational data and finds a *Markov equivalence class* (MEC) of causal graphs that describe the data. In general, the true graph is not uniquely identifiable from observational data, and the MEC is the most informative structure that can be learned. Under standard assumptions in causal discovery, GES is guaranteed to recover the true MEC in the sample limit. In contrast, many causal discovery algorithms—including prominent examples such as max-min hill-climbing [56] and NoTears [61]—lack such large-sample guarantees. Many variants of GES have been developed, including faster, parallelized implementations [43], restricted search over bounded in-degree graphs [10], and Greedy Interventional Equivalence Search (GIES) [24]. The last of these can exploit interventional data, though is not asymptotically correct [58].

Despite its attractive features, the GES family faces challenges shared across most causal discovery algorithms. For instance, the problem of causal discovery is NP-hard [11], and GES commonly struggles to scale in high-dimensional settings. Moreover, in finite-sample regimes, GES often fails to recover the true MEC. In other words, applying GES in practice is challenging due to both

computational complexity (scaling) and sample complexity (accuracy) issues. We refer readers to [56] for an extensive empirical study of GES performance.

At a high level, GES searches over the space of MECs by inserting and deleting edges to maximize a score reflecting data fit. At each state, it evaluates a set of neighbors—possibly exponentially many—and moves to the *highest-scoring* neighbor that scores more than the current MEC. It continues the search greedily until no higher-scoring neighbors are found. In the sample limit, this strategy is guaranteed to find the global optimum of the score: the true MEC.

In this paper, we question the basic assumption of the GES family that the greedy choice of the highest-scoring neighbor is the best one. We first introduce *Generalized GES* (Alg. 3), which allows moving to any score-increasing neighbor of the current MEC, not necessarily the highest-scoring one. This relaxed strategy still finds the global optimum of the score in the sample limit (Thm. 1). More importantly, it opens the door to more strategic neighbor selection.

While it may seem that choosing the highest-scoring neighbor would yield the best performance in practice, surprisingly, we show that this is not the case; a careful and less greedy choice improves both accuracy and runtime. Specializing GGES, we develop the algorithm *Less Greedy Equivalence Search* (LGES) (Alg. 1), advancing on GES and its relatives in the following ways:

1. **Faster, more accurate structure learning.** In Sec. 3.2, we introduce two novel operator selection strategies, CONSERVATIVEINSERT and SAFEINSERT, which LGES exploits for neighbor selection. Empirically, these procedures yield up to a 10-fold reduction in runtime and 2-fold reduction in structural error relative to GES (Experiment 5.1) and other baselines. LGES with SAFEINSERT asymptotically recovers the true MEC (Prop. 2, Cor. 1).
2. **Robustness to misspecified prior knowledge.** In Sec. 3.3, we propose a method whereby LGES can guide neighbor selection based on prior knowledge (Alg. 6) while remaining asymptotically correct even if the knowledge is misspecified. Accurate knowledge improves LGES’ accuracy and runtime; inaccurate knowledge harms LGES significantly less than it does GES initialized with the knowledge (Experiment 5.2).
3. **Refining the learned MEC with interventional data.** In Sec. 4, we develop \mathcal{I} -ORIENT (Alg. 2, Thm. 2), a score-based procedure that LGES (or any structure learning algorithm) can use to refine an observational MEC with interventional data. To our knowledge, this is the first asymptotically correct score-based procedure for learning from interventional data that can scale to graphs with more than a hundred nodes. LGES with \mathcal{I} -ORIENT is 10x faster than GIES [24] while maintaining competitive accuracy (Experiment 5.3).

Proofs for all results are provided in Appendix C. Additional experiments with synthetic data and real-world protein signalling data [48] are provided in Appendix D.

2 Background

Notation. Capital letters denote variables (V), small letters denote their values (v), and bold letters denote sets of variables (\mathbf{V}) and their values (\mathbf{v}). $P(\mathbf{v})$ denotes a probability distribution over a set of variables \mathbf{V} . For disjoint sets of variables $\mathbf{X}, \mathbf{Y}, \mathbf{Z}$, $\mathbf{X} \perp\!\!\!\perp \mathbf{Y} \mid \mathbf{Z}$ denotes that \mathbf{X} and \mathbf{Y} are conditionally independent given \mathbf{Z} and $\mathbf{X} \perp_d \mathbf{Y} \mid \mathbf{Z}$ denotes that \mathbf{X} and \mathbf{Y} are *d-separated* given \mathbf{Z} in the graph in context.

Causal graphs [3, 42]. A causal graph over variables \mathbf{V} is a directed acyclic graph (DAG) with an edge $X \rightarrow Y$ denoting that X is a possible cause of Y . The parents of a variable X in a graph \mathcal{G} , denoted $\text{Pa}_X^{\mathcal{G}}$, are those variables with a directed edge into X . The non-descendants of a variable X in a graph \mathcal{G} , denoted $\text{Nd}_X^{\mathcal{G}}$, are those variables to which there is no directed path from X (excluding X itself). The superscript will be omitted when clear from context. A given distribution $P(\mathbf{v})$ is said to be *Markov* with respect to a DAG \mathcal{G} if for all disjoint sets $\mathbf{X}, \mathbf{Y}, \mathbf{Z} \subseteq \mathbf{V}$, if $\mathbf{X} \perp_d \mathbf{Y} \mid \mathbf{Z}$ in \mathcal{G} , then $\mathbf{X} \perp\!\!\!\perp \mathbf{Y} \mid \mathbf{Z}$ in $P(\mathbf{v})$. If the converse is also true, $P(\mathbf{v})$ is said to be *faithful* to \mathcal{G} . In this work, like GES [9], we assume

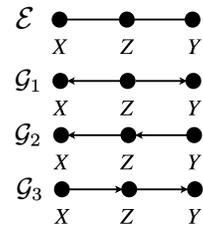


Figure 1: A CPDAG \mathcal{E} and the three DAGs in the MEC it represents, encoding $X \perp_d Y \mid Z$.

that the system of interest is *Markovian*, i.e. it contains no unobserved confounders, and that there exists a DAG \mathcal{G} with respect to which the given $P(\mathbf{v})$ is both Markov and faithful.

Markov equivalence classes [41, 54]. Two causal DAGs \mathcal{G}, \mathcal{H} are said to be Markov equivalent if they encode exactly the same d -separations. The Markov equivalence class (MEC) of a DAG is the set of all graphs that are Markov equivalent to it. A given $P(\mathbf{v})$ may be Markov and faithful to more than one DAG. Hence, the target of causal discovery from observational data is the MEC of DAGs with respect to which $P(\mathbf{v})$ is Markov and faithful. An MEC \mathcal{M} is represented by a unique completed partially directed graph (CPDAG). A CPDAG \mathcal{E} for \mathcal{M} has an undirected edge $X - Y$ if \mathcal{M} contains two DAGs $\mathcal{G}_1, \mathcal{G}_2$ with $X \rightarrow Y$ in \mathcal{G}_1 and $Y \rightarrow X$ in \mathcal{G}_2 . \mathcal{E} has a directed edge $X \rightarrow Y$ if $X \rightarrow Y$ is in every DAG in \mathcal{M} . We frequently refer to an MEC by its representative CPDAG. The adjacencies (neighbours) of a variable X in a CPDAG \mathcal{E} , denoted $\text{Adj}_X^{\mathcal{E}}$ ($\text{Ne}_X^{\mathcal{E}}$), comprise those variables connected by any edge (an undirected edge) to X .

Greedy Equivalence Search [9, 35]. Greedy Equivalence Search (GES) is a score-based algorithm for learning MECs from observational data. It searches for the true MEC by maximizing a scoring criterion given m samples of data $\mathbf{D} \sim P(\mathbf{v})$. For example, a popular choice of scoring criterion is the *Bayesian information criterion* (BIC) [51].

GES assumes that the given scoring criterion is decomposable, consistent, and score-equivalent, so that the score of an MEC is the score of any DAG in that MEC (Defs. A.3, A.4, A.5). BIC satisfies each of these conditions for distributions that are Markov and faithful to some DAG and are curved exponential families, for e.g., linear-Gaussian or multinomial models [9, 21, 23]. Moreover, decomposability and consistency imply local consistency ([9, Lemma 7]), the key property needed for the correctness of GES.

Definition 1 (Locally consistent scoring criterion [9, Def. 6]). Let \mathbf{D} be a dataset consisting of i.i.d. samples from some distribution $P(\mathbf{v})$. Let \mathcal{G} be any DAG, and let \mathcal{G}' be the DAG that results from adding the edge $X \rightarrow Y$ to \mathcal{G} . A scoring criterion S is said to be *locally consistent* if, as the number of samples goes to infinity, the following two properties hold:

1. If $X \not\perp\!\!\!\perp Y \mid \text{Pa}_Y^{\mathcal{G}}$ in $P(\mathbf{v})$ then $S(\mathcal{G}, \mathbf{D}) < S(\mathcal{G}', \mathbf{D})$.
2. If $X \perp\!\!\!\perp Y \mid \text{Pa}_Y^{\mathcal{G}}$ in $P(\mathbf{v})$ then $S(\mathcal{G}, \mathbf{D}) > S(\mathcal{G}', \mathbf{D})$.

Example 1. Consider a distribution $P(\mathbf{v})$ whose true MEC is \mathcal{E} and true DAG is $\mathcal{G}_2 \in \mathcal{E}$ as in Fig. 1. Consider $\mathcal{G}_1 \in \mathcal{E}$ (Fig. 1), and let $\mathcal{G}_1^+ = \mathcal{G}_1 \cup \{X \rightarrow Y\}$ and $\mathcal{G}_1^- = \mathcal{G}_1 \setminus \{Z \rightarrow Y\}$. Since $Y \perp\!\!\!\perp X \mid \text{Pa}_Y^{\mathcal{G}_1^+}$ in $P(\mathbf{v})$, where $\text{Pa}_Y^{\mathcal{G}_1^+} = \{Z\}$, \mathcal{G}_1^+ has a lower score than \mathcal{G}_1 . Since $Y \not\perp\!\!\!\perp Z \mid \text{Pa}_Y^{\mathcal{G}_1^-}$ in $P(\mathbf{v})$, where $\text{Pa}_Y^{\mathcal{G}_1^-} = \emptyset$, \mathcal{G}_1^- has a lower score than \mathcal{G}_1 . \square

Given a scoring criterion satisfying the above conditions and data $\mathbf{D} \sim P(\mathbf{v})$ where $P(\mathbf{v})$ is Markov and faithful to some DAG, GES recovers the true MEC in the sample limit [9, Lemma 10]. The PC algorithm [54], a constraint-based method, has similar asymptotic correctness guarantees, but uses conditional independence (CI) tests instead of a score. PC starts with a fully connected graph and removes edges using CI tests. In contrast, GES starts with a fully disconnected graph and proceeds in two phases. In the forward phase, at each state, GES finds the highest-scoring INSERT operator that results in a score increase, applies it, and repeats until no score-increasing INSERT operator exists. At this point, it has found an MEC \mathcal{E} with respect to which $P(\mathbf{v})$ is Markov.

Definition 2 (INSERT operator, [9, Def. 12]). Given a CPDAG \mathcal{E} , non-adjacent nodes X, Y in \mathcal{E} , and some $\mathbf{T} \subseteq \text{Ne}_Y^{\mathcal{E}} \setminus \text{Adj}_X^{\mathcal{E}}$, the $\text{INSERT}(X, Y, \mathbf{T})$ operator modifies \mathcal{E} by inserting the edge $X \rightarrow Y$ and directing the previously undirected edges $T - Y$ for $T \in \mathbf{T}$ as $T \rightarrow Y$.

Intuitively, an $\text{INSERT}(X, Y, \mathbf{T})$ operator applied to an MEC \mathcal{E} corresponds to choosing a DAG $\mathcal{G} \in \mathcal{E}$ (depending on X, Y , and \mathbf{T}), adding the edge $X \rightarrow Y$ to \mathcal{G} , and computing the MEC of the resulting DAG. Though $P(\mathbf{v})$ is Markov with respect to the MEC \mathcal{E} found in the forward phase, $P(\mathbf{v})$ may not be faithful to \mathcal{E} . In the backward phase, GES starts the search with \mathcal{E} , finds the highest-scoring DELETE(X, Y, \mathbf{H}) operator (Def. A.2) that results in a score increase, applies it, and repeats until no score-increasing DELETE operator exists. At this point, it has found an MEC with respect to which $P(\mathbf{v})$ is both Markov and faithful. An optional turning phase is known to improve performance in practice, but is redundant for asymptotic correctness [24].

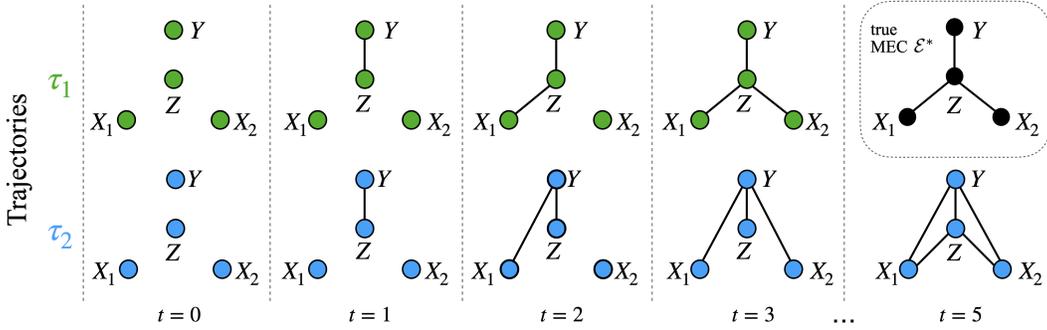


Figure 2: Possible trajectories, τ_1 and τ_2 , that GES may take in the forward phase to obtain an MEC with respect to which a given distribution $P(\mathbf{v})$ is Markov. The true MEC is \mathcal{E}^* (top right). In each trajectory, $\mathcal{E}^{(t+1)}$ results from applying some INSERT operator to $\mathcal{E}^{(t)}$.

3 Less Greedy Equivalence Search

3.1 Generalizing GES

To lay the groundwork for our search strategy, we first introduce Generalized GES (GGES) (Alg. 3), which generalizes GES in three ways. Firstly, GGES allows the search to be initialized from an arbitrary MEC \mathcal{E}_0 , rather than the empty graph. Secondly, GGES is not restricted to the forward-backward-turning phase structure of GES; it can accommodate any order of operators (e.g., deletions before insertions [36]) as specified by the user. Finally, GGES allows the application of *any* valid score-increasing operator, instead of just the highest-scoring one.

At each state \mathcal{E} , GGES calls abstract subroutine GETOPERATOR, which either returns a valid score-increasing operator of the type specified (insertion, deletion, reversal, or either), if one exists, or indicates that there is no such operator. This proceeds until no score-improving operators are found.

Theorem 1 (Correctness of GGES). *Let \mathcal{E} denote the Markov equivalence class that results from GGES (Alg. 3) initialised from an arbitrary MEC \mathcal{E}_0 and let $P(\mathbf{v})$ denote the distribution from which the data \mathbf{D} was generated. Then, as the number of samples goes to infinity, \mathcal{E} is the Markov equivalence class underlying $P(\mathbf{v})$.*

In the next section, we illustrate how GETOPERATOR can be implemented in a way that yields significant improvements in accuracy and runtime relative to GES.

3.2 An improved insertion strategy

In practice, the output of GES is known to include adjacencies between many variables that are non-adjacent in the true MEC [36, 56]. Since these adjacencies are introduced by INSERT operators, this motivates a more careful choice of which INSERT operator to apply. Our approach is grounded in the following observation.

Proposition 1. *Let \mathcal{E} denote an arbitrary CPDAG and let $P(\mathbf{v})$ denote the distribution from which the data \mathbf{D} was generated. Assume, as the number of samples goes to infinity, that there exists a valid score-decreasing $\text{INSERT}(X, Y, \mathbf{T})$ operator for \mathcal{E} . Then, there exists a DAG $\mathcal{G} \in \mathcal{E}$ such that (1) $Y \perp_d X \mid \text{Pa}_Y^{\mathcal{G}}$ and (2) $Y \perp\!\!\!\perp X \mid \text{Pa}_Y^{\mathcal{G}}$ in $P(\mathbf{v})$.*

Then, for variables X and Y , even a single score-decreasing $\text{INSERT}(X, Y, \mathbf{T})$ implies that X and Y are non-adjacent in the true MEC. However, this does not imply that all $\text{INSERT}(X, Y, *)$ are also score-decreasing. GES may thus apply a different $\text{INSERT}(X, Y, \mathbf{T}')$, introducing an adjacency not present in the true MEC. The following example shows how such choices can lead GES to MECs that contain many excess adjacencies.

Example 2. Consider a distribution $P(\mathbf{v})$ over $\mathbf{V} = \{X_1, X_2, Y, Z\}$ whose true MEC is given by \mathcal{E}^* in Fig. 2 (top right). GES starts with the empty graph and successively applies the highest-scoring INSERT operator that it finds. Trajectories τ_1 and τ_2 agree until time $t = 1$. Let $\mathcal{E}^{(1)}$ denote the CPDAG common to τ_1 and τ_2 at $t = 1$. At $t = 1$, GES has many INSERT operators it could apply to

$\mathcal{E}^{(1)}$. Recall that each $\text{INSERT}(\alpha, \beta, \mathbf{T})$ applied to $\mathcal{E}^{(1)}$ corresponds to choosing some DAG \mathcal{G} from $\mathcal{E}^{(1)}$ and adding $\alpha \rightarrow \beta$ to it. The DAG \mathcal{G} is chosen such that for edges $\gamma - \beta$ in $\mathcal{E}^{(1)}$ where α and γ are non-adjacent, \mathcal{G} contains $\gamma \rightarrow \beta$ if $\gamma \in \mathbf{T}$ and $\beta \rightarrow \gamma$ otherwise.

1. $\alpha = X_1, \beta = Z, \mathbf{T} = \emptyset$. This corresponds to choosing $\mathcal{G}_1 \in \mathcal{E}^{(1)}$ (which already has $Z \rightarrow Y$) and adding $X_1 \rightarrow Z$ to it (Fig. 3, left). Since $Z \not\perp\!\!\!\perp X_1 \mid \text{Pa}_Z^{\mathcal{G}_1}$, this edge addition increases the score of \mathcal{G}_1 (by local consistency, Def. 1) and hence of $\mathcal{E}^{(1)}$. This operator is chosen in trajectory τ_1 .
2. $\alpha = X_1, \beta = Y, \mathbf{T} = \{Z\}$. This corresponds to choosing $\mathcal{G}_1 \in \mathcal{E}^{(1)}$ (which already has $Z \rightarrow Y$) and adding $X_1 \rightarrow Y$ to it (Fig. 3, middle). Since $Y \perp\!\!\!\perp X_1 \mid \text{Pa}_Y^{\mathcal{G}_1}$, this edge addition decreases the score of \mathcal{G}_1 and hence of $\mathcal{E}^{(1)}$. This operator is never chosen.
3. $\alpha = X_1, \beta = Y, \mathbf{T} = \emptyset$. This corresponds to choosing $\mathcal{G}_2 \in \mathcal{E}^{(1)}$ (which already has $Y \rightarrow Z$) and adding $X_1 \rightarrow Y$ to it (Fig. 3, right). Since $Y \not\perp\!\!\!\perp X_1 \mid \text{Pa}_Y^{\mathcal{G}_2}$, this edge addition increases the score of \mathcal{G}_2 and hence of $\mathcal{E}^{(1)}$. This operator is chosen in trajectory τ_2 .

In the sample limit, it is unknown whether \mathcal{G}_A or \mathcal{G}_C would score higher. For an extended discussion, see Ex. B.1. ¹ □

The above example shows that GES may insert an edge between non-adjacent variables even in simple settings with a small number of variables. Such choices accumulate in higher-dimensional settings. This motivates avoiding edge insertions for variable pairs (X, Y) for which a score-decreasing INSERT is observed. We hypothesize this has two benefits: (1) *accuracy*: it avoids inserting excess adjacencies that the backward phase may fail to remove, and (2) *efficiency*: it stops the enumeration of (X, Y) insertions when a lower-scoring one is found; moreover, reducing excess adjacencies reduces the number of operators that need to be evaluated in subsequent states.

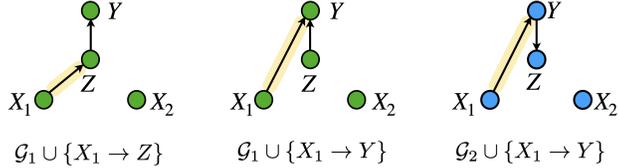


Figure 3: Illustration of some INSERT operators that may be applied to the MEC $\mathcal{E}^{(1)}$ at $t = 1$ in Fig. 2. These operators correspond to various edge additions to the DAGs $\mathcal{G}_1, \mathcal{G}_2 \in \mathcal{E}^{(1)}$, where \mathcal{G}_1 orients $Z - Y$ as $Z \rightarrow Y$ and \mathcal{G}_2 orients $Z - Y$ as $Y \rightarrow Z$.

We now formalize two strategies for avoiding such insertions.

Strategy 1 (CONSERVATIVEINSERT). At a given state with CPDAG \mathcal{E} , for each non-adjacent pair (X, Y) , iterate over valid $\text{INSERT}(X, Y, \mathbf{T})$. If any score-decreasing \mathbf{T} is found, stop, discard all $\text{INSERT}(X, Y, *)$ operators and continue to the next pair. Among all retained candidates, select the highest-scoring operator that results in a score increase, if any.

CONSERVATIVEINSERT avoids inserting edges between any variables (X, Y) for which some conditional independence has been found, as evidenced by a score-decreasing INSERT (Prop. 1). While intuitive, it is unknown if this strategy is guaranteed to find a score-increasing INSERT whenever one exists. To elaborate, the soundness of CONSERVATIVEINSERT rests on the following premise: if $P(\mathbf{v})$ is not Markov with respect to an MEC \mathcal{E} , then there must exist variables X, Y such that for every $\mathcal{G} \in \mathcal{E}$, $X \not\perp\!\!\!\perp Y \mid \text{Pa}_Y^{\mathcal{G}}$. Importantly, the choice of X, Y can not depend on \mathcal{G} . This is the main challenge in proving the soundness of CONSERVATIVEINSERT. Still, we provide partial guarantees in Prop. C.1, C.2, but leave the soundness of CONSERVATIVEINSERT open.

Furthermore, we introduce SAFEINSERT, a relaxation of CONSERVATIVEINSERT that is guaranteed to find a score-increasing INSERT when one exists. The soundness of SAFEINSERT only requires that

¹We can also ask, which of these INSERT operators scores the highest in practice? We generated 100 linear-Gaussian datasets of 100 samples each according to a fixed true DAG in \mathcal{E}^* , following the set-up in Sec. 5.1. Then, we computed the scores of $\mathcal{G}_A : \mathcal{G}_1 \cup \{X_1 \rightarrow Z\}$, $\mathcal{G}_B : \mathcal{G}_1 \cup \{X_1 \rightarrow Y\}$, and $\mathcal{G}_C : \mathcal{G}_2 \cup \{X_1 \rightarrow Y\}$ on each dataset. From the fact that \mathcal{G}_A is closer to the true MEC than \mathcal{G}_C , it may seem that \mathcal{G}_A would always score higher. However, \mathcal{G}_A was the highest-scoring DAG 96% of the time, and \mathcal{G}_C 4% of the time. As expected, \mathcal{G}_B is never the highest-scoring DAG.

Algorithm 1: Less Greedy Equivalence Search (LGES)

Input: Data $\mathbf{D} \sim P(\mathbf{v})$, scoring criterion S , prior assumptions $\mathbf{S} = \langle \mathbf{R}, \mathbf{F} \rangle$, initial MEC \mathcal{E}_0 , insertion strategy $GetInsert$ in $\{\text{GETSAFEINSERT}, \text{GETCONSERVATIVEINSERT}\}$

Output: MEC \mathcal{E} of $P(\mathbf{v})$

```
1  $\mathcal{E} \leftarrow \mathcal{E}_0$ ; // allows initialisation if preferred by user
2 repeat
3   repeat
4      $\mathcal{E} \leftarrow \mathcal{E} +$  the highest-scoring DELETE( $X, Y, \mathbf{T}$ )
5   until no score-increasing deletions exist;
6   repeat
7      $\mathcal{E} \leftarrow \mathcal{E} +$  the highest-scoring TURN( $X, Y, \mathbf{T}$ )
8   until no score-increasing reversals exist;
9    $\mathcal{G} \leftarrow$  some DAG in  $\mathcal{E}$ ;
10   $priorityList \leftarrow \text{GETPRIORITYINSERTS}(\mathcal{E}, \mathcal{G}, \mathbf{S})$ ;
11  foreach candidates in  $priorityList$  do
12     $(X_{max}, Y_{max}, \mathbf{T}_{max}) \leftarrow GetInsert(\mathcal{E}, \mathcal{G}, \mathbf{D}, candidates, S)$ ;
13    if  $(X_{max}, Y_{max}, \mathbf{T}_{max})$  is found then
14       $\mathcal{E} \leftarrow \mathcal{E} + \text{INSERT}(X_{max}, Y_{max}, \mathbf{T}_{max})$ ;
15      break; // no need to check lower priority
16 until no score-increasing operators exist;
17 return  $\mathcal{E}$ 
```

if $P(\mathbf{v})$ is not Markov with respect to \mathcal{E} , then for every $\mathcal{G} \in \mathcal{E}$, there must exist variables X, Y such that $X \not\perp\!\!\!\perp Y \mid \mathbf{Pa}_{\mathcal{G}}^X$. This is unlike CONSERVATIVEINSERT, where X, Y can not depend on \mathcal{G} .

Strategy 2 (SAFEINSERT). At a given state with CPDAG \mathcal{E} , pick an arbitrary DAG $\mathcal{G} \in \mathcal{E}$. For each non-adjacent pair (X, Y) in \mathcal{G} , check if \mathcal{G} has a higher score than $\mathcal{G} \cup \{X \rightarrow Y\}$. If so, discard all $\text{INSERT}(X, Y, *)$ operators and continue to the next pair. Among all retained candidates, select the highest-scoring operator that results in a score increase, if any.

Proposition 2 (Correctness of SAFEINSERT). Let \mathcal{E} denote a Markov equivalence class and let $P(\mathbf{v})$ denote the distribution from which the data \mathbf{D} was generated. Then, as the number of samples goes to infinity, SAFEINSERT returns a valid score-increasing INSERT operator if and only if one exists.

Example 3. (Ex. 2 continued). Let $\mathcal{E}^{(1)}, \mathcal{G}_1$, and \mathcal{G}_2 be as in Ex. 2. Assume GES is at $\mathcal{E}^{(1)}$ and SAFEINSERT picks the DAG $\mathcal{G}_1 \in \mathcal{E}$. Then, $\mathcal{G}_1 \cup \{X_1 \rightarrow Y\}$ has a lower score than \mathcal{G}_1 since $X_1 \perp\!\!\!\perp Y \mid \mathbf{Pa}_{\mathcal{G}_1}^{X_1}$ in $P(\mathbf{v})$, where $\mathbf{Pa}_{\mathcal{G}_1}^{X_1} = \{Z\}$. SAFEINSERT thus does not consider any $\text{INSERT}(X_1, Y, *)$ operators. In contrast, assume SAFEINSERT picks the DAG $\mathcal{G}_2 \in \mathcal{E}$. Then, $\mathcal{G}_2 \cup \{X_1 \rightarrow Y\}$ has a higher score than \mathcal{G}_2 , and SAFEINSERT may still consider $\text{INSERT}(X_1, Y, *)$ operators. However, CONSERVATIVEINSERT will not consider any $\text{INSERT}(X_1, Y, *)$ operators, since $\text{INSERT}(X_1, Y, \{Z\})$, corresponding to $\mathcal{G}_1 \cup \{X_1 \rightarrow Y\}$, results in a lower score than $\mathcal{E}^{(1)}$. \square

Pseudocode for the above insertion strategies are in Algs. 4 and 5. Later, in Sec. 5.1, we compare these strategies, and show how both achieve substantial gains in accuracy and runtime over GES.

3.3 Learning with prior knowledge

In this section, we present another modification of the forward phase of GES: prioritizing edge insertions based on an expert’s prior causal knowledge. The insight that underpins GGES—that we can apply *any* score-increasing insertion—suggests a new way to incorporate such assumptions while still correcting them if contradicted by the data. We assume that we are given a possibly misspecified causal model as a set of required and forbidden edges $\mathbf{S} = \langle \mathbf{R}, \mathbf{F} \rangle$ that may be directed or undirected.

Initialization from prior assumptions. A natural strategy, which we refer to as GES-INIT, initializes the search to an MEC consistent with the assumptions, for e.g., by including all edges in \mathbf{R} , and then proceeds greedily as in standard GES.² This approach, an instantiation of GGES, is sound in the large-sample limit (as a corollary of Thm. 1), even when the assumptions are misspecified.

²This was empirically evaluated in [14]. However, its correctness was not considered.

However, given finite samples, such initialisation may harm both accuracy and runtime. If the expert suggests adjacencies that don't exist in the true MEC, GES-INIT includes them by default in the initialisation and may fail to remove them later. Moreover, such initialisation precludes the use of insertion strategies from Sec. 3.2 that would avoid introducing such excess adjacencies.

Guided search from prior assumptions. We instead propose a strategy that uses prior assumptions to *prioritize* operators, and not to initialize the search. Specifically, for each non-adjacent pair (X, Y) , we rank it into one of four categories based on the constraint set $\mathbf{S} = \langle \mathbf{R}, \mathbf{F} \rangle$ using the procedure GETPRIORITYINSERTS (Alg. 6). Insertions for higher-priority adjacencies are considered first, but only applied if they increase the score. For example, if SAFEINSERT finds no score-increasing insertions for the current MEC, then the remaining expert-provided edges in \mathbf{R} (if any) are not consistent with the data, and will not be inserted. In contrast, GES-INIT inserts all edges by default.

Next, in Sec. 3.4, we incorporate this prioritization scheme into a novel algorithm, combining it with the search strategy of Sec. 3.2 to enable a less greedy search. In Sec. 5.2, we empirically demonstrate the benefit of this prioritization-based strategy.

3.4 The Less Greedy Equivalence Search algorithm

Finally, we introduce the main result of this work: the algorithm Less Greedy Equivalence Search (LGES, Alg. 1). We present three variants:

1. **LGES-0** (Alg. 8), which modifies the insertion step of GES based on our insights in the previous sections, while using the same search strategy as GES in the deletion step,
2. **LGES** (Alg. 1), which is similar to (1) but additionally incorporates the XGES-0 heuristic of prioritizing insertions before deletions [36], and
3. **LGES+** (Alg. 9), which is similar to (2) but additionally incorporates the XGES heuristic of forcing edge deletions and restarting the search [36].³

As a corollary of Thm. 1 and Prop. 2, we can show that each LGES variant with SAFEINSERT recovers the true MEC in the sample limit, even given a misspecified set of prior assumptions.

Corollary 1 (Correctness of LGES). *Let \mathcal{E} denote the Markov equivalence class that results from LGES (Alg. 1) initialized from an arbitrary MEC \mathcal{E}_0 and given prior assumptions $\mathbf{S} = \langle \mathbf{R}, \mathbf{F} \rangle$, and let $P(\mathbf{v})$ denote the distribution from which the data \mathbf{D} was generated. Then, as the number of samples goes to infinity, \mathcal{E} is the Markov equivalence class underlying $P(\mathbf{v})$.*

Remark 1. While we only show that LGES with SAFEINSERT is asymptotically correct, LGES can also be run with CONSERVATIVEINSERT. Since we only have partial guarantees on CONSERVATIVEINSERT (Prop. C.1, C.2), it remains open whether this variant of LGES is asymptotically correct.

4 Score-based learning from interventional data

While the previous sections address causal discovery from observational data, such data alone leaves many edges unoriented. Interventional data can resolve such ambiguities. In this section, we extend the LGES method with \mathcal{I} -ORIENT, a score-based procedure for refining an observational MEC with interventional data. Unlike existing score-based methods, which are asymptotically inconsistent or computationally infeasible even in moderate-dimensional settings [24, 58], our approach scales while preserving soundness.

Background on interventions. Following [24], we assume soft unconditional interventions, including hard (do) interventions as a special case. These set the distribution of a variable X to some fixed $P^*(x)$, thereby removing the influence of its parents. Let \mathcal{I} denote a family of interventional targets, i.e., subsets $\mathbf{I} \subseteq \mathbf{V}$, with the empty intervention $\theta \in \mathcal{I}$ producing the observational distribution. We observe data from distributions $(P_{\mathbf{I}}(\mathbf{v}))_{\mathbf{I} \in \mathcal{I}}$. As in the observational case, we assume there exists a DAG \mathcal{G} such that these distributions are \mathcal{I} -Markov (Def. A.7) and faithful to the corresponding intervention graphs $(\mathcal{G}_{\mathbf{I}})_{\mathbf{I} \in \mathcal{I}}$, obtained by removing edges into any intervened variable $V \in \mathbf{I}$ [16, 42].

³Pseudocode and correctness results for LGES-0 and LGES+ can be found in Appendix C.

Algorithm 2: \mathcal{I} -ORIENT

Input: Intervention targets \mathcal{I} , data $(\mathbf{D}_{\mathbf{I}})_{\mathbf{I} \in \mathcal{I}} \sim (\mathbf{P}_{\mathbf{I}}(\mathbf{v}))_{\mathbf{I} \in \mathcal{I}}$, observational MEC \mathcal{E} , scoring criterion S

Output: \mathcal{I} -MEC \mathcal{E} of $(\mathbf{P}_{\mathbf{I}}(\mathbf{v}))_{\mathbf{I} \in \mathcal{I}}$

```
1 foreach  $X \in \mathcal{E}$  and  $Y \in ne_X^{\mathcal{E}}$  do
2    $\Delta S \leftarrow \sum_{\mathbf{I} \in \mathcal{I}, X \in \mathbf{I}, Y \notin \mathbf{I}} s_{\mathbf{D}_{\mathbf{I}}}(y, x) - s_{\mathbf{D}_{\mathbf{I}}}(y)$ ;
3   if  $\Delta S > 0$  then
4     Orient edge  $X - Y$  as  $X \rightarrow Y$  in  $\mathcal{E}$ ;
5     Apply Meek's rules in  $\mathcal{E}$  to propagate orientations [35];
6   else if  $\Delta S < 0$  then
7     Orient edge  $X - Y$  as  $X \leftarrow Y$  in  $\mathcal{E}$ ;
8     Apply Meek's rules in  $\mathcal{E}$  to propagate orientations [35];
9 return  $\mathcal{E}$ 
```

Just as observational data identifies an MEC, interventional data identifies an \mathcal{I} -MEC, a smaller equivalence class encoding constraints on both observational and interventional data [24, Def. 7].

Orientation procedure. To recover the \mathcal{I} -MEC, we introduce \mathcal{I} -ORIENT (Alg. 2), which orients undirected edges in the observational MEC using scores from interventional data. The underlying idea is simple. Consider some undirected edge $X - Y$ in the observational MEC \mathcal{E} , where the ground truth DAG is \mathcal{G}^* . Say there exists an intervention \mathbf{I} containing X but not Y . When we perform an unconditional intervention on X , we erase the causal influence of $\text{Pa}_X^{\mathcal{G}}$ on X . So, Y is a parent of X in \mathcal{G}^* if and only if X and Y are marginally independent in $P_{\mathbf{I}}(\mathbf{v})$. How can we check for marginal independence using a scoring criterion? The key is to use local consistency (Def. 1): we compare a graph \mathcal{G} where Y has no parents to a graph $\mathcal{G} \cup \{X \rightarrow Y\}$. The former will score higher than the latter if and only if $X \perp\!\!\!\perp Y$ in $P_{\mathbf{I}}(\mathbf{v})$. This is precisely the test in line 2 of \mathcal{I} -ORIENT (Alg. 2).

Theorem 2 (Correctness of \mathcal{I} -ORIENT). *Let \mathcal{E} denote the Markov equivalence class that results from \mathcal{I} -ORIENT (Alg. 2) given an observational MEC \mathcal{E}_0 and interventional targets \mathcal{I} , and let $(\mathbf{P}_{\mathbf{I}}(\mathbf{v}))_{\mathbf{I} \in \mathcal{I}}$ denote the family of distributions from which the data $(\mathbf{D}_{\mathbf{I}})_{\mathbf{I} \in \mathcal{I}}$ was generated. Assume that \mathcal{E}_0 is the MEC underlying $P_{\emptyset}(\mathbf{v})$. Then, as the number of samples goes to infinity for each $\mathbf{I} \in \mathcal{I}$, \mathcal{E} is the \mathcal{I} -Markov equivalence class underlying $(\mathbf{P}_{\mathbf{I}}(\mathbf{v}))_{\mathbf{I} \in \mathcal{I}}$.*

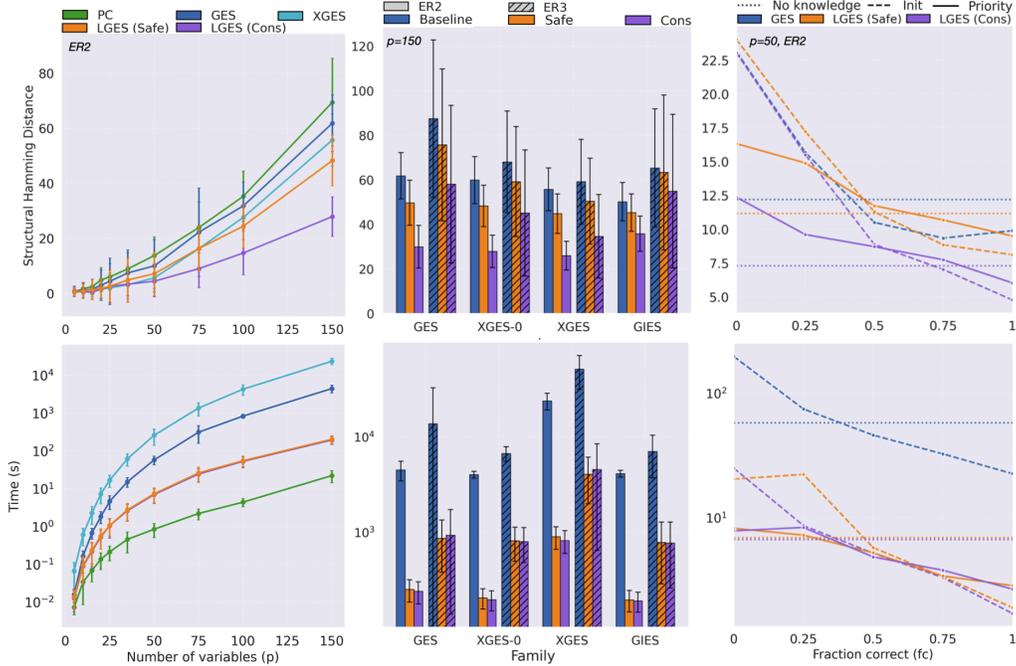
5 Experiments

5.1 Learning from observational data

Synthetic data and baselines. We draw Erdős–Rényi graphs with p variables and $\{1, 2, 3\} \cdot p$ edges in expectation, denoted ER- $\{1, 2, 3\}$ respectively). We run most experiments for p up to 150, with additional experiments for p up to 250 in Sec. D.1.6. For each p , we sample 50 graphs and generate linear-Gaussian data for each graph. Following [37], we draw weights from $\mathcal{U}([-2, -0.5] \cup [0.5, 2])$ and noise variances from $\mathcal{U}([0.1, 0.5])$. We obtain samples of size $n \in \{10^3, 10^4\}$ via `sampler` [20]. We evaluate GES, XGES-0 [36], XGES [36], LGES-0, LGES, LGES+, PC [54], and NoTears [61].⁴

Results. LGES (both Safe and Conservative) significantly outperforms XGES and GES in runtime and accuracy. We measure accuracy by Structural Hamming Distance (SHD) [56] between the estimated and true CPDAGs (Fig. 4a), as well as precision, recall, and F1 score (Fig. D.1.2). Furthermore, the less greedy variants of GES, XGES-0, and XGES all improve on their respective original algorithms (Figs. 4b, D.1.1, D.1.3). Of the GES variants, LGES+ (Conservative) is the most accurate, whereas LGES (Conservative) is the fastest. CONSERVATIVEINSERT outperforms SAFEINSERT in accuracy, but both strategies yield similar runtime.

⁴All implementations in Python. GES variants are implemented by modifying the code in <https://github.com/juangamella/ges>. We use the PC implementation in `causal-learn` [62] and NoTears implementation in `causal-nex` [5]. Our GES implementations share as much code as possible. We do not use the optimized implementation of the operators proposed in [36] across any of our GES variants, to study the effect of the search strategy independent of implementation.



(a) Observational data performance (b) Impact of less greedy insertion (c) Impact of knowledge

Figure 4: Performance of algorithms on 50 simulated datasets from Erdős–Rényi graphs with p variables. **Lower is better** (more accurate / faster) across all plots. **(a)** LGES outperforms baselines in accuracy and runtime on graphs with $2p$ edges in expectation, given $n = 10^4$ observational samples and no prior knowledge. **(b)** Less greedy insertion improves several GES variants on graphs with $p = 150$ variables and $2p$ and $3p$ edges in expectation, given $n = 10^4$ observational samples (and $n = 10^3$ samples per intervention for GIES) and no prior knowledge. LGES-0, LGES+, and LGIES are the less greedy variants of GES, XGES-0, XGES, and GIES respectively. **(c)** Given prior knowledge in the form of $3m/4$ required edges when the true graph contains m edges, LGES’ prioritization strategy is more robust to misspecification in the knowledge than initialization with the same knowledge, given $n = 10^3$ observational samples. See Sec. D for additional results.

To elaborate, LGES (Safe and Conservative) is an order of magnitude faster than GES. LGES (Conservative) is up to 2 times more accurate than GES, for instance, resulting in only ≈ 30 incorrect edges on average in graphs with 150 variables and 300 edges in expectation. These comparisons also hold for LGES-0, the less greedy variant of GES (Figs. 4b, D.1.1, D.1.3). The difference in accuracy is due to excess adjacencies and incorrect orientations; missing adjacencies almost never occur (Fig. D.1.1). PC, though the fastest algorithm, is less accurate than GES for $n = 10^4$, but performs best in the case of many variables and few samples (Fig. D.1.4). NoTears has much worse accuracy than other methods, e.g., average SHD ≈ 125 on graphs with 100 variables, though its runtime appears to scale better (Fig. D.1.1). See Sec. D.1 for additional results and discussion.

5.2 Learning with prior knowledge

Synthetic data and baselines. We study how correctness of prior knowledge affects performance when data is limited ($n = 10^3$) on ER2 graphs with 50 variables, with data generated as in Sec. 5.1. For a true DAG \mathcal{G} with m edges, we generate prior assumptions on $m' \in \{m/2, 3m/4\}$ required edges as follows. We vary the fraction fc of the chosen m' edges that is ‘correct’, with $c \cdot m'$ edges chosen correctly from those in \mathcal{G} and the remaining chosen incorrectly from those not in \mathcal{G} . We compare GES (with and without initialisation) and LGES (without initialisation, with initialisation only, and with priority insertion).

Results. Fig. 4c summarizes results for $m' = 3m/4$, with additional results in Sec. D.2. LGES (Conservative) outperforms GES across all levels of prior correctness in terms of time and SHD.

First, we compare initialization with prioritization. When the knowledge is majority incorrect ($fc \in \{0.5, 0.25, 0.0\}$), the prioritization strategy significantly outperforms initialization in runtime and accuracy. When the knowledge is mostly accurate ($fc \in \{0.75, 1\}$), initialization yields marginally better accuracy than prioritization. Second, we consider when knowledge improves performance. When the knowledge is largely correct ($fc \geq 0.75$), it marginally improves the accuracy of LGES with prioritization, and more visibly improves that of LGES with initialization. However, as long as $fc \geq 0.5$, knowledge significantly improves the runtime of LGES with prioritization. Thus, our prioritization strategy (Sec. 3.3) can leverage knowledge while being robust to misspecification.

5.3 Learning from interventional data

Synthetic data and baselines. We follow a similar set-up as Sec. 5.1 with 10^4 observational samples and ER- $\{2, 3\}$ graphs. For a graph on p variables, we randomly construct $|\mathcal{I}| = p/10$ interventions and generate 10^3 samples for each. We compare GIES [24]; LGES run on observational data followed by \mathcal{I} -ORIENT; LGIES-0; and LGIES. LGIES-0 and LGIES incorporate less greedy insertion into GIES; the latter also prioritizes deletions before insertions.

Results. All less-greedy algorithms are up to $10\times$ faster than GIES (Figs. 4b, 5), with LGIES (Conservative) being the fastest. LGIES (Conservative) is up to $1.5\times$ more accurate than GIES. In general, all less greedy algorithms are more accurate than GIES with the exception of LGES (Safe) + \mathcal{I} -ORIENT, which has competitive accuracy with GIES on ER-3 graphs but performs slightly worse on ER-2 graphs. While being the only combination known to be asymptotically correct, LGES (Safe) + \mathcal{I} -ORIENT is limited by the fact that it uses only observational data to learn the observational MEC, while GIES and LGIES additionally use interventional data to this end.

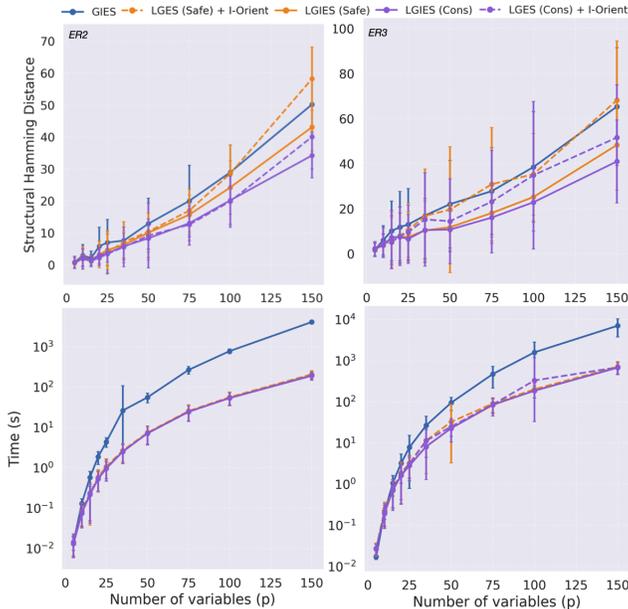


Figure 5: Performance of algorithms on 50 simulated datasets from Erdős-Rényi graphs with p variables and **(left)** $2p$ edges **(right)** $3p$ edges in expectation, given $n = 10^4$ observational samples and $n = 10^3$ samples per intervention (Sec. D.3). **Lower is better** (more accurate / faster) across all plots. LGIES significantly outperforms GIES.

6 Conclusions

In this paper, we introduced a family of less greedy algorithms for score-based causal discovery from observational and interventional data. Our core insight lies in two new operator selection strategies, CONSERVATIVEINSERT and SAFEINSERT, that avoid inserting edges between variables for which the score implies a conditional independence. Building on these ideas, we developed LGES (Alg. 1), which replaces GES’ strictly greedy step with a more careful search. We proved that LGES with SAFEINSERT asymptotically recovers the true Markov equivalence class (Thm. 1, Cor. 1), and showed that LGES can incorporate prior knowledge while remaining robust to misspecification in the knowledge. To extend these advancements beyond purely observational settings, we developed \mathcal{I} -ORIENT (Alg. 2, Thm. 2), a theoretically sound and scalable score-based method for refining an observational MEC using interventional data. Across experiments on random graphs of varying sizes and densities, our less greedy strategies consistently improved both the accuracy and the efficiency existing GES-style algorithms (Sec. 5). The moral, if there is one, is simple: even in causal discovery, temperance is a virtue.

Acknowledgements

This research is supported in part by the NSF, ONR, AFOSR, DoE, Amazon, JP Morgan, and The Alfred P. Sloan Foundation. We thank Jonghwan Kim and the anonymous reviewers for their insightful comments.

References

- [1] Nicos Angelopoulos and James Cussens. Bayesian learning of Bayesian networks with informative priors. *Annals of Mathematics and Artificial Intelligence*, 54(1):53–98, November 2008. ISSN 1573-7470. doi: 10.1007/s10472-009-9133-x. URL <https://doi.org/10.1007/s10472-009-9133-x>.
- [2] E. Bareinboim and J. Pearl. Causal inference and the data-fusion problem. *Proceedings of the National Academy of Sciences*, 113:7345–7352, 2016.
- [3] Elias Bareinboim, Juan D. Correa, Duligur Ibeling, and Thomas Icard. *On Pearl’s Hierarchy and the Foundations of Causal Inference*, page 507–556. Association for Computing Machinery, New York, NY, USA, 1 edition, 2022. ISBN 9781450395861. URL <https://doi.org/10.1145/3501714.3501743>.
- [4] M. J. Bayarri, James O. Berger, Woncheol Jang, Surajit Ray, Luis R. Pericchi, and Ingmar Visser. Prior-based Bayesian information criterion. *Statistical Theory and Related Fields*, 3(1): 2–13, January 2019. ISSN 2475-4269. doi: 10.1080/24754269.2019.1582126. URL <https://doi.org/10.1080/24754269.2019.1582126>. Publisher: Taylor and Francis _eprint: <https://doi.org/10.1080/24754269.2019.1582126>.
- [5] Paul Beaumont, Ben Horsburgh, Philip Pilgerstorfer, Angel Droth, Richard Oentaryo, Steven Ler, Hiep Nguyen, Gabriel Azevedo Ferreira, Zain Patel, and Wesley Leong. CausalNex, October 2021. URL <https://github.com/quantumblacklabs/causalnex>.
- [6] Alexis Bellot, Junzhe Zhang, and Elias Bareinboim. Scores for learning discrete causal graphs with unobserved confounders. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(10):11043–11051, Mar. 2024. doi: 10.1609/aaai.v38i10.28980. URL <https://ojs.aaai.org/index.php/AAAI/article/view/28980>.
- [7] Robert Castelo and Arno Siebes. Priors on network structures. biasing the search for bayesian networks. *International Journal of Approximate Reasoning*, 24(1):39–57, 2000. ISSN 0888-613X. doi: [https://doi.org/10.1016/S0888-613X\(99\)00041-9](https://doi.org/10.1016/S0888-613X(99)00041-9). URL <https://www.sciencedirect.com/science/article/pii/S0888613X99000419>.
- [8] David Maxwell Chickering. A transformational characterization of equivalent bayesian network structures. In *Proceedings of the Eleventh Conference on Uncertainty in Artificial Intelligence, UAI’95*, page 87–98, San Francisco, CA, USA, 1995. Morgan Kaufmann Publishers Inc. ISBN 1558603859.
- [9] David Maxwell Chickering. Optimal structure identification with greedy search. *J. Mach. Learn. Res.*, 3(null):507–554, March 2003. ISSN 1532-4435. doi: 10.1162/153244303321897717. URL <https://doi.org/10.1162/153244303321897717>.
- [10] David Maxwell Chickering and Christopher Meek. Selective greedy equivalence search: finding optimal bayesian networks using a polynomial number of score evaluations. In *Proceedings of the Thirty-First Conference on Uncertainty in Artificial Intelligence, UAI’15*, page 211–219, Arlington, Virginia, USA, 2015. AUAI Press. ISBN 9780996643108.
- [11] David Maxwell Chickering, David Heckerman, and Christopher Meek. Large-sample learning of bayesian networks is np-hard. *J. Mach. Learn. Res.*, 5:1287–1330, December 2004. ISSN 1532-4435.
- [12] Max Chickering. Statistically efficient greedy equivalence search. In Jonas Peters and David Sontag, editors, *Proceedings of the 36th Conference on Uncertainty in Artificial Intelligence (UAI)*, volume 124 of *Proceedings of Machine Learning Research*, pages 241–249. PMLR, 03–06 Aug 2020. URL <https://proceedings.mlr.press/v124/chickering20a.html>.

- [13] Tom Claassen and Ioan G. Bucur. Greedy equivalence search in the presence of latent confounders. In James Cussens and Kun Zhang, editors, *Proceedings of the Thirty-Eighth Conference on Uncertainty in Artificial Intelligence*, volume 180 of *Proceedings of Machine Learning Research*, pages 443–452. PMLR, 01–05 Aug 2022. URL <https://proceedings.mlr.press/v180/claassen22a.html>.
- [14] Anthony C. Constantinou, Zhigao Guo, and Neville K. Kitson. The impact of prior knowledge on causal structure learning. *Knowledge and Information Systems*, 65(8):3385–3434, August 2023. ISSN 0219-3116. doi: 10.1007/s10115-023-01858-x. URL <https://doi.org/10.1007/s10115-023-01858-x>.
- [15] Gregory F. Cooper and Edward Herskovits. A bayesian method for the induction of probabilistic networks from data. *Mach. Learn.*, 9(4):309–347, October 1992. ISSN 0885-6125. doi: 10.1023/A:1022649401552. URL <https://doi.org/10.1023/A:1022649401552>.
- [16] Juan Correa and Elias Bareinboim. A calculus for stochastic interventions:causal effect identification and surrogate experiments. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(06):10093–10100, Apr. 2020. doi: 10.1609/aaai.v34i06.6567. URL <https://ojs.aaai.org/index.php/AAAI/article/view/6567>.
- [17] Atray Dixit, Oren Parnas, Biyu Li, Jenny Chen, Charles P Fulco, Livnat Jerby-Arnon, Nemanja D Marjanovic, Danielle Dionne, Tyler Burks, Raktima Raychowdhury, Britt Adamson, Thomas M Norman, Eric S Lander, Jonathan S Weissman, Nir Friedman, and Aviv Regev. Perturb-seq: Dissecting molecular circuits with scalable single-cell rna profiling of pooled genetic screens. *Cell*, 167(7):1853–1866.e17, December 2016. doi: 10.1016/j.cell.2016.11.038.
- [18] D. Dor and M. Tarsi. A simple algorithm to construct a consistent extension of a partially oriented graph. 1992. URL <https://www.semanticscholar.org/paper/A-simple-algorithm-to-construct-a-consistent-of-a-Dor-Tarsi/3cee18b1965fd94bd98a1e7f2155250276424925>.
- [19] Julien Dubois, Hiroyuki Oya, J. Michael Tyszka, Matthew Howard, Frederick Eberhardt, and Ralph Adolphs. Causal mapping of emotion networks in the human brain: Framework and initial findings. *Neuropsychologia*, 145:106571, 2020. ISSN 0028-3932. doi: <https://doi.org/10.1016/j.neuropsychologia.2017.11.015>. URL <https://www.sciencedirect.com/science/article/pii/S0028393217304281>. The Neural Basis of Emotion.
- [20] Juan L. Gamella, Armeen Taeb, Christina Heinze-Deml, and Peter Bühlmann. Characterization and greedy learning of gaussian structural causal models under unknown interventions. *arXiv preprint arXiv:2211.14897*, 2022.
- [21] Dan Geiger, David Heckerman, Henry King, and Christopher Meek. Stratified exponential families: Graphical models and model selection. *The Annals of Statistics*, 29(2):505–529, 2001. ISSN 00905364, 21688966. URL <http://www.jstor.org/stable/2674112>.
- [22] Uzma Hasan and Md Osman Gani. Optimizing data-driven causal discovery using knowledge-guided search, 2024. URL <https://arxiv.org/abs/2304.05493>.
- [23] Dominique M. A. Haughton. On the choice of a model to fit data from an exponential family. *The Annals of Statistics*, 16(1):342–355, 1988. ISSN 00905364, 21688966. URL <http://www.jstor.org/stable/2241441>.
- [24] Alain Hauser and Peter Bühlmann. Characterization and greedy learning of interventional markov equivalence classes of directed acyclic graphs. *J. Mach. Learn. Res.*, 13(1):2409–2464, August 2012. ISSN 1532-4435.
- [25] Yangbo He, Jinzhu Jia, and Bin Yu. Counting and exploring sizes of markov equivalence classes of directed acyclic graphs. *Journal of Machine Learning Research*, 16(79):2589–2609, 2015. URL <http://jmlr.org/papers/v16/he15a.html>.
- [26] David Heckerman and Dan Geiger. Learning bayesian networks: a unification for discrete and gaussian domains. In *Proceedings of the Eleventh Conference on Uncertainty in Artificial Intelligence*, UAI’95, page 274–284, San Francisco, CA, USA, 1995. Morgan Kaufmann Publishers Inc. ISBN 1558603859.

- [27] A. Jaber, M. Kocaoglu, K. Shanmugam, and E. Bareinboim. Causal discovery from soft interventions with unknown targets: Characterization and learning. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 9551–9561, Vancouver, Canada, Jun 2020. Curran Associates, Inc.
- [28] Hyunchai Jeong, Adiba Ejaz, Jin Tian, and Elias Bareinboim. Testing causal models with hidden variables in polynomial delay via conditional independencies. *Proceedings of the AAAI Conference on Artificial Intelligence*, 39(25):26813–26822, April 2025. ISSN 2159-5399. doi: 10.1609/aaai.v39i25.34885. URL <http://dx.doi.org/10.1609/aaai.v39i25.34885>.
- [29] Murat Kocaoglu, Amin Jaber, Karthikeyan Shanmugam, and Elias Bareinboim. Characterization and learning of causal graphs with latent variables from soft interventions. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper_files/paper/2019/file/c3d96fbd5b1b45096ff04c04038fff5d-Paper.pdf.
- [30] Steffen L Lauritzen, A Philip Dawid, Birgitte N Larsen, and H. G. G Leimer. Independence Properties of Directed Markov Fields. *Networks*, 20(5):491–505, 1990. ISSN 10970037. doi: 10.1002/net.3230200503.
- [31] Thuc Duy Le, Tao Hoang, Jiuyong Li, Lin Liu, Huawen Liu, and Shu Hu. A fast pc algorithm for high dimensional causal discovery with multi-core pcs. *IEEE/ACM Trans. Comput. Biol. Bioinformatics*, 16(5):1483–1495, September 2019. ISSN 1545-5963. doi: 10.1109/TCBB.2016.2591526. URL <https://doi.org/10.1109/TCBB.2016.2591526>.
- [32] Adam Li, Amin Jaber, and Elias Bareinboim. Causal discovery from observational and interventional data across multiple environments. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 16942–16956. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/368cba57d00902c752eaa9e4770bbbbe-Paper-Conference.pdf.
- [33] Stephanie Long, Alexandre Piché, Valentina Zantedeschi, Tibor Schuster, and Alexandre Drouin. Causal discovery with language models as imperfect experts, 2023. URL <https://arxiv.org/abs/2307.02390>.
- [34] Christopher Meek. Causal inference and causal explanation with background knowledge. In *Proceedings of the Eleventh Conference on Uncertainty in Artificial Intelligence, UAI’95*, page 403–410, San Francisco, CA, USA, 1995. Morgan Kaufmann Publishers Inc. ISBN 1558603859.
- [35] Christopher Meek. Graphical Models: Selecting causal and statistical models. 1 1997. doi: 10.1184/R1/22696393.v1. URL https://kilthub.cmu.edu/articles/thesis/Graphical_Models_Selecting_causal_and_statistical_models/22696393.
- [36] Achille Nazaret and David Blei. Extremely greedy equivalence search. In Negar Kiyavash and Joris M. Mooij, editors, *Proceedings of the Fortieth Conference on Uncertainty in Artificial Intelligence*, volume 244 of *Proceedings of Machine Learning Research*, pages 2716–2745. PMLR, 15–19 Jul 2024. URL <https://proceedings.mlr.press/v244/nazaret24a.html>.
- [37] Ignavier Ng, AmirEmad Ghassami, and Kun Zhang. On the role of sparsity and dag constraints for learning linear dags. In *Proceedings of the 34th International Conference on Neural Information Processing Systems, NIPS ’20*, Red Hook, NY, USA, 2020. Curran Associates Inc. ISBN 9781713829546.
- [38] Ignavier Ng, Biwei Huang, and Kun Zhang. Structure learning with continuous optimization: A sober look and beyond. In Francesco Locatello and Vanessa Didelez, editors, *Proceedings of the Third Conference on Causal Learning and Reasoning*, volume 236 of *Proceedings of Machine Learning Research*, pages 71–105. PMLR, 01–03 Apr 2024. URL <https://proceedings.mlr.press/v236/ng24a.html>.

- [39] Chris J. Oates, Jessica Kasza, Julie A. Simpson, and Andrew B. Forbes. Brief report: Repair of partly misspecified causal diagrams. *Epidemiology*, 28(4):pp. 548–552, 2017. ISSN 10443983, 15315487. URL <https://www.jstor.org/stable/26512182>.
- [40] Rodney T. O’Donnell, Lloyd Allison, and Kevin B. Korb. Learning hybrid bayesian networks by mml. In Abdul Sattar and Byeong-ho Kang, editors, *AI 2006: Advances in Artificial Intelligence*, pages 192–203, Berlin, Heidelberg, 2006. Springer Berlin Heidelberg. ISBN 978-3-540-49788-2.
- [41] Judea Pearl. *Probabilistic Reasoning in Intelligent Systems*. Morgan Kaufmann, San Mateo, CA, 1988.
- [42] Judea Pearl. *Causality: Models, Reasoning, and Inference*. Cambridge University Press, New York, 2nd edition, 2009.
- [43] Joseph Ramsey, Madelyn Glymour, Ruben Sanchez-Romero, and Clark Glymour. A million variables and more: the fast greedy equivalence search algorithm for learning high-dimensional graphical causal models, with an application to functional magnetic resonance images. *International Journal of Data Science and Analytics*, 3:121–129, 2017. doi: 10.1007/s41060-016-0032-z. URL <https://doi.org/10.1007/s41060-016-0032-z>.
- [44] Garvesh Raskutti and Caroline Uhler. Learning directed acyclic graph models based on sparsest permutations. *Stat*, 7(1):e183, 2018. doi: <https://doi.org/10.1002/sta4.183>. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/sta4.183>. e183 sta4.183.
- [45] Alexander G. Reisach, Christof Seiler, and Sebastian Weichwald. Beware of the simulated dag! causal discovery benchmarks may be easy to game. In *Proceedings of the 35th International Conference on Neural Information Processing Systems, NIPS ’21*, Red Hook, NY, USA, 2021. Curran Associates Inc. ISBN 9781713845393.
- [46] Thomas Richardson and Peter Spirtes. Ancestral graph Markov models. *The Annals of Statistics*, 30(4):962 – 1030, 2002. doi: 10.1214/aos/1031689015. URL <https://doi.org/10.1214/aos/1031689015>.
- [47] Jakob Runge. Causal network reconstruction from time series: From theoretical assumptions to practical estimation. *Chaos: An Interdisciplinary Journal of Nonlinear Science*, 28(7):075310, 07 2018. ISSN 1054-1500. doi: 10.1063/1.5025050. URL <https://doi.org/10.1063/1.5025050>.
- [48] Karen Sachs, Omar Perez, Dana Pe’er, Douglas A Lauffenburger, and Garry P Nolan. Causal protein-signaling networks derived from multiparameter single-cell data. *Science*, 308(5721): 523–529, 2005.
- [49] Reuben A. Saunders, William E. Allen, Xingjie Pan, Jaspreet Sandhu, Jiaqi Lu, Thomas K. Lau, Karina Smolyar, Zuri A. Sullivan, Catherine Dulac, Jonathan S. Weissman, and Xiaowei Zhuang. A platform for multimodal in vivo pooled genetic screens reveals regulators of liver function. *bioRxiv*, 2024. doi: 10.1101/2024.11.18.624217. URL <https://www.biorxiv.org/content/early/2024/11/21/2024.11.18.624217>.
- [50] Richard Scheines, Peter Spirtes, Clark Glymour, Christopher Meek, and Thomas Richardson. The tetrad project: Constraint based aids to causal model specification. *Multivariate Behavioral Research*, 33(1):65–117, 1998.
- [51] Gideon Schwarz. Estimating the dimension of a model. *The Annals of Statistics*, 6(2):461–464, 1978. ISSN 00905364, 21688966. URL <http://www.jstor.org/stable/2958889>.
- [52] Rajen D. Shah and Jonas Peters. The hardness of conditional independence testing and the generalised covariance measure. *The Annals of Statistics*, 48(3), June 2020. ISSN 0090-5364. doi: 10.1214/19-aos1857. URL <http://dx.doi.org/10.1214/19-AOS1857>.
- [53] Peter Spirtes, Christopher Meek, and Thomas Richardson. Causal inference in the presence of latent variables and selection bias. In *Proceedings of the Eleventh Conference on Uncertainty in Artificial Intelligence, UAI’95*, page 499–506, San Francisco, CA, USA, 1995. Morgan Kaufmann Publishers Inc. ISBN 1558603859.

- [54] Peter Spirtes, Clark N Glymour, and Richard Scheines. *Causation, Prediction, and Search*. MIT Press, Cambridge, MA, 2nd edition, 2000.
- [55] Sofia Triantafillou and Ioannis Tsamardinos. Score based vs constraint based causal learning in the presence of confounders. 2016. URL <http://www.its.caltech.edu/~fehardt/UAI2016WS/papers/Triantafillou.pdf>.
- [56] Ioannis Tsamardinos, Laura E. Brown, and Constantin F. Aliferis. The max-min hill-climbing Bayesian network structure learning algorithm. *Machine Learning*, 65(1):31–78, October 2006. ISSN 1573-0565. doi: 10.1007/s10994-006-6889-7. URL <https://doi.org/10.1007/s10994-006-6889-7>.
- [57] Thomas Verma and Judea Pearl. Equivalence and synthesis of causal models. In *Proceedings of the Sixth Annual Conference on Uncertainty in Artificial Intelligence, UAI '90*, page 255–270, USA, 1990. Elsevier Science Inc. ISBN 0444892648.
- [58] Yuhao Wang, Liam Solus, Karren Dai Yang, and Caroline Uhler. Permutation-based causal inference algorithms with interventions. In *Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS' 17*, page 5824–5833, Red Hook, NY, USA, 2017. Curran Associates Inc. ISBN 9781510860964.
- [59] Marcel Wienöbst, Max Bannach, and Maciej Liśkiewicz. Polynomial-Time Algorithms for Counting and Sampling Markov Equivalent DAGs, December 2020. URL <http://arxiv.org/abs/2012.09679>. arXiv:2012.09679 [cs].
- [60] Jiji Zhang. On the completeness of orientation rules for causal discovery in the presence of latent confounders and selection bias. *Artificial Intelligence*, 172(16):1873–1896, 2008. ISSN 0004-3702. doi: <https://doi.org/10.1016/j.artint.2008.08.001>. URL <https://www.sciencedirect.com/science/article/pii/S0004370208001008>.
- [61] Xun Zheng, Bryon Aragam, Pradeep Ravikumar, and Eric P. Xing. Dags with no tears: continuous optimization for structure learning. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems, NIPS' 18*, page 9492–9503, Red Hook, NY, USA, 2018. Curran Associates Inc.
- [62] 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.

NeurIPS Paper Checklist

The checklist is designed to encourage best practices for responsible machine learning research, addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove the checklist: **The papers not including the checklist will be desk rejected.** The checklist should follow the references and follow the (optional) supplemental material. The checklist does NOT count towards the page limit.

Please read the checklist guidelines carefully for information on how to answer these questions. For each question in the checklist:

- You should answer [Yes], [No], or [NA].
- [NA] means either that the question is Not Applicable for that particular paper or the relevant information is Not Available.
- Please provide a short (1–2 sentence) justification right after your answer (even for NA).

The checklist answers are an integral part of your paper submission. They are visible to the reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it (after eventual revisions) with the final version of your paper, and its final version will be published with the paper.

The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation. While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided a proper justification is given (e.g., "error bars are not reported because it would be too computationally expensive" or "we were unable to find the license for the dataset we used"). In general, answering "[No]" or "[NA]" is not grounds for rejection. While the questions are phrased in a binary way, we acknowledge that the true answer is often more nuanced, so please just use your best judgment and write a justification to elaborate. All supporting evidence can appear either in the main paper or the supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification please point to the section(s) where related material for the question can be found.

IMPORTANT, please:

- **Delete this instruction block, but keep the section heading "NeurIPS Paper Checklist",**
- **Keep the checklist subsection headings, questions/answers and guidelines below.**
- **Do not modify the questions and only use the provided macros for your answers.**

1. Claims

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

Answer: [Yes]

Justification: Yes, we clearly support our theory with assumptions and proofs and our empirical claims with clearly designed experiments.

Guidelines:

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

2. Limitations

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

Answer: [Yes]

Justification: We discuss limitations of our interventional data approach, and demonstrate how our algorithm scales worse than some existing methods.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

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

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

Justification: Clear statements of theoretical results are provided in the main paper. Proofs are provided in the supplement.

Guidelines:

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

4. Experimental result reproducibility

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

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

Justification: Experimental details are provided in the appendix, with the broader experimental design provided in the main paper.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: We are excited to release the code as we want people to use and build on our algorithm. Pseudocode is provided for all algorithms. We will additionally provide anonymized code in the supplementary material.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.

- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

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

Answer: [Yes]

Justification: Yes, these details are outlined in the main paper and provided in more specificity in the appendix.

Guidelines:

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

7. Experiment statistical significance

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

Answer: [Yes]

Justification: Yes, error bars are provided in the main paper; where limited in the main paper for legibility, they are provided in the appendix.

Guidelines:

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

8. Experiments compute resources

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

Answer: [Yes]

Justification: Compute details are provided in the appendix. We only use CPU.

Guidelines:

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

9. Code of ethics

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

Answer: [Yes]

Justification: We have reviewed the code of ethics and do not find concerns related to the three mentioned categories since our work is more in the line of basic research.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

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

Answer: [Yes]

Justification: We present a generic structure learning algorithm without a direct path to particular negative applications. We discuss applications in some problems in biology and neuroscience, though these are still fairly theoretical.

Guidelines:

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

11. Safeguards

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

Answer: [NA]

Justification: We do not use generative models or scraped datasets.

Guidelines:

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

12. Licenses for existing assets

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

Answer: [Yes]

Justification: We credit the implementations we use and build upon.

Guidelines:

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

13. New assets

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

Answer: [Yes]

Justification: Yes, documentation is provided alongside our code which we intend to release publicly.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.

- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

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

Answer: [NA]

Justification: We do not use human subjects in this work.

Guidelines:

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

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

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

Answer: [NA]

Justification: We do not use human subjects in this work.

Guidelines:

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

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigor, or originality of the research, declaration is not required.

Answer: [NA]

Justification: We do not use LLMs as a core method in this work.

Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>) for what should or should not be described.

Appendices

Contents

A Background and related works	23
A.1 Definitions and previous results	23
A.2 Related works	25
A.2.1 Learning from observational data	25
A.2.2 Learning with prior knowledge	26
A.2.3 Learning from interventional data	26
B Discussion and examples	26
B.1 Design of Less Greedy Equivalence Search	26
B.2 Limitations of Less Greedy Equivalence Search	27
B.3 Extended example of GES trajectories	28
C Proofs and pseudocode	28
D Experiments	35
D.1 Learning from observational data	35
D.1.1 Synthetic data details for Experiment 5.1	35
D.1.2 Further baselines and metrics for Experiment 5.1	35
D.1.3 Precision, recall, and F1 score for Experiment 5.1	37
D.1.4 Further experiments with varying edge densities	37
D.1.5 Further experiments with smaller datasets	38
D.1.6 Further experiments with larger graphs	38
D.2 Learning with prior knowledge	39
D.3 Learning from interventional data	39
D.4 Real-world protein signaling data	40
E Frequently Asked Questions	41

A Background and related works

A.1 Definitions and previous results

First, we provide definitions and results used in the main text.

Definition A.1 (*d*-separation [41]). Given a causal DAG \mathcal{G} , a node W on a path π is said to be a collider on π if W has converging arrows into W in π , e.g., $\rightarrow W \leftarrow$ or $\leftrightarrow W \leftarrow$. π is said to be blocked by a set \mathbf{Z} if there exists a node W on π satisfying one of the following two conditions: 1) W is a collider, and neither W nor any of its descendants are in \mathbf{Z} , or 2) W is not a collider, and W is in \mathbf{Z} . Given disjoint sets \mathbf{X} , \mathbf{Y} , and \mathbf{Z} in \mathcal{G} , \mathbf{Z} is said to *d-separate* \mathbf{X} from \mathbf{Y} in \mathcal{G} if \mathbf{Z} blocks every path from a node in \mathbf{X} to a node in \mathbf{Y} according to the *d*-separation criterion.

Definition A.2 (DELETE operator, [9, Def. 13]). For adjacent nodes X, Y in \mathcal{E} connected as either $X \rightarrow Y$ or $X - Y$, and for any $\mathbf{H} \subseteq \text{Ne}_Y^{\mathcal{E}} \cap \text{Adj}_X^{\mathcal{E}}$, the $\text{DELETE}(X, Y, \mathbf{T})$ operator modifies \mathcal{E} by deleting the edge between X and Y , and for each $T \in \mathbf{T}$, directing any undirected edges $X - T$ as $X \rightarrow T$ and any $Y - T$ as $Y \rightarrow T$.

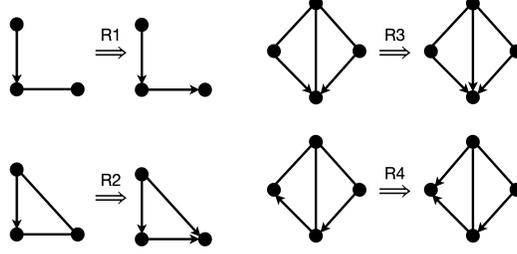


Figure A.1.1: Meek orientation rules for completing partially directed acyclic graphs

The properties and implementation of the TURN operator can be found in [24, Sec. 4.3], though the authors do not provide an exact definition we can reproduce here.

Definition A.3 (Decomposable scoring criterion [9, Sec. 2.3]). Let \mathbf{D} be a set of data consisting of iid samples from some distribution $P(\mathbf{v})$. A scoring criterion S is said to be *decomposable* if it can be written as a sum of measures, each of which is a function of only a single node and its parents, as

$$S(\mathcal{G}, \mathbf{D}) = \sum_{V_i \in \mathbf{V}} s(v_i, pa_i^{\mathcal{G}})$$

Each local score $s(v_i, pa_i^{\mathcal{G}})$ depends only on the values of V_i and \mathbf{Pa}_i in \mathbf{D} .

Definition A.4 (Consistent scoring criterion [9, Def. 5]). Let \mathbf{D} be a set of data consisting of iid samples from some distribution $P(\mathbf{v})$. A scoring criterion S is said to be *consistent* if, as the number of samples goes to infinity, the following two properties hold for any DAGs \mathcal{G}, \mathcal{H} :

1. If $P(\mathbf{v})$ is Markov with respect to \mathcal{G} but not \mathcal{H} , then $S(\mathcal{G}, \mathbf{D}) > S(\mathcal{H}, \mathbf{D})$.
2. If $P(\mathbf{v})$ is Markov with respect to both \mathcal{G} and \mathcal{H} , but \mathcal{G} contains fewer free parameters than \mathcal{H} , then $S(\mathcal{G}, \mathbf{D}) > S(\mathcal{H}, \mathbf{D})$.

Definition A.5 (Score-equivalent scoring criterion [9, Sec 2.3]). Let \mathbf{D} be a set of data consisting of iid samples from some distribution $P(\mathbf{v})$. A scoring criterion S is said to be *score-equivalent* if, as the number of samples goes to infinity, for any two DAGs \mathcal{G}, \mathcal{H} that are Markov equivalent, $S(\mathcal{G}, \mathbf{D}) = S(\mathcal{H}, \mathbf{D})$.

Definition A.6 (Soft unconditional intervention [24, Sec. 2.1]). A soft unconditional intervention on a set of variables \mathbf{X} sets the value of each variable $V_i \in \mathbf{X}$ to an independent random variable U_i from a given set of random variables \mathbf{U} . The resulting distribution is given by

$$P_{\mathbf{X}}(\mathbf{v}) = \prod_{V_i \notin \mathbf{X}} P(v_i | pa_i) \prod_{V_i \in \mathbf{X}} P^*(v_i)$$

where $P^*(v_i)$ denotes the distribution of $U_i \in \mathbf{U}$ corresponding to $V_i \in \mathbf{X}$.

Definition A.7 (\mathcal{I} -Markov property [24, Def. 7]). Let \mathbf{V} be a set of variables, \mathcal{G} a causal DAG over \mathbf{V} , \mathcal{I} a family of interventional targets, and $(\mathbf{P}_{\mathbf{I}}(\mathbf{v}))_{\mathbf{I} \in \mathcal{I}}$ a corresponding family of interventional distributions. We say $(\mathbf{P}_{\mathbf{I}}(\mathbf{v}))_{\mathbf{I} \in \mathcal{I}}$ satisfies the \mathcal{I} -Markov Property of \mathcal{G} if:

1. Each $P_{\mathbf{I}}(\mathbf{v})$ is Markov with respect to the interventional graph $\mathcal{G}_{\mathbf{I}}$, and
2. For interventions $\mathbf{I}, \mathbf{J} \in \mathcal{I}$ and variables $V_i \notin \mathbf{I} \cup \mathbf{J}$, $P_{\mathbf{I}}(v_i | pa_i) = P_{\mathbf{J}}(v_i | pa_i)$.

We let $\mathcal{M}_{\mathcal{I}}(\mathcal{G})$ denote the set of all $(\mathbf{P}_{\mathbf{I}}(\mathbf{v}))_{\mathbf{I} \in \mathcal{I}}$ that are \mathcal{I} -Markov with respect to \mathcal{G} . Two causal DAGs \mathcal{G}, \mathcal{H} are *\mathcal{I} -Markov equivalent* if $\mathcal{M}_{\mathcal{I}}(\mathcal{G}) = \mathcal{M}_{\mathcal{I}}(\mathcal{H})$.

Meek orientation rules. In Fig. A.1.1, we provide Meek's orientation rules used in \mathcal{I} -ORIENT to orient an \mathcal{I} -MEC. These rules provide an algorithm for completing a PDAG to a *completed* PDAG. They are applied repeatedly to a PDAG until no eligible motifs exist.

Next, introduce some additional definitions and results that will be used in Sec. C.

The *skeleton* of a causal DAG \mathcal{G} (denoted $skel(\mathcal{G})$) is the undirected graph that results from ignoring the edge directions of every edge in \mathcal{G} . A triplet of variables (X, Z, Y) in \mathcal{G} is said to be *unshielded*

if (X, Z) and (Y, Z) are adjacent but (X, Y) are not. An unshielded triplet is said to be a *v-structure* (or *unshielded collider*) if it is oriented as $X \rightarrow Z \leftarrow Y$ in \mathcal{G} .

Theorem A.1 (Graphical criterion for Markov equivalence [57, Thm. 1]). *Two DAGs are Markov equivalent if and only if they have the same skeletons and same v-structures.*

Based on the above characterization, to obtain the CPDAG for the MEC corresponding to a DAG \mathcal{G} , one adds an undirected edge for every adjacency in \mathcal{G} ; orients any v-structures according to \mathcal{G} ; then applies Meek’s orientation rules to complete the resulting PDAG to a CPDAG.

Definition A.8 (Global Markov property [41]). A probability distribution $P(\mathbf{v})$ over a set of variables \mathbf{V} is said to satisfy the global Markov property for a causal DAG \mathcal{G} if, for arbitrary disjoint sets $\mathbf{X}, \mathbf{Y}, \mathbf{Z} \subset \mathbf{V}$ with $\mathbf{X}, \mathbf{Y} \neq \emptyset$,

$$\mathbf{X} \perp_d \mathbf{Y} | \mathbf{Z} \implies \mathbf{X} \perp \mathbf{Y} | \mathbf{Z} \text{ in } P(\mathbf{v}).$$

Let $\text{Nd}_X^{\mathcal{G}}$ denote the set of non-descendants of a variable X in \mathcal{G} , i.e. variables in \mathcal{G} (excluding X itself) to which there is no directed path from X .

Definition A.9 (Local Markov property [41]). A probability distribution $P(\mathbf{v})$ over a set of variables \mathbf{V} is said to satisfy the local Markov property for a causal DAG \mathcal{G} if, for every variable $X \in \mathbf{V}$,

$$X \perp \text{Nd}_X^{\mathcal{G}} | \text{Pa}_X^{\mathcal{G}} \text{ in } P(\mathbf{v}).$$

Proposition A.1 (Equivalence of Local and Global Markov Properties [30, Prop. 4]). *Let \mathcal{G} be a causal DAG over variables \mathbf{V} . A probability distribution over \mathbf{V} satisfies the global Markov property for \mathcal{G} if and only if it satisfies the local Markov property for \mathcal{G} .*

Definition A.10 (Covered edge [8, Def. 2]). An edge $X \rightarrow Y$ in a given causal DAG \mathcal{G} is said to be *covered* if $\text{Pa}_Y^{\mathcal{G}} = \text{Pa}_X^{\mathcal{G}} \cup \{X\}$.

Lemma A.1 (Covered edge reversal [8, Lemma. 1]). *Let \mathcal{G} be any causal DAG containing the edge $X \rightarrow Y$ and let \mathcal{G}' be the DAG that is identical to \mathcal{G} except it instead contains the edge $Y \rightarrow X$. Then, \mathcal{G}' is a DAG that is Markov equivalent to \mathcal{G} iff the edge $X \rightarrow Y$ is covered in \mathcal{G} .*

Theorem A.2 (Transformational characterization of Markov equivalent graphs [8, Thm. 2]). *Let $\mathcal{G}, \mathcal{G}'$ be a pair of Markov equivalent causal DAGs. Then, there exists a sequence of covered edge reversals transforming \mathcal{G} to \mathcal{G}' .*

Theorem A.3 (Chickering-Meek theorem [9, Thm. 4]). *Let \mathcal{G} and \mathcal{H} be a pair of causal DAGs such every d -separation that holds in \mathcal{H} also holds in \mathcal{G} . Then, there exists a sequence of edge additions and covered edge reversals transforming \mathcal{G} to \mathcal{H} such that after each reversal and addition, \mathcal{G} is DAG and every d -separation that holds \mathcal{H} also holds in \mathcal{G} .*

A.2 Related works

A.2.1 Learning from observational data

Algorithms for causal discovery fall into three broad categories: constraint-based, score-based, and hybrid. Constraint-based algorithms like PC [54] and the Sparsest Permutations (SP) algorithm [44] learn the true MEC using statistical tests for whether the chosen type of constraints, typically conditional independencies, hold in the data. A challenge to these approaches is improving the accuracy of conditional independence tests, for e.g. in controlling type I error [52]. Score-based algorithms such as GES [9, 34] use a scoring criterion that reflects fit between data and graph, typically in the form of a likelihood plus a complexity penalty. Hybrid algorithms such as max-min hill climbing [56] use a combination of the two approaches, for e.g., first learning the skeleton using a constraint-based method then orienting edges in the skeleton using a score. There is no general claim about the relative accuracy of these methods; we refer readers to [56] for an extensive empirical analysis, who found, for instance, that GES outperforms PC in accuracy across various sample sizes [56, Tables 4, 5].

Commonly, causal discovery algorithms struggle with scaling in high-dimensional settings. This has motivated variants such as Parallel-PC [31] and Fast Greedy Equivalence Search (FGES) [43], which offer faster, parallelized implementations of the algorithms. While FGES offers an additional heuristic over GES, i.e., not adding any edge $X \rightarrow Y$ in the forward phase if X, Y are uncorrelated in the data, this heuristic is not theoretically guaranteed to recover the true MEC. More recent continuous-optimization based approaches such as NoTears [61] are in principle more scalable, but lack theoretical guarantees and show brittle performance even in simulated settings [38, 45].

A.2.2 Learning with prior knowledge

Causal discovery with background knowledge is a well-studied problem [14]. Such background knowledge may be provided by a domain expert or even a large language model [33].

Perfect background knowledge. Many algorithms under this umbrella, such as tiered-FCI [50], the K2 algorithm [15], and Knowledge-Guided GES [22], assume that the background knowledge is correct, i.e., consistent with the ground truth. They use this knowledge to constrain the search, for e.g., by never inserting user-forbidden edges. If the expert’s background knowledge is imperfect, such methods necessarily fail to recover the true MEC. Even if the knowledge is perfect, it is unknown whether these methods will output the true graph; for e.g., in Knowledge-Guided GES [22], restricting insertions no longer guarantees reachability of an MEC with respect to which the given data is Markov in the forward phase, as in GES.

Misspecified background knowledge. The constraint-based approach in [39] allows some misspecification in the expert knowledge in the form of missing or excess adjacencies but not incorrect orientations. However, their approach does not guarantee recovery of the true MEC. Certain score-based approaches treat knowledge about edges a ‘soft’ prior [1, 4, 7, 26, 40] to guide the search, but lack theoretical guarantees on the output graph. One exception with theoretical guarantees is the Sparsest Permutations (SP) algorithm [44], which can initialize the search to an ordering over variables provided by an expert.

A.2.3 Learning from interventional data

Observational learning algorithms can only learn a causal graph up to its observational Markov equivalence class (MEC). The MEC is the limit of what can be identified from observational data without further assumptions. However, MECs can often be large and uninformative for downstream causal tasks [25]. Interventional data can help significantly refine observational MECs [24], and is becoming increasingly available, for e.g., in biological settings due to advances in single-cell technologies [17, 49]. This has motivated the design of algorithms for causal discovery from observational and interventional data such as the score-based Greedy Interventional Equivalence Search (GIES) [24] and the CI-based Interventional Greedy Sparsest Permutations (I-GSP) [58]. However, in [58], it was shown that GIES is inconsistent, i.e., not guaranteed to recover the true interventional MEC in the sample limit. I-GSP is asymptotically consistent, but being permutation-based, struggles to scale in high-dimensional settings.

B Discussion and examples

B.1 Design of Less Greedy Equivalence Search

In this section, we elaborate on certain design choices made in the LGES algorithm.

Choice of DAG in SAFEINSERT. As specified in Def. 2 and Alg. 5, at a given MEC \mathcal{E} , SAFEINSERT chooses some DAG \mathcal{G} from \mathcal{E} to perform its check for conditional independence. In theory, any procedure for choosing a DAG from an MEC suffices. In our implementation, we use Alg. 7, a simple polynomial-time algorithm for converting a CPDAG to DAG presented in [18] and used in the original GES [9]. Another possibility is the polynomial-time algorithm for uniform random sampling of a DAG from an MEC [59]. We leave investigation of DAG choice to future work.

Backward phase of LGES. LGES only modifies the forward phase of GES, while keeping the backward phase intact. Since a score-decreasing INSERT implies independence under some conditioning set (Prop. 1), discarding all $X - Y$ insertions when a single score-decreasing INSERT(X, Y, \mathbf{T}) operator is found promotes more robust edge insertions, i.e., between variables that are not found to be conditionally independent. This significantly reduces false adjacencies in the learned graph (Fig. D.1.1). One may ask, why not use an analogous strategy in the backward phase, discarding all $X - Y$ deletions when a single score-decreasing DELETE(X, Y, \mathbf{T}) operator is found? This is because one score-decreasing deletion implies dependence given one conditioning set, not given any conditioning set. However, a true $X - Y$ adjacency implies dependence given any conditioning set. So, skipping all deletions due to one score decrease would lead to many false adjacencies and

compromise soundness. A true adjacency would imply that all $X - Y$ deletes are score-decreasing, so LGES/GES will not apply any of them. Experiments support this: false non-adjacencies are rare (Fig. D.1.1).

Contrast with PC algorithm. There is a surface-level similarity between LGES and the PC algorithm [54]. The PC algorithm begins with a complete graph, then removes edges between variables for which it finds some separating set. In particular, for adjacent variables X, Y in the current graph \mathcal{G} , it iterates over all subsets $\mathbf{Z} \subseteq \text{Adj}_X^{\mathcal{G}} \cup \text{Adj}_Y^{\mathcal{G}}$ to find some \mathbf{Z} such that $X \perp\!\!\!\perp Y \mid \mathbf{Z}$. If X, Y are non-adjacent in the ground-truth graph, then the PC algorithm is guaranteed to find a \mathbf{Z} which separates them.

LGES, on the other hand, starts with an empty (or user-provided) graph. It then avoids inserting the edge $X \rightarrow Y$ if it finds a set separating X and Y . However, LGES only searches over restricted space of separating sets. CONSERVATIVEINSERT checks if $X \perp\!\!\!\perp Y \mid \text{Pa}_Y^{\mathcal{G}}$ for any \mathcal{G} in the current MEC. SAFEINSERT checks if $X \perp\!\!\!\perp Y \mid \text{Pa}_Y^{\mathcal{G}}$ for a fixed \mathcal{G} in the current MEC. There may exist a set \mathbf{Z} separating X and Y that does not equal $\text{Pa}_Y^{\mathcal{G}}$ for some \mathcal{G} in the current MEC. LGES may not find this \mathbf{Z} , and thus insert an edge between X and Y . In this way, it differs from the PC algorithm.

The advantage of LGES is that it does not require conducting any scoring operations beyond those which GES already conducts. On the other hand, an approach which iterates over all possible separating sets of X and Y would require significantly more scoring operations. Still, future work might consider the theoretical guarantees and empirical performance of such a score-based approach.

B.2 Limitations of Less Greedy Equivalence Search

Causal sufficiency. In this work, we assume that the underlying system is *Markovian*, i.e. no two observed variables have an unobserved common cause. This is also known as the *causal sufficiency* assumption. While this assumption is standard in causal discovery, it can be violated in practice. In settings with unobserved confounders, i.e., *non-Markovian* settings, the equivalence class of graphs that can be identified is typically even larger (and hence more uninformative) than in the Markovian case. One reason for this is that two variables may be non-adjacent in the true graph while still being inseparable by any set due to the existence of an *inducing path* between them [57, Def. 2]. As a result, performing causal inference from equivalence classes of non-Markovian graphs is challenging.

Still, there has been work on causal discovery in non-Markovian settings. Constraint-based approaches include the FCI algorithm [53, 60] and its interventional variants [27, 29, 32], guaranteed to recover the true equivalence class in the sample limit. While there has been progress towards score-based approaches [6, 13, 46, 55], finding algorithms that are asymptotically correct remains an open problem.

There is a fundamental theoretical challenge to generalizing our approach to the non-Markovian setting. GES relies on two useful properties of Markovian DAGs. First is the local Markov property [30] of such DAGs, allowing for a locally consistent scoring criterion (Def. 1). Second are the transformational characterizations of two Markovian causal DAGs \mathcal{G}_1 and \mathcal{G}_2 when (a) \mathcal{G}_1 and \mathcal{G}_2 encode exactly the same d -separations (Thm. A.2) and (b) (a) \mathcal{G}_1 encodes a subset of the d -separations encoded in \mathcal{G}_2 (Thm. A.3). A ‘transformational characterization’ is a procedure whereby \mathcal{G}_1 can be transformed into \mathcal{G}_2 by a sequence of single-edge changes that satisfy certain criteria. The transformational characterization of (a) has been generalized to non-Markovian causal DAGS [60]. Moreover, there exists a local Markov property for non-Markovian causal DAGs, as well as an efficient algorithm to test condition (b) [28]. However, it remains an open problem to fully recover (b) in the non-Markovian setting.

Other assumptions. In this work, we assume we are given a distribution that is Markov and *faithful* with respect to some causal DAG. This is a standard assumption in causal discovery, often justified by the fact that the set of distributions that are Markov but not faithful with respect to a given DAG has Lebesgue measure zero. Moreover, if we assume only that the given distribution is Markov with respect to some causal DAG, we can never rule out the true DAG being the fully connected DAG. Still, there has been work on relaxing the faithfulness assumption, giving rise to the Sparsest Permutations algorithm [44].

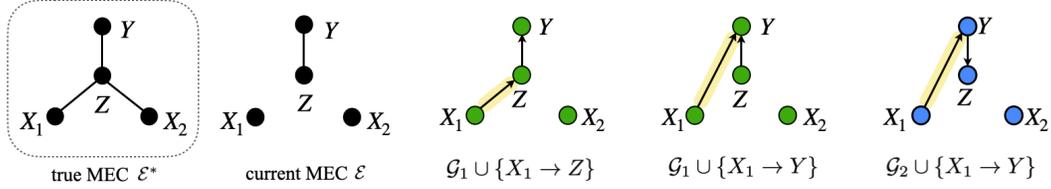


Figure B.3.1: Figs. 2, 3 partially reproduced for convenience. GES is given a distribution $P(\mathbf{v})$ whose true MEC is represented by \mathcal{E}^* (left). GES is currently at MEC \mathcal{E} , evaluating which INSERT operator to apply next. Each INSERT corresponds to picking some $\mathcal{G} \in \mathcal{E}$ and adding some edge to it.

We further assume in this work that we are given a scoring criterion that is decomposable and consistent. This does not strictly mean we make parametric assumptions. Existing scores such as BIC satisfy these criteria as long as the model is a *curved exponential family* [9, 23]. This includes multinomial (discrete) and linear-Gaussian models. For continuous data, the linear-Gaussian assumption can be violated in practice. In this case, one can discretize the data before using it for causal discovery, as in [48]. However, since the parameter space of multinomial models is quite large, and information is lost during discretization, it would be valuable future work to investigate scores for other models of continuous data.

B.3 Extended example of GES trajectories

Finally, we return to our explanation of the two trajectories GES might take in Ex. 2, Fig. 2.

Example B.1. (Ex. 2 continued). Recall that GES is given a distribution $P(\mathbf{v})$ whose true MEC is \mathcal{E}^* (Fig. B.3.1, left). GES is currently at the MEC \mathcal{E} , evaluating which INSERT operator to apply. Each INSERT operator corresponds to picking some DAG $\mathcal{G} \in \mathcal{E}$ and adding some edge to it. One such operator corresponds to picking $\mathcal{G}_1 \in \mathcal{E}$, and adding the edge $X_1 \rightarrow Z$ to it. Another such operator corresponds to picking $\mathcal{G}_2 \in \mathcal{E}$, and adding the edge $X_1 \rightarrow Y$ to it. Both operators, shown in Fig. B.3.1, result in a score increase by local consistency (Def. 1) since $Z \not\perp\!\!\!\perp X_1 \mid \text{Pa}_Z^{\mathcal{G}_1}$ and $Y \not\perp\!\!\!\perp X_1 \mid \text{Pa}_Y^{\mathcal{G}_2}$ in $P(\mathbf{v})$. Although $\mathcal{G}_1 \cup X_1 \rightarrow Z$ looks ‘closer’ to the true MEC \mathcal{E}^* than $\mathcal{G}_2 \cup X_1 \rightarrow Y$, when tested empirically, the latter often scores more than the former (Ex. 2).

Moreover, even in the sample limit, the consistency (Def. A.4) of the scoring criterion does not guarantee that $\mathcal{G}_2 \cup \{X_1 \rightarrow Y\}$ will score lower than $\mathcal{G}_1 \cup \{X_1 \rightarrow Z\}$. Neither the global nor the local consistency of the score provide an immediate guarantee for which will score higher. Global consistency only allows us to compare graphs when $P(\mathbf{v})$ is Markov with respect to at least one of them; however, $P(\mathbf{v})$ is not Markov with respect to $\mathcal{G}_2 \cup \{X_1 \rightarrow Y\}$ or $\mathcal{G}_1 \cup \{X_1 \rightarrow Z\}$. Local consistency (Def. 1) does not let us compare these graphs either, since they differ by more than an edge addition. Therefore, GES may move either to the MEC of $\mathcal{G}_1 \cup \{X_1 \rightarrow Z\}$ (as in τ_1 , Fig. 2) or to the MEC of $\mathcal{G}_2 \cup \{X_1 \rightarrow Y\}$ (as in τ_2 , Fig. 2). It is unknown a priori which operator is the highest-scoring. \square

C Proofs and pseudocode

In this section, we provide proofs, additional results, and pseudocode for the algorithms of the main paper.

Theorem 1 (Correctness of GGES). Let \mathcal{E} denote the Markov equivalence class that results from GGES (Alg. 3) initialized from an arbitrary MEC \mathcal{E}_0 and let $P(\mathbf{v})$ denote the distribution from which the data \mathbf{D} was generated. Then, as the number of samples goes to infinity, \mathcal{E} is the Markov equivalence class underlying $P(\mathbf{v})$.

Proof. The proof is similar to that of [9, Lemma 9, 10]. First, we will prove the correctness of GGES run with the forward-turning-backward phase order.

Forward phase. First, we show that $P(\mathbf{v})$ is Markov with respect to the MEC \mathcal{E}' resulting from the forward phase of GGES. Let \mathcal{G} be any DAG in \mathcal{E}' . By assumption, since GETOPERATOR does not find any valid score-increasing INSERT operators, there exists no such operator. Since there exist no

score-increasing INSERT operators, the local consistency of the scoring criterion (Def. 1) implies that for every $X \in \mathcal{G}$ and $Y \in \text{Nd}_X^{\mathcal{G}}$, $X \perp\!\!\!\perp Y \mid \text{Pa}_X^{\mathcal{G}}$. Otherwise, if we had some $X \in \mathcal{G}$ and $Y \in \text{Nd}_X^{\mathcal{G}}$ such that $X \not\perp\!\!\!\perp Y \mid \text{Pa}_X^{\mathcal{G}}$, the $\text{INSERT}(Y, X, *)$ operator corresponding to $\mathcal{G} \cup \{Y \rightarrow X\}$ would result in a score increase. Since $P(\mathbf{v})$ is faithful to some DAG, it satisfies the *composition axiom* of conditional independence [41]: if $X \perp\!\!\!\perp Y \mid \text{Pa}_X^{\mathcal{G}}$ for every $Y \in \text{Nd}_X^{\mathcal{G}}$, then $X \perp\!\!\!\perp \text{Nd}_X^{\mathcal{G}} \mid \text{Pa}_X^{\mathcal{G}}$. Since this is true for every X , $P(\mathbf{v})$ satisfies the local Markov property of \mathcal{G} . By the equivalence of the global and local Markov properties (Prop. A.1), this means every d-separation in \mathcal{G} implies a corresponding CI in $P(\mathbf{v})$. Thus, $P(\mathbf{v})$ is Markov with respect to \mathcal{G} and hence to \mathcal{E}' .

Turning phase. Next, we show that $P(\mathbf{v})$ is Markov with respect to the MEC \mathcal{E}'' resulting from the turning phase of GGES. The turning phase starts with \mathcal{E}' , output by the forward phase. We have shown that $P(\mathbf{v})$ is Markov with respect to \mathcal{E}' . By construction, each TURN operator applied to the current MEC in the backward phase results in a score increase. If any operator resulted in an \mathcal{E}'' with respect to which $P(\mathbf{v})$ is not Markov, the consistency of the scoring criterion A.4 implies that it would decrease the score. Therefore, $P(\mathbf{v})$ must be Markov with respect to \mathcal{E}'' .

Backward phase. Finally, we show that $P(\mathbf{v})$ is both Markov and faithful to the MEC \mathcal{E} resulting from the backward phase of GGES. The argument that $P(\mathbf{v})$ is Markov with respect to \mathcal{E} is similar to that of the turning phase above. It remains to show that $P(\mathbf{v})$ is faithful to \mathcal{E} . Let \mathcal{E}^* be the true MEC underlying $P(\mathbf{v})$. Since $P(\mathbf{v})$ is both Markov and faithful with respect to \mathcal{E}^* , and $P(\mathbf{v})$ is Markov with respect to \mathcal{E} , every d -separation in \mathcal{E} must also hold in \mathcal{E}^* . Then, by the Chickering-Meek theorem (Thm. A.3), for any $\mathcal{G} \in \mathcal{E}$ and $\mathcal{H} \in \mathcal{E}^*$, there exists a sequence of covered edge reversals and edge additions that transform \mathcal{H} to \mathcal{G} . If this sequence only contains covered edge reversals, then \mathcal{H} and \mathcal{G} are Markov-equivalent (Lemma. A.1), and we are done. Otherwise, let the last edge addition in this sequence add the edge $X \rightarrow Y$, resulting in the DAG \mathcal{G}' . Since \mathcal{G}' can be transformed to \mathcal{G} by a sequence of covered edge reversals, they are Markov equivalent, and we have $\mathcal{G}' \in \mathcal{E}$ (Lemma. A.1). Moreover, since this sequence of transformations includes only covered edge reversals and edge additions, and $P(\mathbf{v})$ is Markov with respect to \mathcal{H} , $P(\mathbf{v})$ is also Markov with respect to $\mathcal{G}' \setminus \{X \rightarrow Y\}$ (by Lemma. A.1, covered edge reversals and additions do not create additional d -separations). By the consistency of the scoring criterion, $\mathcal{G}' \setminus \{X \rightarrow Y\}$ has a higher score than $\mathcal{G}' \in \mathcal{E}$ since the former has fewer parameters. The corresponding DELETE operator thus results in a score increase, and by assumption, GETOPERATOR is guaranteed to find some score-increasing deletion in this case. Thus, we have a contradiction.

Next, we will prove the correctness of GGES without the phase structure, i.e., when GETOPERATOR returns a score-increasing INSERT or DELETE operator if one exists. The proof is similar to the previous case. GGES terminates at an MEC \mathcal{E} only when GETOPERATOR finds no score-increasing INSERT or DELETE operator. By assumption, this implies no such operator exists. If no score-increasing increasing INSERT operator exists, by the same argument as for the forward phase above, the given $P(\mathbf{v})$ must be Markov with respect to \mathcal{E} . Given that $P(\mathbf{v})$ is Markov with respect to \mathcal{E} , and, furthermore, no score-increasing increasing DELETE operator exists, then by the same argument as for the backward phase above, the given $P(\mathbf{v})$ must be both Markov and faithful to \mathcal{E} . \square

Proposition 1. Let \mathcal{E} denote an arbitrary CPDAG and let $P(\mathbf{v})$ denote the distribution from which the data \mathbf{D} was generated. Assume, as the number of samples goes to infinity, that there exists a valid score-decreasing $\text{INSERT}(X, Y, \mathbf{T})$ operator for \mathcal{E} . Then, there exists a DAG $\mathcal{G} \in \mathcal{E}$ such that (1) $Y \perp_d X \mid \text{Pa}_Y^{\mathcal{G}}$ and (2) $Y \perp\!\!\!\perp X \mid \text{Pa}_Y^{\mathcal{G}}$ in $P(\mathbf{v})$.

Proof. The score change of a valid $\text{INSERT}(X, Y, \mathbf{T})$ corresponds to picking a DAG $\mathcal{G} \in \mathcal{E}$ and comparing $S(\mathcal{G}, \mathbf{D})$ with $S(\mathcal{G} \cup \{X \rightarrow Y\}, \mathbf{D})$. By local consistency (Def. 1), if $Y \not\perp\!\!\!\perp X \mid \text{Pa}_Y^{\mathcal{G}}$, then the score must increase. By contrapositive, we have $Y \perp\!\!\!\perp X \mid \text{Pa}_Y^{\mathcal{G}}$ in $P(\mathbf{v})$. Moreover, since X must be a non-descendant of Y for $\mathcal{G} \cup \{X \rightarrow Y\}$ to be a DAG, $Y \perp_d X \mid \text{Pa}_Y^{\mathcal{G}}$ in \mathcal{G} . \square

We provide the following guarantee for LGES Alg. 1 run with CONSERVATIVEINSERT (Strategy 1).

Proposition C.1 (Partial guarantee on CONSERVATIVEINSERT). *Let LGES* denote the variant of LGES (Alg. 1) that uses CONSERVATIVEINSERT instead of SAFEINSERT in the forward phase. Let \mathcal{E} denote the equivalence class that results from the forward phase of LGES* initialized to an arbitrary MEC \mathcal{E}_0 , let $P(\mathbf{v})$ denote the distribution from which the data \mathbf{D} was generated, and let \mathcal{E}^* be the true MEC underlying $P(\mathbf{v})$. Then, as the number of samples goes to infinity,*

1. $skel(\mathcal{E}^*) \subseteq skel(\mathcal{E})$
2. For any unshielded triplet $(X, Z, Y) \in \mathcal{E}^*$, either X, Y are adjacent in \mathcal{E} or (X, Z, Y) is a collider in \mathcal{E}^* if and only if it is a collider in \mathcal{E} .

Proof. For any variables X, Y adjacent in \mathcal{E}^* , since $P(\mathbf{v})$ is faithful to \mathcal{E}^* , X, Y are not independent in the data conditional on any set $\mathbf{Z} \subseteq \mathbf{V}$. Hence, for any $\mathcal{G} \in \mathcal{E}$, $X \not\perp\!\!\!\perp Y \mid \mathbf{Pa}_Y^{\mathcal{G}}$. Therefore, we always have $s(\mathcal{G}) < s(\mathcal{G} \cup \{X \rightarrow Y\})$ for any $\mathcal{G} \in \mathcal{E}$ such that $X \in \mathbf{Nd}_Y^{\mathcal{G}}$, and all valid INSERT($X, Y, *$) operators will result in a score increase. Hence, CONSERVATIVEINSERT will consider all such operators. Since $\hat{\mathcal{E}}$ is a local optimum of the score, any variables that are adjacent in \mathcal{E}^* must also be adjacent in \mathcal{E} . Therefore, $skel(\mathcal{E}^*) \subseteq skel(\mathcal{E})$.

Then, consider some unshielded triplet $(X, Z, Y) \in \mathcal{E}^*$. Since $skel(\mathcal{E}^*) \subseteq skel(\mathcal{E})$, (X, Z) and (Y, Z) must be adjacent in \mathcal{E} . If (X, Y) are also adjacent in \mathcal{E} , we are done. Otherwise, we have an unshielded triplet $(X, Z, Y) \in \mathcal{E}$. Assume (X, Z, Y) is a collider in \mathcal{E}^* . Since CONSERVATIVEINSERT finds no score-increasing INSERT operators for \mathcal{E} , and X, Y are non-adjacent in \mathcal{E} , it must be the case that $\exists \mathcal{G} \in \mathcal{E}$ such that $Y \in \mathbf{Nd}_X^{\mathcal{G}}$ and $X \perp\!\!\!\perp Y \mid \mathbf{Pa}_X^{\mathcal{G}}$ or $X \in \mathbf{Nd}_Y^{\mathcal{G}}$ and $X \perp\!\!\!\perp Y \mid \mathbf{Pa}_Y^{\mathcal{G}}$. Without loss of generality, assume it is the former. Since (X, Z, Y) is a collider in \mathcal{E}^* , $X \not\perp\!\!\!\perp Y \mid \mathbf{Z}$ in $P(\mathbf{v})$ for any set \mathbf{Z} containing Z . Therefore, it must be the case that $Z \notin \mathbf{Pa}_X^{\mathcal{G}}$. Hence, \mathcal{G} contains the edge $X \rightarrow Z$. Since $Y \in \mathbf{Nd}_X^{\mathcal{G}}$, this further implies that \mathcal{G} contains the edge $Y \rightarrow Z$. Therefore, (X, Z, Y) is a collider in \mathcal{G} and hence \mathcal{E}^* . Next, assume that (X, Z, Y) is a collider in \mathcal{E} . Then, $Z \notin \mathbf{Pa}_X^{\mathcal{G}}$ and $Z \notin \mathbf{Pa}_Y^{\mathcal{G}}$ for all $\mathcal{G} \in \mathcal{E}$. As before, since CONSERVATIVEINSERT finds no score-increasing INSERT operators for \mathcal{E} , and X, Y are non-adjacent in \mathcal{E} , it must be the case that $\exists \mathcal{G} \in \mathcal{E}$ such that $Y \in \mathbf{Nd}_X^{\mathcal{G}}$ and $X \perp\!\!\!\perp Y \mid \mathbf{Pa}_X^{\mathcal{G}}$ or $X \in \mathbf{Nd}_Y^{\mathcal{G}}$ and $X \perp\!\!\!\perp Y \mid \mathbf{Pa}_Y^{\mathcal{G}}$. Then, conditioning on Z is not needed to separate X, Y in $P(\mathbf{v})$, which implies that (X, Z, Y) is also a collider in \mathcal{E}^* . □

We give the following condition, sufficient to guarantee that CONSERVATIVEINSERT returns a score-increasing INSERT operator when one exists.

Proposition C.2 (Conditional guarantee on CONSERVATIVEINSERT). *Let \mathcal{E} denote a Markov equivalence class and let $P(\mathbf{v})$ denote the distribution from which the data \mathbf{D} was generated.*

Assume the following holds.

Assumption. *Let \mathcal{G}, \mathcal{H} be two DAGs such that some d -separation encoded in \mathcal{G} does not hold in \mathcal{H} . Then, there exists a pair of variables X, Y non-adjacent in \mathcal{G} with $Y \in \mathbf{Nd}_X^{\mathcal{G}}$ such that for every \mathcal{G}' Markov-equivalent to \mathcal{G} with $Y \in \mathbf{Nd}_X^{\mathcal{G}'}$, $X \not\perp_d Y \mid \mathbf{Pa}_X^{\mathcal{G}'}$ in \mathcal{H} .*

Then, as the number of samples goes to infinity, CONSERVATIVEINSERT returns a valid score-increasing INSERT operator if and only if one exists.

Proof. Let \mathcal{E}^* indicate the true MEC underlying $P(\mathbf{v})$. If there exists a valid score-increasing INSERT operator for the current state \mathcal{E} , then $P(\mathbf{v})$ is not Markov with respect to \mathcal{E} . Since $P(\mathbf{v})$ is faithful to \mathcal{E}^* , this implies that there exists some d -separation encoded in \mathcal{E} that does not hold in \mathcal{E}^* . By the assumption, this implies that there exists X, Y non-adjacent in \mathcal{E} such that for every \mathcal{G} in \mathcal{E} with $Y \in \mathbf{Nd}_X^{\mathcal{G}}$, $X \not\perp_d Y \mid \mathbf{Pa}_X^{\mathcal{G}}$ in \mathcal{E}^* and hence $X \not\perp\!\!\!\perp Y \mid \mathbf{Pa}_X^{\mathcal{G}}$ in $P(\mathbf{v})$. Therefore, every INSERT($Y, X, *$) operator results in a score increase for \mathcal{E} . Then, CONSERVATIVEINSERT is guaranteed to find a score-increasing INSERT. The reverse direction follows by construction, since CONSERVATIVEINSERT enumerates only valid INSERT operators and returns one only if it increases the score. □

As a corollary of the above and Thm. 1, we can also show that LGES with CONSERVATIVEINSERT instead of SAFEINSERT is guaranteed to recover the true MEC in the sample limit, if the assumption in Prop. C.2 holds. We leave the correctness of this assumption open.

Proposition 2 (Correctness of SAFEINSERT). Let \mathcal{E} denote a Markov equivalence class and let $P(\mathbf{v})$ denote the distribution from which the data \mathbf{D} was generated. Then, as the number of samples goes to infinity, SAFEINSERT returns a valid score-increasing INSERT operator if and only if one exists.

Proof. Assume there exists a valid score-increasing INSERT operator for the given MEC \mathcal{E} . Then, $P(\mathbf{v})$ is not Markov with respect to \mathcal{E} . Hence, $P(\mathbf{v})$ is not Markov with respect to the $\mathcal{G} \in \mathcal{E}$ chosen by SAFEINSERT. By the equivalence of the global and local Markov properties (Prop. A.1), this implies that there exists $X \in \mathcal{G}$ such that $X \not\perp\!\!\!\perp \text{Nd}_X^{\mathcal{G}} \mid \text{Pa}_X^{\mathcal{G}}$. Since $P(\mathbf{v})$ is faithful to some DAG, it satisfies the *composition* axiom of conditional independence [41]; hence, there exists some $Y \in \text{Nd}_X^{\mathcal{G}}$ such that $X \not\perp\!\!\!\perp Y \mid \text{Pa}_X^{\mathcal{G}}$. By the local consistency and decomposability of the scoring criterion, we have $s(X, \text{Pa}_X^{\mathcal{G}}) < s(X, \text{Pa}_X^{\mathcal{G}} \cup \{Y\})$. Then, SAFEINSERT will find the valid score-increasing INSERT(Y, X, \mathbf{T}) operator corresponding to $\mathcal{G} \cup \{Y \rightarrow X\}$. The reverse direction is similar. If SAFEINSERT outputs some (X, Y, \mathbf{T}) , this implies it has found some $X \in \mathcal{G}$ and $Y \in \text{Nd}_X^{\mathcal{G}}$ such $s(X, \text{Pa}_X^{\mathcal{G}}) < s(X, \text{Pa}_X^{\mathcal{G}} \cup \{Y\})$, and hence $X \not\perp\!\!\!\perp Y \mid \text{Pa}_X^{\mathcal{G}}$. This implies that $P(\mathbf{v})$ is not Markov with respect to \mathcal{G} and hence to \mathcal{E} , and there exists a valid score-increasing INSERT for \mathcal{E} . The INSERT output by SAFEINSERT is a valid score-increasing operator by construction. \square

We provide pseudocode for the GETSAFEINSERT procedure in Alg. 5. GETSAFEINSERT generalizes SAFEINSERT; instead of searching for a valid INSERT across all non-adjacencies in \mathcal{E} , it searches for a valid INSERT in a subset of the non-adjacencies in \mathcal{E} , given by the *candidates* set. This enables the use of the prioritisation scheme of GETPRIORITYINSERTS.

Proposition C.3 (Correctness of GETPRIORITYINSERTS). *Let $priorityList$ be the list of sets of edges output by GETPRIORITYINSERTS (Alg. 6) given a Markov equivalence class \mathcal{E} and prior assumptions $\mathbf{S} = \langle \mathbf{R}, \mathbf{F} \rangle$. Then, the union of all sets of edges in $priorityList$ is equal to the set of variable pairs (X, Y) that are non-adjacent in \mathcal{E} .*

Proof. This follows from the fact that GETPRIORITYINSERTS loops over all non-adjacencies in \mathcal{E} , and any adjacencies not determined by \mathbf{S} are added to $priorityList[3]$ on line 10. \square

Corollary C.1 (Correctness of LGES-0). *Let \mathcal{E} denote the Markov equivalence class that results from LGES-0 (Alg. 8) with SAFEINSERT, initialised from an arbitrary MEC \mathcal{E}_0 and given prior assumptions $\mathbf{S} = \langle \mathbf{R}, \mathbf{F} \rangle$ using SAFEINSERT, and let $P(\mathbf{v})$ denote the distribution from which the data \mathbf{D} was generated. Then, as the number of samples goes to infinity, \mathcal{E} is the Markov equivalence class underlying $P(\mathbf{v})$.*

Proof. This follows from Prop. 2, Prop. C.3, and Thm. 1. In the forward phase, if $P(\mathbf{v})$ is not Markov with respect to the current MEC \mathcal{E} , SAFEINSERT will find some score-increasing INSERT(X, Y, \mathbf{T})(Prop. 2. Since (X, Y) must be in some set in the priority list returned by GETPRIORITYINSERT (Prop. C.3), some call to GETSAFEINSERT will find a score-increasing INSERT operator. Therefore, each forward step is guaranteed to find a valid score-increasing INSERT operator, if it exists. In the backward phase, since LGES-0 enumerates all valid DELETE operators at each step, it is also guaranteed to find a valid score-increasing DELETE operator when if one exists. Thus, LGES-0 terminates only when there are no score-increasing operators for the current MEC. As such, LGES-0 satisfies the conditions of GGES (Alg. 3) and its correctness follows from Thm. 1. \square

Corollary 1 (Correctness of LGES). Let \mathcal{E} denote the Markov equivalence class that results from LGES (Alg. 1) with SAFEINSERT, initialised from an arbitrary MEC \mathcal{E}_0 and given prior assumptions $\mathbf{S} = \langle \mathbf{R}, \mathbf{F} \rangle$, and let $P(\mathbf{v})$ denote the distribution from which the data \mathbf{D} was generated. Then, as the number of samples goes to infinity, \mathcal{E} is the Markov equivalence class underlying $P(\mathbf{v})$.

Proof. This follows from Prop. 2, Prop. C.3, and Thm. 1. When LGES terminates, there exists no score-increasing deletion for the current MEC \mathcal{E} by construction. By a similar argument to Cor. C.1, we can further show that there exists no score-increasing insertion for the current MEC \mathcal{E} . Thus, LGES satisfies the conditions of GGES (Alg. 3) and its correctness follows from Thm. 1. \square

Corollary C.2 (Correctness of LGES+). *Let \mathcal{E} denote the Markov equivalence class that results from LGES+ (Alg. 8) with SAFEINSERT, initialised from an arbitrary MEC \mathcal{E}_0 and given prior assumptions $\mathbf{S} = \langle \mathbf{R}, \mathbf{F} \rangle$ using SAFEINSERT, and let $P(\mathbf{v})$ denote the distribution from which the*

Algorithm 3: Generalized Greedy Equivalence Search (GGES)

Input: Data $\mathbf{D} \sim \mathbf{P}(\mathbf{v})$, initial MEC \mathcal{E} , scoring criterion S , initial MEC \mathcal{E}_0 , list of phases
phases in $\{[\text{'forward'}, \text{'backward'}], [\text{'any'}]\}$

Output: MEC \mathcal{E} of $\mathbf{P}(\mathbf{v})$

```
1  $\mathcal{E} \leftarrow \mathcal{E}_0$ ;  
2 foreach phase in phases do  
3   repeat  
4      $\text{OPERATE}(X, Y, \mathbf{T}) \leftarrow \text{GETOPERATOR}(\mathcal{E}, \mathbf{D}, S, \textit{phase})$  ;  
5      $\mathcal{E} \leftarrow \mathcal{E} + \text{OPERATE}(X, Y, \mathbf{T})$  ;  
6   until no score-increasing operators exist;  
7 return  $\mathcal{E}$ 
```

data \mathbf{D} was generated. Then, as the number of samples goes to infinity, \mathcal{E} is the Markov equivalence class underlying $\mathbf{P}(\mathbf{v})$.

Proof. The correctness of LGES (Cor. 1) implies that the MEC obtained on line 1 of LGES+ is the highest-scoring MEC. Therefore, the condition on line 7 of LGES+ will never be true, and LGES+ outputs the same MEC as LGES. \square

Theorem 2 (Correctness of \mathcal{I} -ORIENT). Let \mathcal{E} denote the Markov equivalence class that results from \mathcal{I} -ORIENT (Alg. 2) given an observational MEC \mathcal{E}_0 and interventional targets \mathcal{I} , and let $(\mathbf{P}_{\mathbf{I}}(\mathbf{v}))_{\mathbf{I} \in \mathcal{I}}$ denote the family of distributions from which the data $(\mathbf{D}_{\mathbf{I}})_{\mathbf{I} \in \mathcal{I}}$ was generated. Assume that \mathcal{E}_0 is the MEC underlying $P_{\emptyset}(\mathbf{v})$. Then, as the number of samples goes to infinity for each $\mathbf{I} \in \mathcal{I}$, \mathcal{E} is the \mathcal{I} -Markov equivalence class underlying $(\mathbf{P}_{\mathbf{I}}(\mathbf{v}))_{\mathbf{I} \in \mathcal{I}}$.

Proof. Let \mathcal{E}^* denote the true \mathcal{I} -MEC underlying $(\mathbf{P}_{\mathbf{I}}(\mathbf{v}))_{\mathbf{I} \in \mathcal{I}}$. Since \mathcal{E} only orients undirected edges in \mathcal{E}_0 , and \mathcal{E}_0 has the same skeleton and v-structures as \mathcal{E}^* , \mathcal{E} also has the same skeleton and v-structures as \mathcal{E}^* . Next, we show that for every variable pair (X, Y) adjacent in \mathcal{E}^* (and hence \mathcal{E}) for which there exists some $\mathbf{I} \in \mathcal{I}$ with $X \in \mathbf{I}, Y \notin \mathbf{I}$, this edge is directed in both \mathcal{E} and \mathcal{E}^* , and moreover, has the same direction in both.

Consider some edge $(X, Y) \in \mathcal{E}^*$ for which there exists $\mathbf{I} \in \mathcal{I}$ such that $X \in \mathbf{I}, Y \notin \mathbf{I}$. Then, (X, Y) is directed and \mathcal{I} -essential in \mathcal{E} [24, Cor. 13]. We will show that \mathcal{E}^* contains $Y \rightarrow X$ if and only if for every such \mathbf{I} , $s_{\mathbf{D}_{\mathbf{I}}}(y) > s_{\mathbf{D}_{\mathbf{I}}}(y, x)$.

\implies Assume \mathcal{E}^* contains $Y \rightarrow X$. If $X \rightarrow Y \in \mathcal{E}^*$, then $X \rightarrow Y$ in every DAG $\mathcal{G} \in \mathcal{E}^*$. This implies that in every $\mathcal{G} \in \mathcal{E}^*$, there are no directed paths from $Y \rightarrow X$ in \mathcal{G} and hence $\mathcal{G}_{\overline{\mathbf{I}}}$. Moreover, since all edges into X are removed in $\mathcal{G}_{\overline{\mathbf{I}}}$, there are no directed paths from $X \rightarrow Y$ in $\mathcal{G}_{\overline{\mathbf{I}}}$. Since $P_{\mathbf{I}}(\mathbf{v})$ is Markov with respect to $\mathcal{G}_{\overline{\mathbf{I}}}$, this implies $X \perp\!\!\!\perp Y$ in $P_{\mathbf{I}}(\mathbf{v})$. Let \mathcal{H} denote the empty graph over variables \mathbf{V} . Since $X \perp\!\!\!\perp Y$ in $P_{\mathbf{I}}(\mathbf{v})$, by the local consistency of the scoring criterion, \mathcal{H} has a higher score than $\mathcal{H} \cup \{X \rightarrow Y\}$. By the decomposability of the scoring criterion, this implies $s_{\mathbf{D}_{\mathbf{I}}}(y) > s_{\mathbf{D}_{\mathbf{I}}}(y, x)$. Since \mathbf{I} was arbitrary, this must be true for each $\mathbf{I} \in \mathcal{I}$ such that $X \in \mathbf{I}, Y \notin \mathbf{I}$.

\impliedby Assume that $s_{\mathbf{D}_{\mathbf{I}}}(y) > s_{\mathbf{D}_{\mathbf{I}}}(y, x)$ for some $\mathbf{I} \in \mathcal{I}$. Let \mathcal{H} denote the empty graph over variables \mathbf{V} . Then, since $s_{\mathbf{D}_{\mathbf{I}}}(y) > s_{\mathbf{D}_{\mathbf{I}}}(y, x)$, the decomposability of the scoring criterion implies that \mathcal{H} has a higher score than $\mathcal{H} \cup \{X \rightarrow Y\}$. If $Y \not\perp\!\!\!\perp X$ in $P_{\mathbf{I}}(\mathbf{v})$, the local consistency of the scoring criterion would imply that \mathcal{H} has a lower score than $\mathcal{H} \cup \{X \rightarrow Y\}$. By contrapositive, it must be true that $Y \perp\!\!\!\perp X$ in $P_{\mathbf{I}}(\mathbf{v})$. Since $P_{\mathbf{I}}(\mathbf{v})$ is faithful to $\mathcal{G}_{\overline{\mathbf{I}}}$ for some $\mathcal{G} \in \mathcal{E}^*$, this must imply that X, Y are non-adjacent in $\mathcal{G}_{\overline{\mathbf{I}}}$. Since X, Y are adjacent in \mathcal{E}^* , and $X \in \mathbf{I}, Y \notin \mathbf{I}$, this implies that \mathcal{G} and \mathcal{E}^* contain $Y \rightarrow X$. This further implies that if the supposition is true for some $\mathbf{I} \in \mathcal{I}$ with $X \in \mathbf{I}, Y \notin \mathbf{I}$, it must be true for all of them.

The argument to show that \mathcal{E}^* contains $X \rightarrow Y$ if and only if for every \mathbf{I} with $X \in \mathbf{I}, Y \notin \mathbf{I}$, $s_{\mathbf{D}_{\mathbf{I}}}(y) < s_{\mathbf{D}_{\mathbf{I}}}(y, x)$ is analogous.

Moreover, since these statements are true for each $\mathbf{I} \in \mathcal{I}$, they are also true when comparing the sum over sum over all such \mathbf{I} : i.e., $\sum_{\mathbf{I} \in \mathcal{I}, X \in \mathbf{I}, Y \notin \mathbf{I}} s_{\mathbf{D}_{\mathbf{I}}}(y)$ vs $\sum_{\mathbf{I} \in \mathcal{I}, X \in \mathbf{I}, Y \notin \mathbf{I}} s_{\mathbf{D}_{\mathbf{I}}}(y, x)$.

Any edge that is directed in \mathcal{E} is either (a) already directed in \mathcal{E}_0 , in which case it is similarly directed in \mathcal{E}^* , (b) oriented on lines 4 or 7 of \mathcal{I} -ORIENT, in which case it is similarly directed in \mathcal{E}^* by the above argument, or (c) oriented by the Meek rules on lines 5 or 8, in which case it is a consequence of edges directed due to (a) and (b), in which case it is also similarly directed in \mathcal{E}^* . Moreover, the edges directed in \mathcal{E}^* are also due to (a) their being directed in \mathcal{E}_0 , (b) there existing some $\mathbf{I} \in \mathcal{I}$ which contains exactly one endpoint of that edge, or (c) their being a consequence by the Meek rules of these two edge types. Therefore, edges directed in \mathcal{E}^* are also similarly directed in \mathcal{E} , since \mathcal{I} -ORIENT directs each such edge type. We thus have that $\mathcal{E} = \mathcal{E}^*$. \square

Algorithm 4: GETCONSERVATIVEINSERT

Input: MEC \mathcal{E} , DAG $\mathcal{G} \in \mathcal{E}$, data $\mathbf{D} \sim P(\mathbf{v})$, edge insertion candidates *candidates*, scoring criterion S

Output: A valid score-increasing INSERT operator for \mathcal{E} from the adjacencies in *candidates*, or \emptyset if none is found.

```

1  $\Delta S_{max} \leftarrow -\infty;$ 
2  $(X_{max}, Y_{max}, \mathbf{T}_{max}) \leftarrow \emptyset;$ 
3 foreach  $(X, Y)$  in candidates do
4   if  $X \in \text{Nd}_Y^{\mathcal{G}}$  and  $s(Y, \text{Pa}_Y^{\mathcal{G}}) < s(Y, \text{Pa}_Y^{\mathcal{G}} \cup \{X\})$  then
5      $\Delta S_{xy} \leftarrow -\infty;$ 
6      $(\hat{X}, \hat{Y}, \hat{\mathbf{T}}) \leftarrow \emptyset;$ 
7     foreach valid  $\mathbf{T} \subseteq \text{Ne}_Y^{\mathcal{E}} \setminus \text{Adj}_X^{\mathcal{E}}$  do
8        $\Delta S \leftarrow s(Y, (\text{Ne}_Y^{\mathcal{E}} \cap \text{Adj}_X^{\mathcal{E}}) \cup \mathbf{T} \cup \text{Pa}_Y^{\mathcal{E}} \cup \{X\}) - s(Y, (\text{Ne}_Y^{\mathcal{E}} \cap \text{Adj}_X^{\mathcal{E}}) \cup \mathbf{T} \cup \text{Pa}_Y^{\mathcal{E}});$ 
9       if  $\Delta S > \Delta S_{xy}$  then
10          $\Delta S_{xy} \leftarrow \Delta S;$ 
11          $(\hat{X}, \hat{Y}, \hat{\mathbf{T}}) \leftarrow (X, Y, \mathbf{T});$ 
12       else if  $\Delta S < 0$  then
13          $\Delta S_{xy} \leftarrow -\infty;$ 
14         break;
15       if  $\Delta S_{xy} > \Delta S_{max}$  then
16          $\Delta S_{max} \leftarrow \Delta S_{xy};$ 
17          $(X_{max}, Y_{max}, \mathbf{T}_{max}) \leftarrow (\hat{X}, \hat{Y}, \hat{\mathbf{T}});$ 
18 return  $(X_{max}, Y_{max}, \mathbf{T}_{max})$ 

```

Algorithm 5: GETSAFEINSERT

Input: MEC \mathcal{E} , DAG $\mathcal{G} \in \mathcal{E}$, data $\mathbf{D} \sim P(\mathbf{v})$, edge insertion candidates *candidates*, scoring criterion S

Output: A valid score-increasing INSERT operator for \mathcal{E} from the adjacencies in *candidates*, or \emptyset if none is found.

```

1  $\Delta S_{max} \leftarrow -\infty;$ 
2  $(X_{max}, Y_{max}, \mathbf{T}_{max}) \leftarrow \emptyset;$ 
3 foreach  $(X, Y)$  in candidates do
4   if  $X \in \text{Nd}_Y^{\mathcal{G}}$  and  $s(Y, \text{Pa}_Y^{\mathcal{G}}) < s(Y, \text{Pa}_Y^{\mathcal{G}} \cup \{X\})$  then
5     foreach valid  $\mathbf{T} \subseteq \text{Ne}_Y^{\mathcal{E}} \setminus \text{Adj}_X^{\mathcal{E}}$  do
6        $\Delta S \leftarrow s(Y, (\text{Ne}_Y^{\mathcal{E}} \cap \text{Adj}_X^{\mathcal{E}}) \cup \mathbf{T} \cup \text{Pa}_Y^{\mathcal{E}} \cup \{X\}) - s(Y, (\text{Ne}_Y^{\mathcal{E}} \cap \text{Adj}_X^{\mathcal{E}}) \cup \mathbf{T} \cup \text{Pa}_Y^{\mathcal{E}});$ 
7       if  $\Delta S > \Delta S_{max}$  then
8          $\Delta S_{max} \leftarrow \Delta S;$ 
9          $(X_{max}, Y_{max}, \mathbf{T}_{max}) \leftarrow (X, Y, \mathbf{T});$ 
10 return  $(X_{max}, Y_{max}, \mathbf{T}_{max})$ 

```

Algorithm 6: GETPRIORITYINSERTS

Input: MEC \mathcal{E} , prior assumptions $\mathbf{S} = \langle \mathbf{R}, \mathbf{F} \rangle$

```
1 priorityList  $\leftarrow [\{\} \times 4]$ ;
2 foreach  $(X, Y)$  non-adjacent in  $\mathcal{E}$  do
3   if  $X - *Y \in \mathbf{R}$  then
4     | Add  $(X, Y)$  to priorityList[1]; // required
5   else if  $Y \rightarrow X \in \mathbf{R}$  then
6     | Add  $(X, Y)$  to priorityList[2]; // weakly required
7   else if  $X \rightarrow Y \in \mathbf{F}$  or  $X - Y \in \mathbf{F}$  then
8     | Add  $(X, Y)$  to priorityList[4]; // forbidden
9   else
10    | Add  $(X, Y)$  to priorityList[3]; // ambivalent
11 return priorityList
```

Algorithm 7: PDAGTODAG

Input: MEC \mathcal{E}

```
1  $\mathcal{E}' \leftarrow \mathcal{E}$ ;
2 while there exists an undirected edge in  $\mathcal{E}$  do
3   | Choose a vertex  $V$  in  $\mathcal{E}'$  such that (i)  $V$  is a sink node in  $\mathcal{E}'$  (no outgoing directed edges), and
4     | (ii) all undirected neighbors of  $V$  are pairwise adjacent in  $\mathcal{E}'$  (form a clique);
5   foreach undirected edge  $(U, V)$  in  $\mathcal{E}$  do
6     | Orient  $(U, V)$  as  $U \rightarrow V$  in  $\mathcal{E}$ ;
7   Remove  $V$  from  $\mathcal{E}'$ ;
8 return  $\mathcal{E}$ 
```

Algorithm 8: Less Greedy Equivalence Search-0 (LGES-0)

Input: Data $\mathbf{D} \sim \mathbf{P}(\mathbf{v})$, scoring criterion S , prior assumptions $\mathbf{S} = \langle \mathbf{R}, \mathbf{F} \rangle$, initial MEC \mathcal{E}_0 , insertion strategy $GetInsert$ in $\{\text{GETSAFEINSERT}, \text{GETCONSERVATIVEINSERT}\}$

Output: MEC \mathcal{E} of $\mathbf{P}(\mathbf{v})$

```
1  $\mathcal{E} \leftarrow \mathcal{E}_0$ ; // allows initialisation if preferred by user
2 repeat
3   |  $\mathcal{G} \leftarrow$  some DAG in  $\mathcal{E}$ ; // forward phase
4   | priorityList  $\leftarrow$  GETPRIORITYINSERTS( $\mathcal{E}, \mathcal{G}, \mathbf{S}$ );
5   | foreach candidates in priorityList do
6     |  $(X_{max}, Y_{max}, \mathbf{T}_{max}) \leftarrow GetInsert(\mathcal{E}, \mathcal{G}, \mathbf{D}, candidates, S)$ ;
7     | if  $(X_{max}, Y_{max}, \mathbf{T}_{max})$  is found then
8       | |  $\mathcal{E} \leftarrow \mathcal{E} + \text{INSERT}(X_{max}, Y_{max}, \mathbf{T}_{max})$ ;
9       | | break; // no need to check lower priority
10  until no improving insertions exist;
11 repeat
12  |  $\mathcal{E} \leftarrow \mathcal{E} +$  highest-scoring TURN( $X, Y, \mathbf{T}$ ); // turning phase
13  until no score-increasing reversals exist;
14 repeat
15  |  $\mathcal{E} \leftarrow \mathcal{E} +$  highest-scoring DELETE( $X, Y, \mathbf{T}$ ); // backward phase
16  until no score-increasing deletions exist;
17 return  $\mathcal{E}$ 
```

Algorithm 9: Less Greedy Equivalence Search+ (LGES+)

Input: Data $\mathbf{D} \sim \mathbf{P}(\mathbf{v})$, scoring criterion S , prior assumptions $\mathbf{S} = \langle \mathbf{R}, \mathbf{F} \rangle$, initial MEC \mathcal{E}_0 , insertion strategy $GetInsert$ in $\{\text{GETSAFEINSERT}, \text{GETCONSERVATIVEINSERT}\}$

Output: MEC \mathcal{E} of $\mathbf{P}(\mathbf{v})$

```

1  $\mathcal{E} \leftarrow \text{LGES}(\mathbf{D}, S, \mathbf{S}, \mathcal{E}_0, GetInsert)$ ;
2  $\mathcal{D} \leftarrow$  all deletions valid for  $\mathcal{E}$ ;
3 while  $|\mathcal{D}| > 0$  do
4    $(X_{max}, Y_{max}, \mathbf{T}_{max}) \leftarrow$  highest-scoring deletion in  $\mathcal{D}$ ;
5    $\mathcal{E}' \leftarrow \mathcal{E} + \text{DELETE}(X_{max}, Y_{max}, \mathbf{T}_{max})$ ;
6    $\mathcal{E}' \leftarrow \text{LGES}^*(\mathbf{D}, S, \mathbf{S}, \mathcal{E}_0, GetInsert, X_{max}, Y_{max})$ ;
   // LGES* is a modified version of LGES that excludes edge insertions
   // between the given  $X, Y$ 
7   if  $S(\mathcal{E}', \mathbf{D}) > S(\mathcal{E}, \mathbf{D})$  then
8      $\mathcal{E} \leftarrow \mathcal{E}'$ ;
9      $\mathcal{D} \leftarrow$  all deletions valid for  $\mathcal{E}$ ;
10  else
11     $\mathcal{D} \leftarrow \mathcal{D} \setminus \{(X_{max}, Y_{max}, \mathbf{T}_{max})\}$ ;
12 return  $\mathcal{E}$ 

```

D Experiments

Compute details. All experiments were run on a shared compute cluster with 2x Intel Xeon Platinum 8480+ CPUs (112 cores total, 224 threads) at up to 3.8 GHz, and 210 MiB L3 cache.

D.1 Learning from observational data

Baseline details We ran the PC algorithm using significance level $\alpha = 0.05$ for conditional independence tests, with the null hypothesis of independence. Since NoTears often outputs cyclic graphs, we post-processed the output graph by greedily removing the lowest-weight edges until it was acyclic, following [37]. This was done so that we could convert the output to a valid CPDAG for comparison with other algorithms. We ran NoTears with default parameters from the `causalnex` library, which uses a weight threshold of $w = 0$; note that performance may vary depending on the choice of parameters, particularly the weight threshold and the acyclicity penalty parameter. However, since the implementation of NoTears is very compute-heavy, we were not able to tune these parameters. For fair comparison, we implemented all GES variants, including XGES, LGES, and GIES, by modifying the code provided by [20].⁵ Discrepancies in runtime of XGES between our results and the results of [36] may result from the latter’s use of an optimized C++ implementation, which includes an implementation of the search operators that they used in XGES but not in GES.

D.1.1 Synthetic data details for Experiment 5.1

For the results shown in Fig. 4, we draw Erdős–Rényi graphs with p variables and $\{2p$ edges in expectation (ER2), for $p \in \{5, 10, 15, 20, 25, 35, 50, 75, 100, 150\}$. For each p , we sample 50 graphs and generate linear-Gaussian data for each graph. Following [37], we draw weights from $\mathcal{U}([-2, -0.5] \cup [0.5, 2])$. In some cases, the resulting weight matrix was (almost) singular, in which case each row of weights was ℓ_1 normalized. We draw noise means from $\mathcal{N}(0, 1)$ and noise variances from $\mathcal{U}([0.1, 0.5])$. We obtain $n = 10^4$ samples per dataset via `sempler` [20].

D.1.2 Further baselines and metrics for Experiment 5.1

In Fig. D.1.1, we present additional baselines and accuracy metrics for the setting in Sec. 5.1, including the particular types of structural errors (excess adjacencies, missing adjacencies, incorrect orientations). As in the case of SHD, LGES outperforms other algorithms across these metrics, with CONSERVATIVEINSERT outperforming SAFEINSERT.

⁵<https://github.com/juangame11a/ges>

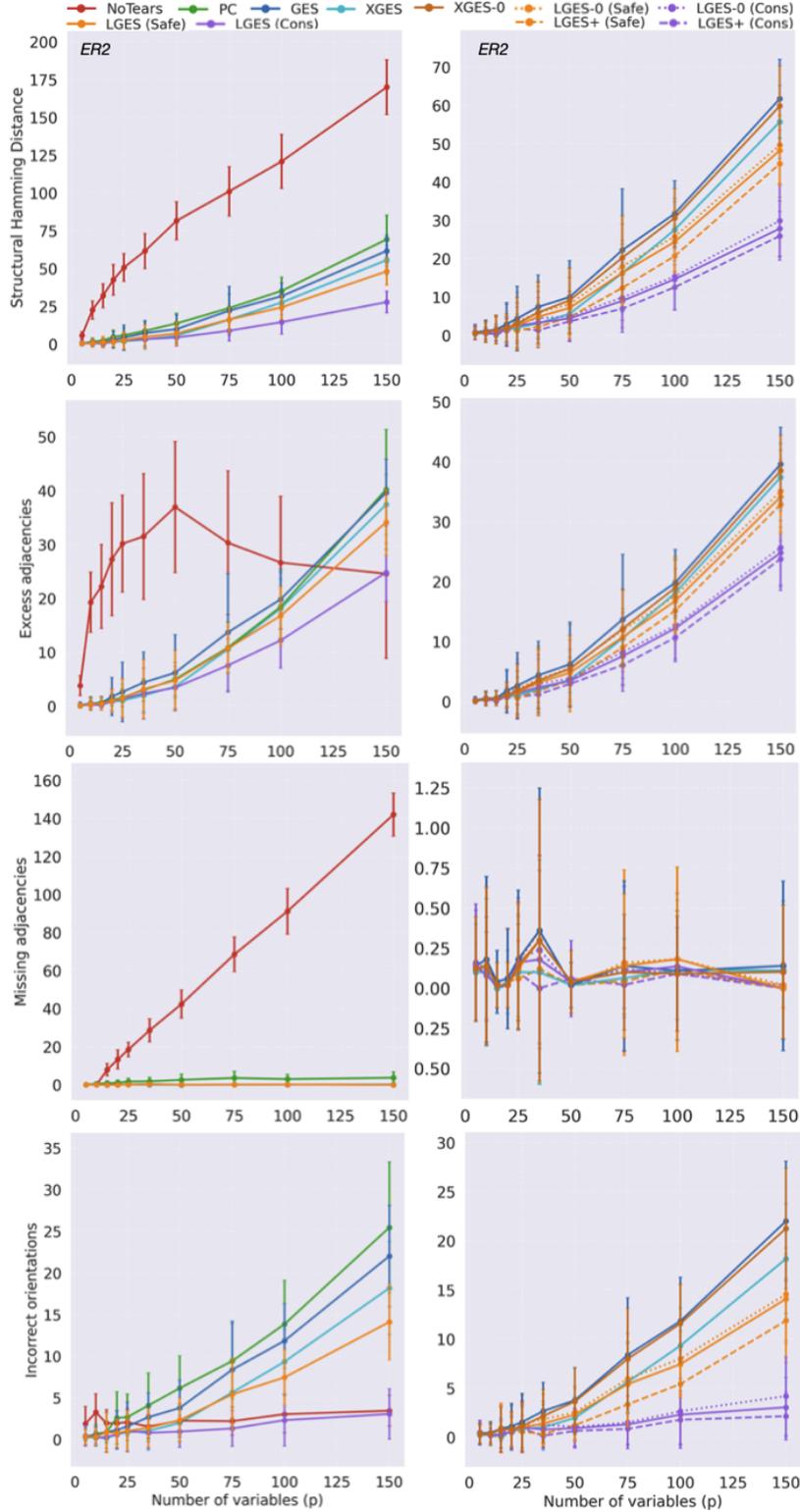


Figure D.1.1: Structural error of algorithms on 50 simulated observational datasets from Erdős–Rényi graphs with p variables and $2p$ edges in expectation, given $n = 10^4$ samples and no prior knowledge (Sec. D.1.2). **Lower is better** (more accurate / faster) across all plots. **(Left)** Common baselines. GES-style algorithms outperform PC and NoTears in accuracy. **(Right)** GES variants. LGES-0, LGES, and LGES+ are the less greedy variants of GES, XGES-0, and XGES respectively. Each less greedy algorithm improves on its greedy counterpart in accuracy.

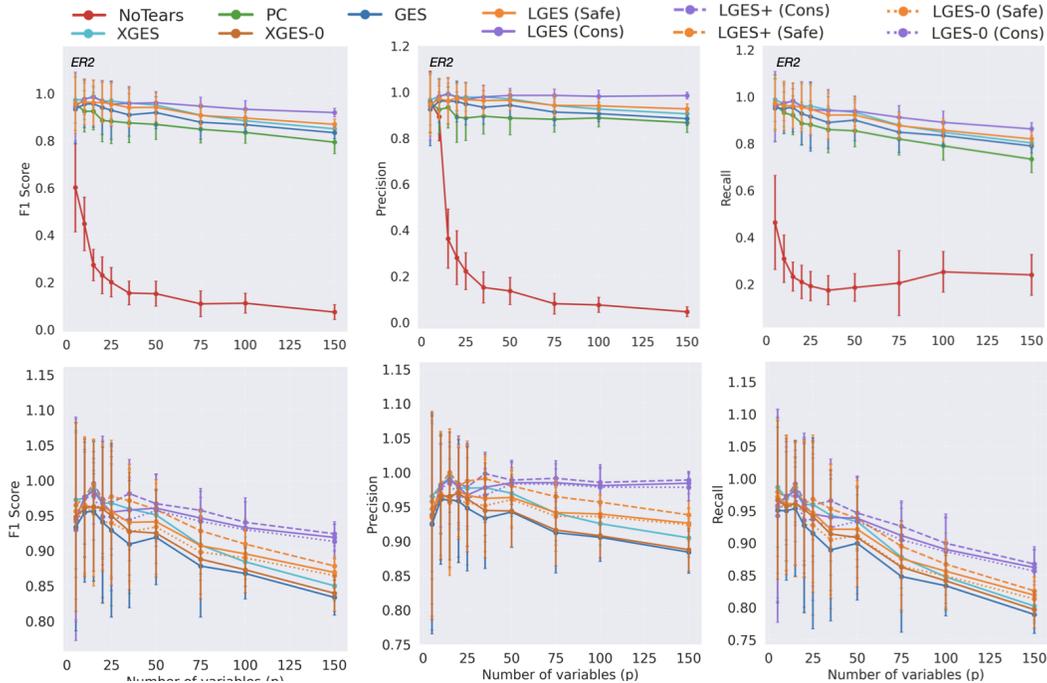


Figure D.1.2: Binary classification accuracy of algorithms on 50 simulated observational datasets from Erdős–Rényi graphs with p variables and $2p$ edges in expectation, given $n = 10^4$ samples and no prior knowledge (Sec. D.1.3). **Higher is better** across all plots. An MEC is said to contain an edge (X, Y) if it contains either $X - Y$ or $X \rightarrow Y$ (but not $Y \rightarrow X$). **(Top)** Common baselines. GES-style algorithms outperform PC and NoTears in accuracy. **(Bottom)** GES variants. LGES-0, LGES, and LGES+ are the less greedy variants of GES, XGES-0, and XGES respectively. Each less greedy variant improves on its greedy counterpart in accuracy.

The behaviour of NoTears is more variable. It misses significantly more adjacencies than all other methods—approximately linearly many in the number of variables. It also includes many more excess adjacencies than the other algorithms up to $p = 50$, after which the number of excess adjacencies begins to decline, ultimately approaching that of LGES for $p = 150$; however, this could be explained by the increasing number of adjacencies that are also missed by NoTears. We show in the next section how, as a result, the F1 score of NoTears is low across all graph sizes (Fig. D.1.2).

D.1.3 Precision, recall, and F1 score for Experiment 5.1

Next, in addition to structural error, we evaluate accuracy by considering causal discovery as a binary classification task. We use the task definition provided in [36]: an MEC \mathcal{E} is said to contain an edge (X, Y) if it contains either $X \rightarrow Y$ or $X - Y$ (but not $Y \rightarrow X$). The results are shown in Fig. D.1.2. For graphs with $p < 20$ nodes, we find that GES and LGES have similar F1 scores. They both outperform PC, which in turn substantially outperforms NoTears. For $p > 20$, LGES+ (Conservative) dominates, followed by LGES (Conservative). For large p , LGES (Conservative) achieves a substantially higher F1 score than GES and other algorithms except LGES+. For instance, for $p = 150$, LGES (Conservative) achieves an F1 score just under 0.95, whereas GES’s F1 score is slightly under 0.85.

D.1.4 Further experiments with varying edge densities

In Fig. D.1.3, we provide results for select baselines on Erdős–Rényi graphs with p variables and $\{1, 3\} \cdot p$ edges in expectation following the set-up in Experiment 5.1. The relative performance of algorithms is similar to the ER2 case (Fig. D.1.1), though PC and GES achieve similar accuracy on ER3 graphs. LGES+ (Conservative) is the most accurate overall, achieving $2\times$ less structural error than GES on ER1 and ER3 graphs. For instance, graphs with 150 variables and 450 edges in expectation, LGES (Conservative) achieves SHD ≈ 40 , and LGES+ (Conservative) ≈ 30 . All less

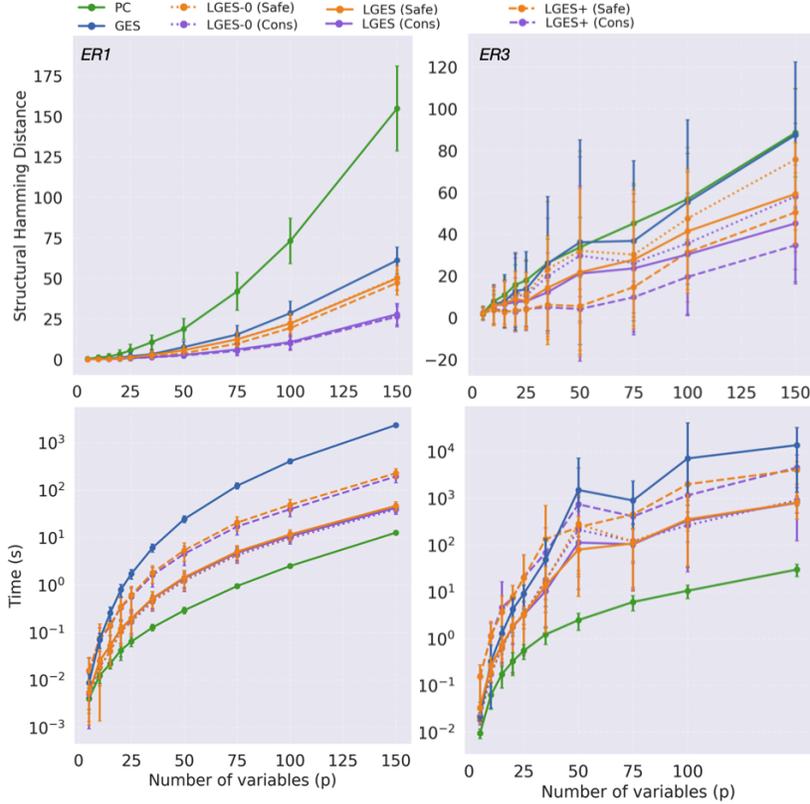


Figure D.1.3: Performance of algorithms on 50 simulated observational datasets from Erdős–Rényi graphs with p variables and **(left)** p edges and **(right)** $3p$ edges in expectation, given $n = 10^4$ samples and no prior knowledge (Sec. D.1.4). **Lower is better** (more accurate / faster) across all plots. See Fig. 4 for graphs with $2p$ edges in expectation. LGES (Conservative) is the fastest GES variant across edge densities, whereas LGES+ (Conservative) is the most accurate overall algorithm across edge densities.

greedy algorithms are faster and more accurate than GES, with LGES-0 and LGES being up to $10\times$ faster.

D.1.5 Further experiments with smaller datasets

To investigate performance in the small-sample setting, we conduct experiments with $n = 10^3$ for ER2 graphs with $p \in \{10, 25, 50, 100\}$ variables. The results are summarized in Fig. D.1.4.

As in previous experiments, LGES (Safe and Conservative) outperforms GES in runtime and accuracy across all graph sizes, with CONSERVATIVEINSERT yielding higher accuracy than SAFEINSERT. Interestingly, the PC algorithm outperforms LGES for the case of $p = 100$, suggesting PC’s usefulness for settings where the number of samples is small compared with the number of variables.

D.1.6 Further experiments with larger graphs

We scaled up Experiment 5.1 to graphs with $p \in \{175, 250\}$ variables. We ran LGES (Safe and Conservative) and PC. We were unable to continue running XGES, GES, and NoTears due to time and compute constraints; for instance, GES took over 10^4 seconds without terminating for $p = 175$ for a single trial. The results are summarized in Table D.1.1.

LGES both with SAFEINSERT and with CONSERVATIVEINSERT is substantially more accurate (in terms of SHD and F1 score) than PC, with CONSERVATIVEINSERT achieving less than half the structural error of PC. While PC is faster than LGES, LGES is much faster than GES, taking only ≈ 5 -6 minutes for $p = 175$ and ≈ 15 -20 minutes for $p = 250$. Recall that for $p = 150$, GES was already approaching a runtime of 10^4 seconds (>2 hours) (Fig. 4a).

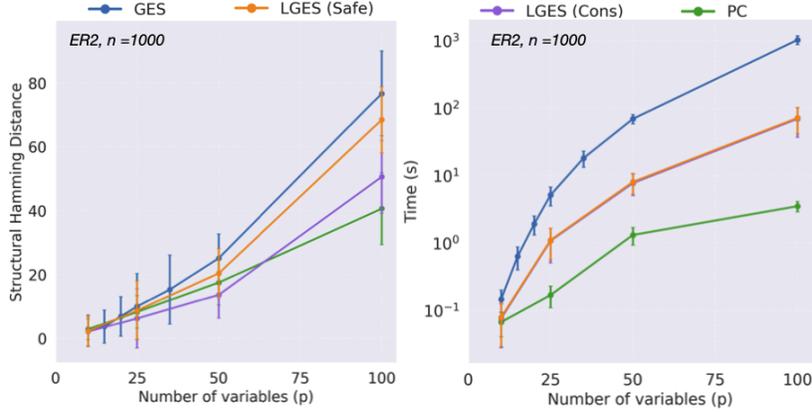


Figure D.1.4: Performance of algorithms on 50 simulated observational datasets from Erdős-Rényi graphs with p variables and $2p$ edges in expectation, given the smaller dataset of 10^3 samples and no prior knowledge (Sec D.1.5). In terms of accuracy, LGES (Conservative) outperforms other algorithms except for the case of $p = 100$, when PC is more accurate.

Metric	Method	$p = 175$	$p = 250$
SHD	LGES-0 (Conservative)	42.04 \pm 11.51	83.06 \pm 14.42
	LGES (Conservative)	39.457 \pm 10.44	82.80 \pm 14.73
	LGES-0 (Safe)	70.62 \pm 11.68	135.24 \pm 14.31
	LGES (Safe)	67.65 \pm 12.31	133.22 \pm 14.07
	PC	91.34 \pm 17.17	166.20 \pm 24.34
Runtime	LGES-0 (Conservative)	315.83 \pm 73.23	989.37 \pm 205.53
	LGES (Conservative)	326.44 \pm 73.166	1019.13 \pm 198.12
	LGES-0 (Safe)	337.45 \pm 78.43	1058.11 \pm 221.32
	LGES (Safe)	348.23 \pm 80.57	1095.71 \pm 209.10
	PC	29.20 \pm 9.50	85.67 \pm 23.52

Table D.1.1: Performance of algorithms (mean \pm std) on 50 simulated observational datasets from large Erdős-Rényi graphs with p variables and $2p$ edges in expectation, given $n = 10^4$ samples and no prior knowledge (Sec. D.1.6). **Lower is better** (more accurate / faster) across all metrics. Other algorithms are omitted as they exceeded the 10^4 second bound on runtime for these graph sizes, prohibiting repeated trials. Among these, PC is the fastest, but makes twice as many structural errors as LGES (Conservative).

D.2 Learning with prior knowledge

In Fig. D.2.1, we present more detailed plots of runtime and SHD for the setting in Sec. 5.2, including background knowledge comprising $m' = m/2$ and $m' = 3m/4$ edges for a ground truth graph with m edges. Results for $m' = m/2$ follow a very similar trend to those with $m' = 3m/4$ discussed in Sec. 5.2. Initialisation outperforms prioritization in LGES when the expert knowledge is mostly correct ($f_c \geq 0.75$), but prioritization does better otherwise.

D.3 Learning from interventional data

In this section, we further discuss the performance of algorithms given interventional data in Sec. 5.3. We observed that LGES (Safe and Conservative) followed by \mathcal{I} -ORIENT has higher SHD from the ground truth than LGIES, Safe and Conservative respectively. We hypothesize that this is because, in our experiments, LGES only uses 10^4 observational samples during the MEC learning phase. It uses the interventional samples only to orient edges using \mathcal{I} -ORIENT. In contrast, LGIES uses $10^4 + k \cdot 10^3$ samples in the MEC learning phase given k interventions, since it uses interventional

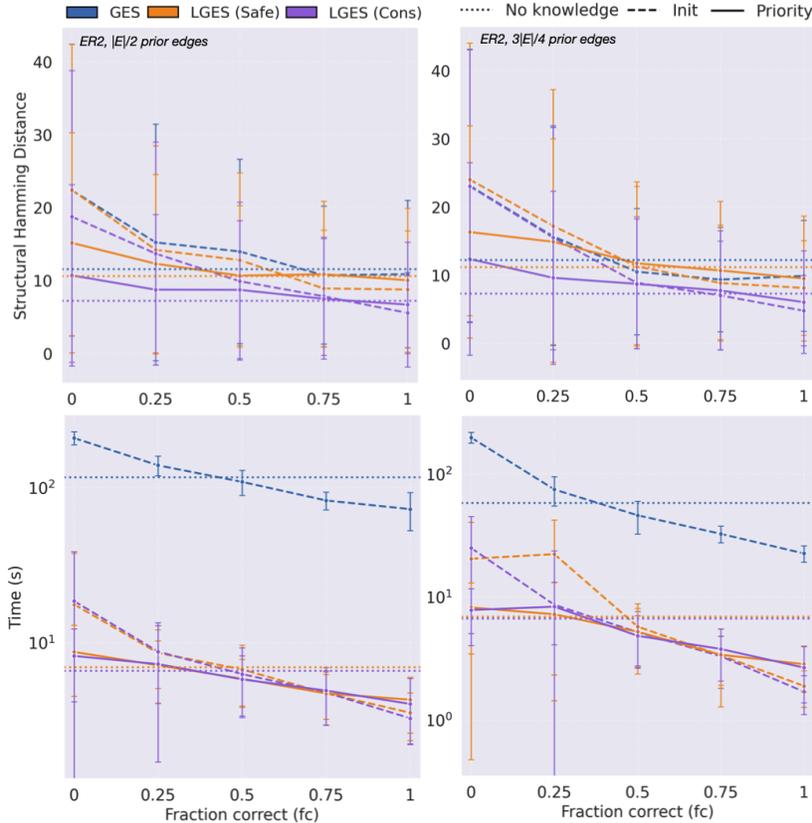


Figure D.2.1: Performance of algorithms on 50 simulated observational datasets from Erdős–Rényi graphs with 50 variables and 100 edges in expectation given $n = 10^3$ samples (Experiment 5.2, Sec. D.2). **Lower is better** (more accurate / faster) across all plots. We vary the correctness of the prior knowledge on the x-axis; higher fc indicates a higher fraction of knowledge that is correct. **(Left)** Algorithms are given $m' = m/2$ required edges as knowledge, for a true graph containing m edges. **(Right)** Similar, but with $m' = 3m/4$. LGES’ prioritization strategy is more robust to misspecification in the knowledge ($\leq 75\%$ correct) than initialization with the same knowledge

data throughout learning, and we generate 10^3 interventional samples per target. Moreover, since we choose $k = p/10$ (where p is the number of variables), LGIES is making use of much more data than LGES for learning the MEC. Although GIES is known not to be asymptotically correct [58], this further motivates research into an asymptotically correct causal discovery algorithm that can make use of interventional data throughout learning.

D.4 Real-world protein signaling data

Sachs dataset. We compare GES and LGES on a real-world protein signaling dataset [48]. The observational dataset consists of 853 measurements of 11 phospholipids and phosphoproteins. We compare the output of our methods with the gold-standard inferred graph [48, Fig. 3]), containing 11 variables and 17 edges. The dataset is continuous but violates the linear-Gaussian assumption.⁶ We test our methods both on the original continuous dataset and on a discretized version (3 categories per variable corresponding to low, medium, and high concentration) from the bnlearn repository.⁷

Results. The learned graphs provided in Fig. D.4.1. All algorithms output the same MEC and thus have the same accuracy in both settings. With discrete data, each algorithm had an SHD of 9 edges from the reference MEC, all of which were missing adjacencies. With continuous data, each

⁶<https://www.bnlearn.com/research/sachs05/>

⁷<https://www.bnlearn.com/book-crc-2ed/>

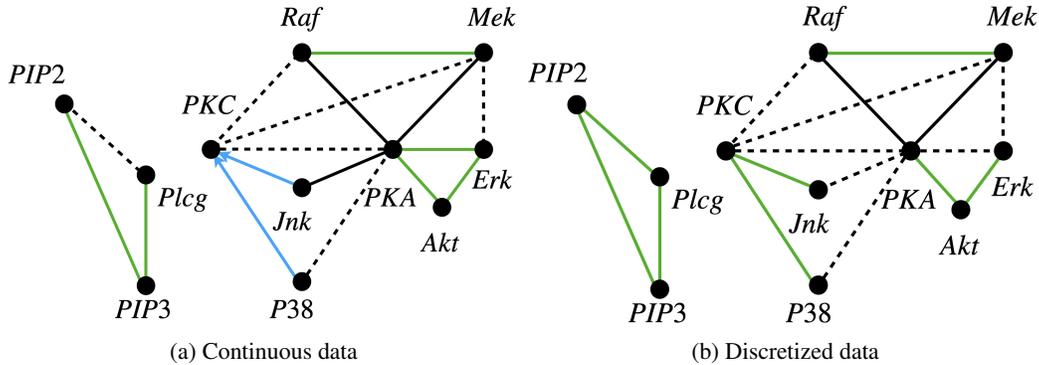


Figure D.4.1: Comparison between the reference MEC and the learned MEC for $n = 853$ observational samples from the Sachs protein-signaling dataset [48]. For both continuous and discretized data, the three algorithms (GES, LGES (SAFEINSERT) and LGES (CONSERVATIVEINSERT)) return the same MEC, so only one learned graph is shown per panel. Green solid lines indicate edges correctly recovered by the algorithms. Blue solid lines indicate edges misoriented by the algorithms. Black dashed lines indicate edges missed by the algorithms. **(a)** Continuous data: nine edges are missed and two are misoriented: $Jnk \rightarrow Pkc$ and $P38 \rightarrow PKC$, both undirected in the reference MEC. **(b)** Discretised data: nine edges are missed and none are misoriented.

algorithm had an SHD of 11 edges from the reference MEC, 9 of which were missing adjacencies and 2 of which were incorrect orientations.

E Frequently Asked Questions

Q1. What is the difference between score-based and constraint-based causal discovery?

Answer. Constraint-based and score-based approaches to causal discovery solve the same problem but in different ways. Constraint-based approaches such as PC [54] and Sparsest Permutations [44] use statistical tests, usually for conditional independence, to learn a Markov equivalence class from data. Score-based approaches such as Greedy Equivalence Search [9, 35] instead attempt to maximize a score (for e.g., the Bayesian Information Criterion or BIC [51]) that reflects the fit between graph and data. There is no general claim about which of these methods is superior; we refer readers to [56] for an extensive empirical analysis, who found, for instance, that GES outperforms PC in accuracy across various sample sizes [56, Tables 4, 5].

Q2. I thought causal discovery was only one problem, but it seems the paper claims to be solving three different tasks: observational observational learning, interventional learning, learning with prior knowledge. Can you elaborate on these tasks (why they are different, how they relate, etc.)?

Answer. The most well-studied problem in causal discovery is that of learning causal graphs from only observational data. However, algorithms for this task have a few limitations. Firstly, they are computationally expensive and often fail to produce quality estimates of the true Markov equivalence class from finite samples. This motivates using prior knowledge in the search to produce better quality estimates of the true Markov equivalence class and do so faster. Secondly, algorithms for learning from observational data only identify a Markov equivalence of graphs, since this is the most informative structure that can be learned from observational data. However, MECs can be quite large, and thus uninformative for downstream tasks such as causal inference. When interventional data is available, more edges in the true graph become identifiable. This motivates using interventional data to identify a smaller and more informative interventional Markov equivalence class—the task of interventional learning. When no prior assumptions or interventional data are available, these tasks collapse to observational learning.

Q3. If GES is asymptotically consistent, why bother to consider LGES?

Answer. While GES is guaranteed to recover the true Markov equivalence class given infinite samples, it faces two challenges. First, computational tractability: structure learning

is an NP-hard problem, and GES commonly struggles to scale in high-dimensional settings. Second, data is often limited in practice, which results in GES failing to recover the true MEC. LGES improves on GES in both of these aspects; it is up to 10 times faster and 2 times more accurate (Experiment 5.1). Moreover, while GES and LGES can both incorporate prior assumptions to guide the search, LGES is more robust to misspecification in the assumptions (Experiment 5.2).

Q4. How well does LGES scale?

Answer. LGES can scale to graphs with hundreds of variables and is up to 10 times faster than GES (Sec. D.1). Both variants of LGES (Safe and Conservative) terminate in less than 20 minutes hours on graphs with 250 variables (Sec. D.1.6) and have substantially better accuracy than other baselines (including PC and NoTears). D.1.6). LGES also outperforms these baselines in settings with dense graphs (Sec. D.1.4).

Q5. LGES may be asymptotically correct, but how well does it perform with finite samples?

Answer. We conduct an extensive empirical analysis of how LGES performs compared with baselines in finite-sample settings. Both variants of LGES (SAFEINSERT and CONSERVATIVEINSERT) have substantially better accuracy than GES, PC, and NoTears in experiments with $n = 10^4$ samples and up to 500 variables (Experiments 5.1, D.1). For instance, in graphs with 150 variables and 300 edges in expectation, LGES with CONSERVATIVEINSERT only makes ≈ 30 structural errors on average, which is twice as accurate as GES, which makes ≈ 60 structural errors. LGES also outperforms GES in smaller sample settings ($n \in \{500, 1000\}$) (Experiment D.1.5).

Q6. What is the difference between SAFEINSERT and CONSERVATIVEINSERT?

Answer. LGES can use either the SAFEINSERT or the CONSERVATIVEINSERT strategy to select INSERT operators in the forward phase; both result in improved accuracy and runtime relative to GES. We show that LGES with SAFEINSERT is asymptotically guaranteed to recover the true MEC (Cor. 1). However, it remains open whether the same is true of LGES with CONSERVATIVEINSERT, though we provide partial guarantees (Prop. C.1, C.2). Both strategies result in similar runtime, though CONSERVATIVEINSERT consistently has greater accuracy than SAFEINSERT across our experiments (Sec. 5, D).

Q7. How is causal discovery with prior knowledge different from initializing a causal discovery task to the hypothesized model?

Answer. Causal discovery with knowledge is the more general problem of using possibly misspecified prior knowledge in the process of causal discovery to aid the search. Initializing the search to a tentative model is one way to achieve this. However, initialization is not robust to misspecification in the assumptions and can result in worse runtime and accuracy than our approach of guiding the search using prior knowledge (Sec. 3.3, Experiment 5.2).

Q8. What is the difference between Greedy Interventional Equivalence Search (GIES) [24] and LGES + \mathcal{I} -ORIENT?

Answer. GIES and LGES are both score-based algorithms for learning from a combination of observational and interventional data. However, LGES has a few primary advantages over GIES. First, LGES with SAFEINSERT is guaranteed to recover the true interventional MEC (Cor. 1, Thm. 2) in the sample limit whereas GIES is not [58]. Second, both variants of LGES are up to 10 times faster than GIES, and LGES with CONSERVATIVEINSERT has accuracy competitive with GIES (Experiment 5.3). However, LGES uses interventional data to only to orient edges in a learned Markov equivalence class (in the \mathcal{I} -ORIENT procedure, Alg. 2), whereas GIES uses a combination of observational and experimental data throughout. We additionally introduce LGIES, which incorporates less greedy insertion into GIES and is thus both faster and more accurate than GIES (Experiment 5.3).

Q9. Can your work be combined with other causal discovery algorithms?

Answer. Several components of our approach are modular and can be combined with other causal discovery algorithms. For one, other algorithms in the GES family like FGES [43] and SGES [12] can be easily modified to use the SAFEINSERT or CONSERVATIVEINSERT strategies that we introduce; a simple extension would investigate the resulting changes in accuracy and runtime. For another, the \mathcal{I} -ORIENT procedure can be used to refine the observational MEC output by any causal discovery algorithm (not necessarily LGES) using interventional data.