
Causal Discovery in Heterogeneous Environments Under the Sparse Mechanism Shift Hypothesis

Ronan Perry

MPI for Intelligent Systems, Tübingen
rflperry@uw.edu

Julius von Kügelgen*

MPI for Intelligent Systems, Tübingen
University of Cambridge
jvk@tue.mpg.de

Bernhard Schölkopf*

MPI for Intelligent Systems, Tübingen
bs@tue.mpg.de

Abstract

Machine learning approaches commonly rely on the assumption of independent and identically distributed (*i.i.d.*) data. In reality, however, this assumption is almost always violated due to distribution shifts between environments. Although valuable learning signals can be provided by heterogeneous data from changing distributions, it is also known that learning under arbitrary (adversarial) changes is impossible. Causality provides a useful framework for modeling distribution shifts, since causal models encode both observational and interventional distributions. In this work, we explore the *sparse mechanism shift hypothesis*, which posits that distribution shifts occur due to a *small* number of changing causal conditionals. Motivated by this idea, we apply it to learning causal structure from heterogeneous environments, where *i.i.d.* data only allows for learning an equivalence class of graphs without restrictive assumptions. We propose the *Mechanism Shift Score* (MSS), a score-based approach amenable to various empirical estimators, which provably identifies the entire causal structure with high probability if the sparse mechanism shift hypothesis holds. Empirically, we verify behavior predicted by the theory and compare multiple estimators and score functions to identify the best approaches in practice. Compared to other methods, we show how MSS bridges a gap by both being nonparametric as well as explicitly leveraging sparse changes.

1 Introduction

Classical machine learning methods and theory presume data to be independently and identically distributed (*i.i.d.*). Although there has been huge success under this assumption, research on topics including adversarial examples [52, 14], distribution shifts [41, 45, 43, 6], and “spurious” correlations [2] has highlighted its fragility. Open questions remain as to how we can relax the *i.i.d.* assumption and still learn useful models, since learning under unrestricted adversarial distribution shifts seems infeasible [4]. Causal models naturally provide structure to a distribution via a factorization into causal *mechanisms*, the processes by which variables are dependent on their direct causes; Hence, they are a natural basis for studying distribution shifts. Based on the idea of the independence of causal mechanisms [36, 45], the *sparse mechanism shift hypothesis* [46] posits that distribution shifts are the result of changes in only a subset of the causal model’s mechanisms. This presents a promising relaxation with many potential applications throughout machine learning [32, 46, 55, 31].

*Shared last author.

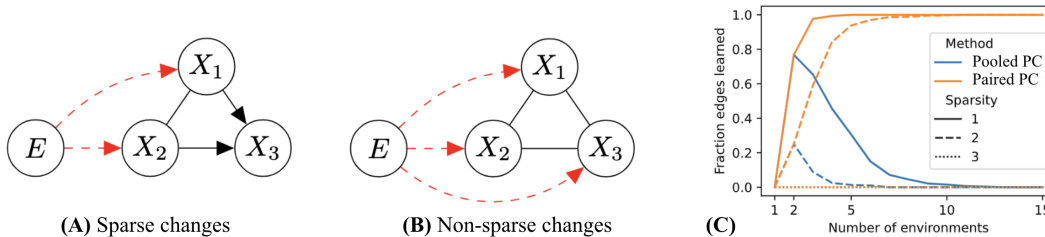


Figure 1: **Sparse shifts yield structure identifiability.** In a causal model with a fully connected DAG over $\mathbf{X} := \{X_1, X_2, X_3\}$, no edge directions can be learned from i.i.d. observational data. Given multiple datasets on \mathbf{X} across different environments with possible distribution shifts, the PC algorithm [51] can be applied to the pooled data augmented by the environment index variable E [25]. (A) A mechanism shift $E \rightarrow X_i$ can allow orientation of some edges. (B) Yet dense shifts prohibit any orientations [25]. (C) Pooling all environments leads to dense shifts, even if pairs of environments differ sparsely. Combining structure learned across *pairs* of pooled environments leads to identifiability under *sparse* shifts; under *dense* shifts, there is no difference.

In particular, causal discovery—the inference of qualitative structure encoded by a graph and underlying causal relationships between variables—is a fundamental scientific problem with broad applications for which distribution shifts have been successfully leveraged. Classical approaches, which assume i.i.d. data from a single environment or domain, can be broadly categorized into three classes. *Constraint-based* methods, such as the PC [51] and FCI [37] algorithms, perform a series of statistical independence tests to identify the existence of causal relationships under various assumptions (e.g., causal faithfulness) and then use certain rules to infer causal directions. *Score-based* methods, such as GES [3], optimize some score function, such as the penalized likelihood BIC [47], over the set of possible graphs. These methods all come with asymptotic *structure identifiability* guarantees of recovering the (Markov) equivalence class [12], but they will rarely uniquely recover the true causal graph. Further work on *functional-form-based* methods, beginning with non-Gaussian assumptions [50], provides various ways to recover a specific DAG by assuming a certain functional form of the causal mechanisms [23, 56, 38, 26], at the cost of potential misspecification.

Although identification is possible through actively specified interventions [8, 7, 18, 17, 49], natural distribution shifts on the same variables across different environments allow for promising approaches under a relaxation of the i.i.d. assumption [5, 45, 30, 39, 20, 19, 13, 24, 10, 11, 25, 33, 15, 29]. Such *multi-environment* methods can learn the direct causes of a specific variable [39, 20] or the entire causal structure [10, 11, 25] without requiring knowledge of which causal mechanisms change. Yet, as noted in these work but not studied, performance is dependent on which and how many changes occur.

Overview and contributions. After a review of causal graphical models and formalization of our setting (§ 2), we discuss key related work on multi-environment causal discovery (§ 3). We demonstrate that sparse distribution changes provide structure identifiability in the bivariate case, and useful information in multivariate settings (§ 4). Based on the observation that pairwise comparisons of environments better leverage sparse changes (see Fig. 1), we propose the *Mechanism Shift Score*, a score-based causal discovery algorithm with theoretical guarantees (§ 5). Empirical results confirm the theory (§ 6). In summary, the present work makes the following contributions:

- We prove that by relaxing the i.i.d. assumption via the sparse mechanism shift assumption, bivariate causal structure is identifiable without parametric assumptions (Cor. 4.2).
- We introduce the *Mechanism Shift Score (MSS)*, defined as the number of conditional distributions implied by a graph which change across all pairs of environments (§ 5). We prove that the true causal (multivariate) graph minimizes the MSS over possible graphs (Prop. 5.1).²
- We provide rates of convergence showing that with a sufficient number of sparsely changing environments, the causal graph *uniquely* minimizes the score function with high probability (Cor. 5.4). Our rates readily apply to existing literature on learning individual mechanisms [39, 20] and the entire graph [11, 25], where a study of the role of sparsity was previously missing.
- We demonstrate empirically that sparsity and pairwise comparisons are useful and show how the MSS accommodates various parametric [39, 10, 11] and nonparametric [20, 25] estimators (§ 6).

²While scores are typically *maximized*, it is more natural to *minimize* distribution shifts and hence the MSS.

2 Problem setting and notation

We start by building up the causal framework needed to understand our work and related literature. It relies on common graph-theoretic terminology which we review for completeness in Appx. A.

Causal terminology. Causal relationships between variables are encoded in a causal graphical model (CGM) which links graphical and distributional properties via certain assumptions.

Definition 2.1 (Causal Graphical Model (CGM)). A causal graphical model (CGM) $\mathcal{M} = (G, \mathbb{P}_{\mathbf{X}})$ over d random variables $\mathbf{X} = \{X_1, \dots, X_d\}$ consists of (i) a *directed acyclic graph* (DAG) G with vertices \mathbf{X} and edges $X_i \rightarrow X_j$ iff X_i is a direct cause of X_j , (ii) and a joint distribution $\mathbb{P}_{\mathbf{X}}$ which factorizes (is *Markovian*) over G .³ Formally, we have the following *Markov* or *causal factorization*:

$$\mathbb{P}_{\mathbf{X}}(X_1, \dots, X_d) = \prod_{j=1}^d \mathbb{P}_{\mathbf{X}}(X_j | \mathbf{PA}_j), \quad (2.1)$$

where \mathbf{PA}_j are the parents (direct causes) of X_j in G and $\mathbb{P}_{\mathbf{X}}(X_j | \mathbf{PA}_j)$ is the *causal mechanism* of X_j .

Implicit in the definition is the assumption that there is *no hidden confounding*, i.e., that any common cause of two or more observed variables is included in \mathbf{X} ; for this reason it is also referred to as *causal sufficiency*. This is a strong assumption to make and important to keep in mind in applications.

The Markov factorization (2.1) of the CGM encodes various conditional independences between variables. The *Markov equivalence class* (MEC) is the set of DAGs which share the same set of conditional independence relations; graphically, it is the set of DAGs which share the same *skeleton* (set of edges regardless of direction) and *v-structures* ($X_i \rightarrow X_j \leftarrow X_k$, but $X_i \not\leftrightarrow X_k$) [40, Lem. 6.25]. A set of DAGs such as the MEC are commonly represented as a *completed partially directed acyclic graph* (CPDAG) in which edges are directed if directed in all DAGs in the set, and otherwise left undirected [16]. Incorrect DAGs in the MEC induce *non-causal factorizations* containing *non-causal conditionals* which differ from the true causal mechanisms. Furthermore, in a CGM (2.1) is equivalent to the *global Markov condition* which states that for disjoint vertex sets $\mathbf{A}, \mathbf{B}, \mathbf{Z}$ in G :

$$\mathbf{A} \perp_G \mathbf{B} \mid \mathbf{Z} \implies \mathbf{A} \perp \mathbf{B} \mid \mathbf{Z} \quad (2.2)$$

where \perp_G denotes d-separation in G (Appx. A) [40, Thm. 6.22]. While the Markov property allows us to derive distributional properties from the DAG, causal discovery concerns deriving graphical properties from distributional properties. This requires the *causal faithfulness assumption*, effectively assuming that variables are not statistically independent unless implied by the graph.

Assumption 2.2 (Causal faithfulness). The observational distribution $\mathbb{P}_{\mathbf{X}}$ is said to be *faithful* to the causal graph G if every conditional independence relationship in $\mathbb{P}_{\mathbf{X}}$ is implied by d-separation in G (i.e., statistical dependence implies d-connection (2.2) [40, Def. 6.33]).

Multi-environment data. We assume that we observe a collection \mathcal{D} of datasets from a set \mathcal{E} of (possibly different) environments, where each dataset \mathcal{D}^e from environment e contains n_e independent and identically distributed (i.i.d.) observations from some joint distribution $\mathbb{P}_{\mathbf{X}}^e$,

$$\mathcal{D} = \{\mathcal{D}^1, \dots, \mathcal{D}^{n_{\mathcal{E}}}\} \quad \text{where} \quad \mathcal{D}^e = \{\mathbf{X}^{e,1}, \dots, \mathbf{X}^{e,n_e}\} \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}_{\mathbf{X}}^e \quad \text{for} \quad e \in \mathcal{E} = \{1, \dots, n_{\mathcal{E}}\}.$$

Environments can encapsulate experimental settings, or broader contexts such as climate or time [33]. A key assumption of our setting is that environments arise from different “soft” interventions on an underlying shared CGM \mathcal{M} . Specifically, the CGM \mathcal{M} entails a set of *interventional distributions*, resulting from changing a subset of the *mechanisms* $\mathbb{P}_{\mathbf{X}}(X_j | \mathbf{PA}_j)$ to some different $\tilde{\mathbb{P}}_{\mathbf{X}}(X_j | \mathbf{PA}_j)$.⁴

Assumption 2.3 (Shared mechanisms). Each environment e independently results from \mathcal{M} by intervening on an (unknown) subset $\mathcal{I}^e \subseteq [d]$ of mechanisms, i.e., $\tilde{\mathbb{P}}_{\mathbf{X}}^e$ can be written as

$$\tilde{\mathbb{P}}_{\mathbf{X}}^e(X_1, \dots, X_d) = \left(\prod_{j \in \mathcal{I}^e} \tilde{\mathbb{P}}_{\mathbf{X}}^e(X_j | \mathbf{PA}_j) \right) \prod_{j \in [d] \setminus \mathcal{I}^e} \mathbb{P}_{\mathbf{X}}(X_j | \mathbf{PA}_j). \quad (2.3)$$

³Throughout, we assume the existence of densities with respect to the Lebesgue measure.

⁴Note that this includes “hard” interventions $\tilde{\mathbb{P}}_{\mathbf{X}}(X_j | \mathbf{PA}_j) = \delta(X_j - x)$ which fix X_j to a specific value x , and also “soft” interventions which merely alter the functional relationship.

We make the following common assumption about how these changes arise.

Assumption 2.4 (Independent causal mechanisms (ICM) [36, 45]). A change in $\mathbb{P}(X_j | \mathbf{PA}_j)$ has no effect on and provides no information on $\mathbb{P}(X_k | \mathbf{PA}_j)$ for any $k \neq j$.

Since causal mechanisms are fixed within an environment and potentially vary across them, we can create an *augmented CGM* to unify CGMs from different environments. Note that in Asm. 2.3 and the following definition, the same causal parents and DAG are preserved over different environments; the distribution changes are limited to soft interventions which do not change a variable’s causal parents.

Definition 2.5 (Augmented CGM [25]). Let $\{(G, \mathbb{P}_{\mathbf{X}}^e)\}_{e=1}^{n_{\mathcal{E}}}$ be a collection of CGMs over the DAG $G = (\mathbf{X}, T)$ from multiple environments. The augmented CGM is defined as $\mathcal{M}' := (G', \mathbb{P}_{\mathbf{X} \cup E})$ where (i) E is a random environment indicator variable with support \mathcal{E} , and (ii) the augmented DAG G' has vertices $\mathbf{X} \cup E$ and edge set $T \cup \{(E, X_j) : \exists e, e' \in \mathcal{E} \text{ s.t. } \mathbb{P}_{\mathbf{X}}^e(X_j | \mathbf{PA}_j^G) \neq \mathbb{P}_{\mathbf{X}}^{e'}(X_j | \mathbf{PA}_j^G)\}$.

Note that $\mathbb{P}_{\mathbf{X} \cup E}$ is Markovian to G' , inheriting the factorization from the underlying CGMs along with the added dependence on E . With respect to the augmented DAG, Asm. 2.4 implies that existence of the edge $E \rightarrow X_i$ provides no information on the existence of an edge $E \rightarrow X_j$ for $j \neq i$. As discussed by Huang et al. [25], since E can be a common cause of variables in \mathbf{X} , causal sufficiency in the original CGM over \mathbf{X} is violated. We instead must assume *pseudo-causal sufficiency*. See Appx. C for further discussion of this assumption and the ICM principle.

Assumption 2.6 (Pseudo causal sufficiency [25]). Any unobserved confounders of variables in \mathbf{X} can be written solely as functions of E . Thus, within any given environment e , all unobserved confounders are fixed and causal sufficiency holds.

Sparse mechanism shift (SMS) hypothesis. Another key assumption, supported Asm. 2.4’s implication that a change to one mechanism does not imply changes to others, is the following:

Assumption 2.7 (SMS [46]). Changes in mechanisms between observed environments are sparse:

$$0 < |\mathcal{I}^e| < d \tag{2.4}$$

The value of this assumption when met is illustrated in Fig. 1 and will be elaborated upon in § 5.

3 Related work on causal discovery from multiple environments

Causal discovery from changing distributions and causal mechanisms has a long history, going back to Simon’s *invariance criterion*, stating that the true causal order is the one that is invariant under the right sort of intervention [21, 22]. Tian and Pearl [53] infer a causal order by testing which marginal distributions change under a single intervention. Invariant causal prediction (ICP) [39] can identify the causal parents of a target variable under the assumption that the target’s causal mechanism is invariant across environments [45, 39, 20]. However, applying ICP to each variable in order to learn all sets of causal parents and hence the causal graph is not immediately admissible: the invariance assumption would imply that all variables are invariant, and thus there are no mechanism changes to learn from.

Learning the MEC from i.i.d. data is a well-studied problem, and under certain assumptions, essentially solved [37, 51, 3]. However, it is still an open question how best to learn the true causal DAG from the MEC. Ghassami et al. [10] apply ideas from linear ICP to identify the causal DAG. Assuming *linear* causal mechanisms in which only the noise distributions change, they compare pairs of environments and orient edges to meet this requirement. Ghassami et al. [11] allow for any kind of change within a *linear* model. Based on the ICM principle, they demonstrate that non-causal DAGs induce a larger number of changes in the linear model parameters than the causal DAG. By counting parameter changes across pairwise environments, they can determine the causal order of variables. Huang et al. [25] remove any functional restrictions through a two-stage approach. First, they use the PC algorithm on the augmented CGM, pooling all the data; this both identifies the MEC but more importantly can also orient additional edges. Since not all edges are guaranteed to be oriented, they propose a second stage, relying on a novel measure of mechanism dependence to individually orient remaining edges.

A consistent yet sparingly explored theme across all of these methods is the impact of the sparsity of changes across environments. In ICP, it is briefly noted that the method is applicable when the target

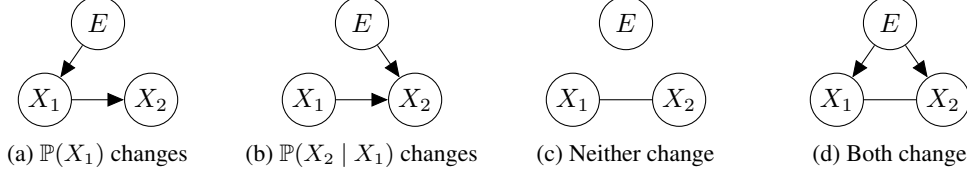


Figure 2: (a, b) Visualizations of the two cases explored in the proof of Prop. 4.1. In both cases, $X_2 \not\perp\!\!\!\perp E$ unconditionally and $X_1 \not\perp\!\!\!\perp E | X_2$ directly or by conditioning on the potential collider X_2 . (c, d) The two other possible situations; neither allow us to orient the edge between X_1 and X_2 .

variable experiences sparse shifts, as long as we assume a maximum degree of sparsity [39, §6.2]. Although not discussed, the performance of methods from Ghassami et al. [10, 11] also depend on sparsity. As mentioned by Huang et al. [25], identifiability requires the invariance of some mechanisms, and their pooled PC approach cannot identify edge directions when both adjacent variables change.

4 Leveraging sparse mechanism changes for causal discovery

The augmented DAG is a powerful tool for understanding the implications of changing causal mechanisms. We build up intuition in the bivariate setting and then consider the general multivariate setting. From now on, we assume causal faithfulness on the augmented CGM.

Bivariate case. Given a causal model composed of two associated variables, identifiability of the causal direction requires assumptions such as the functional form or existence of targeted interventions [40, §4.1]. We show how sparsely changing mechanisms also provide identifiability.

Proposition 4.1 (Both non-causal conditionals change). *Consider the bivariate setting $X_1 \rightarrow X_2$. If either $\mathbb{P}(X_2 | X_1)$ or $\mathbb{P}(X_1)$ change, then both $\mathbb{P}(X_1 | X_2)$ and $\mathbb{P}(X_2)$ change.*

Corollary 4.2. (Bivariate identifiability) *In the setting of Prop. 4.1, if only one mechanism changes (sparsity), then the bivariate causal structure is identifiable.*

Proof. (Prop. 4.1) We consider each case separately and use a proof by contradiction (of faithfulness).

(i) If $\mathbb{P}(X_1)$ changes (see Fig. 2a), then $G_{\mathbf{X} \cup E}$ contains the edge $E \rightarrow X_1$. If $\mathbb{P}(X_2)$ remained invariant, then $X_2 \perp\!\!\!\perp E$ (unfaithful due to the unblocked path $E \rightarrow X_1 \rightarrow X_2$). If $\mathbb{P}(X_1 | X_2)$ remained invariant, then $X_1 \perp\!\!\!\perp E | X_2$ (unfaithful due to the direct path $E \rightarrow X_1$).

(ii) If $\mathbb{P}(X_2 | X_1)$ changes (see Fig. 2b), then $G_{\mathbf{X} \cup E}$ contains the edge $E \rightarrow X_2$. If $\mathbb{P}(X_2)$ remained invariant, then $X_2 \perp\!\!\!\perp E$ (unfaithful due to the direct path $E \rightarrow X_2$). If $\mathbb{P}(X_1 | X_2)$ remained constant, then $X_1 \perp\!\!\!\perp E | X_2$ (unfaithful due to the unblocked collider path $E \rightarrow X_2 \leftarrow X_1$). \square

Proof. (Cor. 4.2) If either causal mechanism changes, by Prop. 4.1 both conditionals in the non-causal factorization change. Hence, the causal structure is the one with only one mechanism change. \square

Multivariate case. Non-causal conditional distributions of X_j may change across environments even if the causal mechanism $\mathbb{P}(X_j | \mathbf{PA}_j)$ does not change. This occurs if the conditioning set leaves open a dependence between E and X_j in $G_{\mathbf{X} \cup E}$.

Lemma 4.3. *For any $X_j \in \mathbf{X}$ and set $\mathbf{Z} \subseteq \mathbf{X} \setminus \{X_j\}$, the conditional distribution $\mathbb{P}(X_j | \mathbf{Z})$ changes if and only if the following d -connectedness relationship holds:*

$$X_j \not\perp\!\!\!\perp_{G_{\mathbf{X} \cup E}} E | \mathbf{Z}.$$

The result follows from the Markov property and faithfulness; for all complete proofs, see Appx. B.

Since the Markov equivalence class is relatively easily available and thus often the starting point of open questions in causal discovery, we specify the implications of Lemma 4.3 in this setting. Note that due to the shared skeleton of all DAGs in the equivalence class, the conditioning set for any X_j in any DAG in the MEC is a subset of X_j 's true parents and children $\mathbf{PA}_j^G \cup \mathbf{CH}_j^G$ [51, 29].

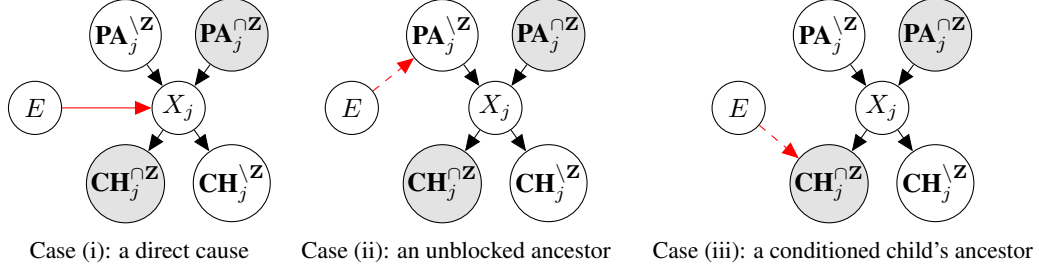


Figure 3: The three possible cases from Cor. 4.4 with a d-connecting path from E to X_j , conditioned on a subset of neighbors (colored in grey) in the MEC; these neighbors are a subset of the true parents and children.

Corollary 4.4. For any variable $X_j \in \mathbf{X}$ and set $\mathbf{Z} \subseteq (\mathbf{PA}_j^G \cup \mathbf{CH}_j^G)$ in the augmented graph, the conditional distribution $\mathbb{P}(X_j | \mathbf{Z})$ changes if and only if at least one of the following holds:

- (i) $E \rightarrow X_j$ [a direct cause].
- (ii) $\exists W_{\mathbf{PA}} \in \mathbf{PA}_j^G \setminus \mathbf{Z}$ such that $W_{\mathbf{PA}} \not\perp_{G_{\mathbf{X} \cup E}} E | \mathbf{Z}$ [unblocked path to unconditioned parent].
- (iii) $\exists W_{\mathbf{CH}} \in \mathbf{CH}_j^G \cap \mathbf{Z}$ such that $W_{\mathbf{CH}} \not\perp_{G_{\mathbf{X} \cup E}} E | \mathbf{Z} \setminus W_{\mathbf{CH}}$ [unblocked path to conditioned child].

Proof sketch. These cases are visualized in Fig. 3. The global Markov property and faithfulness assumption allow us to interchange d-connection and a change of mechanism. In the forward direction, a changing mechanism implies a d-connecting path which is necessarily captured in one of the three cases. In the reverse direction, each case opens a d-connecting path between X_j and E . \square

A direct result of changing causal mechanisms and conditional distributions follows.

Corollary 4.5. Let G^* denote the true (unaugmented) DAG and let G be any other DAG over the same variables. For any $X_j \in \mathbf{X}$, a change in $\mathbb{P}(X_j | \mathbf{PA}_j^{G^*})$ implies a change in $\mathbb{P}(X_j | \mathbf{PA}_j^G)$.

Proof. Under the true causal parents given by G^* , a change in mechanism occurs if and only if $E \rightarrow X_j$ in $G_{\mathbf{X} \cup E}^*$. So, X_j is d-connected to E , no matter the conditioning set $\mathbf{Z} \subset \mathbf{X} \setminus \{X_j\}$. Thus, by Lemma 4.3, $\mathbb{P}(X_j | \mathbf{PA}_j^G)$ necessarily changes. \square

5 Causal discovery via the Mechanism Shift Score (MSS)

We have established that changing mechanisms provide useful information and by Cor. 4.2 can provide identifiability in the bivariate case under faithfulness and sparse changes. We now study identifiability in more general graphs along with approaches to developing practical estimators. Motivated by the novel discovery of the value of comparing pairwise environments if changes are sparse, we propose the *Mechanism Shift Score (MSS)* estimand with useful theoretical guarantees over DAGs in the MEC.

The MSS estimand. Given a set \mathcal{G} of candidate DAGs over the same d variables $\mathbf{X} = \{X_1, \dots, X_d\}$ with data sampled from n_ε distributions $\mathcal{P} := \{\mathbb{P}_{\mathbf{X}}^e\}_{e=1}^{n_\varepsilon}$, we count the number of changing conditional distributions in each graph by defining

$$\text{MSS}_j(G; \mathcal{P}) = \sum_{e' > e}^{n_\varepsilon} \mathbb{I} \left[\mathbb{P}^e(X_j | \mathbf{PA}_j^G) \neq \mathbb{P}^{e'}(X_j | \mathbf{PA}_j^G) \right] \quad \text{and} \quad \text{MSS}(G; \mathcal{P}) = \sum_{j=1}^d \text{MSS}_j(G; \mathcal{P}).$$

$\text{MSS}_j(G; \mathcal{P})$ is the number of pairs of environments in which the conditional distribution of X_j in G changes; $\text{MSS}(G; \mathcal{P})$ is the total number of changes across all variables and pairs of environments according to the Markov factorization implied by G . It follows from Cor. 4.5 that the true DAG G^* minimizes (not necessarily uniquely) the number of changing conditionals among all DAGs.

Proposition 5.1. Let G^* be the true DAG in the set \mathcal{G} of DAGs. Then for all $G \in \mathcal{G}$ and $j \in \{1, \dots, d\}$

$$\text{MSS}_j(G^*; \mathcal{P}) \leq \text{MSS}_j(G; \mathcal{P}) \quad \text{and thus} \quad \text{MSS}(G^*; \mathcal{P}) \leq \text{MSS}(G; \mathcal{P}).$$

Proof. By Cor. 4.5, any change in $\mathbb{P}(X_j|\mathbf{PA}_j^{G^*})$ implies a change in $\mathbb{P}(X_j|\mathbf{PA}_j^G)$ for any other DAG G . Thus, any change counted by MSS_j on the true DAG will also be detected in every other DAG and so both lower bounds hold. \square

Prop. 5.1 can be viewed as the generalization of the Principle of Minimal Changes [11], allowing for mechanism changes beyond the parametric restrictions of a linear model. Identifiability, however, requires us to establish a discerning aspect of the true structure. Recall that the Markov equivalence class \mathcal{G}_{MEC} is identifiable. We define the subset

$$\mathcal{G}_{\text{MEC}}^{\min} := \arg \min_{G \in \mathcal{G}_{\text{MEC}}} \text{MSS}(G; \mathcal{P})$$

of DAGs with minimum MSS, here defined as the number of mechanism shifts across environments. In practice, we may employ any generic conditional test for change in mechanism or choose to use a “softer” score (e.g., based on p -values) to quantify changes along a continuous spectrum.

Prop. 5.1 implies that $G^* \in \mathcal{G}_{\text{MEC}}^{\min}$. Using probabilistic assumptions based on the idea of sparse changes, we show that, given enough environments, the causal parents and full DAG are identifiable.

Lemma 5.2 (Identifiability of causal parents). *Let G^* be the true DAG in the MEC \mathcal{G}_{MEC} and ρ_i the probability that the causal mechanism of X_i is different across any two environments. Under Asms. 2.2 to 2.4 and 2.6, for any $j \in \{1, \dots, d\}$, graph $G \in \mathcal{G}_{\text{MEC}}$ such that $\mathbf{PA}_j^{G^*} \neq \mathbf{PA}_j^G$, and lower and upper bounds on the shift probabilities $\rho_i^{\text{LB}} \leq \rho_i \leq \rho_i^{\text{UB}}$ for all i , we have that*

$$\Pr[\text{MSS}_j(G^*; \mathcal{P}) < \text{MSS}_j(G; \mathcal{P})] \geq 1 - \left(1 - (1 - \rho_j^{\text{UB}}) \min_i \rho_i^{\text{LB}}\right)^{\lfloor n\varepsilon/2 \rfloor}.$$

Proof sketch. From Cor. 4.5, we know: $\text{MSS}_j(G^*; \mathcal{P}) \leq \text{MSS}_j(G; \mathcal{P})$. We use the shared skeleton property of all DAGs in a MEC and Lemma 4.3, for which only case (i) admits a changing mechanism in the true DAG. Based on the ICM principle, we examine sufficient conditions in a pair of environments and establish a probabilistic upper bound on all pairs. \square

Note that this bound is independent of the other DAG G , so long as G is in the MEC and has a different set of causal parents. As a special case of Lemma 5.2, invariance of the mechanism of X_j implies $\rho_j^{\text{UB}} = 0$ and hence provides a bound relevant to invariant causal prediction [39, 20]. Building off of Lemma 5.2, we can provide a probabilistic bound on identifiability of the whole graph.

Theorem 5.3 (Identifiability of the graph). *Let G^* be the true DAG in the MEC \mathcal{G}_{MEC} and ρ_j the probability that the causal mechanism of X_j is different across any two environments. Under assumptions 2.2, 2.3, 2.4, and 2.6, and bounds $\rho_i^{\text{LB}} \leq \rho_i \leq \rho_i^{\text{UB}}$ for all i , we have that*

$$\Pr[\mathcal{G}_{\text{MEC}}^{\min} = \{G^*\}] \geq 1 - |\mathcal{G}_{\text{MEC}}| \left(1 - (1 - \min_i \rho_i^{\text{UB}}) \min_i \rho_i^{\text{LB}}\right)^{\lfloor n\varepsilon/2 \rfloor}.$$

Proof sketch. From Prop. 5.1, G^* is always in $\mathcal{G}_{\text{MEC}}^{\min}$. For each DAG, we use Lemma 5.2 to bound the probability that all mechanisms exhibit the same number of changes. Then we apply the union bound to establish an upper bound across all DAGs. \square

Corollary 5.4. *If ρ_i is bounded away from 0 and 1 for all i , (a probabilistic form of Asm. 2.7),*

$$\Pr[\mathcal{G}_{\text{MEC}}^{\min} = \{G^*\}] \rightarrow 1 \quad \text{as} \quad n\varepsilon \rightarrow \infty$$

That is, with enough environments we can recover the true DAG from the Markov equivalence class.

Proof. The assumption of bounded probability implies that $\rho_i^{\text{UB}} < 1$ and $\rho_i^{\text{LB}} > 0$ for all i . Hence, by the rate established in Thm. 5.3, identifiability is achieved in the limit. \square

The MSS estimator. An MSS estimator of the proposed estimand requires us to be able to test if two conditional distributions change across two environments. This can be done using conditional independence tests or equality of distribution tests [34]. Under parametric assumptions, models may be fit for each mechanism, and these parameters can then be tested across environments [10, 11, 39]. Heinze-Deml et al. [20] provided a comprehensive study of such tests and their power for ICP. More recent but less studied work by Park et al. [35] has provided a kernel-based approach with strong guarantees. In practice care must be taken when using equality of conditional distribution tests, especially if any of their assumptions are violated [48].

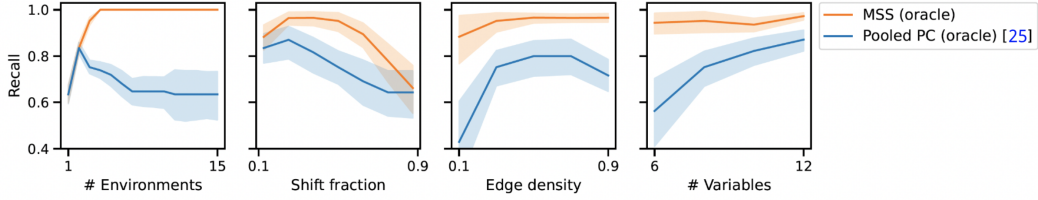


Figure 4: **Oracle rates match the theory.** From left to right (a-d): (a) With sufficiently many environments, our MSS approach learns the true DAG while pooled PC recovers only the original MEC. (b) Pairwise comparisons are less beneficial when shifts are extremely sparse or dense. Differences across the two methods are particularly pronounced both in (c) sparse and dense DAGs, as well as in (d) smaller DAGs. Shaded regions denote 95% confidence intervals calculated from bootstrapped data.

Computational complexity. The score function $\text{MSS}(G; \mathcal{P})$ is *decomposable* [16] in that it is the sum of local scores $\text{MSS}_j(G; \mathcal{P})$. In the most naïve approach, each mechanism in each of $|\mathcal{G}_{\text{MEC}}|$ DAGs must be tested across all pairwise environments, on the overall order of $O(|\mathcal{G}_{\text{MEC}}|dn_{\mathcal{E}}^2)$, without accounting for the complexity of the statistical test. This can be slow, but if the hypothesis test scales with sample size n faster than $O(n^2)$, then pairwise tests may actually be faster than pooling the data. In practice, the decomposable property permits a speedup: since many mechanisms will be shared across DAGs, we can test each unique mechanism and then select the results relevant for each DAG. Experimentally, we find the main bottlenecks to be large sample sizes and numbers of environments.

6 Structure learning experiments

Having established the value of the MSS estimand through theory, we now seek to understand empirically through simulations: (i) how the oracle MSS and pooled PC approaches compare across experimental settings, (ii) which possible MSS estimators perform best, and (iii) how the MSS compares to related approaches. In Appx. D we present a case study application to the real cytometry dataset [44].⁵ For ease of comparison, we adapt the simulation setup of Huang et al. [25]. Specifically, random DAGs are sampled using an Erdős-Rényi model [9] in which each edge has some fixed probability of existing or not (the *density of edges*). In each environment, a random set of variables experience a change in mechanism according to a fixed number or fraction of sparse shifts. Given a DAG, each variable j in environment e has a randomly sampled mechanism

$$X_j^e := \sum_{i \in \text{PA}_j} b_{ji}^e f_{ji}(X_i^e) + \sigma_j^e \epsilon_j^e \quad (6.1)$$

where $b_{ji} \sim \mathcal{U}(0.5, 2.5)$, $\sigma_j^e \sim \mathcal{U}(1, 3)$, and $\epsilon_j^e \sim \mathcal{N}(0, 1)$ or $\mathcal{U}(1, 3)$ with equal probability. The functions f_{ji} are selected uniformly at random from $\{x^2, x^3, \tanh, \text{sinc}\}$. Mechanisms in an unobserved baseline environment are sampled and $\sigma_{ji} = 1$ is fixed. Each observed environment inherits the baseline distributions and mechanisms shifts are resampled per (6.1).

We evaluate the quality of an estimated CPDAG against the true DAG via *precision* and *recall* [10, 11, 25]. Precision is the fraction of *directed* edges in the CPDAG which are correctly oriented. Recall is the fraction of *all* edges in the CPDAG which are oriented. Thus, the true DAG has perfect precision and recall. Since all methods start from the MEC, there are no incorrect edges, only incorrect orientations.

Oracle MSS rates match the theory. Having theory on learning rates bounds under sparsity, we now seek to understand how the empirical performance of MSS and pooled PC depends on a variety of graph and sparsity settings, both which our theory does and does not address. We consider random DAGs over 6 variables with edge density 0.3. Five environments are sampled, in each of which half of the mechanisms shift. In Fig. 4, we hold all of these settings fixed and vary one at a time across 50 repetitions, comparing the recall of the two methods.⁶ Precision is always perfect under the oracle test.

In the first two plots, the empirical results match what the theory predicts. First, per Cor. 5.4, with more environments MSS learns the entire graph while pooled PC learns nothing but the original

⁵All code and experiments are available at https://github.com/rflperry/sparse_shift

⁶Oracle methods used code from the *causaldag* package [3-Clause BSD license].

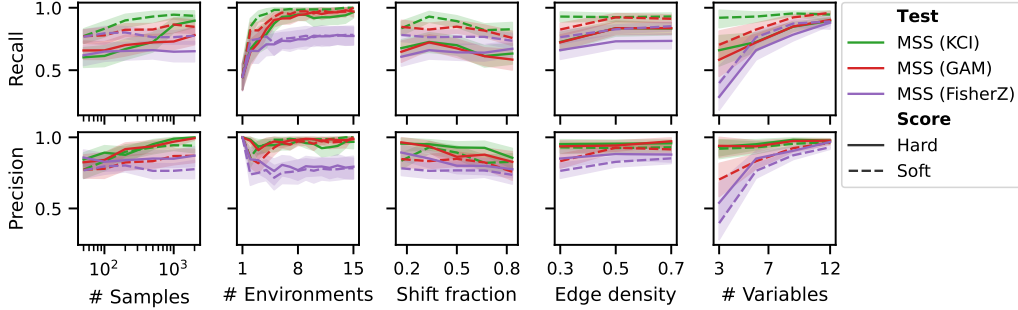


Figure 5: Nonparametric hypothesis tests perform well in nonlinear simulations, and soft scores succeed. Notably, recall converges with increasing environments. KCI appears to best balance high recall and precision.

MEC. Second, per Thm. 5.3, the learning rate decreases when shifts are either uncommon or frequent. Furthermore, we see that differences between the two approaches are accentuated in sparse and dense DAGs, as well as in smaller DAGs. Note that sparsity is a fixed fraction of the variables, and hence larger DAGs experience more shifts in absolute terms. See Appx. D for additional oracle experiments.

MSS performance depends on the chosen estimator. Next, we study how the choice of estimator and type of score affects performance. We use two popular conditional independence tests, the Fisher-Z partial correlation test [27] and the Kernel Conditional Independence (KCI) test [57], as well as the invariant residual test using a generalized additive model (GAM), a top-performing ICP method [20].⁷ Each detects a change if the test p -value is less than $\alpha := 0.05/d$, a Bonferroni correction to bound the false positive rate for each scored DAG. Although the theory pertains to counting shifts, in settings where a “hard” hypothesis test has low power, it may be of more value to use a “softer” score (e.g., based on p -values) to quantify changes along a continuous spectrum. We propose the following “soft” score (see Appx. C for further details):

$$\widehat{\text{MSS}}_j(G; \mathcal{D}) = \sum_{e=1, e' > e}^{n\mathcal{E}} \left[1 - p\text{-value} \left(\mathbb{P}^e(X_j | \mathbf{PA}_j^G) \neq \mathbb{P}^{e'}(X_j | \mathbf{PA}_j^G) \right) \right].$$

DAGs are generated on six variables with edge density 0.3. Three environments are sampled, with 500 samples and two mechanism shifts per environment. We vary each variable while holding the others fixed and compare results across 50 repetitions.

As seen in Fig. 5, the Fisher-Z test performs noticeably worse, presumably due to the unmet parametric assumptions, while the two nonparametric approaches do well, noticeably so at higher sample sizes. As with the oracle, the true DAG is recovered with enough environments. The “soft” versions achieves higher recall at the cost of worse precision, since there won’t be ties between candidate mechanism scores but the score is noisier. The “soft” KCI seems to be best, suggesting it fits the data the best. In practice, it is crucial for a method to model the data well in order for the p -value to be valid [48].

MSS compares favorably against other methods. We now wish to compare the MSS to relevant existing methods. Specifically, the *Minimal changes (MC)* approach [11] tests pairs of environments for changes in the parameters of a linear model. Huang et al. [25] provides a nonparametric version: a two-stage approach which first uses PC with the KCI test on the pooled data. We investigate if pooling data loses information under empirical tests, and how the nonparametric pairwise MSS test combines the best of both approaches.

In Fig. 6, we compare these approaches in the same experimental setting as previously studied. Pooled PC has quite high precision and yet suffers from lower recall, especially with more environments at which point the recall is no better than the base MEC. In contrast, “hard” KCI has much higher recall at the cost of some precision, although this goes away at larger samples sizes. The parametric MC works surprisingly well and yet slightly worse than KCI. Overall, we see that MSS combines the value of the pairwise comparisons from the MC approach with the flexibility of incorporating various nonparametric estimators. Additional experiments in Appx. D confirm this in the bivariate case.

⁷KCI and Fisher-Z are implemented by the *causal-learn* package [GNU General Public License].

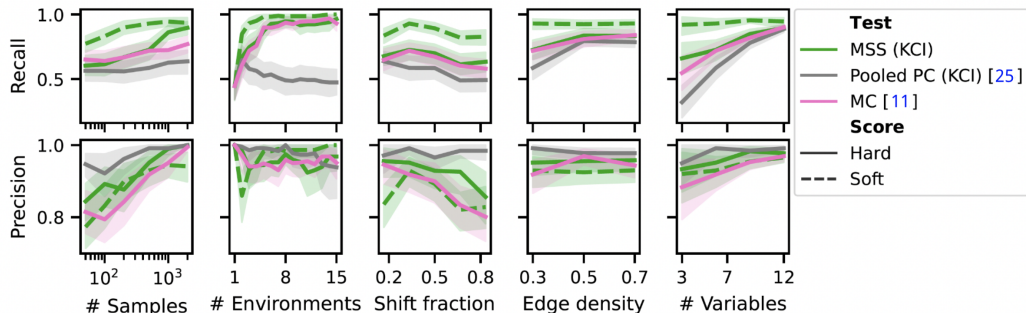


Figure 6: Pairwise approaches improve with more environments, unlike Pooled PC. Although the parametric MC works surprisingly well across settings, the nonparametric MSS with KCI appears superior.

7 Discussion

Sparse shifts as a relaxation of the i.i.d. assumption. Distribution shifts are a common violation of the i.i.d. assumption and a problematic source of error in practice. It has been argued that the issue of robustness to natural shifts is connected to causality [45, 1, 58, 39, 54, 28, 55]. We have shown that viewing a shift in distribution through the causal factorization permits a useful relaxation of the i.i.d. assumption, facilitating causal discovery: if there are no shifts (i.i.d.), we can only identify equivalence classes. If shifts are unrestricted, we cannot meaningfully transfer across distributions [45], and ideas such as Simon’s invariance criterion [22] and ICP [39] do not help. However, if shifts occur *sparsely*—as we formalize—we can provably use this as a learning signal to infer the causal structure.

The Mechanism Shift Score (MSS) and prior methods. The MSS framework extends previous causal discovery work limited to linear mechanisms [11], just as nonlinear ICP [20] extended the initially limited ICP approach [39] beyond linear models. Additionally, we have provided a graph-theoretic analysis proving why pairwise comparisons are actually useful; the learning rates we have established apply to these previous works and provide insight on the role of sparsity. Although Huang et al. [25] provide a thorough nonparametric approach and analysis, we have demonstrated that the first stage of their two-stage-approach is not suited for sparsely-detectable changes since it pools all the data. The MSS is both nonparametric and explicitly leverages sparsity through pairwise comparisons.

Beyond causal discovery. Once the causal graph is known, conditional distributions and hence an entire causal graphical model can be learned. This is harder than learning a statistical model, but has various advantages, especially when distributions shift, as they do in reality. E.g., we may be able to use such a model for causal reasoning, i.e., estimating a certain causal effect. We also expect that MSS can serve as a useful inductive bias for causal representation learning, similar to how invariant prediction [45, 39] inspired invariant risk minimization [28]; recent work has started to explore this [29, 31].

Empirical performance and the validity of the sparse mechanism shift hypothesis. We infer causal structure through a flexible score-based method and, as empirically demonstrated, strong results can be obtained by multiple estimators and when the assumption of sparsity is met. We conjecture that under an oracle, the “hard” MSS is equivalent to the PC algorithm pooled pairwise across environments. It is worth noting that the “hard” approach may still be useful under dense changes if only a sparse number of them are large enough to be *empirically discernible*. Thus the empirical method can actually outperform the oracle baseline and be useful even if the assumption of sparsity is unmet.

Outlook and conclusion. Imagining causal models on an axis of complexity, from the microscopic physical laws of nature to a simplified set of variables and relationships, we transition from a system with no mechanism shifts (effectively a dynamical system) to a system in which all mechanisms shift (due to many unmeasured causes). In the middle, we posit only a sparse number of shifts to be empirically discernible. While *all (causal) models are wrong*, the one which is most invariant to shifts may be the best candidate for supporting robust and transferable inference.

Acknowledgments and Disclosure of Funding

We thank Jonas Kübler, Junhyung Park, Krikamol Muandet, and the Tübingen causality team for helpful discussions. This work was supported by the German Federal Ministry of Education and Research (BMBF): Tübingen AI Center, FKZ: 01IS18039A, 01IS18039B; and by the Machine Learning Cluster of Excellence, EXC number 2064/1 – Project number 390727645. Ronan Perry was supported by a Fulbright Germany research fellowship.

References

- [1] Elias Bareinboim and Judea Pearl. Transportability from multiple environments with limited experiments: Completeness results. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 280–288. Curran Associates, Inc., 2014. [10](#)
- [2] Sara Beery, Grant Van Horn, and Pietro Perona. Recognition in terra incognita. In *Proceedings of the European conference on computer vision (ECCV)*, pages 456–473, 2018. [1](#)
- [3] David Maxwell Chickering. Optimal Structure Identification With Greedy Search. *Journal of Machine Learning Research*, 3(Nov):507–554, 2002. ISSN ISSN 1533-7928. [2](#), [4](#)
- [4] Shai Ben David, Tyler Lu, Teresa Luu, and David Pal. Impossibility Theorems for Domain Adaptation. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pages 129–136. JMLR Workshop and Conference Proceedings, March 2010. ISSN: 1938-7228. [1](#)
- [5] Daniel Eaton and Kevin Murphy. Exact bayesian structure learning from uncertain interventions. In *Artificial intelligence and statistics*, pages 107–114, 2007. [2](#), [22](#), [23](#)
- [6] Frederick Eberhardt. A Sufficient Condition for Pooling Data. *Synthese*, 163(3):433–442, 2008. ISSN 0039-7857. Publisher: Springer. [1](#)
- [7] Frederick Eberhardt and Richard Scheines. Interventions and Causal Inference. *Philosophy of Science*, 74(5):981–995, 2007. ISSN 0031-8248. Publisher: [The University of Chicago Press, Philosophy of Science Association]. [2](#), [19](#)
- [8] Frederick Eberhardt, Clark Glymour, and Richard Scheines. On the Number of Experiments Sufficient and in the Worst Case Necessary to Identify All Causal Relations Among N Variables. In *Proceedings of the 21st Annual Conference on Uncertainty in Artificial Intelligence (UAI)*, pages 178–184, 2005. [2](#)
- [9] Paul Erdős and Alfréd Rényi. On the evolution of random graphs. *Publ. Math. Inst. Hung. Acad. Sci.*, 5(1):17–60, 1960. [8](#)
- [10] Amir Emad Ghassami, Saber Salehkaleybar, Negar Kiyavash, and Kun Zhang. Learning causal structures using regression invariance. In *Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS’17*, page 3015–3025, Red Hook, NY, USA, 2017. Curran Associates Inc. ISBN 9781510860964. [2](#), [4](#), [5](#), [7](#), [8](#)
- [11] Amir Emad Ghassami, Negar Kiyavash, Biwei Huang, and Kun Zhang. Multi-domain Causal Structure Learning in Linear Systems. In *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018. [2](#), [4](#), [5](#), [7](#), [8](#), [9](#), [10](#)
- [12] Clark Glymour, Kun Zhang, and Peter Spirtes. Review of causal discovery methods based on graphical models. *Frontiers in Genetics*, 10, 2019. [2](#)
- [13] M. Gong, K. Zhang, T. Liu, D. Tao, C. Glymour, and B. Schölkopf. Domain adaptation with conditional transferable components. In *Proceedings of the 33rd International Conference on Machine Learning (ICML)*, volume 48, pages 2839–2848, 2016. [2](#)
- [14] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *International Conference on Learning Representations*, 2015. [1](#)

- [15] Siyuan Guo, Viktor Tóth, Bernhard Schölkopf, and Ferenc Huszár. Causal de Finetti: On the Identification of Invariant Causal Structure in Exchangeable Data. *arXiv:2203.15756 [cs, math, stat]*, March 2022. arXiv: 2203.15756. 2, 19
- [16] Alain Hauser and Peter Bühlmann. Characterization and greedy learning of interventional Markov equivalence classes of directed acyclic graphs. *J. Mach. Learn. Res.*, 13(1):2409–2464, August 2012. ISSN 1532-4435. 3, 8
- [17] 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, June 2014. ISSN 0888613X. arXiv: 1205.4174. 2
- [18] Yang-Bo He and Zhi Geng. Active Learning of Causal Networks with Intervention Experiments and Optimal Designs. *Journal of Machine Learning Research*, 9(84):2523–2547, 2008. ISSN 1533-7928. 2
- [19] Yango He and Zhi Geng. Causal Network Learning from Multiple Interventions of Unknown Manipulated Targets. *arXiv:1610.08611*, October 2016. 2
- [20] Christina Heinze-Deml, Nicolai Meinshausen, Jonas Peters, et al. Invariant causal prediction for nonlinear models. *Journal of Causal Inference*, 6(2):1–35, 2018. 2, 4, 7, 9, 10
- [21] K. Hoover. The logic of causal inference. *Economics and Philosophy*, 6:207–234, 1990. 4
- [22] K. D. Hoover. Causality in economics and econometrics. In S. N. Durlauf and L. E. Blume, editors, *The New Palgrave Dictionary of Economics*. Palgrave Macmillan, Basingstoke, UK, 2nd edition, 2008. 4, 10
- [23] Patrik Hoyer, Dominik Janzing, Joris M Mooij, Jonas Peters, and Bernhard Schölkopf. Nonlinear causal discovery with additive noise models. In *Advances in Neural Information Processing Systems*, volume 21. Curran Associates, Inc., 2009. 2
- [24] Biwei Huang, Kun Zhang, Jiji Zhang, Ruben Sanchez-Romero, Clark Glymour, and Bernhard Schölkopf. Behind distribution shift: Mining driving forces of changes and causal arrows. In *2017 IEEE International Conference on Data Mining (ICDM)*, pages 913–918. IEEE, 2017. 2
- [25] Biwei Huang, Kun Zhang, Jiji Zhang, Joseph Ramsey, Ruben Sanchez-Romero, Clark Glymour, and Bernhard Schölkopf. Causal discovery from heterogeneous/nonstationary data. *Journal of Machine Learning Research*, 21(89):1–53, 2020. 2, 4, 5, 8, 9, 10, 19
- [26] Dominik Janzing, Joris Mooij, Kun Zhang, Jan Lemeire, Jakob Zscheischler, Povilas Daniušis, Bastian Steudel, and Bernhard Schölkopf. Information-geometric approach to inferring causal directions. *Artificial Intelligence*, 182-183:1–31, May 2012. ISSN 0004-3702. 2
- [27] Markus Kalisch and Peter Bühlmann. Estimating High-Dimensional Directed Acyclic Graphs with the PC-Algorithm. *Journal of Machine Learning Research*, 8(22):613–636, 2007. 9
- [28] Pritish Kamath, Akilesh Tangella, Danica Sutherland, and Nathan Srebro. Does invariant risk minimization capture invariance? In Arindam Banerjee and Kenji Fukumizu, editors, *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics*, volume 130 of *Proceedings of Machine Learning Research*, pages 4069–4077. PMLR, 13–15 Apr 2021. 10
- [29] Sébastien Lachapelle, Pau Rodríguez López, Yash Sharma, Katie Everett, Rémi Le Priol, Alexandre Lacoste, and Simon Lacoste-Julien. Disentanglement via Mechanism Sparsity Regularization: A New Principle for Nonlinear ICA. *arXiv:2107.10098 [cs, stat]*, February 2022. arXiv: 2107.10098. 2, 5, 10
- [30] Vincenzo Lagani, Ioannis Tsamardinos, and Sofia Triantafillou. Learning from Mixture of Experimental Data: A Constraint-Based Approach. In *Artificial Intelligence: Theories and Applications*, volume 7297, pages 124–131. Springer Berlin Heidelberg, 2012. ISBN 978-3-642-30447-7 978-3-642-30448-4. 2

- [31] Phillip Lippe, Sara Magliacane, Sindy Löwe, Yuki M. Asano, Taco Cohen, and Efstratios Gavves. CITRIS: Causal Identifiability from Temporal Intervened Sequences. *arXiv:2202.03169 [cs, stat]*, February 2022. arXiv: 2202.03169. [1](#), [10](#)
- [32] F. Locatello, B. Poole, G. Rätsch, B. Schölkopf, O. Bachem, and M. Tschannen. Weakly-supervised disentanglement without compromises. In *Proceedings of the 37th International Conference on Machine Learning (ICML)*, volume 119 of *Proceedings of Machine Learning Research*, pages 6348–6359. PMLR, July 2020. [1](#)
- [33] Joris M. Mooij, Sara Magliacane, and Tom Claassen. Joint Causal Inference from Multiple Contexts. *Journal of Machine Learning Research*, 21(99):1–108, 2020. ISSN 1533-7928. [2](#), [3](#), [19](#), [22](#), [23](#)
- [34] Sambit Panda, Cencheng Shen, Ronan Perry, Jelle Zorn, Antoine Lutz, Carey E Priebe, and Joshua T Vogelstein. Nonpar manova via independence testing. *arXiv preprint arXiv:1910.08883*, 2021. [7](#)
- [35] J. Park, U. Shalit, B. Schölkopf, and K. Muandet. Conditional distributional treatment effect with kernel conditional mean embeddings and u-statistic regression. In *Proceedings of 38th International Conference on Machine Learning (ICML)*, volume 139 of *Proceedings of Machine Learning Research*, pages 8401–8412. PMLR, July 2021. [7](#)
- [36] Judea Pearl. *Causality: Models, Reasoning and Inference*. Cambridge University Press, 2000. [1](#), [4](#)
- [37] Peter Spirtes, Christopher Meek, and Thomas S. Richardson. An Algorithm for Causal Inference in the Presence of Latent Variables and Selection Bias. In Clark Glymour and Gregory F. Cooper, editors, *Computation, Causation and Discovery, chapter 6*, pages 211–252. The MIT Press, May 1999. [2](#), [4](#)
- [38] Jonas Peters, Joris M. Mooij, Dominik Janzing, and Bernhard Schölkopf. Identifiability of causal graphs using functional models. In *Proceedings of the Twenty-Seventh Conference on Uncertainty in Artificial Intelligence, UAI’11*, page 589–598, Arlington, Virginia, USA, 2011. AUAI Press. ISBN 9780974903972. [2](#)
- [39] Jonas Peters, Peter Bühlmann, and Nicolai Meinshausen. Causal inference by using invariant prediction: identification and confidence intervals. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 5(78):947–1012, 2016. [2](#), [4](#), [5](#), [7](#), [10](#)
- [40] Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. *Elements of causal inference*. The MIT Press, 2017. [3](#), [5](#), [15](#)
- [41] Joaquin Quiñero-Candela, Masashi Sugiyama, Anton Schwaighofer, and Neil D Lawrence. *Dataset shift in machine learning*. Mit Press, 2008. [1](#)
- [42] Joseph Ramsey and Bryan Andrews. FASK with Interventional Knowledge Recovers Edges from the Sachs Model. *arXiv:1805.03108 [cs, q-bio]*, May 2018. arXiv: 1805.03108. [22](#), [23](#)
- [43] Mateo Rojas-Carulla, Bernhard Schölkopf, Richard Turner, and Jonas Peters. Invariant models for causal transfer learning. *The Journal of Machine Learning Research*, 19(1):1309–1342, 2018. [1](#)
- [44] Karen Sachs, Omar Perez, Dana Pe’er, Douglas A Lauffenburger, and Garry P Nolan. Causal protein-signaling networks derived from multiparameter single-cell data. *Science*, 308(5721): 523–529, 2005. [8](#), [22](#), [23](#)
- [45] B. Schölkopf, D. Janzing, J. Peters, E. Sgouritsa, K. Zhang, and J. Mooij. On causal and anticausal learning. In *Proceedings of the 29th International Conference on Machine Learning*, pages 1255–1262, New York, NY, USA, 2012. Omnipress. [1](#), [2](#), [4](#), [10](#)
- [46] B. Schölkopf, F. Locatello, S. Bauer, N. R. Ke, N. Kalchbrenner, A. Goyal, and Y. Bengio. Toward causal representation learning. *Proceedings of the IEEE - Advances in Machine Learning and Deep Neural Networks*, 109(5):612–634, 2021. [1](#), [4](#)

- [47] Gideon Schwarz. Estimating the Dimension of a Model. *The Annals of Statistics*, 6(2):461 – 464, 1978. [2](#)
- [48] Rajen D. Shah and Jonas Peters. The hardness of conditional independence testing and the generalised covariance measure. *The Annals of Statistics*, 48(3):1514 – 1538, 2020. [7](#), [9](#), [19](#)
- [49] Karthikeyan Shanmugam, Murat Kocaoglu, Alexandros G Dimakis, and Sriram Vishwanath. Learning causal graphs with small interventions. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc., 2015. [2](#)
- [50] 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. [2](#)
- [51] P. Spirtes, C. Glymour, and R. Scheines. *Causation, Prediction, and Search*. MIT press, 2001. [2](#), [4](#), [5](#)
- [52] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian J. Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In *International Conference on Learning Representations*, 2014. [1](#)
- [53] Jin Tian and Judea Pearl. Causal discovery from changes. In *Proceedings of the Seventeenth conference on Uncertainty in artificial intelligence*, pages 512–521, 2001. [4](#)
- [54] Victor Veitch, Alexander D’Amour, Steve Yadlowsky, and Jacob Eisenstein. Counterfactual Invariance to Spurious Correlations: Why and How to Pass Stress Tests. *Advances in Neural Information Processing Systems*, 34, 2021. [10](#)
- [55] J. von Kügelgen, Y. Sharma, L. Gresele, W. Brendel, B. Schölkopf, M. Besserve, and F. Locatello. Self-supervised learning with data augmentations provably isolates content from style. In *Advances in Neural Information Processing Systems 34 (NeurIPS 2021)*, pages 16451–16467. Curran Associates, Inc., December 2021. [1](#), [10](#)
- [56] K Zhang and A Hyvärinen. On the identifiability of the post-nonlinear causal model. In *25th Conference on Uncertainty in Artificial Intelligence (UAI 2009)*, pages 647–655. AUAI Press, 2009. [2](#)
- [57] K. Zhang, J. Peters, D. Janzing, and B. Schölkopf. Kernel-based conditional independence test and application in causal discovery. pages 804–813. AUAI Press, July 2011. [9](#)
- [58] K. Zhang, M. Gong, and B. Schölkopf. Multi-source domain adaptation: A causal view. In *Proceedings of the 29th AAAI Conference on Artificial Intelligence*, pages 3150–3157, 2015. [10](#)

Appendices

Table of Contents

A Graph terminology	15
B Full proofs	16
B.1 Proof of Lemma 4.3	16
B.2 Proof of Corollary 4.4	16
B.3 Proof of Lemma 5.2	17
B.4 Proof of Theorem 5.3	18
C Details of assumptions and methods	19
C.1 Pseudo causal sufficiency and the Independent Causal Mechanisms (ICM) assumption	19
C.2 The p -values “soft” score	19
D Supporting experiments	20
D.1 Additional simulations	20
D.2 Application to real-world cytometry data	21

A Graph terminology

A directed graph $G = (V, T)$ is an object consisting of a set of vertices V and a set of ordered pairs of vertices $T \subset V \times V$ corresponding to directed edges in G . A *path* is a sequence of vertices $(V_{i_1}, \dots, V_{i_n})$ with $n \geq 2$ such that $V_{i_k} \rightarrow V_{i_{k+1}}$ or $V_{i_k} \leftarrow V_{i_{k+1}}$ and in a *directed path* $V_{i_k} \rightarrow V_{i_{k+1}}$ for all k . In graph G : the *children* \mathbf{CH}_j^G of V_j are all V_m such that $V_j \rightarrow V_m$, the *parents* \mathbf{PA}_j^G of V_j are all V_m such that $V_m \rightarrow V_j$, the *ancestors* \mathbf{AN}_j^G of V_j are all V_m such that there exists a directed path (V_m, \dots, V_j) , and the *descendants* \mathbf{DE}_j^G of V_j are V_j and all V_m such that there exists a directed path (V_j, \dots, V_m) . The graph superscript will be omitted unless needed. A *cycle* is a path such that $V_{i_1} = V_{i_n}$ and a *directed acyclic graph (DAG)* is a directed graph with no directed cycles.

On a path $(V_{i_1}, \dots, V_{i_k}, \dots, V_{i_n})$, we say variable V_{i_k} is a *collider* if $V_{i_{k-1}} \rightarrow V_{i_k}$ and $V_{i_k} \leftarrow V_{i_{k+1}}$. A subset $\mathbf{Z} \in \mathbf{V} \setminus \{V_{i_1}, V_{i_n}\}$ *blocks* the path if either (i) \mathbf{Z} contains at least one non-collider vertex on the path or (ii) the path contains a collider with no descendants in \mathbf{Z} (this includes the collider itself by the descendant definition). With this terminology, we say that on the disjoint variable sets \mathbf{A} , \mathbf{B} , and \mathbf{Z} , \mathbf{A} is *d-separated* from \mathbf{B} by \mathbf{Z} iff every path between \mathbf{A} and \mathbf{B} is blocked by \mathbf{Z} [40, Def. 6.1]. This is denoted as $\mathbf{A} \perp_G \mathbf{B} \mid \mathbf{Z}$. If \mathbf{A} and \mathbf{B} are not d-separated, and hence there exists an unblocked path, we say that they are *d-connected*.

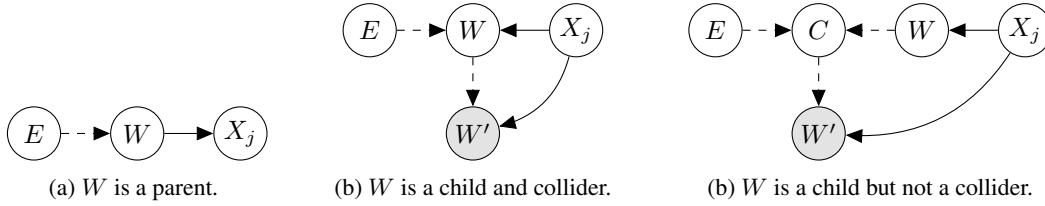


Figure 7: Cases from the proof of Cor. 4.4. Case (a) iff case (ii). Case (b) induces two subcases either of which occur iff case (iii).

B Full proofs

B.1 Proof of Lemma 4.3

Lemma 4.3. *For any $X_j \in \mathbf{X}$ and set $\mathbf{Z} \subseteq \mathbf{X} \setminus \{X_j\}$, the conditional distribution $\mathbb{P}(X_j \mid \mathbf{Z})$ changes if and only if the following d-connectedness relationship holds:*

$$X_j \not\perp_{G_{\mathbf{X} \cup E}} E \mid \mathbf{Z}.$$

Proof. (\Rightarrow) If $\mathbb{P}(X_j \mid \mathbf{Z})$ changes across environments E , then $X_j \not\perp E \mid \mathbf{Z}$. The global Markov property of the CGM states that d-separation implies conditional independence, and thus by the contra-positive the d-connectedness relationship follows.

(\Leftarrow) d-connectedness implies conditional dependence by faithfulness, and thus a change across environments. \square

B.2 Proof of Corollary 4.4

Corollary 4.4. *For any variable $X_j \in \mathbf{X}$ and set $\mathbf{Z} \subseteq (\mathbf{PA}_j^G \cup \mathbf{CH}_j^G)$ in the augmented graph, the conditional distribution $\mathbb{P}(X_j \mid \mathbf{Z})$ changes if and only if at least one of the following holds:*

- (i) $E \rightarrow X_j$ [a direct cause].
- (ii) $\exists W_{\mathbf{PA}} \in \mathbf{PA}_j^G \setminus \mathbf{Z}$ such that $W_{\mathbf{PA}} \not\perp_{G_{\mathbf{X} \cup E}} E \mid \mathbf{Z}$ [unblocked path to unconditioned parent].
- (iii) $\exists W_{\mathbf{CH}} \in \mathbf{CH}_j^G \cap \mathbf{Z}$ such that $W_{\mathbf{CH}} \not\perp_{G_{\mathbf{X} \cup E}} E \mid \mathbf{Z} \setminus W_{\mathbf{CH}}$ [unblocked path to conditioned child].

Proof. By Lemma 4.3, $\mathbb{P}(X_j \mid \mathbf{Z})$ changes iff $X_j \not\perp_{G_{\mathbf{X} \cup E}} E \mid \mathbf{Z}$ (d-connection) and equivalently, iff there is an unblocked path (E, \dots, X_j) in $G_{\mathbf{X} \cup E}$. We assume a generic path and through casework establish that it is unblocked if and only if one of cases (i), (ii), or (iii) holds. The casework is visualized in Figure 7.

Either (E, \dots, X_j) is just $E \rightarrow X_j$ [case (i)] or there must exist W such that (E, \dots, W, X_j) and W is either (a) a collider or (b) not a collider.

(a) If W is a collider, it necessarily a child of X_j . The collided path is unblocked iff $W \not\perp_{G_{\mathbf{X} \cup E}} E \mid \mathbf{Z} \setminus W$ and some descendant $W' \in \mathbf{DE}_W^G$ of W is conditioned on. Thus $W' \in \mathbf{Z} \subset \mathbf{PA}_j^G \cup \mathbf{CH}_j^G$. Without loss of generality, assume W' is the closest descendant to W and hence the path (W, \dots, W') is unblocked by \mathbf{Z} . W' cannot be a parent of X_j , else induce the cycle (X_j, W, \dots, W', X_j) , and so must be a child and case (iii) holds. Specifically, there exists W' such that $W' \not\perp_{G_{\mathbf{X} \cup E}} E \mid \mathbf{Z} \setminus W'$ and W' is a child in \mathbf{Z} .

(b) If W is not a collider, by definition the path is unblocked iff $W \not\perp_{G_{\mathbf{X} \cup E}} E \mid \mathbf{Z}$ and $W \notin \mathbf{Z}$. If W is a parent of X_j (since they are adjacent), case (ii) holds. If W is a child of X_j , because E has an outgoing edge there must exist some collider C on the path such that $(E, \dots, C, \dots, W, X_j)$ and the subpath from W to C is directed into C . The condition $W \not\perp_{G_{\mathbf{X} \cup E}} E \mid \mathbf{Z}$ holds iff some descendant W' of C is in \mathbf{Z} . As before, W' cannot be a parent of X_j or else induce a cycle, and so it must be a child and case (iii) holds. \square

B.3 Proof of Lemma 5.2

Lemma 5.2 (Identifiability of causal parents). *Let G^* be the true DAG in the MEC \mathcal{G}_{MEC} and ρ_i the probability that the causal mechanism of X_i is different across any two environments. Under Asms. 2.2 to 2.4 and 2.6, for any $j \in \{1, \dots, d\}$, graph $G \in \mathcal{G}_{\text{MEC}}$ such that $\mathbf{PA}_j^{G^*} \neq \mathbf{PA}_j^G$, and lower and upper bounds on the shift probabilities $\rho_i^{\text{LB}} \leq \rho_i \leq \rho_i^{\text{UB}}$ for all i , we have that*

$$\Pr[\text{MSS}_j(G^*; \mathcal{P}) < \text{MSS}_j(G; \mathcal{P})] \geq 1 - \left(1 - (1 - \rho_j^{\text{UB}}) \min_i \rho_i^{\text{LB}}\right)^{\lfloor n_{\mathcal{E}}/2 \rfloor}.$$

Proof. By Assumption 2.3, the distribution $\mathbb{P}_{\mathbf{X}}^e$ in each environment $e \in \{1, \dots, n_{\mathcal{E}}\}$ is the result of changing mechanisms from some underlying yet unknown distribution $\mathbb{P}_{\mathbf{X}}$. Let $\Delta^{e,e'}(X_j)$ denote the event $\mathbb{I}[\mathbb{P}_{\mathbf{X}}^e(X_j | \mathbf{PA}_j^{G^*}) \neq \mathbb{P}_{\mathbf{X}}^{e'}(X_j | \mathbf{PA}_j^{G^*})]$ that the mechanism of variable X_j , with respect to the true graph G^* , changes across environments e and e' . Abbreviate $\rho_j^{e,e'} := \Pr[\Delta^{e,e'}(X_j) = 1]$.

Since $\mathbf{PA}_j^{G^*} \neq \mathbf{PA}_j^G$ and G shares the same skeleton as G^* , at least one edge must be oriented incorrectly in G . In the conditioning set \mathbf{PA}_j^G according to the incorrect graph G , there thus exists either an unconditioned true parent $Z \in \mathbf{PA}_j^{G^*} \setminus \mathbf{PA}_j^G$ or a conditioned-upon true child $Z \in \mathbf{CH}_j^{G^*} \cap \mathbf{PA}_j^G$. By Cor. 4.4, we know that if Z is not d-separated from E in the augmented graph, then the conditional $\mathbb{P}(X_j | \mathbf{PA}_j^G)$ changes across E . This occurs at least if the mechanism of Z directly changes, e.g. there is the edge $E \rightarrow Z$ in the augmented graph.

Consider first the case of two environments. We know from Prop. 5.1 that $\text{MSS}_j(G^*; \mathcal{P})$ cannot be greater than $\text{MSS}_j(G; \mathcal{P})$, and will be less if the mechanism of X_j remains invariant while the mechanism of Z changes. By the assumption of independent changing mechanisms,

$$\begin{aligned} \Pr[\text{MSS}_j(G^*; \{\mathcal{D}^1, \mathcal{D}^2\}) = \text{MSS}_j(G; \{\mathcal{D}^1, \mathcal{D}^2\})] \\ &= 1 - \Pr[\text{MSS}_j(G^*; \{\mathcal{D}^1, \mathcal{D}^2\}) < \text{MSS}_j(G; \{\mathcal{D}^1, \mathcal{D}^2\})] \\ &\leq 1 - \Pr[\Delta^{1,2}(X_j) = 0, \Delta^{1,2}(Z) = 1] \\ &= 1 - \Pr[\Delta^{1,2}(X_j) = 0] \Pr[\Delta^{1,2}(Z) = 1] \\ &= 1 - (1 - \rho_j^{1,2}) \rho_Z^{1,2} \end{aligned}$$

Given $n_{\mathcal{E}} > 2$ environments, it follows that

$$\begin{aligned} \Pr[\text{MSS}_j(G^*; \mathcal{P}) = \text{MSS}_j(G; \mathcal{P})] \\ &= \Pr \left[\bigcap_{e,e' > e} \text{MSS}_j(G^*, \{\mathcal{D}^e, \mathcal{D}^{e'}\}) = \text{MSS}_j(G, \{\mathcal{D}^e, \mathcal{D}^{e'}\}) \right] \\ &\leq \Pr \left[\bigcap_{e \in \{1, \dots, \lfloor \mathcal{E}/2 \rfloor\}} \text{MSS}_j(G^*, \{\mathcal{D}^{2e-1}, \mathcal{D}^{2e}\}) = \text{MSS}_j(G, \{\mathcal{D}^{2e-1}, \mathcal{D}^{2e}\}) \right] \\ &= \prod_{e \in \{1, \dots, \lfloor \mathcal{E}/2 \rfloor\}} \Pr [\text{MSS}_j(G^*, \{\mathcal{D}^{2e-1}, \mathcal{D}^{2e}\}) = \text{MSS}_j(G, \{\mathcal{D}^{2e-1}, \mathcal{D}^{2e}\})] \\ &\leq \prod_{e \in \{1, \dots, \lfloor \mathcal{E}/2 \rfloor\}} \left(1 - (1 - \rho_j^{2e-1, 2e}) \rho_Z^{2e-1, 2e}\right) \\ &\leq \left(1 - \min_{e \in \{1, \dots, \lfloor \mathcal{E}/2 \rfloor\}} (1 - \rho_j^{2e-1, 2e}) \rho_Z^{2e-1, 2e}\right)^{\lfloor n_{\mathcal{E}}/2 \rfloor}. \end{aligned}$$

Since Z is arbitrary, we construct an upper bound using the worst case, in which a variable frequently or rarely changes.

$$\begin{aligned} 1 - \min_{e \in \{1, \dots, \lfloor \mathcal{E}/2 \rfloor\}} (1 - \rho_j^{2e-1, 2e}) \rho_Z^{2e-1, 2e} &\leq 1 - (1 - \max_{e \in \{1, \dots, \lfloor \mathcal{E}/2 \rfloor\}} \rho_j^{2e-1, 2e}) \min_{e \in \{1, \dots, \lfloor \mathcal{E}/2 \rfloor\}} \rho_Z^{2e-1, 2e} \\ &\leq 1 - (1 - \max_{e, e' \neq e} \rho_j^{e, e'}) \min_{e, e' \neq e} \rho_Z^{e, e'} \end{aligned}$$

and so to acquire the final bound with simplified notation, for any variable X_i denote the minima and maxima of $\rho_i^{e, e'}$ across any two environments with ρ_i^{LB} and ρ_i^{UB} , respectively. \square

B.4 Proof of Theorem 5.3

Theorem 5.3 (Identifiability of the graph). *Let G^* be the true DAG in the MEC \mathcal{G}_{MEC} and ρ_j the probability that the causal mechanism of X_j is different across any two environments. Under assumptions 2.2, 2.3, 2.4, and 2.6, and bounds $\rho_i^{\text{LB}} \leq \rho_i \leq \rho_i^{\text{UB}}$ for all i , we have that*

$$\Pr[\mathcal{G}_{\text{MEC}}^{\min} = \{G^*\}] \geq 1 - |\mathcal{G}_{\text{MEC}}| \left(1 - (1 - \min_i \rho_i^{\text{UB}}) \min_i \rho_i^{\text{LB}}\right)^{\lfloor n\varepsilon/2 \rfloor}.$$

Proof. Since $\Pr[\mathcal{G}_{\text{MEC}}^{\min} = \{G^*\}] = 1 - \Pr[\mathcal{G}_{\text{MEC}}^{\min} \neq \{G^*\}]$ and by Lemma 5.2,

$$\begin{aligned} \Pr[\mathcal{G}_{\text{MEC}}^{\min} \neq \{G^*\}] &= \Pr \left[\bigcup_{G \in \mathcal{G}_{\text{MEC}} \setminus \{G^*\}} \text{MSS}(G^*; \mathcal{P}) = \text{MSS}(G; \mathcal{P}) \right] \\ &\leq \sum_{G \in \mathcal{G}_{\text{MEC}}} \Pr[\text{MSS}(G^*; \mathcal{P}) = \text{MSS}(G; \mathcal{D})] \\ &\leq \sum_{G \in \mathcal{G}_{\text{MEC}}} \Pr \left[\sum_j \text{MSS}_j(G^*; \mathcal{P}) = \sum_j \text{MSS}_j(G; \mathcal{D}) \right] \\ &= \sum_{G \in \mathcal{G}_{\text{MEC}}} \Pr \left[\sum_j \text{MSS}_j(G^*; \mathcal{P}) = \sum_j \text{MSS}_j(G; \mathcal{D}) \right] \\ &= \sum_{G \in \mathcal{G}_{\text{MEC}}} \Pr \left[\bigcap_j \text{MSS}_j(G^*; \mathcal{P}) = \text{MSS}_j(G; \mathcal{D}) \right] \\ &\leq \sum_{G \in \mathcal{G}_{\text{MEC}}} \min_j \Pr[\text{MSS}_j(G^*; \mathcal{P}) = \text{MSS}_j(G; \mathcal{D})] \\ &\leq \sum_{G \in \mathcal{G}_{\text{MEC}}} \min_j \left(1 - (1 - \rho_j^{\text{UB}}) \min_i \rho_i^{\text{LB}}\right)^{\lfloor n\varepsilon/2 \rfloor} \\ &\leq \sum_{G \in \mathcal{G}_{\text{MEC}}} \left(1 - (1 - \min_j \rho_j^{\text{UB}}) \min_i \rho_i^{\text{LB}}\right)^{\lfloor n\varepsilon/2 \rfloor} \\ &= |\mathcal{G}_{\text{MEC}}| \left(1 - (1 - \min_j \rho_j^{\text{UB}}) \min_i \rho_i^{\text{LB}}\right)^{\lfloor n\varepsilon/2 \rfloor}. \end{aligned}$$

\square

C Details of assumptions and methods

C.1 Pseudo causal sufficiency and the Independent Causal Mechanisms (ICM) assumption

Huang et al. [25] introduced the idea of psuedo-causal sufficiency (Asm. 2.6) and provide a useful discussion on its relation to results on soft interventions by Eberhardt and Scheines [7]. We expand on their discussion in the context of how it relates to our work.

Guo et al. [15] provide a useful formalization of multi-environment data, specifically through a plate-notation representation. An environment e specifies parameters of the causal mechanisms in the CGM over \mathbf{X} ; we can think of environments as encapsulating specific experimental settings, or broad contexts such as climate or time [33]. Under the context of e , there is some distribution $\mathbb{P}_{\mathbf{X}}^e$ and we observe a dataset sampled i.i.d. The ICM assumption tells us that the parameters for each causal mechanism in an environment are chosen or sampled independently, and thus in the augmented CGM the edges from E appear independently.

Within each environment, i.e., when we condition on E , the environmental parameters are fixed; thus we are in the typical i.i.d. setting and causal sufficiency is implied by the CGM. However, without conditioning on E , the environmental parameters are not fixed and across two samples either all remain the same (if the samples are in the same environment) or some change. This dependence between samples through the parameters defined by E is the result of E being a confounder; thus causal sufficiency cannot hold over \mathbf{X} without conditioning on E . Because E is not necessarily a true causal variable but rather an environment encoding a fixed set of unmeasured variables, Huang et al. [25] call it a *pseudo-confounder*. It is worth noting that the second stage of the approach of Huang et al. [25] relies on a novel kernel-based test, which computes a measure of mechanism dependence across all samples. They correctly compare the test statistics rather than examine p -values, because the dependence between the raw samples would lead to a rejected p -value even if the mechanisms were independent.

C.2 The p -values “soft” score

We provide further details on the p -value “soft” score. Recall the modified score definition to be

$$\widehat{\text{MSS}}_j(G; \mathcal{D}) = \sum_{e=1, e' > e}^{n\varepsilon} \left[1 - p\text{-value} \left(\mathbb{P}^e(X_j | \mathbf{PA}_j^G) \neq \mathbb{P}^{e'}(X_j | \mathbf{PA}_j^G) \right) \right].$$

Using a test of equality of distribution, we calculate a test statistic; at a pre-specified level α , if the test is well specified [48], the one-sided p -value is valid and corresponds to the probability under the null hypothesis $H_0 : \mathbb{P}^e(X_j | \mathbf{PA}_j^G) = \mathbb{P}^{e'}(X_j | \mathbf{PA}_j^G)$ of a test statistic as large or larger than the observed test statistic.

If a mechanism changes, a powerful test should yield a small p -value and thus a term close to 1 in the summation, similar to the “hard” score. If a mechanism doesn’t change, since p -values are uniformly distributed in $[0, 1]$ under the null hypothesis, the term in the sum would be uniformly distributed. With enough variables and environments, the variance of the randomness will decrease and the behavior of the score will be dominated by the p -values under the alternatives. It must be noted that the p -values are not independent, as some will use data from the same environments.

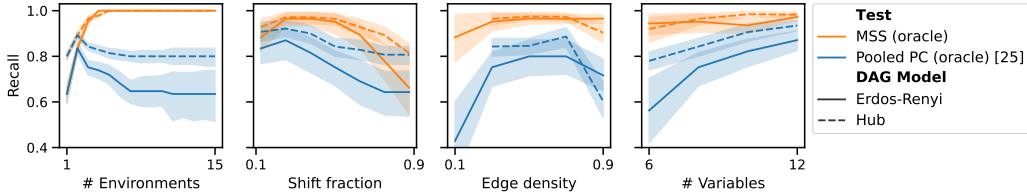


Figure 8: Oracle MSS improves upon pooled PC across both random and hub random graph models.

D Supporting experiments

D.1 Additional simulations

MSS improves upon pooled PC across random graph models. Previously in Fig. 4, we demonstrated that the MSS improves upon pooled PC under an oracle test in simulation settings where DAGs were sampled according to the Erdős-Rényi (ER) random DAG model; in the ER model, each edge is sampled i.i.d. with a fixed probability. Here, we expand upon that simulation by further comparing rates under the Barabasi-Alberts (Hub) scale-free random DAG model; in the Hub model, vertices are sequentially added to the DAG and edges are connected to previous vertices with probability proportional to their existing number of edges.

As seen in Fig. 8, we vary the same parameters as before but compare the two oracle methods across both DAG models this time. First, note that at 1 environment the Hub model exhibits greater recall, indicating that the observational MEC of the Hub graph has fewer unoriented edges than that of the ER graph. Thus, the gap between the methods is lessened in a Hub model as compared to an ER model. Otherwise, the qualitative trends between the two methods are almost identical across the two random graph models. The MSS appears at least mildly robust to the graph structure.

Differences between oracle MSS and pooled PC are most pronounced on sparser and smaller DAGs. Although Fig. 4 highlighted the most important trends of oracle methods in certain fixed settings, for completeness we examine rates of recall across additional fixed settings. As before, we sample DAGs from an Erdős-Rényi distribution and in five environments vary the DAG density, shift fraction, and number of variables. The set of experimental results shown in Fig. 9 convey broader trends in oracle recall rates as multiple variables change across row, column, and the x-axis. We do not vary the number of environments as we can only visualize three variables through our plot and the trend across environments is best understood from the theory. Differences in oracle recall rates are less pronounced on graphs with more variables and when the density of edges is large. Note that we only compare five environments here and that with more environments, differences will again be more pronounced; with enough environments, pooled PC cannot learn more than the MEC.

KCI-based approaches perform the best on bivariate CGMs. We previously examined the empirical rates of recall and precision across various simulated settings, highlighting when methods succeed and fail. Due to the size and complexity of those studied DAGs, not all results are fully interpretable. We seek to further understand empirical performance through the simple bivariate DAG, which contains no indirect effects and few possible interventions to analyze. Specifically, on the DAG $X_1 \rightarrow X_2$, shifts can occur to either $\mathbb{P}(X_1)$, $\mathbb{P}(X_2|X_1)$, neither mechanism, or both mechanisms; the first two shifts are sparse and provide oracle identifiability of the true DAG. Following the simulation setup described by eq. (6.1) on the DAG $X_1 \rightarrow X_2$, we simulate data from one base environment and from one interventional environment subject to one of the four possible shifts. Each environment has 500 samples. We compare the four different MSS methods using parametric and nonlinear equality of distribution tests and conditional independence tests. We also compare the pooled PC and MC approaches. Since only two environments are compared, we conjecture MSS and pooled PC to be equivalent under an oracle test.

Results are shown in Fig. 10. For reference, an oracle method would have recall 1 in the first two (sparse shift) columns and 0 in the other two columns. Although Fisher-Z has high precision when $\mathbb{P}(X_1)$ shifts, it has chance precision when $\mathbb{P}(X_2|X_1)$ shifts. The KCI methods maintain high precision while the precision of other methods is comparable or noticeably lower. With respect to recall,

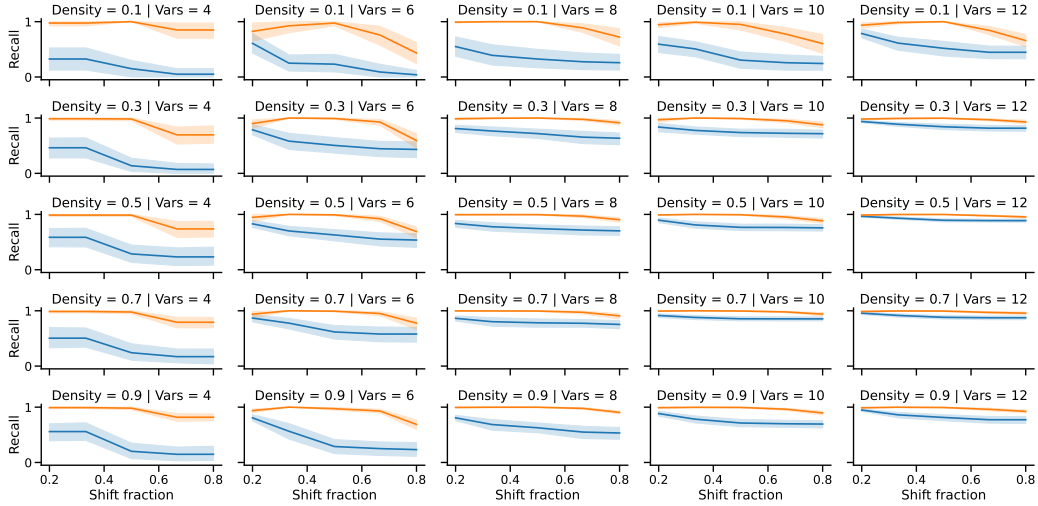


Figure 9: Differences between oracle MSS and pooled PC on 5 environments are most pronounced on smaller and sparse DAGs. For readability, the legend is omitted but we refer back to the same legend in Fig. 4. Specifically, the orange line corresponds to the MSS while the blue line corresponds to pooled PC. Only five environments are sampled, but differences would be exacerbated with additional environments.

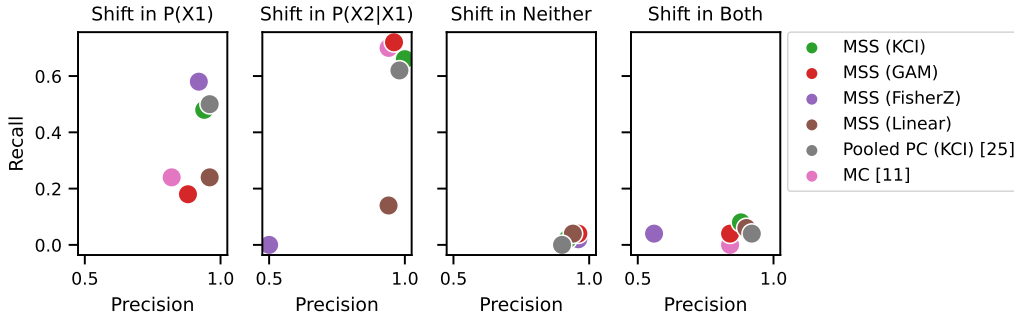


Figure 10: KCI-based tests perform well for causal identification in a bivariate CGM. 500 samples are drawn from a base environment, and a second environment subject to one of four shifts given by the columns; the first two columns are sparse shift settings where we have identifiability. The precision and recall are plotted for each of the methods. It appears that the two KCI-based methods (MSS and pooled PC) achieve the best balance of high power in both sparse shift settings while maintaining high precision in both non-sparse settings. Other methods either have drastically lower recall or precision close to 0.5, indicating random guessing.

when neither or both mechanisms change and thus the DAG is not identifiable, all methods correctly have low recall. However, when just the marginal $\mathbb{P}(X_1)$ changes, the KCI methods dominate in recall whereas the linear MSS MS, and GAM approaches have lower recall, implying they are less often able to detect a change in the reverse conditional $\mathbb{P}(X_1 | X_2)$. When the mechanism $\mathbb{P}(X_2 | X_1)$ shifts, all methods have high recall. Notably, the linear MSS performs much worse than MC. The only difference between them, however, is that MC explicitly counts how many parameters change while the linear MSS simply tests if there is a change; this does come at a slight cost in precision for MS though.

D.2 Application to real-world cytometry data

Although simulations with known ground truth provide useful reference points for comparing methods and evaluating empirical performance, in practice we are interested in studying real data with no known truth and additional challenges such as violated assumptions. To illustrate how one may apply

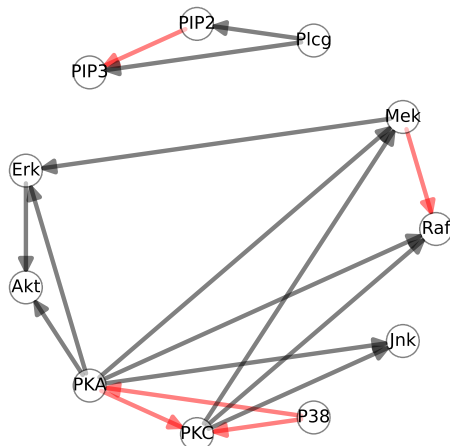


Figure 11: MSS (KCI) edge orientations mostly match the Sachs network [44]. Non-matching edges are posited to be involved in cycles and of ambiguous orientation in the literature, as we discuss and elaborate on.

our method in practice, and to analyze empirical performance on real data, we conduct a case-study application of MSS for causal discovery on a well-studied cytometry dataset [44].

D.2.1 Background

Sachs et al. [44] present a detailed study of the application of Bayesian discovery approaches to learning a causal DAG among protein concentration levels in human immune system cells. In each of 9 experimental environments subject to different perturbations, approximately 700-900 sample measurements were collected; each sample is the concentration levels of 11 proteins from a cell. The learned *Sachs network* is a proposed DAG among the variables, which the authors discuss and contrast with a domain-expert network from the “biologist’s view”. This cytometry data has subsequently been studied in further detail [5, 42, 33]. As is often pointed out, various assumptions may be violated, including the acyclicity assumption, since protein networks contain strong feedback loops [33]. As such, it is not necessarily useful to treat the Sachs network as a ground truth and there are numerous relationships and orientations which should rightfully be questioned [33]. Results must be considered in the context of domain-knowledge and various existing studies in the literature.

D.2.2 Experimental setup

In order to focus on learning edge orientations of undirected edges in the MEC, rather than learning the MEC, we start from the Sachs network despite the potential caveats. The Sachs network from Sachs et al. [44] is a DAG on the 11 variables with 17 edges. We compute the *Sachs MEC* which contains all DAGs which are Markov equivalent to the Sachs network. The Sachs MEC has no directed edges, and thus is simply the undirected skeleton of the graph in Fig. 11. Starting from the Sachs MEC makes our results more interpretable in light of previous works and saves costly computation of the MEC. In practice, we would advise starting from the pooled PC MEC; based on the number of environments and observed density of changes, we would not expect this to orient any edges beyond the observational MEC.

Starting from the Sachs MEC, we apply the MSS using the KCI test, which appears to perform the best among plug-in estimators for MSS in our simulations. Since the feature distributions are heavily skewed, we preprocess them by taking their natural logarithm [42]. Among all DAGs in the Sachs MEC, the DAG with the uniquely minimal number of shifts exhibits approximately 8.9 shifts per pair of environments; this is relatively high but satisfies the assumption of sparse shifts. Violations to assumptions may lead to more shifts than expected.

D.2.3 Results and comparison to related works

The DAG which minimizes the MSS is the unique minimizer and is visualized in Fig. 11. An edge in black is oriented in the same direction as in the Sachs network, while an edge in red is oriented in the opposite direction. Overall, the majority of edges match the Sachs network. The edges which do not match, however, reflect ambiguities and flawed assumptions. We list each edge which does not match the Sachs network and discuss why this might be the case in light of existing work.

- PIP2 \rightarrow PIP3: As illustrated in Sachs et al. [44], these two proteins are actually cyclically related through bi-directed edges in the accepted “biologist’s view”. Indeed, PIP2 \rightarrow PIP3 was similarly recovered by an analysis of Ramsey and Andrews [42], detailed in their Figure 11.
- Mek \rightarrow Raf: that this edge does not match the Sachs network is discussed heavily by Mooij et al. [33] who point to it as a fundamental flaw of the Sachs network. The Mek \rightarrow Raf edge is indeed found by many other methods [5, 42, 33].
- The PKA, PKC, P38 triangle: although there is not a detailed discussion of these variables in other studies, there is strong ambiguity in the edge directions among approaches. Notably, Mooij et al. [33] similarly find strong evidence in their approach for the edge P38 \rightarrow PKC while Ramsey and Andrews [42] and Eaton and Murphy [5] find evidence for PKA \rightarrow PKC. However, all other approaches agree that the edge P38 \rightarrow PKA is incorrect. Although we do not explore further, it is worth noting that the 3rd minimal MSS DAG (not shown) is the same as the one shown, except it contains the presumed correct edge PKA \rightarrow P38.

As an additional note, we see in Fig. 11 that the edge Mek \rightarrow Erk is correctly recovered. Sachs et al. [44] similarly recover this well-known connection and point to it as strong evidence of success. In contrast, neither Eaton and Murphy [5], Ramsey and Andrews [42], nor Mooij et al. [33] recover the edge with their methods. Indeed, Ramsey and Andrews [42] specifically discuss how their approach incorrectly missed this edge, potentially the result of signal being lost when all the data is pooled. Pooled PC would face a similar issue, exacerbated by additional environments, while the pairwise comparisons by the MSS help to avoid this issue.