
Learning Linear Causal Representations from General Environments: Identifiability and Intrinsic Ambiguity

Jikai Jin

Institute for Computational and Mathematical Engineering
Stanford University
Stanford, CA 94305
jkjin@stanford.edu

Vasilis Syrgkanis

Management Science and Engineering
Stanford University
Stanford, CA 94305
vsyrgk@stanford.edu

Abstract

We study causal representation learning, the task of recovering high-level latent variables and their causal relationships in the form of a causal graph from low-level observed data (such as text and images), assuming access to observations generated from multiple environments. Prior results on the identifiability of causal representations typically assume access to single-node interventions which is rather unrealistic in practice, since the latent variables are unknown in the first place. In this work, we consider the task of learning causal representation learning with data collected from *general environments*. We show that even when the causal model and the mixing function are both linear, there exists a *surrounded-node ambiguity* (SNA) [46] which is basically unavoidable in our setting. On the other hand, in the same linear case, we show that identification up to SNA is possible under mild conditions, and propose an algorithm, LiNGCR_eL which provably achieves such identifiability guarantee. We conduct extensive experiments on synthetic data and demonstrate the effectiveness of LiNGCR_eL in the finite-sample regime.

1 Introduction

Artificial intelligence (AI) has achieved tremendous success in various domains in the past decade [4, 40, 6]. However, current approaches are largely based on learning the *statistical* structures and relationships in the data that we observe. As a result, it is not surprising that these approaches often capture spurious statistical dependencies between different features, resulting in poor performance in the presence of test distribution shift [30, 22] or adversarial attacks [3, 50].

In view of these pitfalls, a recent line of work has explored the problem of *causal representation learning* (CRL) [34], the task of learning the causal relationships between high-level latent variables underlying our low-level observations. Notably, it is widely believed in cognitive psychology that humans take a causal approach to distill information from the world and make decisions to achieve their goals [37, 12, 19]. As a result, there is reason to believe that learning causal representations has the potential to significantly improve the power of AI, especially on tasks where performance lags far behind human level [17].

Despite such promise, a crucial challenge in CRL is the *identifiability* of the data generating process; in other words, given the data that we observe, can we uniquely identify the underlying causal model. It has been shown that given observational data (*i.e.*, i.i.d. data generated from a single environment), the model is already non-identifiable in strictly simpler settings where the latent variables are known to be independent [25, 26], or where there is no mixing function and one directly observes the latent variables [39]. As a result, existing algorithms for CRL with observational data [52, 53, 11] typically require additional assumptions on the structure of the underlying causal graph. A natural question that arises is what types of data do we need to acquire to make identification possible in the general case.

One line of works assumes access to counterfactual data [27, 48, 5], where some form of *weak supervision* is typically required. A common assumption here is that one observes data in *pairs*, where each pair of data is related via sharing part of the latent representation. However, such data is hard to acquire since it requires direct control on the latent representation.

Another line of works [1, 49, 7, 47] instead considers an interventional setting, where the learner observes data generated from multiple different environments. This is arguably a much more realistic setup and reflects common practices in robotics [24] and genomics [28, 43] applications. However, a vast majority of identifiability guarantees assume that each environment corresponds to *single-node, hard* interventions, which is defined as interventions that isolate a single latent variable from its causal parents. Again, this is quite a restrictive assumption because of two reasons. *First*, since the latent variables are unknown and need to be learned from data, it is unclear how to perform interventions that only affect one variable. *Second*, even if one can perform single-node interventions, it may not be feasible to artificially remove causal effects in the data generating processes. This issue is ubiquitous in real-world applications as pointed out in Campbell [8], Eberhardt [14], Eronen [15]. Motivated by these challenges, we make the following contributions in this paper:

- Assuming access to data collected from multiple environments, but not necessarily from single-node, hard interventions, we identify an intrinsic surrounded-node ambiguity (SNA) in learning the underlying causal representations. We show in [Theorem 3](#) that SNA is unavoidable even if (1) both the mixing function and the causal model are known to be linear and (2) one has access to single-node, soft interventions. This highlights a remarkable difference with existing literature which shows that perfect identification can be achieved with hard interventions.
- When the causal model and the mixing function are both linear, we prove in [Theorem 1](#) that identification up to SNA is achievable with $\mathcal{O}(d)$ diverse environments ([Assumption 4](#)), where d is the size of the latent causal graph. To the best of our knowledge, this is the first identification guarantee that applies to fully general environments and makes no assumption on their relationship or similarity. Interestingly, we also show in [Theorem 2](#) that one would require $\Omega(d^2)$ single-node soft interventions to achieve the same identification guarantee, indicating the benefit of learning from diverse environments.
- We propose an algorithm, LiNGCR_{EL}, in [Section 5](#) that provably recovers the ground-truth model up to SNA ([Theorem 4](#)) in the setting of [Theorem 1](#) when perfect information of the observation distributions is available. To demonstrate the effectiveness of LiNGCR_{EL} in finite-sample regime, we conduct extensive experiments on synthetic data, and our results reported in [Section 6](#) show that LiNGCR_{EL} is capable of recovering the true causal model up to SNA with high accuracy.

Due to space limit, proofs of all our statements and additional theoretical results are given in the appendix.

2 Preliminaries

We consider the standard setup of CRL from multiple environments $E \in \mathcal{E}$. Let $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ be the ground-truth causal graph which is directed and acyclic (DAG), where $\mathcal{V} = [d]$ and \mathcal{E} describes the causal relationship between different nodes. Each node corresponds to a latent variable $z_i \in \mathbb{R}$.

For any node i , we let $\text{pa}_{\mathcal{G}}(i)$, $\text{ch}_{\mathcal{G}}(i)$, $\text{ans}_{\mathcal{G}}(i)$ and $\text{nd}_{\mathcal{G}}(i)$ to be the set of all parents, children, ancestors and non-descendants of i in \mathcal{G} respectively. We also define $\overline{\text{pa}}_{\mathcal{G}}(i) = \text{pa}_{\mathcal{G}}(i) \cup \{i\}$ and similarly for $\overline{\text{ch}}_{\mathcal{G}}(i)$, $\overline{\text{ans}}_{\mathcal{G}}(i)$ and $\overline{\text{nd}}_{\mathcal{G}}(i)$. Assuming that all probability distributions have continuous

densities, the joint density of the latent variables \mathbf{z} can then be written as

$$p_E(\mathbf{z}) = \prod_{i=1}^d p_i^E(\mathbf{z}_i \mid \mathbf{z}_{\text{pa}_{\mathcal{G}}(i)}). \quad (1)$$

where p_i^E is the (unknown) latent generating distribution from environment E at node i . Here for a given vector \mathbf{v} , we write $\mathbf{v}_i = \mathbf{e}_i^\top \mathbf{v}$, and let $\mathbf{v}_S = (\mathbf{v}_i : i \in S) \in \mathbb{R}^{|S|}$.

The causal graph model with density given by (1) necessarily enjoys the following property:

Definition 1 (Causal Markov Condition). *For any node i , conditioned on $\mathbf{z}_{\text{pa}_{\mathcal{G}}(i)}$, \mathbf{z}_i is independent of $\mathbf{z}_{\text{nd}_{\mathcal{G}}(i)}$. As a consequence, for any node $i, j \in [d]$ and $S \subseteq [d]$, if S d -separates i from j (cf. Definition 7), then $\mathbf{z}_i \perp\!\!\!\perp \mathbf{z}_j \mid \mathbf{z}_S$.*

The latent variables \mathbf{z} are unknown to the learner. Instead, the learner has access to observations $\mathbf{x} \in \mathbb{R}^n$ ($n \geq d$) from all environments $E \in \mathfrak{E}$ that are related to the latent \mathbf{z} via an injective mixing function \mathbf{g} :

$$\mathbf{x} = \mathbf{g}(\mathbf{z}). \quad (2)$$

The main assumption here that the mixing function is the same across all environments:

Assumption 1. *All environments $E \in \mathfrak{E}$ share the same diffeomorphic mixing function $\mathbf{g} : \mathbb{R}^d \mapsto \mathbb{R}^n$.*

In CRL, the goal of the learner is to 1) recover the inverse of the mixing function $\mathbf{h} = \mathbf{g}^{-1}$ (often called the *unmixing* function) which allows recovering the latent variables given any observations, and, 2) recover the underlying causal graph \mathcal{G} . In the remaining part of this paper, we refer to $(\mathbf{h}, \mathcal{G})$ as the causal model to be learned. Obviously, there would be some ambiguities in learning $(\mathbf{h}, \mathcal{G})$. For example, choosing a different permutation of the nodes in the causal graph would lead to a different model, and so does element-wise transformations on each component \mathbf{h}_i of \mathbf{h} .

A line of recent works show that the ground-truth model can be identified up to these ambiguities in various settings, assuming access to single-node hard interventions [36, 49, 47]. On the other hand, some weaker notions of identifiability have also been proposed and studied in the literature [36, 46, 23] for single-node soft interventions. Here, we provide a generic definition of single-node soft interventions that we will rely on in this paper.

Definition 2. *We say that a collection of environments $\hat{\mathfrak{E}}$ is a set of (soft) interventions on a subset of latent variables $\{\mathbf{z}_j, j \in S\}$ if for any $i \in [d]$ and any $E_1, E_2 \in \hat{\mathfrak{E}}, E_1 \neq E_2$, we have $p_i^{E_1} = p_i^{E_2}$ if and only if $i \notin S$ (the notation p_i^E comes from (1)). Equivalently, we write $\mathcal{I}_{\mathbf{z}}^{\hat{\mathfrak{E}}} = S$.*

We note that soft interventions are very different from hard interventions, since they do not remove causal relationships between latent variables. The goal of this paper is to address the following question:

What is the best-achievable identification guarantee when hard interventions are not available, and what are the intrinsic ambiguities?

3 The surrounding set and a notion of identifiability

One may expect that identifiability with soft interventions is not much different from hard interventions, since soft interventions can approximate hard interventions with arbitrary accuracy. However, we will show that this is not the case. At a high level, hard intervention is more powerful than soft intervention because it is capable of isolating a latent variable from its direct cause while soft interventions is not, so soft interventions can sometimes fail to identify the true causal relationship from a mixture of causal effects.

To quantify what kind of ambiguities may arise, we can define the surrounding set for each node in a causal graph \mathcal{G} as follows:

Definition 3. (46, Definition 3) *For two nodes $i, j \in [d]$ in \mathcal{G} , we say that j is surrounded by i , or $i \in \text{sur}_{\mathcal{G}}(j)$ if $i \in \text{pa}_{\mathcal{G}}(j)$, and $\text{ch}_{\mathcal{G}}(j) \subseteq \text{ch}_{\mathcal{G}}(i)$. Moreover, we define $\overline{\text{sur}}_{\mathcal{G}}(j) = \text{sur}_{\mathcal{G}}(j) \cup \{j\}$.*

Intuitively, if there exists some $i \in \text{sur}_{\mathcal{G}}(j)$, then ambiguities may arise for the causal variable at node j , since any effect of j on any of its child k can also be interpreted as a mixture of the effect of i

and j . In [Appendix E](#) we discuss an example with three causal variables to further illustrate such ambiguities.

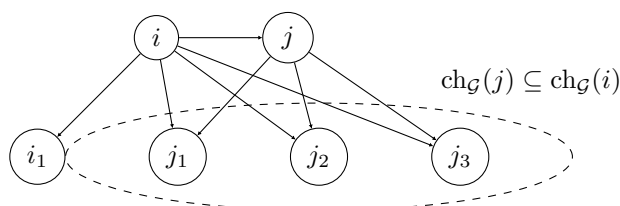


Figure 1: An illustration of [Definition 3](#); here $i \in \text{sur}_G(j)$.

[Definition 3](#) naturally induces the following relationship between causal models:

Definition 4. Using the notations in [Definition 10](#), we write $(\mathbf{h}, \mathcal{G}) \sim_{\text{sur}} (\hat{\mathbf{h}}, \hat{\mathcal{G}})$ if there exists a permutations π on $[d]$, and a diffeomorphism $\psi : \mathbb{R}^d \mapsto \mathbb{R}^d$ where the j -th component of ψ , denoted by $\psi_j(\mathbf{z})$, is a function of $\mathbf{z}_{\text{sur}_G(j)}$ for $\forall j \in [d]$, such that the following holds:

- For any $i, j \in [d]$, $i \in \text{pa}_G(j)$ if and only if $\pi(i) \in \text{pa}_{\hat{\mathcal{G}}}(\pi(j))$, and
- $\mathbf{P}_\pi \circ \hat{\mathbf{h}} = \psi \circ \mathbf{h}$, where \mathbf{P}_π is a permutation matrix satisfying $(\mathbf{P}_\pi)_{ij} = 1$ if $j = \pi(i)$ and $(\mathbf{P}_\pi)_{ij} = 0$ otherwise.

In other words, \sim_{sur} requires that the causal graph to be exactly the same up to some permutation of nodes, but allows each latent variable v_i to be a mixture of $\mathbf{z}_{\text{sur}_G(i)}$. Although not obvious from definition, one can actually check that \sim_{sur} defines an *equivalence relation* (see [Lemma 11](#)). Moreover, we will show in the following section that \sim_{sur} is in general the best that we can hope for in our problem setting.

4 Identifiability theory for linear CRL with general environments

In this section, we consider learning causal models from *general* environments. Specifically, we assume that the environments $E_k, k \in [K]$ share the same causal graph, but the dependencies between connected nodes (latent variables) are completely unknown, and, in contrast with existing literature on single-node interventions, we impose no similarity constraints on the environments. We begin our investigation of identifiability in this setting in the context of linear causal models with a linear mixing function.

4.1 Problem setup

Formally, we assume the following generative model in K distinct environments $\mathfrak{E} = \{E_k : k \in [K]\}$ with data generating process

$$\mathbf{z} = \mathbf{A}_k \mathbf{z} + \Omega_k^{\frac{1}{2}} \epsilon, \quad \mathbf{x} = \mathbf{G} \mathbf{z} \quad k \in [K], \quad (3)$$

where the matrix \mathbf{A}_k satisfies $(\mathbf{A}_k)_{ij} \neq 0$ if and only if $j \rightarrow i$ in \mathcal{G} . We refer to (\mathbf{A}_k, Ω_k) as the weight matrices of latent variables \mathbf{z} in the environment E_k . It is easy to see that [Assumption 1](#) in our general setup translates into the following assumption:

Assumption 2. The mixing matrix $\mathbf{G} \in \mathbb{R}^{n \times d}$ has full column rank. Equivalently, the unmixing matrix $\mathbf{H} = \mathbf{G}^\dagger$ has full row rank.

Let $\mathbf{B}_k = \Omega_k^{-\frac{1}{2}}(\mathbf{I} - \mathbf{A}_k), k \in [K]$, then we have $\epsilon = \mathbf{B}_k \mathbf{z} = \mathbf{B}_k \mathbf{H} \mathbf{x}$. Since in the linear case, there is an easy to see one-to-one correspondence between the matrix \mathbf{H} and the un-mixing function $\mathbf{x} \mapsto \mathbf{H} \mathbf{x}$, we abuse the notation and write $(\mathbf{H}, \mathcal{G})$ to represent the model instead of $(\mathbf{h}, \mathcal{G})$. Using \mathbf{h}_i to denote the i -th row of \mathbf{H} , the following lemma translates [Definition 4](#) the the linear setting:

Lemma 1. According to [Definition 4](#), $(\mathbf{H}, \mathcal{G}) \sim_{\text{sur}} (\hat{\mathbf{H}}, \hat{\mathcal{G}})$ if and only if there exists a permutation π on $[d]$, such that the following statements hold:

1. For all $i, j \in [d]$, $i \in \text{pa}_{\mathcal{G}}(j)$ if and only if $\pi(i) \in \text{pa}_{\hat{\mathcal{G}}}(\pi(j))$, and
2. For all $i \in [d]$, $\hat{\mathbf{h}}_i \in \text{span} \langle \mathbf{h}_j : \pi(j) \in \overline{\text{sur}}_{\mathcal{G}}(i) \rangle$.

We also need to make the following assumption on noise.

Assumption 3. The noise vector $\epsilon \in \mathbb{R}^d$ has independent components, at most one component is Gaussian distributed, and any two components have different distribution.

The non-gaussianity of the noise vectors is a typical assumption in causal discovery within linear models [9, 39] and is always assumed in the LinGAM setting [38]. The assumption that all components have a different distribution is not so standard, but is quite natural in real-world scenarios.

4.2 Identifiability guarantee

For each node $i \in [d]$ of \mathcal{G} , we use $\mathbf{w}_k(i)$ to be the *weight vector* of environment E_k at node i , i.e., $\mathbf{w}_k(i) = ((\mathbf{A}_k)_{ij} : j \in \text{pa}_{\mathcal{G}}(i)) \in \mathbb{R}^{|\text{pa}_{\mathcal{G}}(i)|}$. In other words, the structural equation for node i in environment k is of the form:

$$z_i = \mathbf{w}_k(i)^\top z_{\text{pa}_{\mathcal{G}}(i)} + \sqrt{\omega_{k,i,i}} \epsilon_i \quad (4)$$

To obtain our identifiability result, the main assumption we need to make is the non-degeneracy of the weights at each node:

Assumption 4. For each node $i \in [d]$ of \mathcal{G} , we have $\text{aff}(\mathbf{w}_k(i) : k \in [K]) = \mathbb{R}^{|\text{pa}_{\mathcal{G}}(i)|}$ where $\text{aff}(\cdot)$ denotes the affine hull. Equivalently, the weights $\mathbf{w}_k(i), k = 1, 2, \dots, K$ do not lie in a $(|\text{pa}_{\mathcal{G}}(i)| - 1)$ -dimensional hyperplane of $\mathbb{R}^{|\text{pa}_{\mathcal{G}}(i)|}$.

This assumption is quite mild since it only requires the weight vectors to be in general positions, and it holds with probability 1 if the weights at each node are sampled from continuous distributions. Moreover, as shown in Lemma 5, it is equivalent to the following assumption.

Assumption 5 (Node-level non-degeneracy). We say that the matrices $\{\mathbf{B}_k\}_{k=1}^K$ are node-level non-degenerate if for all node $i \in [d]$, we have $\dim \text{span} \langle (\mathbf{B}_k)_i : k \in [K] \rangle = |\text{pa}_{\mathcal{G}}(i)| + 1$, where $(\mathbf{B}_k)_i$ is the i -th row of \mathbf{B}_k .

In the following, we state our main result in this section, which shows that $K = d$ non-degenerate environments suffices for the model to be identifiable up to \sim_{sur} .

Theorem 1. Suppose that $K \geq d$ and we have access to observations generated from the linear causal model $(\mathbf{H}, \mathcal{G})$ across multiple environments $\mathfrak{E} = \{E_k : k \in [K]\}$ with observation distributions $\{\mathbb{P}_{\mathbf{x}}^E\}_{E \in \mathfrak{E}}$, and the data generating processes are given by (3). Let $(\hat{\mathbf{H}}, \hat{\mathcal{G}})$ be any candidate solution with the hypothetical data generating process

$$\mathbf{v} = \hat{\mathbf{A}}_k \mathbf{v} + \hat{\Omega}_k^{\frac{1}{2}} \hat{\epsilon}, \quad \mathbf{x} = \hat{\mathbf{H}}^\top \mathbf{v} \quad \text{in the environment } E_k$$

where $\hat{\mathbf{H}}$ has full row rank, such that

- (i) the observation distribution that this hypothetical model generates in E_k is exactly $\mathbb{P}_{\mathbf{x}}^{E_k}$;
- (ii) all environments share the same causal graph: $\forall k \in [K]$ and $i, j \in [d]$, $(\mathbf{A}_k)_{ij} \neq 0 \Leftrightarrow j \in \text{pa}_{\mathcal{G}}(i)$, $(\hat{\mathbf{A}}_k)_{ij} \neq 0 \Leftrightarrow j \in \text{pa}_{\hat{\mathcal{G}}}(i)$ and $\Omega_k, \hat{\Omega}_k$ are diagonal matrices with positive entries;
- (iii) $\{\mathbf{B}_k\}_{k=1}^K$ and $\left\{ \hat{\mathbf{B}}_k = \hat{\Omega}_k^{-\frac{1}{2}} (\mathbf{I} - \hat{\mathbf{A}}_k) \right\}_{k=1}^K$ are non-degenerate in the sense of Assumption 5;
- (iv) the noise variables ϵ and $\hat{\epsilon}$ satisfy Assumption 3.

Then we must have $(\mathbf{H}, \mathcal{G}) \sim_{\text{sur}} (\hat{\mathbf{H}}, \hat{\mathcal{G}})$.

The proof of Theorem 1 is given in Appendix H.1. In the next section, we will introduce an algorithm, LiNGReL, that provably recovers the ground-truth up to \sim_{sur} .

To the best of our knowledge, this is the first identifiability guarantee in the literature for CRL from general environments, even for the linear case. Our result is closely related but fundamentally different from Xie et al. [52, 53], Dong et al. [11] that consider the task of linear CRL using *observational data*. As discussed before, with observational data the causal graph can at best be identified up to Markov equivalence. As a result, one typically requires additional assumptions on the structure of the causal graph to obtain stronger guarantees. In contrast, we show that with data from multiple environments, exact recovery of the causal graph is possible without any structural assumptions.

Interestingly, while the fact that existing works focus on single-node interventions seem to suggest that learning from diverse environments is hard, it turns out that such diversity is actually helpful. Specifically, we show that in the worst case, $\Theta(d^2)$ interventions are required for identifying the ground-truth model under \sim_{sur} :

Theorem 2 (informal version of [Theorem 6](#)). *There exists a causal graph \mathcal{G} with $\Theta(d^2)$ edges, such that for any unmixing matrix $\mathbf{H} \in \mathbb{R}^{d \times n}$ with full row rank, any independent noise variables ϵ , and any $0 < s_i \leq |\text{pa}_{\mathcal{G}}(i)|$, $i \in [d]$, the ground-truth model $(\mathbf{H}, \mathcal{G})$ is non-identifiable up to \sim_{sur} with s_i soft interventions for node i , unless the (ground-truth and intervened) weights of the causal model lie in a null set (w.r.t the Lebesgue measure).*

A formal version and the proof of [Theorem 2](#) can be found in [Appendix H.2](#). On the other hand, by having d single-node interventions per node, [Assumption 5](#) can be satisfied as long as the weights are in general positions, so in this case we have $(\mathbf{H}, \mathcal{G}) \sim_{\text{sur}} (\hat{\mathbf{H}}, \hat{\mathcal{G}})$ by [Theorem 1](#). Therefore, [Theorems 1](#) and [6](#) together imply that $\Theta(d^2)$ single-node interventions are necessary and sufficient for identification up to \sim_{sur} .

Given that [Theorem 1](#) only guarantees identification up to \sim_{sur} that is strictly weaker than full identification, one might naturally ask whether [Theorem 1](#) can be further improved. Our last theorem in this section indicates that \sim_{sur} is indeed a fundamental barrier that exists even when we access to *single node, soft interventions*.

Theorem 3 (Counterpart to [Theorem 1](#), informal version of [Theorem 9](#)). *For any linear causal model $(\mathbf{H}, \mathcal{G})$ and any set of environments $\mathfrak{E} = \{E_k : k \in [K]\}$ such that all conditions in [Theorem 1](#) are satisfied, there must exist a candidate solution $(\hat{\mathbf{H}}, \hat{\mathcal{G}})$ and a hypothetical data generating process that satisfy the same set of conditions, but*

$$\frac{\partial v_i}{\partial z_j} \neq 0, \quad \forall j \in \overline{\text{sur}}_{\mathcal{G}}(i).$$

Moreover, if we additionally assume that the environments are groups of single-node soft interventions, then we can guarantee the existence of $(\hat{\mathbf{H}}, \hat{\mathcal{G}})$ and weight matrices which, besides the properties listed above, are also groups of single-node soft interventions.

5 LinGCR_EL: Algorithm for linear non-Gaussian causal representation learning

In this section, we introduce Linear Non-Gaussian Causal Representation Learning (LinGCR_EL), an algorithm that provably recovers the underlying causal graph and latent variables up to \sim_{sur} in the infinite-sample limit. At this point, it is instructive to recall the celebrated LiNGAM algorithm [38] for linear causal graph discovery. Different from their setting, we only observe some unknown linear mixture of the latent variables. Hence, running linear ICA as in LiNGAM only gives us $M_k = B_k \mathbf{H}$ rather than the weight matrix B_k itself.

The key idea in our approach is an effect cancellation scheme that allows us to determine the “remaining degree of freedom” (RDF) of any node (*a.k.a.* latent variable) given any subset of its ancestors. This scheme allows us to not only find a topological order of the nodes, but also figure out direct causes by tracking the changes of the RDF. In the following, we present the main steps of LinGCR_EL in more details.

Suppose that we are given samples of observations $\mathbf{X}^{(k)} = \{\mathbf{x}_i^{(k)}\}_{i=1}^N$, $k \in [K]$ where $\mathbf{x}_i^{(k)}$ is the i -th sample from the k -th environment.

Step 1. Recover the matrices $M_k = B_k \mathbf{H}$ Since $\epsilon = B_k \mathbf{z} = B_k \mathbf{H} \mathbf{x}$ in the k -th environment, so we can use any identification algorithm for linear ICA to recover the matrix M_k . Then we properly

rearrange the rows of M_k so that all $M_k \mathbf{x}$, $k = 1, 2, \dots, K$ correspond to the same permutation of noise variables. This step is quite standard and details can be found in [Appendix B.1](#).

Step 2. CRL based on M_k Now we have obtained $M_k = B_k H$, but the unmixing matrix H is still unknown. We propose [Algorithm 3](#) to learn H and the causal graph \mathcal{G} . The main part of [Algorithm 3](#) contains a loop that maintains a node set S which, we will show later, is ancestral, *i.e.*, $i \in S \Rightarrow \text{ans}_{\mathcal{G}}(i) \subseteq S$. In each round the algorithm finds a new node $i \notin S$ such that $\text{ans}_{\mathcal{G}}(i) \subseteq S$, and a subroutine [Identify-Parents](#) ([Algorithm 2](#)) is used to find all parents of i . After that, we append i into S and continue until S contains all nodes in \mathcal{G} . Finally, the rows of the mixing matrix H is obtained by intersections of properly-chosen row spaces of M_k .

Both [Algorithm 2](#) and [Algorithm 3](#) include a crucial step, which we call it *orthogonal projection*, as described in [Algorithm 1](#). At a high level, it helps determine the minimal RDF for z_i after fixing the latent variables z_S , and this exactly corresponds to the number of parents of z_i that are not in z_S . We provide a simple example in [Appendix E.2](#) to illustrate why this approach works.

The following result states that [Algorithm 3](#) can recover the ground-truth causal model up to \sim_{sur} :

Theorem 4. Suppose that $M_k, k \in [K]$ are perfectly identified in [Step 1](#). Let $(\hat{H}, \hat{\mathcal{G}})$ be the solution returned by [Algorithm 3](#), then we must have $(H, \mathcal{G}) \sim_{\text{sur}} (\hat{H}, \hat{\mathcal{G}})$.

The full proof of [Theorem 4](#) is given in [Appendix H.3](#). It crucially relies on the following two propositions that reveal how [Algorithm 3](#) and the subroutine [Algorithm 2](#) work.

Algorithm 1 Orthogonal-projections

```

1: Input: Ordered set  $S = \{s_1, s_2, \dots, s_m\} \subseteq [d]$ , index  $i \notin S$ , matrices  $M_k \in \mathbb{R}^{d \times n}$ ,  $k \in [K]$ 
2: Output: Set of vectors  $\{p_k\}_{k=1}^K$ 
3: for  $k \leftarrow 1$  to  $K$  do
4:    $W \leftarrow \text{span}(\{(M_k)_s : s \in S\})$   $\triangleright (M_k)_s$  is the  $s$ -th row of  $M_k$ 
5:    $p_k \leftarrow \text{proj}_{W^\perp}((M_k)_i)$ 
6: end for

```

Proposition 1. The following two propositions hold for [Algorithm 3](#):

- $\text{ans}_{\mathcal{G}}(i) \subseteq S \Leftrightarrow$ the *if* condition in line 8 of [Algorithm 3](#) is fulfilled;
- the set S maintained in [Algorithm 3](#) is always an ancestral set, in the sense that $j \in S \Rightarrow \text{ans}_{\mathcal{G}}(j) \subseteq S$.

Proposition 2. Given any ordered ancestral set S that contains $\text{pa}_{\mathcal{G}}(i)$ for some $i \notin S$, [Algorithm 2](#) returns a set $P_i \subseteq S$ that is exactly $\text{pa}_{\mathcal{G}}(i)$.

Algorithm 2 Identify-Parents

```

1: Input: An ordered set  $S = \{s_1, s_2, \dots, s_m\} \subseteq [d]$ , a node  $i \notin S$  and matrices  $M_k, k \in [K]$ 
2: Output: The parent set  $P_i$  of node  $i$ 
3:  $P_i \leftarrow \emptyset$ 
4: for  $m' \leftarrow 0$  to  $m$  do
5:    $\{p_k\}_{k=1}^K \leftarrow \text{Orthogonal-projections}(\{s_j : j \leq m'\}, i, \{M_k\}_{k \in [K]})$ 
6:    $r_{m'} \leftarrow \dim \text{span}(\{p_k : k \in [K]\})$ 
7:   if  $m' \geq 1$  and  $r_{m'} = r_{m'-1} - 1$  then
8:      $P_i \leftarrow P_i \cup \{m'\}$ 
9:   end if
10: end for

```

6 Experiments

In this section, we present our experimental setup and results for LiNGCReL. Note that LiNGCReL as described in the previous section only works in the population regime. When the number of samples is limited, two main challenges in implementing LiNGCReL are to accurately compute the dimension

Algorithm 3 Learn-Causal-Model

```
1: Input: Matrices  $\mathbf{M}_k, k \in [K]$ 
2: Output: The edge set  $\mathcal{E}$  on the vertex set  $[d]$  and the mixing matrix  $\hat{\mathbf{H}}$ 
3:  $S \leftarrow \emptyset;$   $\triangleright S$  is an ordered set of nodes
4:  $\mathcal{E} \leftarrow \emptyset;$   $\triangleright \mathcal{E}$  is the edge set
5: while  $|S| < d$  do
6:   for  $i \notin S$  do
7:      $\{\mathbf{p}_k\}_{k=1}^K \leftarrow \text{Orthogonal-projections}(S, i, \{\mathbf{M}_k\}_{k \in [K]})$ 
8:     if  $\dim \text{span} \langle \mathbf{q}_k : k \in [K] \rangle = 1$  then
9:       break  $\triangleright$  Proposition 1 guarantees that such an  $i$  must exist
10:    end if
11:  end for
12:   $P_i \leftarrow \text{Identify-Parents}(S, i)$ 
13:   $S \leftarrow S \cup \{i\}$ 
14:   $\mathcal{E} \leftarrow \mathcal{E} \cup \{(j, i) : j \in P_i\}$ 
15: end while
16: for  $i = 1$  to  $d$  do
17:    $E_i \leftarrow \text{span} \langle (\mathbf{M}_k)_i : k \in [K] \rangle$ 
18: end for
19: for  $i = 1$  to  $d$  do
20:    $\hat{\mathbf{h}}_i \leftarrow$  any non-zero vector in  $(\cap_{j:(i,j) \in \mathcal{E}} E_j) \cap E_i$ 
21: end for
22:  $\hat{\mathbf{H}} \leftarrow [\hat{\mathbf{h}}_1^\top, \hat{\mathbf{h}}_2^\top, \dots, \hat{\mathbf{h}}_d^\top]^\top$ 
```

of a subspace (line 6 of Algorithm 2 and line 8 of Algorithm 3), and to find a vector in the intersection of multiple subspaces (line 20, Algorithm 3). Due to space limit, the implementation details are described in Appendix B.2.

Experimental setup. We generate the independent noise variables from generalized Gaussian distributions $p_\beta(x) \propto \exp(-|x|^\beta)$ with parameters $\beta_k = 0.2k^2, k = 1, 2, \dots, d$, multiplied by normalization constants to make their variances equal to 1. The ground-truth causal graph is generated by first fixing a total order of the vertices, say $1, 2, \dots, d$, then add directed edges $i \rightarrow j (i < j)$ according to i.i.d. Bernoulli(p) distributions, where $p \in (0, 1)$. The non-zero entries of matrices \mathbf{B}_k and \mathbf{H} are all generated independently from Gaussian distributions. For simplicity, we focus on the case $n = d$ since recovery of the latent graphs only requires information from d components of \mathbf{x} .

Metrics of estimation error. Since CRL seeks to learn both the causal graphs and the latent variables, for each output of our algorithm we first check if it exactly recovers the ground-truth causal graph. Then, recall that the latent variables and the observations are related by $\mathbf{z} = \mathbf{H}\mathbf{x}$, given any output unmixing matrix $\hat{\mathbf{H}}$ from Algorithm 3, we define the relative estimation error Δ_i for \mathbf{z}_i as the solution of the following optimization problem:

$$\min \|\Delta\|_\infty \quad s.t. \Delta_i = \frac{\left\| \text{proj}_{\text{span} \langle \mathbf{h}_j : j \in \overline{\text{sur}}_{\mathcal{G}}(i)}(\hat{\mathbf{h}}_i) \right\|_2}{\left\| \hat{\mathbf{h}}_i \right\|_2}, \quad (5)$$

$$\hat{\mathbf{H}} = \mathbf{P}\hat{\mathbf{H}} \text{ for some signed permutation matrix } \mathbf{P}.$$

where signed permutation is allowed here since the noise distribution in our experiments is symmetric and the order of latent variables $\mathbf{z}_i, i = 1, 2, \dots, d$ does not matter. We refer to the errors Δ_i defined in (5) as the *SNA error*. The SNA error measures how much of the row $\hat{\mathbf{h}}_i$ that we learn is contained in the span of the ground-truth rows $\mathbf{h}_j, j \in \overline{\text{sur}}_{\mathcal{G}}(i)$. Indeed, recall that given any observation \mathbf{x} , the ground-truth latent variable is $\mathbf{z} = \mathbf{H}\mathbf{x}$ while our algorithm outputs $\hat{\mathbf{v}}_i = \hat{\mathbf{h}}_i^\top \mathbf{x}$, so the SNA error essentially captures whether the recovered latent variable is close to some linear mixture of latent variables in the effect-dominating set of i . When the SNA error is zero for some node i , we know that the recovered latent variable at node i is exactly a linear mixture of the ground-truth latent variables in $\overline{\text{sur}}_{\mathcal{G}}(i)$, according to Lemma 1.

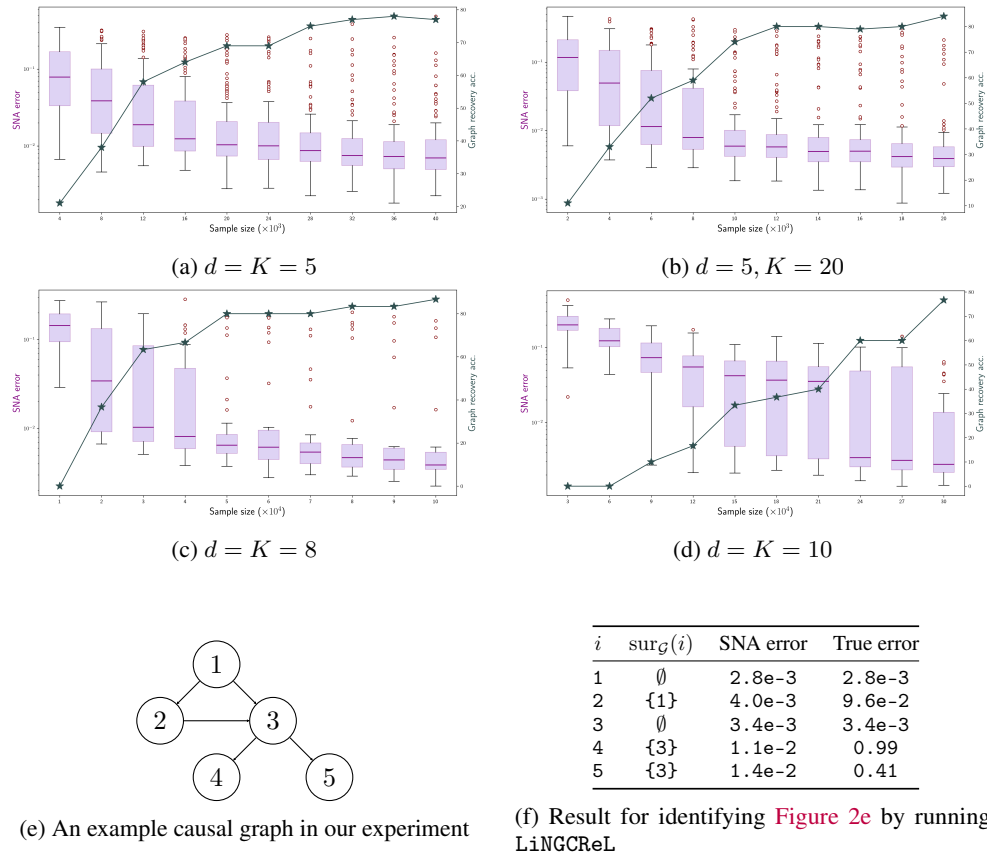


Figure 2: *First two rows*: plots of SNA Error and graph recovery accuracy achieved by LiNGCReL as functions of sample size (per environment) for different choices of graph size d and number of environments K . *Third row*: an example of causal graph generated in our experiments, and the estimation error of LiNGCReL for each node.

We also define the *true error* for estimating each latent variable. Formally, let $\hat{\mathbf{H}}$ be the unmixing matrix that corresponds to the solution of (5), then we define the true estimation error $\tilde{\Delta}_i$ of \mathbf{z}_i as

$$\tilde{\Delta}_i = \left\| (\mathbf{I} - \mathbf{h}_i \mathbf{h}_i^\top) \hat{\mathbf{h}}_i \right\|_2. \quad (6)$$

Results. We randomly sample 100 causal models with size $d = 5$, 30 causal models with size $d = 8$ and 30 causal models of size $d = 10$. In light of Theorem 1, for each $d \in \{5, 8, 10\}$, we sample data from $K = d$ randomly chosen environments; for $d = 5$ we also consider $K = 20$ to study how different choices of K can affect the result. We run LiNGCReL for each model with different sample sizes, compute the SNA error and true error of the obtained solution from (5) and (6) respectively for each latent variable, and check whether the ground-truth causal graph is exactly recovered.

Figure 2 shows how the average SNA error (over all latent variables) and the accuracy of graph recovery changes when sample size grows. We can see LiNGCReL successfully recovers about 80% of all models within each category, and the median of the average SNA error is smaller than 1%. Moreover, by comparing Figure 2a with Figure 2b, one can observe that if we fix the total number of samples but choose a larger K (i.e., fewer samples per environment), LiNGCReL can still achieve the same level of performance compared with the choice $K = d$. Intuitively, this is because $K \gg d$ vectors sampled from an r ($r \leq d$) dimensional subspace are unlikely to approximately lie in an $(r - 1)$ -dimensional subspace, so that the calculation of line 6 of Algorithm 2 and line 8 of Algorithm 3 can be more accurate. We leave a better and quantitative understanding of the trade-off between d and K to future work.

SNA error v.s. true error. To understand the implication of our theory, we dive deeper by looking into the learning outcome of LiNGCReL on a specific model, of which the causal graph is shown in

Figure 2e. In **Figure 2f**, we list the surrounding set of each node and the corresponding SNA error and true error. We can see that if $\text{sur}_G(i) = \emptyset$, the two errors equal and both are small, but if $\text{sur}_G(i) \neq \emptyset$, the true error is much larger than the SNA error. This indicates that LiNGCReL indeed learns the ground-truth model up to \sim_{sur} , as **Theorem 1** predicts.

7 Conclusions

This paper studies the limit of learning identifiable causal representations using data from multiple environments. When hard interventions are not available, we provide theory and algorithm for identification up to SNA, and also show that SNA is an intrinsic ambiguity in our setting.

It is interesting to further investigate the setting where we do not assume that the causal model is linear. Moreover, it is important to understand the concrete form of available interventions in real-world applications. For instance, it is suggested that for single-cell genomics, the intervention is sometimes a "mixture" of hard and soft interventions, and sometimes can even reverse the direction of an edge [43]. Modelling such more complicated interventions appears to be crucial to reveal the underlying causal mechanisms in real-world problems.

Acknowledgments and Disclosure of Funding

VS is supported by NSF Award IIS-2337916 and a 2023 Google Research Scholar Award.

References

- [1] Kartik Ahuja, Divyat Mahajan, Yixin Wang, and Yoshua Bengio. Interventional causal representation learning. In *International Conference on Machine Learning*, pages 372–407. PMLR, 2023.
- [2] Kartik Ahuja, Amin Mansouri, and Yixin Wang. Multi-domain causal representation learning via weak distributional invariances. *arXiv preprint arXiv:2310.02854*, 2023.
- [3] Naveed Akhtar and Ajmal Mian. Threat of adversarial attacks on deep learning in computer vision: A survey. *Ieee Access*, 6:14410–14430, 2018.
- [4] Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8): 1798–1828, 2013.
- [5] Johann Brehmer, Pim De Haan, Phillip Lippe, and Taco S Cohen. Weakly supervised causal representation learning. *Advances in Neural Information Processing Systems*, 35:38319–38331, 2022.
- [6] Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*, 2023.
- [7] Simon Buchholz, Goutham Rajendran, Elan Rosenfeld, Bryon Aragam, Bernhard Schölkopf, and Pradeep Ravikumar. Learning linear causal representations from interventions under general nonlinear mixing. *arXiv preprint arXiv:2306.02235*, 2023.
- [8] John Campbell. An interventionist approach to causation in psychology. *Causal learning: Psychology, philosophy, and computation*, pages 58–66, 2007.
- [9] Pierre Comon. Independent component analysis, a new concept? *Signal processing*, 36(3): 287–314, 1994.
- [10] Gregory F Cooper and Changwon Yoo. Causal discovery from a mixture of experimental and observational data. In *Proceedings of the Fifteenth conference on Uncertainty in artificial intelligence*, pages 116–125, 1999.

- [11] Xinshuai Dong, Biwei Huang, Ignavier Ng, Xiangchen Song, Yujia Zheng, Songyao Jin, Roberto Legaspi, Peter Spirtes, and Kun Zhang. A versatile causal discovery framework to allow causally-related hidden variables. In *The Twelfth International Conference on Learning Representations*, 2023.
- [12] Kevin N Dunbar and Jonathan A Fugelsang. Causal thinking in science: How scientists and students interpret the unexpected. In *Scientific and technological thinking*, pages 57–79. Psychology Press, 2004.
- [13] Frederick Eberhardt. Almost optimal intervention sets for causal discovery. In *Proceedings of the Twenty-Fourth Conference on Uncertainty in Artificial Intelligence*, pages 161–168, 2008.
- [14] Frederick Eberhardt. Direct causes and the trouble with soft interventions. *Erkenntnis*, 79: 755–777, 2014.
- [15] Markus I Eronen. Causal discovery and the problem of psychological interventions. *New Ideas in Psychology*, 59:100785, 2020.
- [16] Ronald Aylmer Fisher et al. The design of experiments. *The design of experiments.*, (7th Ed), 1960.
- [17] Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias Bethge, and Felix A Wichmann. Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2(11):665–673, 2020.
- [18] Alain Hauser and Peter Bühlmann. Two optimal strategies for active learning of causal models from interventional data. *International Journal of Approximate Reasoning*, 55(4):926–939, 2014.
- [19] Keith J Holyoak and Patricia W Cheng. Causal learning and inference as a rational process: The new synthesis. *Annual review of psychology*, 62:135–163, 2011.
- [20] Antti Hyttinen, Frederick Eberhardt, and Patrik O Hoyer. Experiment selection for causal discovery. *Journal of Machine Learning Research*, 14:3041–3071, 2013.
- [21] Ilyes Khemakhem, Diederik Kingma, Ricardo Monti, and Aapo Hyvarinen. Variational autoencoders and nonlinear ica: A unifying framework. In *International Conference on Artificial Intelligence and Statistics*, pages 2207–2217. PMLR, 2020.
- [22] Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Lanus Phillips, Irena Gao, et al. Wilds: A benchmark of in-the-wild distribution shifts. In *International Conference on Machine Learning*, pages 5637–5664. PMLR, 2021.
- [23] Wendong Liang, Armin Kekić, Julius von Kügelgen, Simon Buchholz, Michel Besserve, Luigi Gresele, and Bernhard Schölkopf. Causal component analysis. *arXiv preprint arXiv:2305.17225*, 2023.
- [24] Phillip Lippe, Sara Magliacane, Sindy Löwe, Yuki M Asano, Taco Cohen, and Efstratios Gavves. Biscuit: Causal representation learning from binary interactions. *arXiv preprint arXiv:2306.09643*, 2023.
- [25] Francesco Locatello, Stefan Bauer, Mario Lucic, Gunnar Raetsch, Sylvain Gelly, Bernhard Schölkopf, and Olivier Bachem. Challenging common assumptions in the unsupervised learning of disentangled representations. In *international conference on machine learning*, pages 4114–4124. PMLR, 2019.
- [26] Francesco Locatello, Stefan Bauer, Mario Lucic, Gunnar Rätsch, Sylvain Gelly, Bernhard Schölkopf, and Olivier Bachem. A sober look at the unsupervised learning of disentangled representations and their evaluation. *The Journal of Machine Learning Research*, 21(1):8629–8690, 2020.
- [27] Francesco Locatello, Ben Poole, Gunnar Rätsch, Bernhard Schölkopf, Olivier Bachem, and Michael Tschannen. Weakly-supervised disentanglement without compromises. In *International Conference on Machine Learning*, pages 6348–6359. PMLR, 2020.

- [28] Romain Lopez, Natasa Tagasovska, Stephen Ra, Kyunghyun Cho, Jonathan Pritchard, and Aviv Regev. Learning causal representations of single cells via sparse mechanism shift modeling. In *Conference on Causal Learning and Reasoning*, pages 662–691. PMLR, 2023.
- [29] Chaochao Lu, Yuhuai Wu, José Miguel Hernández-Lobato, and Bernhard Schölkopf. Invariant causal representation learning for out-of-distribution generalization. In *International Conference on Learning Representations*, 2021.
- [30] Yaniv Ovadia, Emily Fertig, Jie Ren, Zachary Nado, David Sculley, Sebastian Nowozin, Joshua Dillon, Balaji Lakshminarayanan, and Jasper Snoek. Can you trust your model’s uncertainty? evaluating predictive uncertainty under dataset shift. *Advances in neural information processing systems*, 32, 2019.
- [31] Judea Pearl. *Causality*. Cambridge university press, 2009.
- [32] Geoffrey Roeder, Luke Metz, and Durk Kingma. On linear identifiability of learned representations. In *International Conference on Machine Learning*, pages 9030–9039. PMLR, 2021.
- [33] Mark Rudelson and Roman Vershynin. Smallest singular value of a random rectangular matrix. *Communications on Pure and Applied Mathematics: A Journal Issued by the Courant Institute of Mathematical Sciences*, 62(12):1707–1739, 2009.
- [34] Bernhard Schölkopf, Francesco Locatello, Stefan Bauer, Nan Rosemary Ke, Nal Kalchbrenner, Anirudh Goyal, and Yoshua Bengio. Toward causal representation learning. *Proceedings of the IEEE*, 109(5):612–634, 2021.
- [35] J Schwartz. The formula for change in variables in a multiple integral. *The American Mathematical Monthly*, 61(2):81–85, 1954.
- [36] Anna Seigal, Chandler Squires, and Caroline Uhler. Linear causal disentanglement via interventions. *arXiv preprint arXiv:2211.16467*, 2022.
- [37] David R Shanks and Anthony Dickinson. Associative accounts of causality judgment. In *Psychology of learning and motivation*, volume 21, pages 229–261. Elsevier, 1988.
- [38] Shohei Shimizu, Patrik O Hoyer, Aapo Hyvärinen, Antti Kerminen, and Michael Jordan. A linear non-gaussian acyclic model for causal discovery. *Journal of Machine Learning Research*, 7(10), 2006.
- [39] Ricardo Silva, Richard Scheines, Clark Glymour, Peter Spirtes, and David Maxwell Chickering. Learning the structure of linear latent variable models. *Journal of Machine Learning Research*, 7(2), 2006.
- [40] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.
- [41] Peter Spirtes, Clark N Glymour, and Richard Scheines. *Causation, prediction, and search*. 2000.
- [42] Michael Strevens. Review of woodward," making things happen", 2007.
- [43] Alejandro Tejada-Lapuerta, Paul Bertin, Stefan Bauer, Hananeh Aliee, Yoshua Bengio, and Fabian J Theis. Causal machine learning for single-cell genomics. *arXiv preprint arXiv:2310.14935*, 2023.
- [44] Robert E Tillman and Frederick Eberhardt. Learning causal structure from multiple datasets with similar variable sets. *Behaviormetrika*, 41(1):41–64, 2014.
- [45] Simon Tong and Daphne Koller. Active learning for structure in bayesian networks. In *International joint conference on artificial intelligence*, volume 17, pages 863–869. Citeseer, 2001.

- [46] Burak Varici, Emre Acarturk, Karthikeyan Shanmugam, Abhishek Kumar, and Ali Tajer. Score-based causal representation learning with interventions. *arXiv preprint arXiv:2301.08230*, 2023.
- [47] Burak Varıcı, Emre Acartürk, Karthikeyan Shanmugam, and Ali Tajer. General identifiability and achievability for causal representation learning. *arXiv preprint arXiv:2310.15450*, 2023.
- [48] Julius Von Kügelgen, Yash Sharma, Luigi Gresele, Wieland Brendel, Bernhard Schölkopf, Michel Besserve, and Francesco Locatello. Self-supervised learning with data augmentations provably isolates content from style. *Advances in neural information processing systems*, 34: 16451–16467, 2021.
- [49] Julius von Kügelgen, Michel Besserve, Wendong Liang, Luigi Gresele, Armin Kekić, Elias Bareinboim, David M Blei, and Bernhard Schölkopf. Nonparametric identifiability of causal representations from unknown interventions. *arXiv preprint arXiv:2306.00542*, 2023.
- [50] Tony Tong Wang, Adam Gleave, Tom Tseng, Kellin Pelrine, Nora Belrose, Joseph Miller, Michael D Dennis, Yawen Duan, Viktor Pogrebniak, Sergey Levine, et al. Adversarial policies beat superhuman go ais. 2023.
- [51] James Woodward. *Making things happen: A theory of causal explanation*. Oxford university press, 2005.
- [52] Feng Xie, Ruichu Cai, Biwei Huang, Clark Glymour, Zhifeng Hao, and Kun Zhang. Generalized independent noise condition for estimating latent variable causal graphs. *Advances in neural information processing systems*, 33:14891–14902, 2020.
- [53] Feng Xie, Biwei Huang, Zhengming Chen, Yangbo He, Zhi Geng, and Kun Zhang. Identification of linear non-gaussian latent hierarchical structure. In *International Conference on Machine Learning*, pages 24370–24387. PMLR, 2022.
- [54] Jiaqi Zhang, Chandler Squires, Kristjan Greenewald, Akash Srivastava, Karthikeyan Shanmugam, and Caroline Uhler. Identifiability guarantees for causal disentanglement from soft interventions. *arXiv preprint arXiv:2307.06250*, 2023.

A Related works

The interventionist approach to causation For the problem of causal graph discovery, it is well-known that the underlying causal structure is non-identifiable given only “passively observed” (equivalently, *i.i.d.*) data alone. As a result, randomized controlled experiments [16] is often used to infer causality. These experiments typically take the form of interventions [41, 31], *i.e.*, manipulations on the “natural state” of the system of interest. Early works [51, 42] define the “hard” (also called “surgical” or “arrow-breaking”) interventions in which the value of the intervened variable is entirely determined by the experimenter, thereby removing the dependence of this variable on its direct causes. This type of intervention is arguably the most natural one to consider, and following this definition, a line of works explore sufficient conditions for designing experiments that guarantee identifiability of the causal model in various settings [10, 45, 13, 20, 18].

Intervention *v.s.* passive observation While extensive works demonstrate the success of the interventionist approach, it faces several key challenges that significantly limit its applicability. First, Eberhardt [14] finds that in the presence of unobserved variables, certain causal structures are indistinguishable if we only perform hard interventions. This issue can be resolved by performing soft interventions *i.e.*, interventions that do not remove the dependency on direct causes but only changes the conditional distribution. Second, as pointed out in [44], interventions — whether hard or soft — are often expensive or even infeasible to perform in practice. For example, a psychological intervention is likely to affect multiple psychological variables simultaneously Eronen [15]. As a result, [44] returns to the “passive observation” setting but with multiple datasets with overlapping latent variables.

Interventional causal representation learning Motivated by the interventionist literature in causal graph discovery, a recent line of works [1, 36, 46, 49, 7, 54, 47] consider performing interventions to resolve the non-identifiability issue in causal representation learning [25]. Roughly speaking, these result indicate that identification (possibly with some ambiguities) is possible if one can perform intervention on every latent variable. However, it is unclear how to perform such interventions in practice, given that the underlying latent variables are unknown. Khemakhem et al. [21], Lu et al. [29], Roeder et al. [32] do not require single-node interventions to achieve identifiability, but assumes that the joint distribution of latent variables in each environment lie in a certain exponential family. This assumption can be understood as a prior on the latent variables, but it is unclear when or why it is reasonable to make in reality. Recently, Ahuja et al. [2] considers learning causal representations from multiple domains that relate to each other via an invariance constraint on the subset \mathcal{S} of *stable* latent variables, and they prove identification up to affine mixtures within \mathcal{S} .

B Experiment details for Section 6

B.1 Details for step 1 in Section 5

Since $\epsilon = B_k z = B_k H x$ in the k -th environment, so we can use any identification algorithm for linear ICA to recover the matrix M_k . Note that while standard linear ICA algorithms only apply to the case where $n = d$, for $n > d$ we can arbitrarily choose d principal components of x to reduce it to the $n = d$ case. This is without loss of generality, since when $n > d$ there is redundant information in x .

After recovering M_k for each k by running linear ICA, we still do not know whether each $M_k x$ corresponds to the same permutation of the ground-truth noise variables ϵ . To resolve this issue, we choose test function Ψ mapping any distribution on \mathbb{R} to a deterministic real value, which we expect to take different values for different ϵ_i ’s. We choose $\Psi(\mathbb{P}) = \mathbb{P}[|X| \leq 1]$ in our experiments. For all $k \geq 2$, we calculate the Ψ value of each component of the d -dimensional empirical distribution $\hat{\mathbb{P}}_k = \frac{1}{N} \sum_{i=1}^N \mathbb{1}_{M_k x_i^{(k)}}$, and choose a permutation π_k to rearrange them in increasing order. Then, we rearrange the columns of M_k using the same permutation π_k . This procedure would asymptotically produce correct alignments as long as $\Psi(\epsilon_i), i \in [d]$ are different, and we find that it empirically works well.

Alternatively, this alignment step can be done as follows: for each pair of environments (E_1, E_t) , and for each pair of nodes (i, j) , we calculate the distribution distance between ϵ_i in environment E_1 and ϵ_j in environment E_t , based on some notion of distribution distance (*e.g.* kernel maximum

mean discrepancy). Then we find the min-cost perfect matching, where the cost of an edge is the distribution distance.

B.2 Details for the implementation of LiNGCReL in the finite-sample regime

Although LiNGCReL provably works in the population regime, it faces several challenges when there is only a finite number of samples:

- First, since rank is not a continuous function, it is sensitive to finite-sample estimation errors. In our implementation of [Algorithm 3](#), in each iteration we instead choose $i \notin S$ that has the largest ratio between the first and second singular values of $[q_1, q_2, \dots, q_K]$. And in line 6 of [Algorithm 2](#), we introduce a hyper-parameter $\mathfrak{t}1$ such that the matrix $[q_1, q_2, \dots, q_K]$ is considered to have rank $r_{m'}-1$ if its $r_{m'}$ -th singular value is smaller than $\mathfrak{t}1$. Since the smallest singular value of a random matrix $A \in \mathbb{R}^{K \times m}$ ($K \geq m$) is at the order of $\sqrt{K} - \sqrt{m-1}$ with high probability [33], when $K = d$ one shall choose $\mathfrak{t}1 \sim \sqrt{d} - \sqrt{d-1} = \mathcal{O}\left(\frac{1}{\sqrt{d}}\right)$. On the other hand, for larger K we can correspondingly choose a larger $\mathfrak{t}1$. Note that a small $\mathfrak{t}1$ potentially has the risk of being dominated the noise in the estimation, which means that we need more samples per environment to reduce the noise. In contrast, for larger $\mathfrak{t}1$ the estimation is more robust to noise and we can use fewer samples.
- Second, finite-sample estimation errors of M_k make it harder to obtain h_i in [Algorithm 3](#) of [Algorithm 3](#). We implement this step in the following way: first let Q_j be the orthogonal projection matrix onto E_j^\perp i.e., $Q_j^\top x = \text{proj}_{E_j^\perp}(x)$, then choose h_i to be the singular vector of $\sum_{j:(j,i) \in \mathcal{E} \text{ or } j=i} Q_j^\top Q_j$ that corresponds to the smallest singular value (including zero). Indeed, in the noiseless case we would have $\left(\sum_{j:(j,i) \in \mathcal{E} \text{ or } j=i} Q_j^\top Q_j\right) h_i = 0$ if and only if $h_i \in \left(\cap_{j:(i,j) \in \mathcal{E}} E_j\right) \cap E_i$.

C Further experiment results

SNA error v.s. true error We plot the SNA error v.s. true error achieved by LiNGCReL in [Figure 3](#). We observe that

- For most nodes, SNA error is exactly equal to the true error and both errors are small, indicating that the corresponding latent variables have been successfully learned by LiNGCReL.
- The remaining nodes typically have true error much larger than SNA error. This indicates that there exists some ambiguities at these nodes in the sense that $\text{sur}_G(i) \neq \emptyset$. Note that the true error for many nodes are close to 1; one possible reason is that one selects the wrong singular vector in the second part of [Appendix B.2](#), so that it is orthogonal to the ground-truth vector.

Sensitivity of LiNGCReL to the hyperparameter $\mathfrak{t}1$ We examine how different choices of $\mathfrak{t}1$ would affect the performance of LiNGCReL. Specifically, we run LiNGCReL on the 100 models with size $d = 5$ and number of environments $K = 5$ sampled in [Section 6](#) with $\mathfrak{t}1 \in \{0.1, 0.15, 0.2, 0.25, 0.3\}$ and the results are reported in [Figure 4](#). We can see that the performance is actually quite sensitive to $\mathfrak{t}1$.

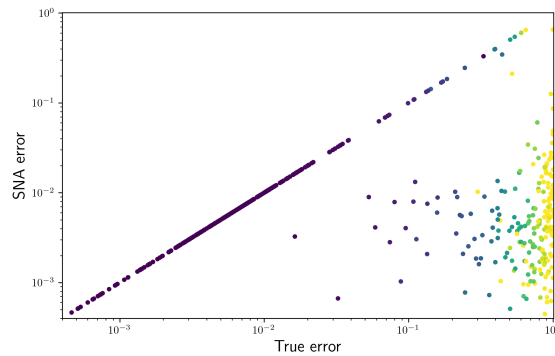


Figure 3: Comparing SNA error with true error for the 500 latent variables in the 100 graphs of size $d = 5$ that we sample in Section 6.

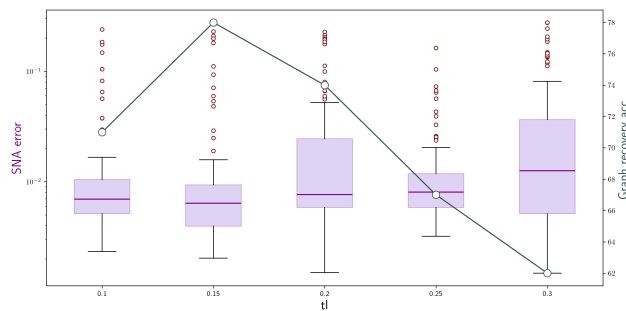


Figure 4: Performance of LiNGReL as a function of τ_1 . $\tau_1 = 0.15$ achieves the best performance in terms of both SNA error and graph recovery accuracy.

D Background on causal representation learning

It is common to assume some axioms on what kind of (conditional) dependency information is encoded in a causal graph (see 41, Section 3.4 for a detailed discussion). The most natural one is the Causal Markov Condition introduced in Definition 1 that gives sufficient conditions for conditional independence via d -separation. We introduce the formal definition of d -separation below:

Definition 5 (paths and colliders). Let i, j be two nodes of a DAG \mathcal{G} , a path is a sequence of nodes $i_0 = i, i_1, \dots, i_k = j$ such that there is an edge (in either direction) between i_j and i_{j+1} , $j = 0, 1, \dots, k-1$. A node i_j is called a collider on this path if $i_j \in \text{ch}_{\mathcal{G}}(i_{j-1}) \cap \text{ch}_{\mathcal{G}}(i_{j+1})$.

Definition 6 (blocked path). A path in a DAG \mathcal{G} between node i and node j is said to be blocked by a node set S if either of the following holds:

- there exists a node v on the path that is in S but not a collider; or
- there exists a node v on the path that is a collider, but none of its descendants (including itself) are in S .

Definition 7 (d -separation). For a DAG \mathcal{G} with node set $[d]$, any two nodes $i \neq j$ are said to be d -separated by a set $S \subset [d] \setminus \{i, j\}$ if all paths from i to j are blocked by S .

The minimality condition states that there is no redundant edges in the causal graph, and is a natural consequence of the Occam's Razor Principle.

Assumption 6 (Causal minimality, 41, Section 3.4.2). For latent variables \mathbf{z} , removing any edge from \mathcal{G} would render violation of the causal Markov condition Definition 1. In other words, let \mathcal{G}_1 be the graph obtained by removing any single edge from \mathcal{G} , then there must exist $i \in [d]$ such that $\mathbf{z}_i \not\perp\!\!\!\perp \mathbf{z}_{\text{nd}_{\mathcal{G}_1}(i)} \mid \mathbf{z}_{\text{pa}_{\mathcal{G}_1}(i)}$.

The *faithfulness* condition states that the Causal Markov Condition actually entails all (conditional) independence in the latent variables.

Assumption 7 (Faithfulness, 41, Section 3.4.3). *Every (conditional) independence in the latent variables \mathbf{z} is entailed by the Causal Markov Condition applied to \mathcal{G} . In other words, $\mathbf{z}_i \perp \mathbf{z}_j \mid \mathbf{z}_S \Leftrightarrow i, j$ are d -separated by S .*

Existing works have explored different notions of identifiability. For observational data, it is well known that Markov equivalence of graphs is an intrinsic ambiguity that one cannot resolve:

Definition 8 (Markov equivalence/Faithful Indistinguishability, 41, Section 4.2). *If two DAGs encode the same set of dependency relations, we say that they are Markov equivalent.*

Any DAG \mathcal{G} induces a partial order on its nodes which we denote by $\prec_{\mathcal{G}}$. In the special case when for all $i \neq j$, either $i \prec_{\mathcal{G}} j$ or $j \prec_{\mathcal{G}} i$ holds, we say that $\prec_{\mathcal{G}}$ is a total order. This partial order is equivalent to the transitive closure of the graph, as defined below:

Definition 9 (Transitional closure). *Given any DAG \mathcal{G} , its transitional closure $\bar{\mathcal{G}}$ is defined to be the graph obtained by connecting all edges $i \rightarrow j$ where i is an ancestor of j in \mathcal{G} .*

When $\prec_{\mathcal{G}}$ is a total order, each pair of nodes are connected by a directed edge in its transitive closure $\bar{\mathcal{G}}$. Such $\bar{\mathcal{G}}$ is often called a *tournament* in graph theory.

In the following, we list different forms of identifiability that appear in the literature:

Definition 10 (different notions of identifiability). *Let $\mathcal{H} : \mathbb{R}^n \supseteq \mathcal{X} \mapsto \mathbb{R}^d$ be the space of diffeomorphic mappings from observation to latent, and \mathfrak{G} be the space of all DAGs with d nodes, then for $h, \hat{h} \in \mathcal{H}$ and $\mathcal{G}, \hat{\mathcal{G}} \in \mathfrak{G}$, we write*

- (i) [36, 23] $(h, \mathcal{G}) \stackrel{T}{\sim}_G (\hat{h}, \hat{\mathcal{G}})$ if there exists a permutation π on $[d]$ such that $\pi(\mathcal{G})$ and $\hat{\mathcal{G}}$ have the same transitional closure;
- (ii) [49, 47] $(h, \mathcal{G}) \sim_{\text{CRL}} (\hat{h}, \hat{\mathcal{G}})$ if we actually have $\mathcal{G} = \hat{\mathcal{G}}$ for the ϕ defined above.

Given an equivalence relation \sim on $\mathcal{H} \times \mathfrak{G}$, we say that a causal model (h, \mathcal{G}) is identifiable under \sim if any candidate solution $(\hat{h}, \hat{\mathcal{G}})$ satisfies $(\hat{h}, \hat{\mathcal{G}}) \sim (h, \mathcal{G})$. The notion of identification up to $\stackrel{T}{\sim}_G$, as shown in Seigal et al. [36] with single-node soft interventions on linear causal models, is highly related to this paper. Compared with their result, our \sim_{sur} guarantee is much stronger, since not only the causal graph can be fully recovered, but the latent variables can be identified up to mixtures of the effect-dominating sets as well.

E Illustrating examples for our theory and algorithm

E.1 An example for understanding the SNA ambiguity

We provide a simple example below to illustrate the SNA ambiguity discussed in Section 3.

Example 1. *Let G be a causal graph with $d = 3$ nodes and edges $1 \rightarrow 2$ and $2 \rightarrow 3$. We have access to observations from a set of environments \mathfrak{E} . It turns out that there is no way to distinguish between the following two structural equation models:*

$$\begin{aligned} \mathbf{z}_1 &= \epsilon_1^E & \mathbf{v}_1 &= \epsilon_1^E \\ \mathbf{z}_2 &= f_2^E(\mathbf{z}_1, \epsilon_2^E) & \mathbf{v}_2 &= f_2^E(\mathbf{v}_1, \epsilon_2^E) \\ \mathbf{z}_3 &= f_3^E(\mathbf{z}_2, \epsilon_3^E) & \mathbf{v}_3 &= \mathbf{v}_2 + f_3^E(\mathbf{v}_2, \epsilon_3^E) \\ \mathbf{x} = \mathbf{z} &= (\mathbf{z}_1, \mathbf{z}_2, \mathbf{z}_3)^\top & \mathbf{x} &= (\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3 - \mathbf{v}_2)^\top \end{aligned}$$

where $\epsilon_i^E, i = 1, 2, 3$ are independent noise variables, if we do not change the causal graph \mathcal{G} , no matter what environment E that we have.

This issue does not exist when we assume access to hard interventions on node 3, which effectively removes the edge $2 \rightarrow 3$. Specifically, with hard intervention on \mathbf{z}_3 , the variables \mathbf{z}_2 and \mathbf{z}_3 become independent. But by definition, $\mathbf{v}_2 = \mathbf{z}_2$ and $\mathbf{v}_3 = \mathbf{z}_2 + \mathbf{z}_3$ must be dependent, so this intervention

cannot be realized by any hard intervention on v_3 , thereby providing a way to distinguish between the above models.

Without node 3, the same ambiguity would arise on node 2. However, node 3 can help us to overcome this ambiguity, thanks to the fact that node 2 is the only causal parent of node 3. Suppose for example that $v_2 = m(z_1, z_2)$ is some mixture of z_1 and z_2 , then $v_3 = \hat{f}_3^E(v_2, \epsilon_3^E) = \hat{f}_3^E(m(z_1, z_2), \epsilon_3^E)$. Since all environments share the same mixing function, v_3 must be some deterministic function $\psi_3(z)$ of z , where ψ_3 is the same across all environment E . Hence, we have

$$\hat{f}_3^E(m(z_1, z_2), \epsilon_3^E) = \psi_3(z_1, z_2, \hat{f}_3^E(z_2, \epsilon_3^E)) \quad (7)$$

Now we note that the dependencies of LHS on z_1 and z_2 are through a single scalar-valued function m , but since we would have different \hat{f}_3^E 's in different environments, this in general does not hold for the RHS. Therefore, any causal model with latent variable v_2 as a mixture of z_1 and z_2 cannot be equivalent to the ground-truth model.

According to Definition 3, in Example 1 we have $\text{sur}_{\mathcal{G}}(1) = \text{sur}_{\mathcal{G}}(2) = \emptyset$ but $\text{sur}_{\mathcal{G}}(3) = \{2\}$.

E.2 An example for the main idea behind LiNGCReL

To illustrate our main algorithm on how we can recover the graph \mathcal{G} and the matrix H , we first provide some intuition using a simple three-node example:

Example 2. Let \mathcal{G} be the graph with $d = 3$ nodes and edges $1 \rightarrow 2, 1 \rightarrow 3$ and $2 \rightarrow 3$, so that each B_k is of form

$$B_k = \begin{pmatrix} \times & 0 & 0 \\ \times & \times & 0 \\ \times & \times & \times \end{pmatrix} \rightsquigarrow \begin{matrix} \mathbf{b}_{k1} \\ \mathbf{b}_{k2} \\ \mathbf{b}_{k3} \end{matrix} \quad (8)$$

We can identify the graph as follows: first, for $i \in \{1, 2, 3\}$, look at the space \mathbf{W}_i spanned by the rows $(\mathbf{M}_k)_i, k \in [K]$. If $\dim \mathbf{W}_i = 1$, we know that i is a source node (i.e., $\text{pa}_{\mathcal{G}}(i) = \emptyset$) in \mathcal{G} . Otherwise it is not, due to Assumption 5. Hence we can know that node 1 is a source node.

In our example, there is no other node that satisfies this requirement. We then proceed to search for some $i \neq 1$ such that the projection of \mathbf{W}_i onto \mathbf{W}_1^\perp has dimension 1. If this holds, then one can show that $\text{pa}_{\mathcal{G}}(i) = \{1\}$. Otherwise, i must have parents other than 1.

It turns this requirement is satisfied for node 2 since $\dim(\text{proj}_{\mathbf{h}_1} \text{span}(\mathbf{h}_1, \mathbf{h}_2)) = 1$, but is not satisfied for node 3 since $\dim(\text{proj}_{\mathbf{h}_1} \text{span}(\mathbf{h}_1, \mathbf{h}_2, \mathbf{h}_3)) \geq 2$ (by Lemma 4). Hence we know that $\text{pa}_{\mathcal{G}}(2) = \{1\}$.

Finally, it remains to determine $\text{pa}_{\mathcal{G}}(3)$. To do this, we first note that $\dim \mathbf{W}_3 = 3$. Then we project \mathbf{W}_3 onto \mathbf{W}_1^\perp and \mathbf{W}_2^\perp respectively, and the resulting dimensions are 2 and 1. As we rigorously show in Proposition 2, a decrease of the dimension exactly indicates finding a new parent. Thus we have $\text{pa}_{\mathcal{G}}(3) = \{1, 2\}$, completing the recovery of the graph.

Finally, we recover the unmixing matrix H (and thus the latent variables) by noticing that $\mathbf{h}_1 \in \mathbf{W}_1$, $\mathbf{h}_2 \in \mathbf{W}_2 \cap \mathbf{W}_3$ and $\mathbf{h}_3 \in \mathbf{W}_3$. Ambiguities would arise at nodes 2 and 3, which are exactly the nodes that have non-empty effect-dominating sets.

F Auxiliary lemmas

Lemma 2. For any family of m -dimensional vectors $\{\mathbf{v}_k\}_{k=1}^K$ and $\{\mathbf{z}_k\}_{k=1}^K$ if $\mathbf{v}_k = \mathbf{z}_k \mathbf{T}$ and $\mathbf{T} \in \mathbb{R}^{m \times m}$ is invertible, then

$$\dim \text{span} \langle \mathbf{v}_k : k \in [K] \rangle = \dim \text{span} \langle \mathbf{z}_k : k \in [K] \rangle$$

Theorem 5 (Darmois-Skitovic Theorem). Let $\epsilon_i, i \in [d]$ be independent random variables and $X = \sum_{i=1}^d \alpha_i \epsilon_i, Y = \sum_{i=1}^d \beta_i \epsilon_i$. If $X \perp\!\!\!\perp Y$, then for $\forall i \in [d], \alpha_i \beta_i \neq 0 \Rightarrow \epsilon_i$ is Gaussian distributed.

Lemma 3. Suppose that $\epsilon = (\epsilon_1, \dots, \epsilon_d)$ is a d -dimensional random vector with independent components such that $\text{Var}(\epsilon_i) = 1, \forall i \in [d]$, and there exists an invertible and non-diagonal matrix M such that $M\epsilon \stackrel{d}{=} \epsilon$, then at least one of the following statements must hold:

(1) there exists at least two Gaussian variables in $\epsilon_1, \dots, \epsilon_d$;

(2) M is a permutation matrix and there exists $1 \leq i < j \leq d$ such that $\epsilon_i \stackrel{d}{=} \epsilon_j$.

Proof. Suppose that (1) does not hold, then there is at most one Gaussian variable in $\epsilon_1, \dots, \epsilon_d$. We assume WLOG that $\epsilon_1, \dots, \epsilon_{d-1}$ are all non-Gaussian. Then by the Darmois-Skitovic Theorem, we know that for $\forall 1 \leq j < k \leq [d]$ and $i \in [d-1]$, $M_{ji} \cdot M_{ki} = 0 \Rightarrow$ there is at most one non-zero entry in each of the first $d-1$ columns of M .

Assume that $M_{k_i, i} \neq 0$, $i \in [d-1]$. Since M is invertible, we know that $k_i, i \in [d-1]$ must be different. Let k_d be the remaining element in $[d]$ that does not appear in $k_i, i < d$, then $(M\epsilon)_{k_d} = M_{k_d, d}\epsilon_d$, while $(M\epsilon)_{k_i} = M_{k_i, i}\epsilon_i + M_{k_i, d}\epsilon_d$. Since the components of $M\epsilon$ are independent, it is easy to see that $M_{id} \neq 0, \forall i \neq k_d$. In other words, M only has non-zero entries at $(k_i, i), i \in [d]$.

Since $\text{Var}(\epsilon_i) = 1$, we know that M must be a signed permutation matrix. Finally, let π be the permutation on $[d]$ such that $M_{i, \pi(i)} \neq 0$. Since M is not diagonal, π must have a cycle (i_1, i_2, \dots, i_k) with length $k \geq 2$, so that $\epsilon_{i_1}, \dots, \epsilon_{i_k}$ all have the same distribution, which implies that (2) holds, as desired. \square

Lemma 4. Let V_1, V_2 be two subspaces of \mathbb{R}^d such that $V_1 \cap V_2 = \{0\}$, and $P_{V_1^\perp}$ be the orthogonal projection onto V_1^\perp , then we have that $\dim(V_2) = \dim(P_{V_1^\perp} V_2)$.

Proof. Obviously we have $\dim(V_2) \geq \dim(P_{V_1^\perp} V_2)$. On the other hand, let u_1, u_2, \dots, u_m be a basis of V_2 , then $w_i = P_{V_1^\perp} u_i, i = 1, 2, \dots, m$ are also independent. Indeed, suppose that $\lambda_i, i = 1, 2, \dots, m$ satisfy $\sum_{i=1}^m \lambda_i w_i = 0$, then $P_{V_1^\perp}(\sum_{i=1}^m \lambda_i u_i) = 0$, implying that $\sum_{i=1}^m \lambda_i u_i \in V_1$. However, we know that $V_1 \cap V_2 = \{0\}$, so $\lambda_1 = \dots = \lambda_m = 0$. This concludes the proof. \square

Lemma 5. Assumption 4 is equivalent to Assumption 5.

Proof. The main observation is that for each $k \in [K]$, $(B_k)_i$ only has non-zero entries at the j -th coordinate where $j \in \overline{\text{pa}}_{\mathcal{G}}(i)$. Moreover, let $\tilde{w}_k(i)$ be the vector consisting of these entries, then $\tilde{w}_k(i) = (\Omega_k)_{ii}^{-\frac{1}{2}}(-w_k(i), 1)$. Hence,

$$\dim \text{span} \langle (B_k)_i : k \in [K] \rangle = \dim \text{span} \langle (-w_k(i), 1) : k \in [K] \rangle.$$

Suppose that Assumption 4 holds, then for $\forall x \in \mathbb{R}^{|\text{pa}_{\mathcal{G}}(i)|}$, there exists $\lambda_k \in \mathbb{R}, 1 \leq k \leq |\text{pa}_{\mathcal{G}}(i)|$ such that $\sum_k \lambda_k = 1$ and $\sum_k \lambda_k w_k(i) = x$. Hence,

$$(x, 1) = \sum_k \lambda_k \tilde{w}_k(i) \in \text{span} \langle (B_k)_i : k \in [K] \rangle.$$

This immediately implies that $\text{span} \langle (B_k)_i : k \in [K] \rangle = \mathbb{R}^{|\text{pa}_{\mathcal{G}}(i)|+1}$, so that Assumption 5 holds.

Conversely, suppose that Assumption 5 holds, then for $\forall x \in \mathbb{R}^{|\text{pa}_{\mathcal{G}}(i)|}$, there exists $\lambda_k \in \mathbb{R}, 1 \leq k \leq |\text{pa}_{\mathcal{G}}(i)|$ such that $\sum_k \lambda_k \tilde{w}_k(i) = (x, 1)$. Hence we have $\sum_k \lambda_k w_k(i) = x$ and $\sum_k \lambda_k = 1$, implying Assumption 4. \square

G Properties of effect-domination sets

Lemma 6. • $j \in \text{sur}_{\mathcal{G}}(i)$ if and only if $\overline{\text{ch}}_{\mathcal{G}}(i) \subseteq \text{ch}_{\mathcal{G}}(j)$;

• when $i \neq j, j \in \text{sur}_{\mathcal{G}}(i)$ if and only if $\overline{\text{ch}}_{\mathcal{G}}(i) \subseteq \overline{\text{ch}}_{\mathcal{G}}(j)$.

Proof. If $j \in \text{sur}_{\mathcal{G}}(i)$, by definition $i \in \text{ch}_{\mathcal{G}}(j)$ and $\text{ch}_{\mathcal{G}}(i) \subseteq \text{ch}_{\mathcal{G}}(j)$, so that $\overline{\text{ch}}_{\mathcal{G}}(i) \subseteq \text{ch}_{\mathcal{G}}(j)$. Conversely, $\overline{\text{ch}}_{\mathcal{G}}(i) \subseteq \text{ch}_{\mathcal{G}}(j)$ implies that $i \in \text{ch}_{\mathcal{G}}(j)$ and $\text{ch}_{\mathcal{G}}(i) \subseteq \text{ch}_{\mathcal{G}}(j)$, so $j \in \text{sur}_{\mathcal{G}}(i)$. This proves the first claim.

To prove the second claim, assume that $\overline{\text{ch}}_{\mathcal{G}}(i) \subseteq \overline{\text{ch}}_{\mathcal{G}}(j)$ holds but $\overline{\text{ch}}_{\mathcal{G}}(i) \subseteq \text{ch}_{\mathcal{G}}(j)$ does not hold, then we must have $j \in \overline{\text{ch}}_{\mathcal{G}}(i)$. since $j \neq i$, we have $j \in \text{ch}_{\mathcal{G}}(i)$, but then $i \notin \overline{\text{ch}}_{\mathcal{G}}(j)$, which is a contradiction. Hence $\overline{\text{ch}}_{\mathcal{G}}(i) \subseteq \text{ch}_{\mathcal{G}}(j)$ and the conclusion follows from the first claim. \square

Lemma 7. Let \mathcal{G} be a DAG and i be its node, then for $\forall j \in \text{pa}_{\mathcal{G}}(i)$, we have $\text{sur}_{\mathcal{G}}(j) \subseteq \text{pa}_{\mathcal{G}}(i)$.

Proof. Let $k \in \text{sur}_{\mathcal{G}}(j)$, then by definition we have $\text{ch}_{\mathcal{G}}(j) \subseteq \text{ch}_{\mathcal{G}}(k)$. In particular, we have $i \in \text{ch}_{\mathcal{G}}(k) \Rightarrow k \in \text{pa}_{\mathcal{G}}(i)$. \square

Lemma 8. Let \mathcal{G} be a DAG and i be its node, then for $\forall j \in \text{sur}_{\mathcal{G}}(i)$, we have $\text{sur}_{\mathcal{G}}(j) \subseteq \text{sur}_{\mathcal{G}}(i)$.

Proof. Let $k \in \text{sur}_{\mathcal{G}}(j)$, then by definition we have $\overline{\text{ch}}_{\mathcal{G}}(j) \subset \overline{\text{ch}}_{\mathcal{G}}(k)$. We also know that $\overline{\text{ch}}_{\mathcal{G}}(i) \subset \overline{\text{ch}}_{\mathcal{G}}(j)$, so $\overline{\text{ch}}_{\mathcal{G}}(i) \subset \overline{\text{ch}}_{\mathcal{G}}(k)$, implying that $k \in \text{sur}_{\mathcal{G}}(i)$. \square

Lemma 9. If $M \in \mathcal{M}_{\text{sur}}^0(\mathcal{G})$, then $M^{-1} \in \mathcal{M}_{\text{sur}}^0(\mathcal{G})$.

Proof. Assume WLOG that the nodes of \mathcal{G} satisfy $i \in \text{pa}_{\mathcal{G}}(j) \Rightarrow i < j$ (otherwise we can choose a different index of the nodes and correspondingly swap some rows and columns of M). Since $i \in \text{sur}_{\mathcal{G}}(j) \Rightarrow i \in \text{pa}_{\mathcal{G}}(j)$, it follows that M must be lower triangular and the diagonal entries are nonzero.

Let $N = M^{-1}$, then for $\forall i \in [d]$, we have

$$\sum_{j=1}^d N_{ij} M_{j\ell} = 0, \quad \forall \ell \notin \overline{\text{sur}}_{\mathcal{G}}(i). \quad (9)$$

Since $M \in \mathcal{M}_{\text{sur}}^0(\mathcal{G})$, we have $M_{j\ell} = 0$ for $\forall j$ such that $\ell \notin \overline{\text{sur}}_{\mathcal{G}}(j)$. By Lemma 8, if $j \in \overline{\text{sur}}_{\mathcal{G}}(i)$, then $\ell \notin \overline{\text{sur}}_{\mathcal{G}}(i)$ necessarily implies that $\ell \notin \overline{\text{sur}}_{\mathcal{G}}(j)$. Hence the left hand side of (9) is essentially a sum over $j \notin \overline{\text{sur}}_{\mathcal{G}}(i)$, i.e.,

$$\sum_{j \notin \overline{\text{sur}}_{\mathcal{G}}(i)} N_{ij} M_{j\ell} = 0, \quad \forall \ell \notin \overline{\text{sur}}_{\mathcal{G}}(i).$$

Viewing the above as a system of linear equations in $N_{ij}, j \notin \overline{\text{sur}}_{\mathcal{G}}(i)$, the coefficient matrix $(M_{j\ell})_{j, \ell \notin \overline{\text{sur}}_{\mathcal{G}}(i)}$ must be invertible since it is a sub-matrix of the invertible lower-triangular matrix M . As a result, we necessary have $N_{ij} = 0, \forall j \notin \overline{\text{sur}}_{\mathcal{G}}(i)$. Finally, $N = M^{-1}$ must be invertible, so $N \in \mathcal{M}_{\text{sur}}^0(\mathcal{G})$ as desired. \square

Lemma 10. Suppose that $\psi : \mathbb{R}^d \mapsto \mathbb{R}^d$ is a diffeomorphism and \mathcal{G} be a DAG, such that for $\forall i \in [d]$, $\psi_i(\mathbf{z})$ is a function of $\mathbf{z}_{\overline{\text{sur}}_{\mathcal{G}}(i)}$. Then for $\forall j \in [d]$, $(\psi^{-1})_j(\mathbf{v})$ is a function of $\mathbf{v}_{\overline{\text{sur}}_{\mathcal{G}}(j)}$.

Proof. Let $\mathbf{J}_{\mathbf{z}} = \mathbf{J}_{\psi}(\mathbf{z})$ be the Jacobian matrix of ψ . Since ψ is a diffeomorphism, $\mathbf{J}_{\mathbf{z}}$ is invertible for any $\mathbf{z} \in \mathbb{R}^d$. Moreover, our assumption implies that $(\mathbf{J}_{\mathbf{z}})_{ij} = 0, \forall j \notin \overline{\text{sur}}_{\mathcal{G}}(i)$, so $\mathbf{J}_{\mathbf{z}} \in \mathcal{M}_{\text{sur}}^0(\mathcal{G})$. By Lemma 9, $\mathbf{J}_{\mathbf{z}}^{-1} \in \mathcal{M}_{\text{sur}}^0(\mathcal{G})$. But $\mathbf{J}_{\mathbf{z}}^{-1}$ is exactly the Jacobian matrix of ψ^{-1} at $\mathbf{v} = \psi(\mathbf{z})$, hence it follows that $(\psi^{-1})_j(\mathbf{v})$ is only a function of $\mathbf{v}_{\overline{\text{sur}}_{\mathcal{G}}(j)}$, as desired. \square

Lemma 11. The binary relation \sim_{sur} defined in Definition 4 is an equivalence relation.

Proof. It is obvious that $(\mathbf{h}, \mathcal{G}) \sim_{\text{sur}} (\mathbf{h}, \mathcal{G})$ holds for any model $(\mathbf{h}, \mathcal{G})$.

Suppose that $(\mathbf{h}_1, \mathcal{G}_1) \sim_{\text{sur}} (\mathbf{h}_2, \mathcal{G}_2)$, then there exists a permutation π on $[d]$ and a diffeomorphism $\psi : \mathbb{R}^d \mapsto \mathbb{R}^d$ where $\psi_i(\mathbf{z})$ is a function of $\mathbf{z}_{\overline{\text{sur}}_{\mathcal{G}_1}(i)}$, such that $i \in \text{pa}_{\mathcal{G}_1}(j) \Leftrightarrow \pi(i) \in \text{pa}_{\mathcal{G}_2}(\pi(j))$ and $\mathbf{P}_{\pi} \circ \mathbf{h}_2 = \psi \circ \mathbf{h}_1$. Then we can write $\mathbf{P}_{\pi}^{-1} \circ \mathbf{h}_1 = \hat{\psi} \circ \mathbf{h}_2$ where $\hat{\psi} = \mathbf{P}_{\pi}^{-1} \circ \psi^{-1} \circ \mathbf{P}_{\pi}$. By Lemma 10, we know that $(\psi^{-1})_j(\mathbf{v})$ is a function of $\mathbf{v}_{\overline{\text{sur}}_{\mathcal{G}_1}(j)}$, so $(\hat{\psi})_j$ is a function of $\mathbf{v}_{\pi(\overline{\text{sur}}_{\mathcal{G}_1}(j))} = \mathbf{v}_{\overline{\text{sur}}_{\mathcal{G}_2}(j)}$, implying that $(\mathbf{h}_2, \mathcal{G}_2) \sim_{\text{sur}} (\mathbf{h}_1, \mathcal{G}_1)$.

Finally, let $(\mathbf{h}_1, \mathcal{G}_1) \sim_{\text{sur}} (\mathbf{h}_2, \mathcal{G}_2)$ and $(\mathbf{h}_2, \mathcal{G}_2) \sim_{\text{sur}} (\mathbf{h}_3, \mathcal{G}_3)$, then we can write

$$\mathbf{P}_\pi \circ \mathbf{h}_2 = \psi \circ \mathbf{h}_1 \quad \text{and} \quad \mathbf{P}_{\hat{\pi}} \circ \mathbf{h}_3 = \hat{\psi} \circ \mathbf{h}_2$$

where: for $\forall i \in [d]$, $\psi_i(\mathbf{z})$ is a function of $\mathbf{z}_{\overline{\text{sur}}_{\mathcal{G}_1}(i)}$, $\hat{\psi}_i(\mathbf{z})$ is a function of $\mathbf{z}_{\overline{\text{sur}}_{\mathcal{G}_2}(i)}$, $i \in \text{pa}_{\mathcal{G}_1}(j) \Leftrightarrow \pi(i) \in \text{pa}_{\mathcal{G}_2}(\pi(j))$ and $i \in \text{pa}_{\mathcal{G}_2}(j) \Leftrightarrow \hat{\pi}(i) \in \text{pa}_{\mathcal{G}_3}(\hat{\pi}(j))$. Then, we can write

$$\mathbf{P}_\pi \circ \mathbf{P}_{\hat{\pi}} \circ \mathbf{h}_3 = \mathbf{P}_\pi \circ \hat{\psi} \circ \mathbf{P}_\pi^{-1} \circ \psi \circ \mathbf{h}_1.$$

Since $\hat{\psi}_i(\mathbf{z})$ is a function of $\mathbf{z}_{\overline{\text{sur}}_{\mathcal{G}_2}(i)}$, we deduce that $\left(\mathbf{P}_\pi \circ \hat{\psi} \circ \mathbf{P}_\pi^{-1}\right)_i(\mathbf{z})$ is a function of $\mathbf{z}_{\overline{\text{sur}}_{\mathcal{G}_1}(i)}$. Hence, $\left(\mathbf{P}_\pi \circ \hat{\psi} \circ \mathbf{P}_\pi^{-1} \circ \psi\right)_i(\mathbf{z}) = \left(\mathbf{P}_\pi \circ \hat{\psi} \circ \mathbf{P}_\pi^{-1}\right)_i(\psi(\mathbf{z}))$ is a function of $\psi_{\overline{\text{sur}}_{\mathcal{G}_1}(i)}(\mathbf{z})$. The definition of ψ implies that for each $j \in \overline{\text{sur}}_{\mathcal{G}_1}(i)$, $\psi_j(\mathbf{z})$ is a function of $\mathbf{z}_{\overline{\text{sur}}_{\mathcal{G}_1}(j)}$. By **Lemma 8**, we have $\cup_{j \in \overline{\text{sur}}_{\mathcal{G}_1}(i)} \overline{\text{sur}}_{\mathcal{G}_1}(j) \subseteq \overline{\text{sur}}_{\mathcal{G}_1}(i)$. Hence $\left(\mathbf{P}_\pi \circ \hat{\psi} \circ \mathbf{P}_\pi^{-1} \circ \psi\right)_i(\mathbf{z})$ is still a function of $\mathbf{z}_{\overline{\text{sur}}_{\mathcal{G}_1}(i)}$. Moreover, we also have $i \in \text{pa}_{\mathcal{G}_1}(j) \Leftrightarrow \pi(i) \in \text{pa}_{\mathcal{G}_2}(\pi(j)) \Leftrightarrow \hat{\pi} \circ \pi(i) \in \text{pa}_{\mathcal{G}_3}(\hat{\pi} \circ \pi(j))$, so by definition, $(\mathbf{h}_1, \mathcal{G}_1) \sim_{\text{sur}} (\mathbf{h}_3, \mathcal{G}_3)$, as desired. \square

H Omitted proofs from Section 4 and Section 5

H.1 Proof of Theorem 1

According to the assumption, we have that $\epsilon = \mathbf{B}_k \mathbf{H} \mathbf{x}$ and $\hat{\epsilon} = \hat{\mathbf{B}}_k \hat{\mathbf{H}} \mathbf{x}$, so that $\epsilon = \mathbf{B}_k \mathbf{H} (\hat{\mathbf{B}}_k \hat{\mathbf{H}})^\dagger \hat{\epsilon}, \forall k \in [K]$. By **Lemma 3**, we know that for each k , $\mathbf{P}_k := \mathbf{B}_k \mathbf{H} (\hat{\mathbf{B}}_k \hat{\mathbf{H}})^\dagger$ is a signed permutation matrix, so that $\epsilon = \mathbf{P}_k \hat{\epsilon}$. Since for any $i \neq j$, $\hat{\epsilon}_i \neq \hat{\epsilon}_j$, we must have $|\mathbf{P}|_1 = |\mathbf{P}|_2 = \dots = |\mathbf{P}|_K =: \mathbf{P}$ and $\epsilon = \mathbf{P} \hat{\epsilon}$, where $|\mathbf{P}|$ denotes the resulting matrix by taking the absolute value of all entries in \mathbf{P} . Thus, we can WLOG assume that $\epsilon = \hat{\epsilon}$, since otherwise we can permute the noise variables $\hat{\epsilon}$, and also permute the rows of \mathbf{B}_k correspondingly. In other words, suppose that the permutation matrix $|\mathbf{P}|$ has $|\mathbf{P}|_{k_i, i} = 1, i \in [d]$, then we can assign to each node i in $\hat{\mathcal{G}}$ a new index k_i and work with the new indices.

In this case, by **Lemma 3** we have $\mathbf{B}_k \mathbf{H} = \Sigma_k \hat{\mathbf{B}}_k \hat{\mathbf{H}}, \forall k \in [K]$ or equivalently $\Sigma_k \hat{\mathbf{B}}_k = \mathbf{B}_k \mathbf{T}$, where $\mathbf{T} = \mathbf{H} \hat{\mathbf{H}}^\dagger \in \mathbb{R}^{d \times d}$, and Σ_k is a diagonal matrix with diagonal entries in $\{+1, -1\}$. Let $\hat{\hat{\mathbf{B}}}_k = \Sigma_k \hat{\mathbf{B}}_k$, then the rows of $\hat{\hat{\mathbf{B}}}_k$ equals (up to sign) to the rows of $\hat{\mathbf{B}}_k$.

To summarize, we now know that i) $\hat{\hat{\mathbf{B}}}_k = \mathbf{B}_k \mathbf{T}, k \in [K]$, ii) $(\mathbf{B}_k)_{ij} \neq 0 \Leftrightarrow j \in \overline{\text{pa}}_{\mathcal{G}}(i)$, and similarly, $(\hat{\hat{\mathbf{B}}}_k)_{ij} \neq 0 \Leftrightarrow j \in \overline{\text{pa}}_{\hat{\mathcal{G}}}(i)$, and iii) Both $\{\mathbf{B}_k\}$ and $\{\hat{\hat{\mathbf{B}}}_k\}$ satisfy the node-level non-degeneracy assumption **Assumption 5**. For any two such matrices that satisfy such a set of conditions, it must necessarily be true that $\mathcal{G} = \hat{\mathcal{G}}$.

Lemma 12 (Graph Identifiability). *Consider any two sets matrices $\{\hat{\mathbf{B}}_k\}_{k \in [K]}$ and $\{\mathbf{B}_k\}_{k \in [K]}$ and associated graphs $\mathcal{G}, \hat{\mathcal{G}}$. If these sets and graphs satisfy that:*

1. $\hat{\hat{\mathbf{B}}}_k = \mathbf{B}_k \mathbf{T}, k \in [K]$;
2. $(\mathbf{B}_k)_{ij} \neq 0 \Leftrightarrow j \in \overline{\text{pa}}_{\mathcal{G}}(i)$, and similarly, $(\hat{\hat{\mathbf{B}}}_k)_{ij} \neq 0 \Leftrightarrow j \in \overline{\text{pa}}_{\hat{\mathcal{G}}}(i)$.
3. Both $\{\mathbf{B}_k\}$ and $\{\hat{\hat{\mathbf{B}}}_k\}$ satisfy the node-level non-degeneracy assumption **Assumption 5**.

then it must hold that $\mathcal{G} = \hat{\mathcal{G}}$.

Proof. We prove this via induction on the size of the graph d . Note that here $\mathcal{G} = \hat{\mathcal{G}}$ is not up to permutation and our statement is equivalent to $\text{pa}_{\mathcal{G}}(i) = \text{pa}_{\hat{\mathcal{G}}}(i), \forall i \in [d]$.

If $d = 1$, i.e., $\mathcal{G} = \hat{\mathcal{G}}$ obviously holds since both are graphs with only 1 node.

Suppose that for all graphs \mathcal{G} of size $d - 1$, the graph $\hat{\mathcal{G}}$ satisfying all given assumptions must necessarily be equal to \mathcal{G} . Now, we consider the case that \mathcal{G} has d nodes. WLOG we can assume that the nodes of \mathcal{G} are properly indexed such that $i \in \text{pa}_{\mathcal{G}}(j) \Rightarrow i < j$, so $\mathbf{B}_k, k \in [K]$ are lower-triangular matrices. (However, it is currently unknown whether $\hat{\mathbf{B}}_k$ are also lower-triangular.)

By our assumption that $i \in \text{pa}_{\mathcal{G}}(j) \Rightarrow i < j$, the node d in \mathcal{G} has no child. Thus we can write

$$\mathbf{B}_k = \begin{pmatrix} \mathbf{B}_k^- & \mathbf{0} \\ \mathbf{b}_k & c_k \end{pmatrix}, \mathbf{T} = \begin{pmatrix} \mathbf{T}^- & \times \\ \times & \times \end{pmatrix} \text{ and } \hat{\mathbf{B}}_k = \mathbf{B}_k \mathbf{T} = \begin{pmatrix} \hat{\mathbf{B}}_k^- & \times \\ \times & \times \end{pmatrix}$$

where $\mathbf{B}_k^-, \mathbf{T}, \hat{\mathbf{B}}_k^- = \mathbf{B}_k^- \mathbf{T}^- \in \mathbb{R}^{(d-1) \times (d-1)}, \mathbf{b}_k \in \mathbb{R}^{d-1}, c_k \in \mathbb{R}$ and \times denotes irrelevant entries.

Let $\mathbf{A}_k^-, \hat{\mathbf{A}}_k^-, \mathbf{\Omega}_k^-$ and $\hat{\mathbf{\Omega}}_k^-$ be the top-left $(d-1) \times (d-1)$ sub-matrices of $\mathbf{A}_k, \hat{\mathbf{A}}, \mathbf{\Omega}_k$ and $\hat{\mathbf{\Omega}}_k$ respectively, and \mathcal{G}^- and $\hat{\mathcal{G}}^-$ are graphs obtained by deleting node d and all related edges from \mathcal{G} and $\hat{\mathcal{G}}$. Then it is easy to see that

$$(\mathbf{A}_k^-)_{ij} \neq 0 \Leftrightarrow j \in \text{pa}_{\mathcal{G}^-}(i) \quad \text{and} \quad (\hat{\mathbf{A}}_k^-)_{ij} \neq 0 \Leftrightarrow j \in \text{pa}_{\hat{\mathcal{G}}^-}(i). \quad (10)$$

Moreover,

$$\begin{pmatrix} \mathbf{B}_k^- & \mathbf{0} \\ \mathbf{b}_k & c_k \end{pmatrix} = \mathbf{B}_k = \mathbf{\Omega}_k^{-\frac{1}{2}} (\mathbf{I} - \mathbf{A}_k) = \begin{pmatrix} (\mathbf{\Omega}_k^-)^{-\frac{1}{2}} & \mathbf{0} \\ \mathbf{0} & \times \end{pmatrix} \begin{pmatrix} \mathbf{I} - \mathbf{A}_k^- & \times \\ \times & \times \end{pmatrix} = \begin{pmatrix} (\mathbf{\Omega}_k^-)^{-\frac{1}{2}} (\mathbf{I} - \mathbf{A}_k^-) & \times \\ \times & \times \end{pmatrix}$$

so that $\mathbf{B}_k^- = (\mathbf{\Omega}_k^-)^{-\frac{1}{2}} (\mathbf{I} - \mathbf{A}_k^-)$. Similarly, we have $\hat{\mathbf{B}}_k^- = (\hat{\mathbf{\Omega}}_k^-)^{-\frac{1}{2}} (\mathbf{I} - \hat{\mathbf{A}}_k^-)$.

We can also verify that $\{\mathbf{B}_k^-\}_{k=1}^K$ and $\{\hat{\mathbf{B}}_k^-\}_{k=1}^K$ are node-level independent in the sense of **Assumption 5**. We only prove this for $\{\hat{\mathbf{B}}_k^-\}_{k=1}^K$; the arguments used for $\{\mathbf{B}_k^-\}_{k=1}^K$ are exactly the same as the first case considered below. Now for each $i \in [d-1]$, let $\mathbf{R}_i \in \mathbb{R}^{K \times d}$ be the matrix whose k -th row is the i -th row of $\hat{\mathbf{B}}_k$, and $\mathbf{R}_i^- \in \mathbb{R}^{K \times (d-1)}$ be the matrix whose k -th row is the i -th row of $\hat{\mathbf{B}}_k^-$, then obviously \mathbf{R}_i is of form $[\mathbf{R}_i^-, \mathbf{r}_i]$. We consider two cases:

- **Case 1.** $d \notin \text{pa}_{\hat{\mathcal{G}}}(i)$ This means that the last entry of the i -th row of $\hat{\mathbf{B}}_k$ is zero. Thus $\mathbf{r}_i = \mathbf{0}$, and $\text{rank}(\mathbf{R}_i^-) = \text{rank}(\mathbf{R}_i) = |\overline{\text{pa}}_{\hat{\mathcal{G}}}(i)| = |\overline{\text{pa}}_{\hat{\mathcal{G}}^-}(i)|$, where the second equality follows from **Assumption 5**.
- **Case 2.** $d \in \text{pa}_{\hat{\mathcal{G}}}(i)$ In this case we have $\text{rank}(\mathbf{R}_i^-) \geq \text{rank}(\mathbf{R}_i) - 1 = |\overline{\text{pa}}_{\hat{\mathcal{G}}}(i)| - 1 = |\overline{\text{pa}}_{\hat{\mathcal{G}}^-}(i)|$. Due to our assumption on $\hat{\mathbf{A}}_k$ and the relationship $\hat{\mathbf{B}}_k^- = (\hat{\mathbf{\Omega}}_k^-)^{-\frac{1}{2}} (\mathbf{I} - \hat{\mathbf{A}}_k^-)$, we know that each row of \mathbf{R}_i^- , namely the i -th row of some $\hat{\mathbf{B}}_k$ only has $|\overline{\text{pa}}_{\hat{\mathcal{G}}}(i)| - 1 = |\overline{\text{pa}}_{\hat{\mathcal{G}}^-}(i)|$ non-zero entries, so that $\text{rank}(\mathbf{R}_i^-) = |\overline{\text{pa}}_{\hat{\mathcal{G}}^-}(i)|$ holds.

Since we have shown that the matrices \mathbf{B}_k^- and $\hat{\mathbf{B}}_k^-$ satisfy the three properties that we assume for induction with \mathbf{T} replaced by \mathbf{T}^- and $\mathcal{G}, \hat{\mathcal{G}}$ replaced by $\mathcal{G}^-, \hat{\mathcal{G}}^-$ respectively, by induction hypothesis, we can thus deduce that $\mathcal{G}^- = \hat{\mathcal{G}}^-$. To prove $\mathcal{G} = \hat{\mathcal{G}}$ it remains to show that the dependency of node d on the remaining nodes are the same in \mathcal{G} and $\hat{\mathcal{G}}$.

First, we show that $\text{ch}_{\hat{\mathcal{G}}}(d) = \emptyset$. Suppose in contrary that there is some $i \in \text{ch}_{\hat{\mathcal{G}}}(d)$, then $|\text{pa}_{\mathcal{G}}(i)| = |\text{pa}_{\mathcal{G}^-}(i)| = |\text{pa}_{\hat{\mathcal{G}}^-}(i)| = |\text{pa}_{\hat{\mathcal{G}}}(i)| - 1$. Recalling that $(\mathbf{B})_i$ denotes the i -th row of matrix \mathbf{B} , we have

$$\begin{aligned} \dim \left(\text{span} \left\langle (\hat{\mathbf{B}}_k)_i : 1 \leq k \leq K \right\rangle \right) &= \dim \left(\text{span} \left\langle (\mathbf{B}_k)_i : 1 \leq k \leq K \right\rangle \right) \\ &\leq |\text{pa}_{\mathcal{G}}(i)| + 1 < |\text{pa}_{\hat{\mathcal{G}}}(i)| + 1, \end{aligned} \quad (11)$$

where the first inequality follows from $(\hat{\mathbf{B}}_k)_i = (\mathbf{B}_k)_i \mathbf{T}$ and **Lemma 2**, the second holds since each $(\mathbf{B}_k)_i$ has nonzero elements only at coordinates in $j \in \overline{\text{pa}}_{\mathcal{G}}(i)$, and the last one holds since

$|\text{pa}_{\mathcal{G}}(i)| = |\text{pa}_{\hat{\mathcal{G}}}(i)| - 1$. However, (11) contradicts the non-degeneracy condition **Assumption 5** that we assume for matrices $\hat{\mathbf{B}}_k, k \in [K]$ in the statement of the theorem. Therefore we have $\text{ch}_{\hat{\mathcal{G}}}(d) = \emptyset = \text{ch}_{\mathcal{G}}(d)$.

Second, by a similar argument comparing the number of nonzero elements in the last row of \mathbf{B}_k and $\hat{\mathbf{B}}_k$, we can also deduce that

$$|\text{pa}_{\mathcal{G}}(d)| = |\text{pa}_{\hat{\mathcal{G}}}(d)|.$$

Indeed, since $\left(\hat{\mathbf{B}}_k\right)_d = (\mathbf{B}_k)_d \mathbf{T}$, by **Lemma 2** we have

$$\dim \left(\text{span} \left\langle \left(\hat{\mathbf{B}}_k\right)_d : 1 \leq k \leq K \right\rangle \right) = \dim \left(\text{span} \left\langle (\mathbf{B}_k)_d : 1 \leq k \leq K \right\rangle \right)$$

However, since we assume that **Assumption 5** is satisfied for $\{\mathbf{B}_k\}_{k=1}^K$ and $\{\hat{\mathbf{B}}_k\}_{k=1}^K$, we know that the LHS and RHS of the above equation are equal to $|\text{pa}_{\mathcal{G}}(d)| + 1$ and $|\text{pa}_{\hat{\mathcal{G}}}(d)| + 1$ respectively, implying (12).

Third, we show that $\text{pa}_{\mathcal{G}}(d) = \text{pa}_{\hat{\mathcal{G}}}(d)$. Suppose the contrary, let ℓ be the smallest element in $\text{pa}_{\mathcal{G}}(d) \Delta \text{pa}_{\hat{\mathcal{G}}}(d)$, where $A \Delta B := (A \setminus B) \cup (B \setminus A)$. Recall that while \mathcal{G} and $\hat{\mathcal{G}}$ are originally not symmetric as nodes are topologically sorted according to \mathcal{G} , now we have shown that $\mathcal{G}^- \equiv \hat{\mathcal{G}}^-$ and that $\text{ch}_{\mathcal{G}}(d) = \text{ch}_{\hat{\mathcal{G}}}(d) = \emptyset$, so we can assume WLOG that $\ell \in \text{pa}_{\mathcal{G}}(d)$ and $\ell \notin \text{pa}_{\hat{\mathcal{G}}}(d)$, and the other case can be handled symmetrically. Since \mathbf{B}_k is lower triangular and $(\mathbf{B}_k)_{jj} = (\Omega_k)_{jj}^{-\frac{1}{2}} \neq 0, \forall j \in [d]$, the top-left $\ell \times \ell$ sub-matrix of \mathbf{B}_k , which we denote by $[\mathbf{B}_k]_{\ell, \ell}$, must be invertible. This implies that $\left\{ [\mathbf{B}_k]_{\ell, \ell}^\top \boldsymbol{\lambda} : \boldsymbol{\lambda} \in \mathbb{R}^\ell \right\} = \mathbb{R}^\ell$, so we can always find coefficients $\lambda_{kj}, j \in [\ell]$ such that the first ℓ entries of the vector $(\mathbf{B}_k)_d - \sum_{i=1}^\ell \lambda_{ki} (\mathbf{B}_k)_i \in \mathbb{R}^d$ are all zero. Since $\hat{\mathbf{B}}_k = \mathbf{B}_k \mathbf{T}$ and \mathbf{T} is invertible, we have $\left(\hat{\mathbf{B}}_k\right)_d - \sum_{j=1}^\ell \lambda_{kj} \left(\hat{\mathbf{B}}_k\right)_j = \left((\mathbf{B}_k)_d - \sum_{j=1}^\ell \lambda_{kj} (\mathbf{B}_k)_j\right) \mathbf{T}, \forall k \in [K]$ and

$$\begin{aligned} \dim \left(\text{span} \left\langle \left(\hat{\mathbf{B}}_k\right)_d - \sum_{j=1}^\ell \lambda_{kj} \left(\hat{\mathbf{B}}_k\right)_j : k \in [K] \right\rangle \right) &= \dim \left(\text{span} \left\langle (\mathbf{B}_k)_d - \sum_{j=1}^\ell \lambda_{kj} (\mathbf{B}_k)_j : k \in [K] \right\rangle \right) \\ &\leq |\text{pa}_{\mathcal{G}}(d) \setminus [\ell]| + 1. \end{aligned}$$

Here, the inequality holds because for any coordinate $t \in [d]$,

$$\left((\mathbf{B}_k)_d - \sum_{j=1}^\ell \lambda_{kj} (\mathbf{B}_k)_j \right)_t = \begin{cases} 0 & \text{if } t \leq \ell \\ (\mathbf{B}_k)_{d,t} & \text{otherwise} \end{cases} \quad (12)$$

where we note that \mathbf{B}_k is lower-triangular and thus $(\mathbf{B}_k)_{j,t} = 0, \forall j \leq \ell, t > \ell$. This implies that $\left((\mathbf{B}_k)_d - \sum_{j=1}^\ell \lambda_{kj} (\mathbf{B}_k)_j\right)_t$ is nonzero only if $t > \ell$ and $t \in \text{pa}_{\mathcal{G}}(d)$.

On the other hand, let $S = (\text{pa}_{\hat{\mathcal{G}}}(d) \cap [\ell]^c) \cup \{d\}$, then

$$\begin{aligned} \dim \left(\text{span} \left\langle \left(\hat{\mathbf{B}}_k\right)_d - \sum_{j=1}^\ell \lambda_{kj} \left(\hat{\mathbf{B}}_k\right)_j : k \in [K] \right\rangle \right) &\geq \dim \left(\text{span} \left\langle \left(\left(\hat{\mathbf{B}}_k\right)_d - \sum_{j=1}^\ell \lambda_{kj} \left(\hat{\mathbf{B}}_k\right)_j\right)_S : k \in [K] \right\rangle \right) \\ &= \dim \left(\text{span} \left\langle \left(\left(\hat{\mathbf{B}}_k\right)_d\right)_S : k \in [K] \right\rangle \right) = |S|. \end{aligned}$$

where we recall that \mathbf{u}_S denotes the vector $(u_i : i \in S) \in \mathbb{R}^{|S|}$. Here the first equality holds due to the same reason as (12), and the second follows from **Assumption 5**. To see why this is the case, note that **Assumption 5** implies that the $K \times (|\text{pa}_{\mathcal{G}}(d)| + 1)$ having $((\mathbf{B}_k)_d)_{\overline{\text{pa}_{\mathcal{G}}}(d)}$ as the k -th row has full column rank, so that the sub-matrix obtained by extracting columns corresponding to the node set S also has full column rank.

We have shown that $|\overline{\text{pa}}_{\hat{\mathcal{G}}}(d) \cap [\ell]^c| = |S| \leq |\text{pa}_{\mathcal{G}}(d) \cap [\ell]^c| + 1 = |\overline{\text{pa}}_{\mathcal{G}}(d) \cap [\ell]^c|$. On the other hand, recall that by our choice of ℓ , we have $|\overline{\text{pa}}_{\mathcal{G}}(d) \cap [\ell - 1]| = |\overline{\text{pa}}_{\hat{\mathcal{G}}}(d) \cap [\ell - 1]|$ and $\ell \in \overline{\text{pa}}_{\mathcal{G}}(d) \setminus \overline{\text{pa}}_{\hat{\mathcal{G}}}(d)$. Putting these together, we have $|\overline{\text{pa}}_{\mathcal{G}}(d)| > |\overline{\text{pa}}_{\hat{\mathcal{G}}}(d)|$. However, we know from (12) that $|\text{pa}_{\mathcal{G}}(d)| = |\text{pa}_{\hat{\mathcal{G}}}(d)|$, leading to a contradiction. Hence, such ℓ shouldn't exist and we must have $\text{pa}_{\mathcal{G}}(d) = \text{pa}_{\hat{\mathcal{G}}}(d)$, completing the induction step for graphs of size d .

By the principle of induction, we have shown that $\mathcal{G} = \hat{\mathcal{G}}$ holds for any graphs under given assumptions. \square

Now that we have established that $\mathcal{G} = \hat{\mathcal{G}}$, we prove the remaining part of the theorem. Note that for any $i, j \in [d]$ such that $i \notin \overline{\text{pa}}_{\mathcal{G}}(j)$, we have $(B_k)_{ji} = (\hat{B}_k)_{ji} = 0, \forall k \in [K]$. Since $\hat{B}_k = B_k T$, we have

$$\sum_{\ell \in \overline{\text{pa}}_{\mathcal{G}}(j)} (B_k)_{j\ell} T_{\ell i} = 0.$$

By Assumption 5, the above implies that $T_{\ell i} = 0$ for $\forall \ell \in \overline{\text{pa}}_{\mathcal{G}}(j)$. In short, we have argued that if there exists j such that $i \notin \overline{\text{pa}}_{\mathcal{G}}(j)$ and $\ell \in \overline{\text{pa}}_{\mathcal{G}}(j)$, then $T_{\ell i} = 0$.

This implies that $T_{\ell i}$ is non-zero only if $\bar{\text{ch}}_{\mathcal{G}}(\ell) \subseteq \bar{\text{ch}}_{\mathcal{G}}(i)$. Since $v = Tz$, we have $v_{\ell} = \sum_{i=1}^d T_{\ell i} z_i = \sum_{i \in [d]: \bar{\text{ch}}_{\mathcal{G}}(\ell) \subseteq \bar{\text{ch}}_{\mathcal{G}}(i)} T_{\ell i} z_i$. Note that when $i \neq \ell$, $\bar{\text{ch}}_{\mathcal{G}}(\ell) \subseteq \bar{\text{ch}}_{\mathcal{G}}(i)$ is equivalent to $i \in \text{sur}_{\mathcal{G}}(\ell)$, so v_{ℓ} only depends on $z_{\text{sur}_{\mathcal{G}}(\ell)}$ by Lemma 6, as desired.

H.2 Formal version and proof of Theorem 2

In previous works [36, 54], it is common to consider single-node soft interventions in the following sense:

Assumption 8. For $\forall 2 \leq k \leq K$, there exists $i_k \in [d]$, such that the structural equation in environment k satisfies (4) satisfies $w_k(i) = w_1(i)$ and $\omega_{k,i,i} = \omega_{1,i,i}$ for $\forall i \neq i_k$.

Let $S_i = \{k : 2 \leq k \leq K, i_k = i\}, i \in [d]$ and $s_i = |S_i|$. Suppose that \mathcal{G} has $e = \sum_{i=1}^d |\text{pa}_{\mathcal{G}}(i)|$ edges, then we can view the weight vectors $\{(w_k(i), \omega_{k,i,i}) : k = 1 \text{ or } i = i_k\}$ as elements of the Euclidean space $\mathbb{R}^{e + \sum_{k=2}^K |\text{pa}_{\mathcal{G}}(i_k)|} \times \mathbb{R}_+^{d+K-1}$. Under Assumption 8, the models can be fully determined by these weight vectors. The following result states that if we restrict ourselves to single-node interventions, then in the worst case, $\Theta(d^2)$ interventions are required.

Theorem 6. There exists a causal graph \mathcal{G} with $\Theta(d^2)$ edges, such that for any unmixing matrix $H \in \mathbb{R}^{d \times n}$ with full row rank, any independent noise variables ϵ , and any $s_i > 0, i \in [d]$ such that $s_i \leq |\text{pa}_{\mathcal{G}}(i)|$ for some i , the following holds: except from a null set of the weight space $\mathbb{R}^{e + \sum_{k=2}^K |\text{pa}_{\mathcal{G}}(i_k)|} \times \mathbb{R}_+^{d+K-1}$ (w.r.t the Lebesgue measure), there must exist a candidate solution $(\hat{H}, \hat{\mathcal{G}})$ and a hypothetical data generating process

$$\forall k \in [K], \quad v = \hat{A}_k v + \hat{\Omega}_k^{\frac{1}{2}} \epsilon, \quad x = \hat{H}^\dagger v$$

such that

- (i') the unmixing matrix $\hat{H} \in \mathbb{R}^{d \times n}$ has full row rank;
- (ii') $\forall k \in [K]$ and $i, j \in [d]$, $(\hat{A}_k)_{ij} \neq 0 \Leftrightarrow j \in \text{pa}_{\hat{\mathcal{G}}}(i)$ and $\hat{\Omega}_k$ is a diagonal matrix with positive entries;
- (iii') for $\forall 2 \leq k \leq K$, the weight matrices $\hat{A}_k, \hat{\Omega}_k$ of environment E_k are from a single-node soft intervention on E_1 on node i_k , in the sense of Assumption 8,

but \mathcal{G} is non-isomorphic to $\hat{\mathcal{G}}$.

In this subsection we give the full proof of Theorem 6. We say that $S \subseteq \mathbb{R}^m$ is a null set if it has zero Lebesgue measure. Obviously, any hyperplanes in \mathbb{R}^m are null sets. We will also need the following simple lemma:

Lemma 13. Suppose that $m \in \mathbb{Z}_+$ and V is a subspace of \mathbb{R}^m . Then for any set of vectors $\mathbf{u}_i \in \mathbb{R}^m, i = 1, 2, \dots, n$ that does not lie in V , there must exist $\mathbf{v} \in \mathbb{R}^m$ such that $\mathbf{u}_i^\top \mathbf{v} \neq 0, \forall i \in [n]$ but $\mathbf{v} \in V^\perp$, where V^\perp is the orthogonal space of V .

Proof. Let \mathbf{w}_i be the orthogonal projection of \mathbf{u}_i onto V^\perp . Since $\mathbf{u}_i \notin V$, we know that $\mathbf{w}_i \neq \mathbf{0}$. The solution space of each equation $\mathbf{w}_i^\top \mathbf{v} = 0$ in V^\perp must then be a proper subspace of V^\perp . Equipped with the Lebesgue measure, all these spaces are null sets in V^\perp , so one can always choose a $\mathbf{v} \in V^\perp$ that does not lie in any of these solution spaces. Such \mathbf{v} satisfies all the requirements. \square

We choose \mathcal{G} to be the graph with $i \rightarrow j$ for $\forall 1 \leq i < j \leq d$, so that \mathcal{G} has $\frac{d(d-1)}{2}$ edges. Suppose that $i_0 \in [d]$ satisfies $s_i \leq |\text{pa}_{\mathcal{G}}(i)| - 1$, then we must have $i_0 \geq 2$, so there is an edge $1 \rightarrow i_0$ in \mathcal{G} . Let $\hat{\mathcal{G}}$ be the resulting graph obtained via removing the edge $1 \rightarrow i_0$ in \mathcal{G} , then \mathcal{G} and $\hat{\mathcal{G}}$ are clearly non-isomorphic.

Note that the i -th row of \mathbf{B}_k can be written as $\omega_{k,i,i}^{-\frac{1}{2}} (\mathbf{e}_i - (\mathbf{A}_k)_i)$. Let's choose an lower-triangular matrix $\mathbf{T} = (t_{ij})_{i,j=1}^d \in \mathbb{R}^{d \times d}$ with columns $\mathbf{t}_i, i \in [d]$ such that the following holds:

$$(\mathbf{e}_i - (\mathbf{A}_k)_i)^\top \mathbf{t}_j = \begin{cases} = 0, & \forall k \in \{1\} \cup S_{i_0}, j = 1 \text{ and } i = i_0 \\ > 0, & \forall i = j \text{ and } k \in \{1\} \cup S_i \\ \neq 0, & \forall \text{ remaining } (i, j, k) \in \{k = 1, j < i\} \cup \{k \geq 2, i = i_k, j < i\} \end{cases} \quad (13)$$

and

$$t_{ii} \neq 0, \quad \forall i \in [d]. \quad (14)$$

We now show that: except from a null set in the weight space, such \mathbf{T} can always be chosen. To see why this is the case, we first consider all the constraints on \mathbf{t}_1 :

$$(\mathbf{e}_i - (\mathbf{A}_k)_i)^\top \mathbf{t}_1 = \begin{cases} = 0, & \forall k \in \{1\} \cup S_{i_0} \text{ and } i = i_0 \\ > 0, & \forall i = 1 \text{ and } k \in \{1\} \cup S_i \\ \neq 0, & \forall \text{ remaining } (i, k) \in \{k = 1, i > 1\} \cup \{k \geq 2, i = i_k > 1\} \end{cases} \quad (15)$$

Now let $V = \text{span} \langle \mathbf{e}_i - (\mathbf{A}_k)_i : k \in \{1\} \cup S_{i_0} \text{ and } i = i_0 \rangle$ and R be the set of pairs (i, k) specified in the second and third row of (15). For $\forall (i, k)$, let $\mathbf{w}_k(i)$ be the weight vector of node i in the environment k , i.e., the vector of nonzero entries in $(\mathbf{A}_k)_i$. Then for $\forall (i, k) \in R$, the following set (as a subset of the weight space)

$$\bigcup_{k^* \in \{1\} \cup S_{i_0}} \{ \mathbf{e}_{i_0} - \mathbf{w}_{k^*}(i_0) \in \text{span} \langle \mathbf{e}_i - \mathbf{w}_k(i), \mathbf{e}_{i_0} - \mathbf{w}_{k'}(i_0) : k' \in \{1\} \cup S_{i_0} \setminus \{k^*\} \rangle \} \quad (16)$$

must be a null set. Thus

$$\mathbf{E} = \bigcup_{(i,k) \in R} \bigcup_{k^* \in \{1\} \cup S_{i_0}} \{ \mathbf{e}_{i_0} - \mathbf{w}_{k^*}(i_0) \in \text{span} \langle \mathbf{e}_i - \mathbf{w}_k(i), \mathbf{e}_{i_0} - \mathbf{w}_{k'}(i_0) : k' \in \{1\} \cup S_{i_0} \setminus \{k^*\} \rangle \} \quad (17)$$

is also a null set. For any weights that are not in \mathbf{E} , we necessarily have

$$\mathbf{e}_i - \mathbf{w}_k(i) \notin \text{span} \langle \mathbf{e}_i - (\mathbf{A}_k)_i : k \in \{1\} \cup S_{i_0} \text{ and } i = i_0 \rangle = V, \quad (i, k) \in R.$$

Let $U = \{ \mathbf{e}_i - \mathbf{w}_k(i) : (i, k) \in R \}$, then we can apply Lemma 13 to deduce that there exists \mathbf{t}_1 such that

$$(\mathbf{e}_i - (\mathbf{A}_k)_i)^\top \mathbf{t}_1 = \begin{cases} = 0, & \forall k \in \{1\} \cup S_{i_0} \text{ and } i = i_0 \\ \neq 0, & \forall \text{ remaining } (i, k) \in \{k = 1\} \cup \{k \geq 2, i = i_k\} \end{cases} \quad (18)$$

Note that the only difference between (18) and (15) is that the latter one further requires that

$$(\mathbf{e}_1 - (\mathbf{A}_k)_1)^\top \mathbf{t}_1 > 0, \quad \forall k \in \{1\} \cup S_i.$$

while the former only guarantees that these terms are nonzero. However, recall that $(\mathbf{A}_k)_{ij} \neq 0 \Rightarrow j \in \text{pa}_{\mathcal{G}}(i) \Rightarrow j < i$, so the above essentially says that $t_{11} > 0$. This can be easily guaranteed by replacing the solution \mathbf{t}_1 we obtained satisfying (18) with $-\mathbf{t}_1$ if needed.

Assuming that the weights do not lie in the null set \mathcal{E} we have shown that \mathbf{t}_1 can always be chosen to satisfy all constraints imposed on it. We now proceed to choose the remaining entries of \mathbf{T} . The remaining entries in \mathbf{t}_1 can be chosen arbitrarily. For $\mathbf{t}_j, j > 1$, we note that the remaining constraints in (13) that need to be satisfied consist of the "nonzero" part and the "positivity" part. The positivity constraints can always be satisfied by choosing a sufficiently large t_{jj} for $j > 1$.

After choosing the \mathbf{t}_j 's satisfying the positivity constraints, the nonzero constraints along with (14) are easy to fulfill by slightly perturbing \mathbf{t}_j if they are violated; since each of these constraints are only violated in a zero-measure set of the weight space. Hence, we have shown that except a null set \mathcal{E} in the weight space, there always exists some \mathbf{T} satisfying (13). Such \mathbf{T} must be invertible since it is lower-triangular and its diagonal entries are nonzero. Now let $\hat{\mathbf{H}} = \mathbf{T}^{-1}\mathbf{H}$ and $\hat{\Omega}_k$ be the diagonal matrix with entries $\hat{\omega}_{k,i,i} = t_{ii}^{-2} \cdot \omega_{k,i,i}, i \in [d]$ and

$$\hat{\mathbf{A}}_k = \mathbf{I} - \hat{\Omega}_k^{\frac{1}{2}} \Omega_k^{-\frac{1}{2}} (\mathbf{I} - \mathbf{A}_k) \mathbf{T}. \quad (19)$$

First since \mathbf{T} is invertible and \mathbf{H} has full rank, $\hat{\mathbf{H}}$ must also have full row rank. Second,

$$(\hat{\mathbf{A}}_k)_{ij} = \begin{cases} 1 - \hat{\omega}_{k,i,i}^{\frac{1}{2}} \omega_{k,i,i}^{-\frac{1}{2}} t_{ii} = 0 & \text{if } j = i \\ -\hat{\omega}_{k,i,i}^{\frac{1}{2}} \omega_{k,i,i}^{-\frac{1}{2}} (\mathbf{e}_i - (\mathbf{A}_k)_i)^\top \mathbf{t}_j = 0 & \text{if } j > i \\ -\hat{\omega}_{k,i,i}^{\frac{1}{2}} \omega_{k,i,i}^{-\frac{1}{2}} (\mathbf{e}_i - (\mathbf{A}_k)_i)^\top \mathbf{t}_j & \text{if } j < i. \end{cases}$$

where we again recall that both \mathbf{A}_k and \mathbf{T} are lower-triangular. From (13) we can see that

- When $i = i_0$ and $j = 1$, we have
 - $(\hat{\mathbf{A}}_k)_{i_0,1} = 0$ if $k \in \{1\} \cup S_{i_0}$, and
 - $(\hat{\mathbf{A}}_k)_{i_0,1} = (\hat{\mathbf{A}}_1)_{i_0,1} = 0$ if $k \notin \{1\} \cup S_{i_0}$, by definition of S_{i_0} and Assumption 8.
- When $i > j$ and $(i, j) \neq (i_0, 1)$, we have
 - $(\hat{\mathbf{A}}_k)_{ij} \neq 0$ if $k = 1$ or $i = i_k$, which directly follows from (13), and
 - $(\hat{\mathbf{A}}_k)_{ij} = (\hat{\mathbf{A}}_1)_{ij} \neq 0$, by Assumption 8.

To summarize, for each k , $(\mathbf{A}_k)_{ij} \neq 0 \Leftrightarrow j \in \text{pa}_{\mathcal{G}}(i)$ and $(i, j) \neq (i_0, 1)$.

Finally, let $\hat{\mathbf{w}}_k(i)$ be the weight vector of node i in environment k in the hypothetical model *i.e.*, the vector of nonzero entries in $(\mathbf{A}_k)_i$, and \mathbf{T}_S be the submatrix of \mathbf{T} by selecting the rows and columns in the index set S , then by (19) we have that

$$\hat{\omega}_{k,i,i} = t_{ii}^{-2} \cdot \omega_{k,i,i}, \quad \hat{\omega}_{k,i,i}^{\frac{1}{2}} \omega_{k,i,i}^{-\frac{1}{2}} \mathbf{w}_k(i) \mathbf{T}_{\text{pa}_{\mathcal{G}}(i)} = \begin{cases} \hat{\mathbf{w}}_k(i) & \text{if } i \neq i_0 \\ [0, \hat{\mathbf{w}}_k(i)] & \text{if } i = i_0 \end{cases} \quad (20)$$

By our assumption, for $\forall k \geq 2, i \neq i_k \Rightarrow \mathbf{w}_k(i) = \mathbf{w}_1(i)$ and $\omega_{k,i,i} = \omega_{1,i,i}$. Thus (20) imply that $\forall k \geq 2, i \neq i_k \Rightarrow \hat{\mathbf{w}}_k(i) = \hat{\mathbf{w}}_1(i)$ and $\hat{\omega}_{k,i,i} = \hat{\omega}_{1,i,i}$. In other words, a single-node intervention on node i_k in environment k in the ground-truth model corresponds to a single-node intervention on node i_k in environment k in the hypothetical model, thereby completing the proof.

H.3 Proof of Theorem 4

We first prove two lemmas.

Lemma 14. $\forall i \in [d]$, we have $\text{span} \langle (\mathbf{M}_k)_i : k \in [K] \rangle = \text{span} \langle \mathbf{h}_j : j \in \overline{\text{pa}}_{\mathcal{G}}(i) \rangle$.

Proof. Since $(\mathbf{M}_k)_i = (\mathbf{B}_k)_i \mathbf{H}$, and $(\mathbf{B}_k)_{ij} \neq 0 \Leftrightarrow j \in \overline{\text{pa}}_{\mathcal{G}}(i)$, we can see that $(\mathbf{M}_k)_i \in \text{span} \langle \mathbf{h}_j : j \in \overline{\text{pa}}_{\mathcal{G}}(i) \rangle$. On the other hand, since \mathbf{H} is invertible, by Assumption 5 we have $\dim \text{span} \langle (\mathbf{M}_k)_i : k \in [K] \rangle = \dim \text{span} \langle (\mathbf{B}_k)_i : k \in [K] \rangle = |\overline{\text{pa}}_{\mathcal{G}}(i)|$. Thus we must have $\text{span} \langle (\mathbf{M}_k)_i : k \in [K] \rangle = \text{span} \langle \mathbf{h}_j : j \in \overline{\text{pa}}_{\mathcal{G}}(i) \rangle$. \square

Lemma 15. Let \hat{S} be an ancestral set of graph \mathcal{G} and $\hat{\mathbf{V}}_k = \text{span} \langle (\mathbf{M}_k)_s : s \in \hat{S} \rangle, k \in [K]$. Then we have $\mathbf{V}_1 = \mathbf{V}_2 = \dots = \mathbf{V}_K = \text{span} \langle \mathbf{h}_s : s \in \hat{S} \rangle$.

Proof. Recall that $M_k = B_k H$, so for $\forall s \in \hat{S}$, the s -th row of M_k can be written as

$$(M_k)_s = \sum_{t=1}^d (B_k)_{st} \mathbf{h}_t = \sum_{t \in \overline{\text{pa}}_G(s)} (B_k)_{st} \mathbf{h}_t \in \text{span} \langle \mathbf{h}_s : s \in \hat{S} \rangle \quad (21)$$

where the last equation is because \hat{S} is ancestral $\Rightarrow \overline{\text{pa}}_G(s) \subseteq \hat{S}$. Thus, for $\forall k \in [K]$, $\hat{V}_k = \text{span} \langle (M_k)_s : s \in \hat{S} \rangle \subseteq \text{span} \langle \mathbf{h}_s : s \in \hat{S} \rangle$. On the other hand, recall that both B_k and H have full rank, so M_k has full row rank as well, which implies that $\dim V_k = |S| = \dim \text{span} \langle \mathbf{h}_s : s \in \hat{S} \rangle$. Hence, $V_k = \text{span} \langle \mathbf{h}_s : s \in \hat{S} \rangle, \forall k \in [K]$. \square

The following two propositions show that our algorithm always maintain an ancestral set, recursively adds a new node into the set and correctly identifies its parents.

Proposition 3 (Proposition 1 restated). *The following two propositions hold for Algorithm 3:*

- $\text{ans}_G(i) \subseteq S \Leftrightarrow$ the if condition in line 8 of Algorithm 3 is fulfilled;
- the set S maintained in Algorithm 3 is always an ancestral set, in the sense that $j \in S \Rightarrow \text{ans}_G(j) \subseteq S$.

Proof. At the starting point, we have $S = \emptyset$ which is obviously an ancestral set. Now suppose that after the ℓ -th iteration, $S = \{s_1, s_2, \dots, s_\ell\}$ is an ancestral set. In the following, we show that $\text{ans}_G(i) \subseteq S \Leftrightarrow$ the if condition in line 8 is fulfilled. This would immediately imply that there always exists a node i that can be added into S in the $(\ell + 1)$ -th iteration, and that after adding i , S is still an ancestral set.

Suppose that $\text{ans}_G(i) \subseteq S$ for some $i \notin S$, by Lemma 14 we know that $(M_k)_i \in \text{span} \langle \mathbf{h}_j : j \in \overline{\text{pa}}_G(i) \rangle$, so there exists $\alpha_k \in \mathbb{R}$ such that $(M_k)_i - \alpha_k \mathbf{h}_i \in \text{span} \langle \mathbf{h}_j : j \in \text{pa}_G(i) \rangle$. Moreover, since $(M_k)_i = \sum_{j \in \overline{\text{pa}}_G(i)} (B_k)_{ij} \mathbf{h}_j$, $(B_k)_{ii} = \omega_{k,i,i}^{-\frac{1}{2}} \neq 0$ and H has full row rank by assumption, we must have $(M_k)_i \notin \text{span} \langle \mathbf{h}_j : j \in \text{pa}_G(i) \rangle$ and so $\alpha_k \neq 0$. Thus, we have by the linearity of the projection operator

$$\mathbf{q}_k := \text{proj}_{V_k^\perp}((M_k)_i) = \text{proj}_{V_k^\perp}((M_k)_i - \alpha_k \mathbf{h}_i) + \text{proj}_{V_k^\perp}(\alpha_k \mathbf{h}_i) = \alpha_k \text{proj}_{V_k^\perp}(\mathbf{h}_i).$$

Recall that all the V_k 's are the same and equal $\text{span} \langle \mathbf{h}_s : s \in S \rangle$ by Lemma 15. So $\dim \text{span} \langle \mathbf{q}_k : k \in [K] \rangle \leq 1$. Since H has full row rank, we have $\mathbf{h}_i \notin \text{span} \langle \mathbf{h}_s : s \in S \rangle = V_k$, so that $\dim \text{span} \langle \mathbf{q}_k : k \in [K] \rangle = 1$ holds, which is exactly the if condition in line 8.

Conversely, suppose that there is an $i \notin S$ such that $\text{ans}_G(i) \not\subseteq S$ but $\dim \text{span} \langle \mathbf{q}_k : k \in [K] \rangle = 1$ holds. Since S is ancestral, we know that there must be some $j \in \text{pa}_G(i)$ such that $j \notin S$. Since \mathbf{e}_i and \mathbf{e}_j both have support on the coordinates in $\overline{\text{pa}}_G(i)$, by Assumption 5 we know that $\text{span} \langle \mathbf{e}_i, \mathbf{e}_j \rangle \subseteq \text{span} \langle (B_k)_i : k \in [K] \rangle$, so that $\text{span} \langle \mathbf{h}_i, \mathbf{h}_j \rangle = \text{span} \langle \mathbf{e}_i, \mathbf{e}_j \rangle H \subseteq \text{span} \langle (B_k)_i : k \in [K] \rangle H = \text{span} \langle (M_k)_i : k \in [K] \rangle$. Since $\dim \text{span} \langle \mathbf{q}_k : k \in [K] \rangle = 1$, there must exist some vector $\mathbf{u} \in \mathbb{R}^n$ and $\alpha_i, \alpha_j \in \mathbb{R}$ such that $\mathbf{h}_i - \alpha_i \mathbf{u}, \mathbf{h}_j - \alpha_j \mathbf{u} \in V_k = \text{span} \langle \mathbf{h}_s : s \in S \rangle$. Since $i, j \notin S$ and H has full row rank, we can deduce that $\mathbf{h}_i, \mathbf{h}_j \notin \text{span} \langle \mathbf{h}_s : s \in S \rangle$, and so both of α_i and α_j are non-zero. Hence $\alpha_j \mathbf{h}_i - \alpha_i \mathbf{h}_j \in \text{span} \langle \mathbf{h}_s : s \in S \rangle$, which is impossible since we know that H has full row-rank. \square

Proposition 4 (Proposition 2 restated). *Given any ordered ancestral set S that contains $\text{pa}_G(i)$ for some $i \notin S$, Algorithm 2 returns a set $P_i \subseteq S$ that is exactly $\text{pa}_G(i)$.*

Proof. As we have shown in Proposition 1, for each possible input (S, i) to Algorithm 2, both S and $S \cup \{i\}$ are ancestral sets, so that $\text{ans}_G(i) \subseteq S$. Similarly one can see that inside the set $S := \{s_1, s_2, \dots, s_m\}$, all the ancestors of s_j are contained in $\{s_1, s_2, \dots, s_{j-1}\}$. In the following, we show that $\forall m' \in \{0, \dots, m\}, r_{m'} = |\overline{\text{pa}}_G(i) - \{s_j : j \leq m'\}|$ (*).

By [Lemma 15](#) we have $\mathbf{W}_1 = \mathbf{W}_2 = \dots = \mathbf{W}_K = \text{span} \langle \mathbf{h}_{s_j} : j \leq m' \rangle$. Let t_1, t_2, \dots, t_ℓ be elements of $\overline{\text{pa}}_{\mathcal{G}}(i)$ that are not in $\{s_j : j \leq m'\}$, then

$$\begin{aligned} r_{m'} &= \dim \text{span} \langle \mathbf{p}_k : k \in [K] \rangle = \dim \left(\text{proj}_{\text{span} \langle \mathbf{h}_{s_j} : j \leq m' \rangle^\perp} \text{span} \langle (\mathbf{M}_k)_i : k \in [K] \rangle \right) \\ &= \dim \left(\text{proj}_{\text{span} \langle \mathbf{h}_{s_j} : j \leq m' \rangle^\perp} \text{span} \langle \mathbf{h}_j : j \in \overline{\text{pa}}_{\mathcal{G}}(i) \rangle \right) \quad (\text{by } \textcolor{red}{\text{Lemma 14}}) \\ &= \dim \left(\text{proj}_{\text{span} \langle \mathbf{h}_{s_j} : j \leq m' \rangle^\perp} \text{span} \langle \mathbf{h}_{t_1}, \mathbf{h}_{t_2}, \dots, \mathbf{h}_{t_\ell} \rangle \right) \\ &= \ell \quad (\text{by } \textcolor{red}{\text{Lemma 4}} \text{ and non-degeneracy of } \mathbf{H}) \end{aligned}$$

which proves (*). From (*) it is easy to see that $m' \in \overline{\text{pa}}_{\mathcal{G}}(i)$ (and thus in $\text{pa}_{\mathcal{G}}(i)$ since $i \notin S$) if and only if $r_{m'} = r_{m'-1} - 1$. \square

Now we conclude the proof of [Theorem 4](#). [Propositions 1](#) and [2](#) directly imply that [Algorithm 3](#) is able to exactly recover the ground-truth causal graph \mathcal{G} . It remains to show that Line 20 in [Algorithm 3](#) produces the correct $\hat{\mathbf{h}}_i$'s. By [Lemma 14](#) we know that $E_j = \text{span} \langle \mathbf{h}_\ell : \ell \in \overline{\text{pa}}_{\mathcal{G}}(j) \rangle$, so

$$\cap_{j \in \overline{\text{ch}}_{\mathcal{G}}(i)} E_j = \cap_{j \in \overline{\text{ch}}_{\mathcal{G}}(i)} \text{span} \langle \mathbf{h}_\ell : \ell \in \overline{\text{pa}}_{\mathcal{G}}(j) \rangle = \text{span} \langle \mathbf{h}_\ell : \ell \in \overline{\text{sur}}_{\mathcal{G}}(i) \rangle.$$

where the last step holds because \mathbf{H} has full row rank and $\cap_{j \in \overline{\text{ch}}_{\mathcal{G}}(i)} \overline{\text{pa}}_{\mathcal{G}}(j) = \overline{\text{sur}}_{\mathcal{G}}(i)$ by definition. Hence, each $\hat{\mathbf{h}}_i$ is a linear combination of $\mathbf{h}_\ell, \ell \in \overline{\text{sur}}_{\mathcal{G}}(i)$, completing the proof.

I Identification limit of general causal models with soft interventions

While [Theorem 1](#) guarantees identifiability with general environments, it only applies to linear causal models. In this section, we show that if we have access to single-node soft interventions, then we can identify general non-parametric causal models up to \sim_{sur} . To obtain our identifiability result, we also require that the environments are non-degenerate in the following sense:

Definition 11 (Non-degeneracy set of interventions). Let $\hat{p}_k(z_i | z_{\text{pa}_{\mathcal{G}}(i)}), k \in [K_i]$ be conditional probability densities at node i , then $\{\hat{p}_k\}_{k=1}^{K_i}$ is said to be non-degenerate on node i at point $\hat{\mathbf{z}} \in \mathbb{R}^d$ if all these conditional densities are well-defined and positive at $\hat{\mathbf{z}}$, and the matrix

$$\left[\frac{\partial (\hat{p}_1 / \hat{p}_k)}{\partial \mathbf{z}_j} \right]_{2 \leq k \leq K_i, j \in \overline{\text{pa}}_{\mathcal{G}}(i)} \bigg|_{\mathbf{z}=\hat{\mathbf{z}}} \in \mathbb{R}^{(K_i-1) \times (|\text{pa}_{\mathcal{G}}(i)|+1)}$$

has full row rank. Moreover, we say that $\{\hat{p}_k\}_{k=1}^{K_i}$ is non-degenerate in a point set \mathcal{O} if for all $\hat{\mathbf{z}} \in \mathcal{O}$, it is non-degenerate at $\hat{\mathbf{z}}$.

The following lemma shows how [Definition 11](#) is related to [Assumption 5](#) in the linear setting:

Lemma 16. Suppose that $\hat{p}_k(\mathbf{z}) = \prod_{i=1}^d \hat{p}_k(z_i | z_{\text{pa}_{\mathcal{G}}(i)}), k \in [K]$ be probability distributions of latent variables \mathbf{z} generated from the linear causal models [\(3\)](#), such that for $\forall i \in [d]$, $\hat{p}_k(z_i | z_{\text{pa}_{\mathcal{G}}(i)}), k \in [K]$ are non-degenerate on node i in the sense of [Definition 11](#). Then the corresponding matrices $\mathbf{B}_k, k \in [K]$ satisfy [Assumption 5](#).

Now we are ready to state our main result in this section:

Theorem 7. Suppose that we have access to observations generated from multiple environments $\{P_{\mathbf{X}}^E\}_{E \in \mathfrak{E}}$. Let $(\hat{\mathbf{h}}, \hat{\mathcal{G}})$ be any candidate solution with data generated according to [Assumption 1](#) with latent variables $\mathbf{v} = \hat{\mathbf{h}}(\mathbf{x})$ and joint distribution q_E with factors q_i^E . Assuming that

- (i) the joint densities $\{p_E(\mathbf{z})\}_{E \in \mathfrak{E}}$ are continuous differentiable on \mathbb{R}^d with common support $\mathcal{O}_{\mathbf{z}}$, and $\{q_E(\mathbf{v})\}_{E \in \mathfrak{E}}$ are continuous differentiable on \mathbb{R}^d with common support $\mathcal{O}_{\mathbf{v}}$;
- (ii) we have access to multiple single-node soft interventions on each node with unknown targets: there exists a partition $\mathfrak{E} = \cup_{i=1}^d \mathfrak{E}_i$ such that $\mathcal{I}_{\mathbf{z}}^{\mathfrak{E}_i} = \{\pi(i)\}, \mathcal{I}_{\mathbf{v}}^{\mathfrak{E}_i} = \{\pi'(i)\}, \forall i \in [d]$ for some unknown permutations π and π' on $[d]$;

- (iii) the intervention distributions on each node are non-degenerate in the sense of **Definition 11**: there exists $N_z \subseteq O_z$ and $N_v \subseteq O_v$ satisfying $N_z^\circ = N_v^\circ = \emptyset$ where S° denotes the interior of a set S , such that for all $i \in [d]$, $\{p_i^E(\cdot) : E \in \mathfrak{E}_{\pi^{-1}(i)}\}$ (resp. $\{q_i^E(\cdot) : E \in \mathfrak{E}_{\pi'^{-1}(i)}\}$) is non-degenerate on node i in $O_z \setminus N_z$ (resp. $O_v \setminus N_v$).

Then we must have $(\mathbf{h}, \mathcal{G}) \sim_{\text{sur}} (\hat{\mathbf{h}}, \hat{\mathcal{G}})$.

Previous works on the identifiability of non-parametric causal models typically require that all the joint distributions are supported on the whole space \mathbb{R}^d [49, 23, 47]. In contrast, we only assume that the densities have common and unknown support across all interventions.

Theorem 7 can be regarded as a soft-intervention version of 49, Theorem 4.3, which assumes access to hard interventions and only need two paired interventions per node. While they are able to show full identifiability, we show in the following that identifiability up to \sim_{sur} is the best we can hope for with soft interventions.

Theorem 8 (Counterpart to **Theorem 7**, informal version of **Theorem 10**). *For any causal model $(\mathbf{h}, \mathcal{G})$ and any set of environments $\mathfrak{E} = \{E_k : k \in [K]\}$ such that all conditions in **Theorem 7** are satisfied, there must exist a candidate solution $(\hat{\mathbf{h}}, \mathcal{G})$ and a hypothetical data generating process that satisfy the same set of conditions, but*

$$\frac{\partial \mathbf{v}_i}{\partial \mathbf{z}_j} \neq 0, \quad \forall j \in \overline{\text{sur}}_{\mathcal{G}}(i).$$

Finally, the ambiguity still exists if we additionally assume standard axioms such as causal minimality (**Assumption 6**) and faithfulness (**Assumption 7**) on the causal model.

I.1 Proof of **Lemma 16**

Let $\mathbf{w}_k(i) \in \mathbb{R}^{|\text{pa}_{\mathcal{G}}(i)|}$ be the vector obtained by removing all zero entries in the i -th row of \mathbf{A}_k and $\omega_{k,i,i}$ be the i -th diagonal entry in Ω_k , then for the k -th environment we have $\mathbf{z}_i = \mathbf{w}_k(i)^\top \mathbf{z}_{\text{pa}_{\mathcal{G}}(i)} + \omega_{k,i,i}^{-\frac{1}{2}} \epsilon_i$, so that

$$\hat{p}_k(\mathbf{z}_i | \mathbf{z}_{\text{pa}_{\mathcal{G}}(i)}) = \omega_{k,i,i}^{-\frac{1}{2}} p_{\epsilon_i} \left(\omega_{k,i,i}^{-\frac{1}{2}} (\mathbf{z}_i - \langle \mathbf{w}_k(i), \mathbf{z}_{\text{pa}_{\mathcal{G}}(i)} \rangle) \right)$$

where $p_{\epsilon_i}(\cdot)$ is the density of ϵ_i . As a result, we have

$$\begin{aligned} \nabla \frac{\hat{p}_1}{\hat{p}_k}(\mathbf{z}_i | \mathbf{z}_{\text{pa}_{\mathcal{G}}(i)}) &= \frac{\hat{p}_1}{\hat{p}_k}(\mathbf{z}_i | \mathbf{z}_{\text{pa}_{\mathcal{G}}(i)}) \cdot \nabla \log \frac{\hat{p}_1}{\hat{p}_k}(\mathbf{z}_i | \mathbf{z}_{\text{pa}_{\mathcal{G}}(i)}) \\ &= \frac{\hat{p}_1}{\hat{p}_k}(\mathbf{z}_i | \mathbf{z}_{\text{pa}_{\mathcal{G}}(i)}) \cdot [c_{i1}(1, -\mathbf{w}_1(i)) - c_{ik}(1, -\mathbf{w}_k(i))] \end{aligned}$$

where for convenience we use ∇ to denote the gradient with respect to all variables $\mathbf{z}_{\overline{\text{pa}}_{\mathcal{G}}(i)}$, and $c_{ik} = \omega_{k,i,i}^{-\frac{1}{2}} \cdot \frac{p'_{\epsilon_i}}{p_{\epsilon_i}} \left(\omega_{k,i,i}^{-\frac{1}{2}} (\mathbf{z}_i - \langle \mathbf{w}_k(i), \mathbf{z}_{\text{pa}_{\mathcal{G}}(i)} \rangle) \right)$ (we omit the dependency on \mathbf{z} for simplicity).

Definition 11 implies that $\text{span} \langle c_{i1}(1, -\mathbf{w}_1(i)) - c_{ik}(1, -\mathbf{w}_k(i)) : 2 \leq k \leq K \rangle = \mathbb{R}^{|\text{pa}_{\mathcal{G}}(i)|+1}$, thus it holds that $\text{span} \langle (1, -\mathbf{w}_k(i)) : k \in [K] \rangle = \mathbb{R}^{|\text{pa}_{\mathcal{G}}(i)|+1}$ as well. By definition of \mathbf{B}_k , this immediately implies that $\dim(\text{span} \langle (\mathbf{B}_k)_i : k \in [K] \rangle) = |\text{pa}_{\mathcal{G}}(i)| + 1$ as desired.

I.2 Proof of **Theorem 7**

Define $\tau := \hat{\mathbf{h}} \circ \mathbf{h}^{-1} : \mathbb{R}^d \mapsto \mathbb{R}^d$, then we have that $\mathbf{v} = \tau(\mathbf{z})$. Since both \mathbf{h} and $\hat{\mathbf{h}}$ are diffeomorphisms by assumption, so is τ . To avoid confusion, in this section we use \mathbf{z} (resp. \mathbf{v}) to denote random variables while using $\hat{\mathbf{z}}$ (resp. $\hat{\mathbf{v}}$) to denote (deterministic) vectors.

Let $\mathfrak{E}_j = \{E_k^{(j)} : k \in [K_j]\}$ be the j -th collection of environments according to our assumption. We first prove the following lemma:

Lemma 17. $O_v = \tau(O_z)$.

Proof. By the change of variable formula [35], for $\forall \hat{z} \in \mathbb{R}^d$ and $\forall E \in \mathfrak{E}$ we have $p_E(\hat{z}) = q_E(\hat{v}) |\det \mathbf{J}_\tau(\hat{z})|$, where $\hat{v} = \tau(\hat{z})$. Since τ is a diffeomorphism, we must have $|\det \mathbf{J}_\tau(\hat{z})| \neq 0$, so $\hat{z} \in \mathbf{O}_z \Leftrightarrow \hat{v} = \tau(\hat{z}) \in \mathbf{O}_v$, concluding the proof. \square

Lemma 18. Let $\hat{z} \in \mathbf{O}_z$. For $\forall j \in [d]$ and $2 \leq k \leq K_j$, we have

$$\frac{p_j^{E_k^{(j)}}}{p_j^{E_1^{(j)}}}(\hat{z}_j \mid \hat{z}_{\text{pa}_G(j)}) = \frac{q_j^{E_k^{(j)}}}{q_j^{E_1^{(j)}}}(\hat{v}_j \mid \hat{v}_{\text{pa}_G(j)}), \quad (22)$$

where $\hat{v} = \tau(\hat{z}) \in \mathbf{O}_v$.

Proof. Since $v = \tau(z)$, by the change-of-measure formula [35] we have that for $\forall \hat{z} \in \mathbf{O}_z$,

$$\prod_{i=1}^d p_i^E(\hat{z}_i \mid \hat{z}_{\text{pa}_G(i)}) = p_E(\hat{z}) = q_E(\hat{v}) |\det \mathbf{J}_\tau(\hat{z})| = \prod_{i=1}^d q_i^E(\hat{v}_i \mid \hat{v}_{\text{pa}_G(i)}) |\det \mathbf{J}_\tau(\hat{z})| \quad (23)$$

for all $E \in \mathfrak{E}_j$, where $\hat{v} = \tau(\hat{z})$. By Assumption (ii) and Definition 2, we know that $p_i^{E_k^{(1)}} = p_i^{E_1^{(1)}} \Leftrightarrow i \neq 1$ and $q_i^{E_k^{(1)}} = q_i^{E_1^{(1)}} \Leftrightarrow i \neq 1$ for all $k > 1$. Thus, we have that

$$\prod_{i=1}^d \frac{p_i^{E_k^{(j)}}}{p_i^{E_1^{(j)}}}(\hat{z}_i \mid \hat{z}_{\text{pa}_G(i)}) = \frac{p_j^{E_k^{(j)}}}{p_j^{E_1^{(j)}}}(\hat{z}_j \mid \hat{z}_{\text{pa}_G(j)})$$

and

$$\prod_{i=1}^d \frac{q_i^{E_k^{(j)}}}{q_i^{E_1^{(j)}}}(\hat{v}_i \mid \hat{v}_{\text{pa}_G(i)}) = \frac{q_j^{E_k^{(j)}}}{q_j^{E_1^{(j)}}}(\hat{v}_j \mid \hat{v}_{\text{pa}_G(j)}).$$

Since the LHS of the above two equations are the same by (23), the RHS must also be the same, concluding the proof. \square

We assume WLOG that the vertices of \mathcal{G} are labelled such that $i \rightarrow j \Rightarrow i < j$, and that $\pi(i) = i, \forall i \in [d]$. Also we can assume the nodes are fixed and only consider how they are connected, i.e., $\pi'(i) = i, \forall i \in [d]$.¹

Lemma 19. We have $(\tau(N_z))^o = (\tau^{-1}(N_v))^o = \emptyset$.

Proof. The result immediately follows from the assumption that $N_z^o = N_v^o = \emptyset$ and that τ is a diffeomorphism. \square

For any vertex set V , we use \mathcal{G}_V to denote its corresponding induced subgraph of \mathcal{G} . We first prove the following statements by induction on j :

- (1) $\forall i \neq j, i \in \text{pa}_G(j) \Leftrightarrow i \in \text{pa}_{G'}(j)$;
- (2) $\forall j \in [d]$, there exists a continuously differentiable function ϕ_i such that $v_j = \phi_j(z_{\text{pa}_G(j)})$.
Moreover, $\frac{\partial \phi_j}{\partial z_j} \neq 0$ (i.e., not always zero).
- (3) $\forall j \in [d]$, there exists a continuously differentiable function Υ_j such that $v_{\text{pa}_G(j)} = \Upsilon_j(z_{\text{pa}_G(j)})$.

For $j = 1$, by assumption $\text{pa}_G(j) = \emptyset$. Lemma 18 implies that for any $\hat{z} \in \mathbf{O}_z$,

$$\frac{p_1^{E_k^{(1)}}}{p_1^{E_1^{(1)}}}(\hat{z}_1) = \frac{q_1^{E_k^{(1)}}}{q_1^{E_1^{(1)}}}(\hat{v}_1 \mid \hat{v}_{\text{pa}_G(1)}), \forall 2 \leq k \leq K_1. \quad (24)$$

¹This is also WLOG because we now have groups of soft interventions where each group corresponds to a single node, so we can just relabel the node in \mathcal{G} that corresponds to the i -th group as node i .

Then for $\forall i \in \overline{\text{pa}}_{\hat{G}}(1)$, taking the partial derivative w.r.t \mathbf{v}_j gives

$$\frac{\partial}{\partial \hat{\mathbf{v}}_i} \frac{q_1^{E_k^{(1)}}}{q_1^{E_1^{(1)}}} \left(\hat{\mathbf{v}}_1 \mid \hat{\mathbf{v}}_{\text{pa}_{\hat{G}}(1)} \right) = \left(\frac{p_1^{E_k^{(1)}}}{p_1^{E_1^{(1)}}} \right)' (\hat{\mathbf{z}}_1) \cdot \frac{\partial \hat{\mathbf{z}}_1}{\partial \hat{\mathbf{v}}_i} \Rightarrow \nabla_{\mathbf{v}_{\overline{\text{pa}}_{\hat{G}}(1)}} \frac{q_1^{E_k^{(1)}}}{q_1^{E_1^{(1)}}} \left(\hat{\mathbf{v}}_1 \mid \hat{\mathbf{v}}_{\text{pa}_{\hat{G}}(1)} \right) = \left(\frac{p_1^{E_k^{(1)}}}{p_1^{E_1^{(1)}}} \right)' (\hat{\mathbf{z}}_1) \cdot \nabla_{\mathbf{v}_{\overline{\text{pa}}_{\hat{G}}(1)}} \hat{\mathbf{z}}_1.$$

Thus,

$$\text{rank} \left[\nabla_{\mathbf{v}_{\overline{\text{pa}}_{\hat{G}}(1)}} \frac{q_1^{E_k^{(1)}}}{q_1^{E_1^{(1)}}} \left(\hat{\mathbf{v}}_1 \mid \hat{\mathbf{v}}_{\text{pa}_{\hat{G}}(1)} \right) : 2 \leq k \leq K_1 \right] \leq 1.$$

Note that the above inequality holds for $\forall \hat{\mathbf{v}} \in \mathbf{O}_v$. If $\text{pa}_{\hat{G}}(1) \neq \emptyset$, then this would contradict the non-degeneracy assumption (iii) which implies that the above matrix should have $\text{rank} \geq 2$ at some point $\hat{\mathbf{v}} \in \mathbf{O}_v$. Hence we must have $\text{pa}_{\hat{G}}(1) = \emptyset$, implying that (1) holds for $j = 1$.

Taking the derivative of both sides of (24) w.r.t $\mathbf{z}_i, i \geq 2$ implies that $\left(\frac{q_1^{E_k^{(1)}}}{q_1^{E_1^{(1)}}} \right)' (\hat{\mathbf{v}}_1) \cdot \frac{\partial \hat{\mathbf{v}}_1}{\partial \mathbf{z}_i} = 0$. By

our assumption (iii), for $\forall \hat{\mathbf{v}} \in \mathbf{O}_v \setminus \mathbf{N}_v$, there exists $2 \leq k \leq K_1$ such that $\left(\frac{q_1^{E_k^{(1)}}}{q_1^{E_1^{(1)}}} \right)' (\hat{\mathbf{v}}_1) \neq 0$,

and thus we have $\frac{\partial \hat{\mathbf{v}}_1}{\partial \mathbf{z}_i} = 0, \forall \hat{\mathbf{z}} \in \tau^{-1}(\mathbf{O}_v \setminus \mathbf{N}_v)$. Since τ is a diffeomorphism, we can deduce that $\tau^{-1}(\mathbf{O}_v \setminus \mathbf{N}_v) = \mathbf{O}_z \setminus \tau^{-1}(\mathbf{N}_v)$ and $(\tau^{-1}(\mathbf{N}_v))^0 = \emptyset$ by Lemma 19. As a result, we actually have $\frac{\partial \hat{\mathbf{v}}_1}{\partial \mathbf{z}_i} = 0, \forall \hat{\mathbf{z}} \in \mathbf{O}_z$. Hence in \mathbf{O}_z there exists a continuous differentiable function ϕ_1 such that $\mathbf{v}_1 = \phi_1(\mathbf{z}_1)$, proving (2). Finally, (3) directly follows from (2) since $\text{pa}_G(1) = \emptyset$, concluding the proof for $j = 1$.

Now suppose that the statement holds up to $j - 1$, and we need to prove it for j . Again by Lemma 18 we have for $\forall \hat{\mathbf{z}} \in \mathbf{O}_z$ that

$$\frac{p_j^{E_k^{(j)}}}{p_j^{E_1^{(j)}}} (\hat{\mathbf{z}}_j \mid \hat{\mathbf{z}}_{\text{pa}_G(j)}) = \frac{q_j^{E_k^{(j)}}}{q_j^{E_1^{(j)}}} \left(\hat{\mathbf{v}}_j \mid \hat{\mathbf{v}}_{\text{pa}_{\hat{G}}(j)} \right), \quad \forall 2 \leq k \leq K_j. \quad (25)$$

For all $i \notin \overline{\text{pa}}_{\hat{G}}(j)$, taking partial derivative w.r.t. \mathbf{z}_i gives

$$0 = \sum_{\ell \in \overline{\text{pa}}_{\hat{G}}(j)} \frac{\partial}{\partial \hat{\mathbf{v}}_\ell} \frac{q_j^{E_k^{(j)}}}{q_j^{E_1^{(j)}}} \left(\hat{\mathbf{v}}_j \mid \hat{\mathbf{v}}_{\text{pa}_{\hat{G}}(j)} \right) \cdot \frac{\partial \hat{\mathbf{v}}_\ell}{\partial \mathbf{z}_i}, \quad \forall 2 \leq k \leq K_j,$$

i.e.,

$$\left[\nabla_{\mathbf{v}_{\overline{\text{pa}}_{\hat{G}}(j)}} \frac{q_j^{E_k^{(j)}}}{q_j^{E_1^{(j)}}} \left(\hat{\mathbf{v}}_j \mid \hat{\mathbf{v}}_{\text{pa}_{\hat{G}}(j)} \right) : 2 \leq k \leq K_j \right]^\top \frac{\partial \hat{\mathbf{v}}_{\overline{\text{pa}}_{\hat{G}}(j)}}{\partial \mathbf{z}_i} = 0.$$

Similar to the $j = 1$ case, by assumption (iii), we know that the above coefficient matrix has full row rank for $\forall \hat{\mathbf{v}} \in \mathbf{O}_v \setminus \mathbf{N}_v$, so for $\forall \mathbf{z} \in \tau^{-1}(\mathbf{O}_v \setminus \mathbf{N}_v) = \mathbf{O}_z \setminus \tau^{-1}(\mathbf{N}_v)$, we have $\frac{\partial \hat{\mathbf{v}}_{\overline{\text{pa}}_{\hat{G}}(j)}}{\partial \mathbf{z}_i} = 0$. Since $(\tau^{-1}(\mathbf{N}_v))^0 = \emptyset$ by Lemma 19, for all $\hat{\mathbf{z}} \in \mathbf{N}_z$ we can choose a sequence of points $\hat{\mathbf{z}}^{(i)}, i = 1, 2, \dots$ in \mathbf{O}_z such that $\hat{\mathbf{z}}^{(i)} \rightarrow \hat{\mathbf{z}}$. Since τ is a diffeomorphism, its derivatives are continuous and we can deduce that $\frac{\partial \hat{\mathbf{v}}_{\overline{\text{pa}}_{\hat{G}}(j)}}{\partial \mathbf{z}_i} = \lim_{\ell \rightarrow +\infty} \frac{\partial \hat{\mathbf{v}}_{\overline{\text{pa}}_{\hat{G}}(j)}^{(\ell)}}{\partial \hat{\mathbf{z}}_i^{(\ell)}} = 0$. As a result, $\frac{\partial \hat{\mathbf{v}}_{\overline{\text{pa}}_{\hat{G}}(j)}}{\partial \mathbf{z}_i} = 0$ actually holds for all $\mathbf{z} \in \mathbf{O}_z$. Hence, there exists a continuous differentiable function Υ_j such that $\mathbf{v}_{\overline{\text{pa}}_{\hat{G}}(j)} = \Upsilon_j(\mathbf{z}_{\overline{\text{pa}}_{\hat{G}}(j)})$.

By our assumption, $\text{pa}_G(j) \subseteq [j-1]$. Suppose that $\text{pa}_{\hat{G}}(j) \not\subseteq \{i : i < j\}$, let $\ell \in \text{pa}_{\hat{G}}(j) \setminus \{i : i < j\}$, then by induction hypothesis, $\hat{\mathbf{v}}_\ell = \tau_\ell(\hat{\mathbf{z}})$, $\hat{\mathbf{z}} \in \mathbf{O}_z, t = 1, 2, \dots, j, \ell$ are all functions of $\hat{\mathbf{z}}_1, \dots, \hat{\mathbf{z}}_j$. Since τ is a diffeomorphism and \mathbf{O}_z is the support of the distributions $p_E, E \in \mathfrak{E}$, we can deduce that the support of the latent variables $(\mathbf{v}_t : t = 1, 2, \dots, j, \ell)$ lie on a submanifold with dimension $\leq j$, which is impossible since \mathbf{v} is supported on the open set $\mathbf{O}_v \subseteq \mathbb{R}^d$ by assumption (i).

Hence, we must have $\text{pa}_{\hat{\mathcal{G}}}(j) \subseteq \{i : i < j\}$. Furthermore, if there exists $i \in \text{pa}_{\hat{\mathcal{G}}}(j)$ such that $i \notin \text{pa}_{\mathcal{G}}(j)$, then the induction hypothesis implies that $\frac{\partial v_i}{\partial z_i} \neq 0$, but v_i is a function of $\mathbf{z}_{\overline{\text{pa}}_{\mathcal{G}}(j)}$ as previously derived, which is also a contradiction. Thus we actually have $\text{pa}_{\hat{\mathcal{G}}}(j) \subseteq \text{pa}_{\mathcal{G}}(j)$.

In a completely symmetric manner, we can take the derivatives of (25) w.r.t. $v_i, \forall i \in \overline{\text{pa}}_{\hat{\mathcal{G}}}(j)$ and obtain that $\text{pa}_{\mathcal{G}}(j) \subseteq \text{pa}_{\hat{\mathcal{G}}}(j)$. Hence, $\text{pa}_{\hat{\mathcal{G}}}(j) = \text{pa}_{\mathcal{G}}(j)$, completing the proof of (1) and (3) for the j case.

Finally, if $\frac{\partial v_j}{\partial z_j} \equiv 0$, then by (3) and the induction hypothesis, v_1, \dots, v_j are all functions of $\mathbf{z}_{[j-1]}$, which implies that (v_1, \dots, v_j) lies on a submanifold with dimension $\leq j-1$, again contradicting assumption (i). Thus $\frac{\partial v_j}{\partial z_j} \neq 0$. This completes the proof of our inductive step.

To recap, we now know that

- $\mathcal{G} = \hat{\mathcal{G}}$, and
- For $\forall i \in [d]$, there exists a function Υ_i such that $\mathbf{v}_{\overline{\text{pa}}_{\mathcal{G}}(i)} = \Upsilon_i(\mathbf{z}_{\overline{\text{pa}}_{\mathcal{G}}(i)})$.

It remains to show that for $\forall k \in \text{pa}_{\mathcal{G}}(i) \setminus \text{sur}_{\mathcal{G}}(i)$, Υ_i doesn't depend on z_k .

By definition, if $k \in \text{pa}_{\mathcal{G}}(i) \setminus \text{sur}_{\mathcal{G}}(i)$, we know that there exists $j \in \text{ch}_{\mathcal{G}}(i)$ such that $j \notin \text{ch}_{\mathcal{G}}(k)$. We have shown that v_i , as a component of $\mathbf{v}_{\overline{\text{pa}}_{\mathcal{G}}(j)}$, is a function of $\mathbf{z}_{\overline{\text{pa}}_{\mathcal{G}}(j)}$. By the choice of k , we have $k \notin \overline{\text{pa}}_{\mathcal{G}}(j)$, so that v_i does not depend on z_k . The conclusion follows.

J Omitted Proofs for Theorem 3 and Theorem 8

In this section we provide detailed proofs of main ambiguity results.

Definition 12. We say that a matrix $\mathbf{M} \in \mathbb{R}^{d \times d}$ is effect-respecting for a causal graph \mathcal{G} , or $\mathbf{M} \in \mathcal{M}_{\text{sur}}(\mathcal{G})$, if $M_{ij} \neq 0 \Leftrightarrow j \in \overline{\text{sur}}_{\mathcal{G}}(i)$. We also write $\mathbf{M} \in \mathcal{M}_{\text{sur}}^0(\mathcal{G})$ if \mathbf{M} is invertible and $M_{ij} \neq 0 \Rightarrow j \in \overline{\text{sur}}_{\mathcal{G}}(i)$. Finally, we write $\mathbf{M} \in \overline{\mathcal{M}}_{\text{sur}}(\mathcal{G})$ if $M_{ij} \neq 0 \Rightarrow j \in \overline{\text{sur}}_{\mathcal{G}}(i)$.

Remark 1. By definition $\mathcal{M}_{\text{sur}}^0(\mathcal{G})$ is the set of all matrices \mathbf{M} where $M_{ij} \neq 0, \forall j \notin \overline{\text{sur}}_{\mathcal{G}}(i)$, so it can be identified as $\mathbb{R}^{d+d_{\mathcal{G}}}$ where $d_{\mathcal{G}} = \sum_{i=1}^d |\text{sur}_{\mathcal{G}}(i)|$. Equipped with the Lebesgue measure, we have $\mathcal{M}_{\text{sur}}(\mathcal{G}) \subset \mathcal{M}_{\text{sur}}^0(\mathcal{G}) \subset \overline{\mathcal{M}}_{\text{sur}}(\mathcal{G})$ and $\overline{\mathcal{M}}_{\text{sur}}(\mathcal{G}) \setminus \mathcal{M}_{\text{sur}}(\mathcal{G})$ is a null set. In the remaining part of this section, we will use measure-theoretic statement for $\mathbf{M} \in \mathcal{M}_{\text{sur}}(\mathcal{G})$ in the above sense.

We first present a result that serves as a good starting point to understand why this is the case. It states that latent representations that are equivalent under \sim_{sur} are essentially generated from the same causal graph.

Proposition 5. Let \mathbf{M} be an invertible matrix such that $M_{ij} \neq 0 \Rightarrow j \in \overline{\text{sur}}_{\mathcal{G}}(i)$. Suppose that the latent variables $\mathbf{z} \in \mathbb{R}^d$ are generated from any distributions $p_i(\mathbf{z}_i | \mathbf{z}_{\text{pa}_{\mathcal{G}}(i)}), i \in [d]$ with joint density $p(\mathbf{z}) = \prod_{i=1}^d p_i(\mathbf{z}_i | \mathbf{z}_{\text{pa}_{\mathcal{G}}(i)})$, then the joint density of $\mathbf{v} = \mathbf{M}\mathbf{z}$ can be written as $q(\mathbf{v}) = \prod_{i=1}^d q_i(\mathbf{v}_i | \mathbf{v}_{\text{pa}_{\mathcal{G}}(i)})$ for some density functions $q_i, i \in [d]$.

J.1 Proof of Proposition 5

We first prove the following lemma:

Lemma 20. Let $\mathbf{M} \in \mathcal{M}_{\text{sur}}^0(\mathcal{G})$ and latent variables $\mathbf{v} = \mathbf{M}\mathbf{z}$, then for $\forall i \in [d]$, there exists invertible matrices \mathbf{M}_i and \mathbf{M}_i^- such that $\mathbf{v}_{\text{pa}_{\mathcal{G}}(i)} = \mathbf{M}_i^- \mathbf{z}_{\text{pa}_{\mathcal{G}}(i)}$ and $\mathbf{v}_{\overline{\text{pa}}_{\mathcal{G}}(i)} = \mathbf{M}_i \mathbf{z}_{\overline{\text{pa}}_{\mathcal{G}}(i)}$.

Proof. $\forall j \in \overline{\text{pa}}_{\mathcal{G}}(i)$, we know that v_j is a linear function of $\mathbf{z}_{\ell}, \ell \in \overline{\text{sur}}_{\mathcal{G}}(j)$. By Lemma 7, we know that $\overline{\text{sur}}_{\mathcal{G}}(j) \subseteq \overline{\text{pa}}_{\mathcal{G}}(i)$, so each $v_j, j \in \overline{\text{pa}}_{\mathcal{G}}(i)$ is a linear function of $\mathbf{z}_{\overline{\text{pa}}_{\mathcal{G}}(i)}$. Thus we can write $\mathbf{v}_{\overline{\text{pa}}_{\mathcal{G}}(i)} = \mathbf{M}_i \mathbf{z}_{\overline{\text{pa}}_{\mathcal{G}}(i)}$. In the following we argue that \mathbf{M}_i is invertible. Let π be a permutation on $\overline{\text{pa}}_{\mathcal{G}}(i)$ such that $k \in \text{pa}_{\mathcal{G}}(\ell) \Rightarrow \pi(k) < \pi(\ell)$ (such π can always be chosen since \mathcal{G} is acyclic), then we can write

$$(\hat{\mathbf{v}}_{\pi(j)} : j \in \overline{\text{pa}}_{\mathcal{G}}(i))^{\top} = \tilde{\mathbf{M}}_i (\hat{\mathbf{z}}_{\pi(j)} : j \in \overline{\text{pa}}_{\mathcal{G}}(i))^{\top} \quad (26)$$

where \tilde{M}_i is an upper triangular matrix with non-zero diagonal entries by our choice of M . Since M_i can be obtained from \tilde{M}_i by exchanging a few rows and columns, M_i is invertible as well.

Similarly, using the fact that $\forall j \in \text{pa}_{\mathcal{G}}(i), \overline{\text{sur}}_{\mathcal{G}}(j) \subseteq \text{pa}_{\mathcal{G}}(i)$, we can prove the existence of an invertible matrix M_i^- such that $v_{\text{pa}_{\mathcal{G}}(i)} = M_i^- z_{\text{pa}_{\mathcal{G}}(i)}$. \square

Returning to the proof of **Proposition 5**. Assume WLOG that the nodes of \mathcal{G} are ordered in a way such that $i \in \text{pa}_{\mathcal{G}}(j) \Rightarrow i < j$, so that M is a lower-triangular matrix. The joint density of v can be written as

$$q(v) = \prod_{i=1}^d q(v_i | v_1, \dots, v_{i-1}).$$

Since $v = Mz$ and M is lower triangular and invertible (hence, with non-zero diagonals), we know that $(v_1, v_2, \dots, v_{i-1})$ is an invertible linear function of $(z_1, z_2, \dots, z_{i-1})$ and (v_1, v_2, \dots, v_i) is an invertible linear function of (z_1, z_2, \dots, z_i) . Let $\hat{v} = M\hat{z} \in \mathbb{R}^d$, then we have

$$\begin{aligned} q(\hat{v}_i | \hat{v}_1, \dots, \hat{v}_{i-1}) &= \frac{q(\hat{v}_1, \hat{v}_2, \dots, \hat{v}_i)}{q(\hat{v}_1, \hat{v}_2, \dots, \hat{v}_{i-1})} = \frac{p(\hat{z}_1, \hat{z}_2, \dots, \hat{z}_i) \det \hat{M}_{1:i, 1:i}}{p(\hat{z}_1, \hat{z}_2, \dots, \hat{z}_{i-1}) \det \hat{M}_{1:i-1, 1:i-1}} \\ &\propto \frac{p(\hat{z}_1, \hat{z}_2, \dots, \hat{z}_i)}{p(\hat{z}_1, \hat{z}_2, \dots, \hat{z}_{i-1})} = p(\hat{z}_i | \hat{z}_1, \dots, \hat{z}_{i-1}) = p_i(\hat{z}_i | \hat{z}_{\text{pa}_{\mathcal{G}}(i)}), \end{aligned}$$

where $\hat{M}_{1:i, 1:i}$ denotes that top-left submatrix of \hat{M} of size $i \times i$, and the last step follows from the causal Markov condition (**Definition 1**). On the other hand, let $q_i(\hat{v}_i | \hat{v}_{\text{pa}_{\mathcal{G}}(i)})$ be the conditional density of v_i on its parents at $\hat{v} \in \mathbb{R}^d$. For $\forall j \in \text{pa}_{\mathcal{G}}(i)$, from $v = Mz$ we know that v_j is a linear function of $z_{\overline{\text{sur}}_{\mathcal{G}}(j)}$. By **Lemma 20** we know that $\hat{v}_{\text{pa}_{\mathcal{G}}(i)}$ is a linear function of $\hat{z}_{\text{pa}_{\mathcal{G}}(i)}$ and $\hat{v}_{\overline{\text{pa}}_{\mathcal{G}}(i)}$ is a linear function of $\hat{z}_{\overline{\text{pa}}_{\mathcal{G}}(i)}$, so that

$$q(\hat{v}_{\text{pa}_{\mathcal{G}}(i)}) \propto p(\hat{z}_{\text{pa}_{\mathcal{G}}(i)}) \quad \text{and} \quad q(\hat{v}_{\overline{\text{pa}}_{\mathcal{G}}(i)}) \propto p(\hat{z}_{\overline{\text{pa}}_{\mathcal{G}}(i)})$$

and

$$q_i(\hat{v}_i | \hat{v}_{\text{pa}_{\mathcal{G}}(i)}) \propto \frac{p(\hat{z}_{\overline{\text{pa}}_{\mathcal{G}}(i)})}{p(\hat{z}_{\text{pa}_{\mathcal{G}}(i)})} = p_i(\hat{z}_i | \hat{z}_{\text{pa}_{\mathcal{G}}(i)}).$$

Hence, we have $q_i(\hat{v}_i | \hat{v}_{\text{pa}_{\mathcal{G}}(i)}) \propto q(\hat{v}_i | \hat{v}_1, \dots, \hat{v}_{i-1})$, so that

$$q(\hat{v}) = \prod_{i=1}^d q_i(\hat{v}_i | \hat{v}_{\text{pa}_{\mathcal{G}}(i)}) \propto \prod_{i=1}^d q(\hat{v}_i | \hat{v}_{\text{pa}_{\mathcal{G}}(i)}).$$

Since both sides integrate to 1, it turns out that they are equal, as desired.

J.2 Formal version and proof of **Theorem 3**: the linear case

Theorem 9 (Counterpart to **Theorem 1**). *For any causal model (H, \mathcal{G}) and any set of environments $\mathfrak{E} = \{E_k : k \in [K]\}$, suppose that we have observations $\{P_{\mathbf{X}}^E\}_{E \in \mathfrak{E}}$ satisfying **Assumption 1**:*

$$\forall k \in [K], \quad \mathbf{z} = \mathbf{A}_k \mathbf{z} + \mathbf{\Omega}_k^{\frac{1}{2}} \epsilon, \quad \mathbf{x} = \mathbf{H}^\dagger \mathbf{z}$$

such that

- (i) the unmixing matrix $\mathbf{H} \in \mathbb{R}^{d \times n}$ has full row rank;
- (ii) $\forall k \in [K]$ and $i, j \in [d]$, $(\mathbf{A}_k)_{ij} \neq 0 \Leftrightarrow j \in \text{pa}_{\mathcal{G}}(i)$ and $\mathbf{\Omega}_k$ is a diagonal matrix with positive entries;
- (iii) $\left\{ \mathbf{B}_k = \mathbf{\Omega}_k^{-\frac{1}{2}} (\mathbf{I} - \mathbf{A}_k) \right\}_{k=1}^K$ are node level non-degenerate in the sense of **Assumption 5**,

then there must exist a candidate solution (\hat{H}, \mathcal{G}) and a hypothetical data generating process

$$\forall k \in [K], \quad \mathbf{v} = \hat{\mathbf{A}}_k \mathbf{v} + \hat{\mathbf{\Omega}}_k^{\frac{1}{2}} \epsilon, \quad \mathbf{x} = \hat{\mathbf{H}}^\dagger \mathbf{v}$$

such that

- (i') the unmixing matrix $\hat{\mathbf{H}} \in \mathbb{R}^{d \times n}$ has full row rank;
- (ii') $\forall k \in [K]$ and $i, j \in [d]$, $(\hat{\mathbf{A}}_k)_{ij} \neq 0 \Leftrightarrow j \in \text{pa}_{\mathcal{G}}(i)$ and $\hat{\mathbf{\Omega}}_k$ is a diagonal matrix with positive entries;
- (iii') $\left\{ \hat{\mathbf{B}}_k = \hat{\mathbf{\Omega}}_k^{-\frac{1}{2}} (\mathbf{I} - \hat{\mathbf{A}}_k) \right\}_{k=1}^K$ are node level non-degenerate in the sense of [Assumption 5](#),

but

$$\frac{\partial v_i}{\partial z_j} \neq 0, \quad \forall j \in \overline{\text{sur}}_{\mathcal{G}}(i).$$

Finally, if we additionally assume that

- (iii) the environments are groups of single-node interventions: there exists a partition $\mathfrak{E} = \bigcup_{i=1}^d \mathfrak{E}_i$ such that $\mathcal{I}_{\mathbf{z}}^{\mathfrak{E}_i} = \{i\}$ (see [Definition 2](#)),

then we can guarantee the existence of $(\hat{\mathbf{H}}, \mathcal{G})$ and weight matrices which, besides the properties listed above, also satisfy

- (iii') for the same partition $\mathfrak{E} = \bigcup_{i=1}^d \mathfrak{E}_i$, we have $\mathcal{I}_{\mathbf{v}}^{\mathfrak{E}_i} = \{i\}$.

In other words, additionally assuming that the environments are from single-node interventions does not resolve the ambiguity.

Remark 2. Compared with our identifiability guarantee [Theorem 1](#), [Theorem 9](#) actually demonstrates a stronger form of impossibility. Specifically, it states that the SNA cannot be resolved even if both the ground-truth causal graph and the noise variables are known.

We define

$$\mathbf{v} = \mathbf{M}\mathbf{z} \quad (27)$$

where \mathbf{M} is an effect-respecting matrix. At this point we do not make any other restrictions on \mathbf{M} , but we will specify the appropriate choice of \mathbf{M} later.

By assumption, the latent variables in the k -th environment are generated by

$$\mathbf{z} = \mathbf{A}_k \mathbf{z} + \mathbf{\Omega}_k^{\frac{1}{2}} \epsilon,$$

then $\mathbf{v} = \mathbf{M}(\mathbf{I} - \mathbf{A}_k)^{-1} \mathbf{\Omega}_k^{\frac{1}{2}} \epsilon$. Let $\hat{\mathbf{\Omega}}_k$ be the diagonal matrix with entries $M_{ii}^2 \cdot (\mathbf{\Omega}_k)_{ii}$, $i \in [d]$ and $\hat{\mathbf{A}}_k = \mathbf{I} - \hat{\mathbf{\Omega}}_k^{\frac{1}{2}} \mathbf{\Omega}_k^{-\frac{1}{2}} (\mathbf{I} - \mathbf{A}_k) \mathbf{M}^{-1}$, then $\mathbf{v} = \hat{\mathbf{A}}_k \mathbf{v} + \hat{\mathbf{\Omega}}_k^{\frac{1}{2}} \epsilon$. Note that the choice of $\hat{\mathbf{\Omega}}_k$ here is to that the diagonal entries of $\hat{\mathbf{A}}_k$ are zero, as we show below. It remains to show that: for almost all $\mathbf{M} \in \mathcal{M}_{\text{sur}}^0(\mathcal{G})$, it holds for $\forall k \in [K]$ that $(\hat{\mathbf{A}}_k)_{ij} = 0 \Leftrightarrow j \notin \text{pa}_{\mathcal{G}}(i)$.

For the \Leftarrow direction, since $\mathbf{M} \in \mathcal{M}_{\text{sur}}^0(\mathcal{G})$, $\mathbf{M}^{-1} \in \mathcal{M}_{\text{sur}}^0(\mathcal{G})$ as well. Thus, $\forall j \notin \text{pa}_{\mathcal{G}}(i)$ we have

$$\begin{aligned} [(I - A_k)M^{-1}]_{ij} &= \sum_{\ell=1}^d (I - A_k)_{i\ell} \cdot (M^{-1})_{\ell j} = \sum_{\ell \in \overline{\text{pa}}_{\mathcal{G}}(i) \cap \{\ell': j \in \overline{\text{sur}}_{\mathcal{G}}(\ell')\}} (I - A_k)_{i\ell} \cdot (M^{-1})_{\ell j} \\ &= \begin{cases} 0 & \text{if } j \notin \overline{\text{pa}}_{\mathcal{G}}(i) \\ (M^{-1})_{ii} & \text{if } j = i \end{cases} \end{aligned}$$

where the last step holds because $\forall \ell \in [d]$, $\ell \in \overline{\text{pa}}_{\mathcal{G}}(i)$, $j \in \overline{\text{sur}}_{\mathcal{G}}(\ell) \Rightarrow j \in \overline{\text{pa}}_i$, and when $j = i$, the only such ℓ is $\ell = i$. Hence, we can see that our choice of $\hat{\mathbf{A}}_k$ satisfies

$$(\hat{\mathbf{A}}_k)_{ij} = \begin{cases} 0 - 0 = 0 & \text{if } j \notin \overline{\text{pa}}_{\mathcal{G}}(i) \\ 1 - \hat{\omega}_{k,i,i}^{\frac{1}{2}} \omega_{k,i,i}^{-\frac{1}{2}} (M^{-1})_{ii} = 0 & \text{if } j = i, \end{cases}$$

so $(\hat{\mathbf{A}}_k)_{ij} \neq 0 \Rightarrow j \in \text{pa}_{\mathcal{G}}(i)$.

Conversely, for $\forall j \in \text{pa}_{\mathcal{G}}(i)$,

$$(\hat{\mathbf{A}}_k)_{ij} = 0 \Leftrightarrow \sum_{s \in \overline{\text{pa}}_{\mathcal{G}}(i)} (I - A_k)_{is} (M^{-1})_{sj} = 0 \Leftrightarrow \sum_{s \in \overline{\text{pa}}_{\mathcal{G}}(i)} (-1)^s (I - A_k)_{is} \det M_{sj}^{-} = 0 \quad (28)$$

where M_{sj}^- is the $(d-1) \times (d-1)$ matrix obtained by removing the s -th row and j -th column of M , and the second step in the equation above follows from the fact that $M^{-1} = \det(M)^{-1} \text{adj}(M)$, where $\text{adj}(M)$ denotes the adjugate matrix of M whose (i, j) -th entry is $(-1)^{i+j} \det M_{ij}^-$.

(28) holds if only if M takes values on a lower-dimensional algebraic manifold of its embedded space $\mathbb{R}^{d+d_{\mathcal{G}}}$ (see Remark 1). As a result, for almost every $M \in \mathcal{M}_{\text{sur}}^0(\mathcal{G})$, v is generated from a linear causal model with graph \mathcal{G} as defined in (3). Moreover, let $\hat{B}_k = B_k M^{-1}$, $k \in [K]$, so that $\epsilon = \hat{B}_k v$ in the k -th environment. Then for all nodes $i \in [d]$ and $S \subseteq \text{pa}(i) \cup \{i\}$, we have

$$\begin{aligned} \dim \text{span} \left\langle \left(\hat{B}_k^\top e_i \right)_S : k \in [K] \right\rangle &= \dim \text{span} \left\langle M^{-\top} \left((B_k^\top e_i)_S : k \in [K] \right) \right\rangle \\ &= \dim \text{span} \left\langle (B_k^\top e_i)_S : k \in [K] \right\rangle = |\text{pa}_{\mathcal{G}}(i)| + 1, \end{aligned}$$

implying that \hat{B}_k , $k \in [K]$ satisfy Assumption 5.

Now we have shown that for almost every $M \in \mathcal{M}_{\text{sur}}^0(\mathcal{G})$, we can construct a hypothetical data generating process with latent variables $v = Mz$ that satisfies all requirements in Theorem 9. Choose an arbitrary M that is in $\mathcal{M}_{\text{sur}}(\mathcal{G})$, then we have that

$$\frac{\partial v_i}{\partial z_j} \neq 0, \quad j \notin \overline{\text{sur}}_{\mathcal{G}}(i).$$

Finally, if we additionally assume single-node interventions, $\forall k, \ell \in \mathfrak{E}_i$, we have that $(B_k)_j \neq (B_\ell)_j \Leftrightarrow j = i$. For any $M \in \mathcal{M}_{\text{sur}}^0(\mathcal{G})$ (and specifically the M that we have already chosen above), we have $(\hat{B}_k)_j = (B_k)_j M^{-1}$ and $(\hat{B}_\ell)_j = (B_\ell)_j M^{-1}$, $\forall j \in [d]$. Thus, $(\hat{B}_k)_j \neq (\hat{B}_\ell)_j \Leftrightarrow j = i$ as well, implying that \mathfrak{E}_i is also a group of single-node interventions on v , concluding the proof.

J.3 Formal statement and proof of Theorem 10: the non-parametric case

Theorem 10 (Counterpart to Theorem 7). *For any causal model (h, \mathcal{G}) and any set of environments \mathfrak{E} , suppose that we have observations $\{P_X^E\}_{E \in \mathfrak{E}}$ satisfying Assumption 1:*

$$\forall E \in \mathfrak{E}, z \sim p_E(\hat{z}) = \prod_{i=1}^d p_i^E(\hat{z}_i | \hat{z}_{\text{pa}_{\mathcal{G}}(i)}), x = h^{-1}(z)$$

such that

- (i) all densities p_i^E are continuously differentiable and the joint density p_E is positive everywhere;
- (ii) the environments are groups of single-node interventions: there exists a partition $\mathfrak{E} = \cup_{i=1}^d \mathfrak{E}_i$ such that $\mathcal{I}_z^{\mathfrak{E}_i} = \{i\}$;
- (iii) the intervention distributions on each node are non-degenerate: $\forall i \in [d]$, the set of distributions $\{p_i^E : E \in \mathfrak{E}_i\}$ satisfy Definition 11 at any point $\hat{z} \in \mathbb{R}^d$,

then there must exist a candidate solution (\hat{h}, \mathcal{G}) and a hypothetical data generating process

$$\forall E \in \mathfrak{E}, v \sim q_E(\hat{v}) = \prod_{i=1}^d q_i^E(\hat{v}_i | \hat{v}_{\text{pa}_{\mathcal{G}}(i)}), x = \hat{h}^{-1}(v)$$

such that

- (i') all densities q_i^E are continuously differentiable and the joint density q_E is positive everywhere;
- (ii') for the same partition $\mathfrak{E} = \cup_{i=1}^d \mathfrak{E}_i$, we have $\mathcal{I}_v^{\mathfrak{E}_i} = \{i\}$;
- (iii') $\forall i \in [d]$, the set of distributions $\{q_i^E : E \in \mathfrak{E}_i\}$ satisfy Definition 11 at any point $\hat{v} \in \mathbb{R}^d$,

but

$$\frac{\partial \mathbf{v}_i}{\partial \mathbf{z}_j} \neq 0, \quad \forall j \in \overline{\text{sur}}_{\mathcal{G}}(i).$$

Remark 3. Similar to the case of [Theorem 9](#), [Appendix J.3](#) also establishes a stronger form of identifiability. First, it is assumed that the causal graph \mathcal{G} is known. Second, we only focus on a special case of the setting of [Theorem 7](#) by assuming that the support is the whole space, and the non-degeneracy condition [Definition 11](#) holds at any point. Even in this case, we show that our identification guarantee up to SNA cannot be improved.

We state and prove a stronger version of [Theorem 10](#):

Theorem 11. For any causal model $(\mathbf{h}, \mathcal{G})$ and any set of environments \mathfrak{E} , suppose that we have observations $\{P_{\mathbf{X}}^E\}_{E \in \mathfrak{E}}$ satisfying [Assumption 1](#):

$$\forall E \in \mathfrak{E}, \quad \mathbf{z} \sim p_E(\mathbf{z}) = \prod_{i=1}^d p_i^E(z_i | z_{\text{pa}_{\mathcal{G}}(i)}), \quad \mathbf{x} = \mathbf{h}^{-1}(\mathbf{z})$$

such that

- (i) all densities p_i^E are continuously differentiable and the joint density p_E is positive everywhere;
- (ii) the environments are groups of single-node interventions: there exists a partition $\mathfrak{E} = \bigcup_{i=1}^d \mathfrak{E}_i$ such that $\mathcal{I}_{\mathbf{z}}^{\mathfrak{E}_i} = \{i\}$;
- (iii) the intervention distributions on each node are non-degenerate: $\forall i \in [d]$, the set of distributions $\{p_i^E : E \in \mathfrak{E}_i\}$ satisfy [Definition 11](#),

then there must exist a candidate solution $(\hat{\mathbf{h}}, \mathcal{G})$ and a hypothetical data generating process

$$\forall E \in \mathfrak{E}, \quad \mathbf{v} \sim q_E(\mathbf{v}) = \prod_{i=1}^d q_i^E(v_i | v_{\text{pa}_{\mathcal{G}}(i)}), \quad \mathbf{x} = \hat{\mathbf{h}}^{-1}(\mathbf{v})$$

such that

- (i') all densities q_i^E are continuously differentiable and the joint density q_E is positive everywhere;
- (ii') for the same partition $\mathfrak{E} = \bigcup_{i=1}^d \mathfrak{E}_i$, we have $\mathcal{I}_{\mathbf{v}}^{\mathfrak{E}_i} = \{i\}$;
- (iii') $\forall i \in [d]$, the set of distributions $\{q_i^E : E \in \mathfrak{E}_i\}$ satisfy [Definition 11](#),

but

$$\frac{\partial \mathbf{v}_i}{\partial \mathbf{z}_j} \neq 0, \quad \forall j \in \overline{\text{sur}}_{\mathcal{G}}(i).$$

Finally, if we additionally assume minimality ([Assumption 6](#)) and/or faithfulness ([Assumption 7](#)) of all p_E 's, we can guarantee the existence of $(\hat{\mathbf{h}}, \mathcal{G})$ and q_E 's satisfying minimality and/or faithfulness in addition to the properties listed above. In other words, assuming minimality and/or faithfulness does not resolve the ambiguity.

Proof. We define

$$\mathbf{v} = \mathbf{M}\mathbf{z} \tag{29}$$

where \mathbf{M} is an effect-respecting matrix. At this point we do not make any other restrictions on \mathbf{M} , and we will choose appropriate \mathbf{M} later. By [Lemma 20](#), there exists invertible matrices \mathbf{M}_i and \mathbf{M}_i^- such that $\mathbf{v}_{\text{pa}_{\mathcal{G}}(i)} = \mathbf{M}_i^- \mathbf{z}_{\text{pa}_{\mathcal{G}}(i)}$ and $\mathbf{v}_{\overline{\text{pa}}_{\mathcal{G}}(i)} = \mathbf{M}_i \mathbf{z}_{\overline{\text{pa}}_{\mathcal{G}}(i)}$, so for all environment $E \in \mathfrak{E}$ we have

$$q_i^E(\mathbf{v}_{\text{pa}_{\mathcal{G}}(i)}) = p_i^E(\mathbf{z}_{\text{pa}_{\mathcal{G}}(i)}) \cdot |\det(\mathbf{M}_i^-)^{-1}|, \quad q_i^E(\mathbf{v}_{\overline{\text{pa}}_{\mathcal{G}}(i)}) = p_i^E(\mathbf{z}_{\overline{\text{pa}}_{\mathcal{G}}(i)}) \cdot |\det(\mathbf{M}_i)^{-1}|$$

so that

$$q_i^E(\mathbf{v}_i | \mathbf{v}_{\text{pa}_{\mathcal{G}}(i)}) = p_i^E(\mathbf{z}_i | \mathbf{z}_{\text{pa}_{\mathcal{G}}(i)}) \frac{|\det \mathbf{M}_i^{-1}|}{|\det(\mathbf{M}_i^-)^{-1}|}, \quad \forall i \in [d]. \tag{30}$$

In the following, assuming that $(p_i^E : E \in \mathfrak{E})$ satisfies any of the listed assumptions, we show that $(q_i^E : E \in \mathfrak{E})$ satisfies the same assumption as well.

Firstly, (30) immediately implies that the density of \mathbf{v} is continuous differentiable and positive everywhere. Secondly, $\forall k, \ell \in \mathfrak{E}_i$, we have that

$$p_j^{E_k}(\mathbf{z}_j | \mathbf{z}_{\text{pa}_{\mathcal{G}}(j)}) = p_j^{E_\ell}(\mathbf{z}_j | \mathbf{z}_{\text{pa}_{\mathcal{G}}(j)}) \Leftrightarrow j = i.$$

By (30) it is easy to see that

$$q_j^{E_k}(\mathbf{v}_j | \mathbf{v}_{\text{pa}_{\mathcal{G}}(j)}) = q_j^{E_\ell}(\mathbf{v}_j | \mathbf{v}_{\text{pa}_{\mathcal{G}}(j)}) \Leftrightarrow j = i$$

as well, i.e., $q^k, k \in \mathfrak{E}_i$ are single-node interventions on \mathbf{v}_i according to Definition 2.

Thirdly, we verify the non-degeneracy condition for $q_i^{E_1}$'s. Indeed we have for $\forall k \geq 2$ that

$$\nabla_{\mathbf{v}_{\overline{\text{pa}}_{\mathcal{G}}(i)}} \frac{q_i^{E_1}}{q_i^{E_k}}(\mathbf{v}_i | \mathbf{v}_{\text{pa}_{\mathcal{G}}(i)}) = \frac{\partial \mathbf{z}_{\overline{\text{pa}}_{\mathcal{G}}(i)}}{\partial \mathbf{v}_{\overline{\text{pa}}_{\mathcal{G}}(i)}} \nabla_{\mathbf{z}_{\overline{\text{pa}}_{\mathcal{G}}(i)}} \frac{q_i^{E_1}}{q_i^{E_k}}(\mathbf{z}_i | \mathbf{z}_{\text{pa}_{\mathcal{G}}(i)}) = \mathbf{M}_i^{-1} \nabla_{\mathbf{z}_{\overline{\text{pa}}_{\mathcal{G}}(i)}} \frac{q_i^{E_1}}{q_i^{E_k}}(\mathbf{z}_i | \mathbf{z}_{\text{pa}_{\mathcal{G}}(i)}).$$

Since \mathbf{M}_i is invertible, the above equation and the non-degeneracy of $p^{E_k}, k \in [K]$ immediately implies that non-degeneracy of $q^{E_k}, k \in [K]$.

Thus, for arbitrary $\mathbf{M} \in \mathcal{M}_{\text{sur}}(\mathcal{G})$, we have constructed a hypothetical data generating process with latent variable $\mathbf{v} = \mathbf{M}\mathbf{z}$ that satisfies all given conditions. It remains to show that such construction is still possible under additional minimality and faithfulness conditions.

Claim 1. There exists a neighbourhood O of the identity matrix \mathbf{I} in $\overline{\mathcal{M}}_{\text{sur}}(\mathcal{G})$ (in the sense of Remark 1) such that for $\forall \mathbf{M} \in O \cap \mathcal{M}_{\text{sur}}^0(\mathcal{G}), p^{E_k}, k \in [K]$ satisfy Assumption 7 $\Rightarrow q^{E_k}, k \in [K]$ satisfy Assumption 7.

For $\forall i, j$ not d -separated by $S \subseteq [d]$, for all $k \in [K]$ there exists $\hat{\mathbf{z}} \in \mathbb{R}^d$ such that $\Delta_k^{(i,j,S)} = p^{E_k}(\hat{\mathbf{z}}_i, \hat{\mathbf{z}}_j | \hat{\mathbf{z}}_S) - p^{E_k}(\hat{\mathbf{z}}_i | \hat{\mathbf{z}}_S) p^{E_k}(\hat{\mathbf{z}}_j | \hat{\mathbf{z}}_S) \neq 0$. By continuous differentiability of p^{E_k} , we know that there exists $\delta_k^{(i,j,S)} > 0$ such that for all $\mathbf{M} \in \overline{\mathcal{M}}_{\text{sur}}(\mathcal{G})$ such that $\|\mathbf{M} - \mathbf{I}\|_F \leq \delta_k^{(i,j,S)}$, the density of the variable $\mathbf{v} = \mathbf{M}\mathbf{z}$ satisfies $q^{E_k}(\hat{\mathbf{v}}_i, \hat{\mathbf{v}}_j | \hat{\mathbf{v}}_S) \neq q^{E_k}(\hat{\mathbf{v}}_i | \hat{\mathbf{v}}_S) q^{E_k}(\hat{\mathbf{v}}_j | \hat{\mathbf{v}}_S)$ for $\hat{\mathbf{v}} = \mathbf{M}\hat{\mathbf{z}}$, which implies that \mathbf{v}_i and \mathbf{v}_j are dependent given \mathbf{v}_S . Now choose $\delta = \min_{k,i,j,S} \delta_k^{(i,j,S)} > 0$, then for all $\mathbf{M} \in \overline{\mathcal{M}}_{\text{sur}}(\mathcal{G})$ such that $\|\mathbf{M} - \mathbf{I}\|_F \leq \delta$, the resulting distributions $q^{E_k}, k \in [K]$ satisfy assumption Assumption 6.

Claim 2. There exists a neighbourhood O of \mathbf{I} in $\overline{\mathcal{M}}_{\text{sur}}(\mathcal{G})$ (in the sense of Remark 1) such that for almost all $\mathbf{M} \in O \cap \mathcal{M}_{\text{sur}}^0(\mathcal{G}), p^{E_k}, k \in [K]$ satisfies Assumption 6 $\Rightarrow p^{E_k}, k \in [K]$ satisfies Assumption 6.

The proof is similar to the previous statement. Since Assumption 6 causal minimality is satisfied for \mathbf{z} , for $\forall k \in [K], i \in [d]$, let \mathcal{G}_{ij} be the resulting graph obtained by removing the edge $j \rightarrow i$ from \mathcal{G} , then there must exist some $\alpha_{ijk} \in [d]$ such that $\mathbf{z}_{\alpha_{ijk}} \not\perp \mathbf{z}_{\text{nd}_{\mathcal{G}_{ij}}(\alpha_{ijk})} | \mathbf{z}_{\text{pa}_{\mathcal{G}_{ij}}(\alpha_{ijk})}$. Hence, there exists $\hat{\mathbf{z}}^{ijk} \in \mathbb{R}^d$ such that

$$p^{E_k}(\hat{\mathbf{z}}_{\alpha_{ijk}}^{ijk} | \hat{\mathbf{z}}_{\text{pa}_{\mathcal{G}_{ij}}(\alpha_{ijk})}^{ijk}) p^{E_k}(\hat{\mathbf{z}}_{\text{nd}_{\mathcal{G}_{ij}}(\alpha_{ijk})}^{ijk} | \hat{\mathbf{z}}_{\text{pa}_{\mathcal{G}_{ij}}(\alpha_{ijk})}^{ijk}) \neq p^{E_k}(\hat{\mathbf{z}}_{\text{nd}_{\mathcal{G}_{ij}}(\alpha_{ijk})}^{ijk} | \hat{\mathbf{z}}_{\text{pa}_{\mathcal{G}_{ij}}(\alpha_{ijk})}^{ijk}).$$

By continuous differentiability of p^{E_k} , there exists $\delta_k^{(i,j)} > 0$ such that for all $\mathbf{M} \in \overline{\mathcal{M}}_{\text{sur}}(\mathcal{G})$ such that $\|\mathbf{M} - \mathbf{I}\|_F \leq \delta_k^{(i,j)}$, the density $q_{ij}^{E_k}$ of the variable $\hat{\mathbf{v}}^{ijk} = \mathbf{M}\hat{\mathbf{z}}^{ijk}$ satisfies

$$q^{E_k}(\hat{\mathbf{v}}_{\alpha_{ijk}}^{ijk} | \hat{\mathbf{v}}_{\text{pa}_{\mathcal{G}_{ij}}(\alpha_{ijk})}^{ijk}) q^{E_k}(\hat{\mathbf{v}}_{\text{nd}_{\mathcal{G}_{ij}}(\alpha_{ijk})}^{ijk} | \hat{\mathbf{v}}_{\text{pa}_{\mathcal{G}_{ij}}(\alpha_{ijk})}^{ijk}) \neq q^{E_k}(\hat{\mathbf{v}}_{\text{nd}_{\mathcal{G}_{ij}}(\alpha_{ijk})}^{ijk} | \hat{\mathbf{v}}_{\text{pa}_{\mathcal{G}_{ij}}(\alpha_{ijk})}^{ijk}).$$

for $\hat{\mathbf{v}}^{ijk} = \mathbf{M}\hat{\mathbf{z}}^{ijk}$. This implies that removing the edge $j \rightarrow i$ in \mathcal{G} would break the causal Markov condition for q^{E_k} . Now let $\delta = \min_{k,i,j} \delta_k^{(i,j)} > 0$, then for all $\mathbf{M} \in \overline{\mathcal{M}}_{\text{sur}}(\mathcal{G})$ such that $\|\mathbf{M} - \mathbf{I}\|_F \leq \delta$, the resulting distributions $q^{E_k}, k \in [K]$ satisfy assumption Assumption 1.

Combining the above two statements and what we have proven before, it is straightforward to see that one can choose some $\mathbf{M} \in \mathcal{M}_{\text{sur}}(\mathcal{G})$ in a small neighbourhood of \mathbf{I} that satisfies all the requirements, completing the proof. \square

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: Our abstract and introduction provide the readers a sense of our main results.

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 compare with existing works in the introduction and the related work sections.

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: Rigorous proofs are provided in the appendix.

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: We introduce our experimental setup in details.

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: [No]

Justification: Code will be released after review.

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: Details are provided.

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: we run experiments on 100 random causal graphs and report the overall accuracy.

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: [\[No\]](#)

Justification: The experiments do not require huge computational resources and can be run on a local computer.

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: The paper conforms the code of ethics.

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: [NA]

Justification:

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:

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: [NA]

Justification:

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: [NA]

Justification:

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:

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:

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.