
TOWARDS CHARACTERIZING DOMAIN COUNTERFACTUALS FOR INVERTIBLE LATENT CAUSAL MODELS

A PREPRINT

Zeyu Zhou*, Ruqi Bai*, Sean Kulinski*, Murat Kocaoglu, David I. Inouye
 Elmore Family School of Electrical and Computer Engineering
 Purdue University
 {zhou1059, bai116, skulinsk, mkocaoglu, dinouye}@purdue.edu

November 16, 2023

ABSTRACT

Answering counterfactual queries has many important applications such as knowledge discovery and explainability, but is challenging when causal variables are unobserved and we only see a projection onto an observation space, for instance, image pixels. One approach is to recover the latent Structural Causal Model (SCM), but this typically needs unrealistic assumptions, such as linearity of the causal mechanisms. Another approach is to use naïve ML approximations, such as generative models, to generate counterfactual samples; however, these lack guarantees of accuracy. In this work, we strive to strike a balance between practicality and theoretical guarantees by focusing on a specific type of causal query called *domain counterfactuals*, which hypothesizes what a sample would have looked like if it had been generated in a different domain (or environment). Concretely, by only assuming invertibility, sparse domain interventions and access to observational data from different domains, we aim to improve domain counterfactual estimation both theoretically and practically with less restrictive assumptions. We define *domain counterfactually equivalent* models and prove necessary and sufficient properties for equivalent models that provide a tight characterization of the domain counterfactual equivalence classes. Building upon this result, we prove that every equivalence class contains a model where all intervened variables are at the end when topologically sorted by the causal DAG. This surprising result suggests that a model design that only allows intervention in the last k latent variables may improve model estimation for counterfactuals. We then test this model design on extensive simulated and image-based experiments which show the sparse canonical model indeed improves counterfactual estimation over baseline non-sparse models.

1 Introduction

Causal reasoning and machine learning, two fields which historically evolved disconnected from each other, have recently started to merge with several recent results leveraging the available causal knowledge to develop better ML solutions [Kusner et al., 2017, Moraffah et al., 2020, Nemirovsky et al., 2022]. One such setting is causal representation learning [Schölkopf et al., 2021, Brehmer et al., 2022], which aims to take data from a complex observed space (e.g., images) and learn the *latent* causal factors that generate the data. A common scenario is when we have access to multiple datasets from different domains, where from a causal perspective, each domain is generated via an *unknown* intervention on some domain-specific latent causal mechanisms. Schölkopf et al. [2021] introduce the Sparse Mechanism Hypothesis which hypothesizes that real-world domain interventions are sparse in the latent causal space, i.e., they only change a few causal mechanisms, even though the *observed* distribution might be wholly different [Schölkopf et al., 2021]. With this in mind, we focus on a specific causal query called a *domain counterfactual*, which hypothesizes: “What would this sample look like if it had been generated in a different domain (or environment)?” For example, given a patient’s medical imaging from Hospital A, what would it look like if it had been taken at Hospital

*Equal contribution. Listing order is random.

B? Answering this domain counterfactual query could have applications in knowledge discovery, explainability, and model robustness.

If one has access to or can recover the causal structure, it can be used to generate samples from counterfactual queries [Kocaoglu et al., 2018, Sauer and Geiger, 2021, Nemirovsky et al., 2022]. However, most of these existing methods assume that the causal variables are observed and thus are inapplicable in our setting where the causal variables are latent. For example, if the same patient was imaged at different hospitals, the differences in the images would be caused by complex factors such as different radiographer technicians, different calibrations on the imaging equipment, etc.² Recently, there have been results on learning the latent causal structure from observed data under well-defined assumptions [Xie et al., 2023, Yang et al., 2022, Huang et al., 2022, Liu et al., 2022a,b, Xie et al., 2022, Chen et al., 2022, Brehmer et al., 2022, Squires et al., 2023]. However, as seen in Table 1, many of these theoretic works make strong assumptions such as linearity, access to counterfactual pairs, etc., which often do not hold in real-world scenarios. In contrast to causal methods, the naïve ML approach is to simply train generative models such as VAEs to map between the two distributions without any causal assumptions or causal constraints (e.g., [Kulinski and Inouye, 2023]). However, there are many possible ways to map between two distributions and only some of them could produce the counterfactuals equivalent to the ground-truth SCM. Thus, the results can be highly dependent on the inductive bias of the model and could yield poor counterfactuals, as we see in our experiment with unconstrained VAEs.

In this paper, we aim for a practical yet theoretically grounded approach to estimating domain counterfactuals under minimal assumptions about the true model and available data. Concretely, *assuming only (1) continuous and bijective observation function and causal function, (2) sparsity of intervention, and (3) access to domain datasets (which represent unknown interventions), we seek to develop a method that improves domain counterfactual estimation both theoretically and practically over naïve methods.* Given assumption (3) that an algorithm solely has access to domain datasets (i.e., no counterfactual pairs or extra information is available), a learning algorithm can only ensure that the model matches the domain distributions—a property we call distribution equivalence. Thus, we seek the minimal modeling assumptions that are needed to improve domain counterfactuals while still being realistic. Towards this end, we first define the invertible latent domain causal model (ILD), which has a shared invertible observation function and a set of domain-specific invertible latent SCMs. Here we only assume that the latent causal model is invertible and has sparse latent intervention. Given this setup, we summarize our contributions as follows:

- C1** As a step towards theoretic understanding, we show that recovering the true ILD is unnecessary for estimating domain counterfactuals by proving a necessary and sufficient characterization of domain counterfactual equivalence. This allows us to construct counterfactually equivalent models and validate if two ILD models are counterfactually equivalent.
- C2** Given this characterization, we prove that *any* ILD can be written in a *canonical* form where only the last variables are intervened. Theoretically, this result simplifies proofs because a canonical ILD model has a simpler form. Practically, this result means that if we assume an intervention sparsity of k , an algorithm only needs to optimize over one sparsity structure for ILD models with sparsity k rather than searching over all $\binom{m}{k}$ sparsity structures for a general ILD models that have a sparsity of k and latent dimension m .
- C3** Leveraging the canonical ILD theory, we then prove that all ILD models can be split into disjoint equivalence classes with respect to their sparsity k . This result suggests that if the true ILD has an intervention sparsity of k^* , it is advantageous for an algorithm to restrict the sparsity to k^* because it will exclude ILDs with sparsity greater than k^* .
- C4** In light of these theoretic results, we propose an algorithm for estimating domain counterfactuals by searching over *canonical* ILD models (inspired by C2) while restricting intervention sparsity (inspired by C3). We validate our algorithm on both simulated and image-based experiments.

Notation. For a scalar n , $[n]$ denotes the set $\{1, \dots, n\}$. For a vector \mathbf{x} , the i -th entry is denoted by x_i and $[\mathbf{x}]_i$. Similarly, we denote $\mathbf{x}_{\leq i}$ and $[\mathbf{x}]_{\leq i}$ as the vector from the first entry to i -th entry. For a function $f : \mathbb{R}^m \rightarrow \mathbb{R}^m$, $[f(\mathbf{x})]_j \in \mathbb{R}$ refers to the j -th output of $f(\mathbf{x})$ and $[f(\mathbf{x})]_{\leq j} \in \mathbb{R}^j$ refers to the outputs from index 1 to j inclusive). We denote function equality between two functions $f : \mathcal{X} \rightarrow \mathcal{Y}$ and $f' : \mathcal{X} \rightarrow \mathcal{Y}$ as simply $f = f'$, which more formally can be stated as $\forall \mathbf{x} \in \mathcal{X}, f(\mathbf{x}) = f'(\mathbf{x})$. Similarly, $f \neq f'$ means $\exists \mathbf{x} \in \mathcal{X}, f(\mathbf{x}) \neq f'(\mathbf{x})$. We use \circ to denote function composition, e.g., $g(f(\mathbf{x})) = g \circ f(\mathbf{x})$ or simply $h = g \circ f$. We denote Id as the identity function.

²This is different from the selection bias issue which refers to different hospitals having different populations of patients.

Table 1: This table of related causal representation learning works, focuses mostly on works that study learning a *latent* SCM, shows that most prior works in this area aim for identifiability of the (latent) SCM, and thus require strong technical assumptions which may not hold in real-world scenarios (e.g., perfect single-node interventions for each variable).

	SCM	Observation Function	Other Assumptions	Observ. Function Identifiability	Characterization of Counterfactual Equivalence
Nasr-Esfahany et al. [2023]	Invertible observed	N/A (Does not study latent SCM)	1) Access to ground-truth DAG	N/A	Single mechanism counterfactuals under specific contexts
Brehmer et al. [2022]	Invertible latent	Invertible	1) Atomic stochastic hard interv. per node 2) Training set is counterfactuals pairs 3) SCM is faithful DAG	Mixing and elementwise transform	N/A - Counterfactuals are input
Squires et al. [2023]	Linear latent	Linear	1) Atomic hard interv.	Scaling	No
Liu et al. [2022a]	Linear latent	Non-linear	1) Significant causal weights variation	Mixing and scaling	No
Varici et al. [2023]	Latent non-linear	Linear	1) Atomic stochastic hard interv. 2) Each latent variable is intervened on	Mixing or scaling	No
Khemakhem et al. [2021]	Invertible observed (implicit)	Affine	1) Bivariate requirement for identifiability	Full (for bivariate case)	No
Ours	Invertible latent	Invertible	1) Access to domain labels 2) Sparse Mechanism Hypothesis	No	Domain counterfactual

2 Invertible Structural Causal Models

A structural causal model (SCM) considers m endogenous variables z_j and m exogenous noises ϵ_j , $j \in [m]$, where each variable is a deterministic function of its parents and independent exogenous noise. Formally, we denote $\mathbf{z} \in \mathbb{R}^m$ as a vector where its entries are endogenous variables and

$$z_j \triangleq \tilde{f}^{(j)}(\epsilon_j, \mathbf{z}_{\text{Pa}(j)}), \quad (1)$$

where $\tilde{f}^{(j)}(\cdot, \cdot) : \mathbb{R} \times \mathbb{R}^{|\text{Pa}(j)|} \rightarrow \mathbb{R}$, and $\mathbf{z}_{\text{Pa}(j)} \in \mathbb{R}^{|\text{Pa}(j)|}$ denotes the parents of z_j . The deterministic function $\tilde{f}^{(j)}$ is called the *causal mechanism* of the j -th variable. If all the exogenous noise ϵ_j can be recovered from all variables given all the causal mechanism. We say such SCM is *invertible*.

While it may at first seem like we are limiting ourselves by only considering invertible SCMs, the following lemma shows that this constraint does not reduce the expressivity of distributions.

Proposition 1 (Expressivity of Invertible SCM). *Invertible SCMs can model any continuous distribution if the exogenous noise distribution is continuous.*

The full proof is in Appendix A.1. The proof leverages the invertible Rosenblatt transformation [Rosenblatt, 1952, Melchers and Beck, 2018, Chapter B] that can transform any distribution to the uniform distribution or vice versa for its inverse.

2.1 Invertible autoregressive function represents invertible SCM

In this section, we show that all invertible SCM could be uniquely represented by an invertible autoregressive function $f \in \mathcal{F}_{IA}$. We further express the intervention between two SCM using f .

Definition 1 (Autoregressive Function). *A function $f : \mathbb{R}^m \rightarrow \mathbb{R}^m$ is autoregressive, denoted by $f \in \mathcal{F}_A$, if for all i , the i -th output can be written as a function of its corresponding input predecessors, i.e.,*

$$f \in \mathcal{F}_A \Leftrightarrow \forall j, \exists f^{(j)} \text{ s.t. } [f(\epsilon)]_j \triangleq f^{(j)}(\epsilon_{\leq j}), \text{ where } \epsilon \in \mathbb{R}^m. \quad (2)$$

Proposition 2 (Invertible SCM Representation). *An invertible SCM defined by a set of causal mechanisms $\{\tilde{f}^{(j)}(\epsilon_j, \mathbf{z}_{<j})\}_{j=1}^m$ could be uniquely represented by an invertible autoregressive function $f \in \mathcal{F}_{IA}$ and vice versa.*

See Appendix A.2 for proof. Given an Invertible SCM $\{\tilde{f}^{(j)}\}_{j=1}^m$, w.l.o.g., we assume all its variables are topological ordered, i.e., the parents always have smaller index than their children. Proposition 2 ensures there exists $f \in \mathcal{F}_{IA}$ that uniquely represents this SCM. We now define the intervention set between two SCMs.

Definition 2 (Intervention Set). *The intervention set between two SCMs \tilde{f}, \tilde{f}' is the set of indices where the causal mechanism changes, i.e.,*

$$\mathcal{I}(\tilde{f}, \tilde{f}') \triangleq \{j : \tilde{f}^{(j)} \neq \tilde{f}'^{(j)}\}. \quad (3)$$

The intervention could be arbitrary as long as it reserves the invertibility in the new SCM. Proposition 3 shows that the intervention set between f and f' could be directly computed.

Proposition 3 (Representation of Intervention Set Using Autoregressive Functions). *We denote the autoregressive representation f and f' for two SCMs \tilde{f}, \tilde{f}' respectively, then*

$$\mathcal{I}(\tilde{f}, \tilde{f}') = \mathcal{I}(f, f') = \{j : f^{-1} \neq f'^{-1}\}. \quad (4)$$

Please see the proof in Appendix A.3. Equipped with this proposition, we can use the autoregressive invertible function to analyse the latent causal models in the rest of the paper.

2.2 Invertible Latent Domain Causal Model (ILD)

Our goal is to propose a model designed to encompass latent causal mechanisms. In this section, we propose the invertible latent domain causal model to capture multiple latent SCMs that are emerged through intervention. The data generated by one SCM forms a *domain*.

Definition 3 (Invertible Latent Domain Causal Model). *An invertible latent domain causal model (ILD), denoted by (g, \mathcal{F}) , is a shared observation function g and a set of invertible SCM \mathcal{F} with N_d domains (i.e., $|\mathcal{F}| = N_d$) satisfying the following properties:*

1. [Latent Domain-Specific Invertible SCMs] $\mathcal{F} \triangleq \{f_d : \mathbb{R}^m \rightarrow \mathbb{R}^m \in \mathcal{F}_{IA}\}_{d=1}^{N_d}$.
2. [Invertible Observation Function] $g : \mathbb{R}^m \rightarrow \mathbb{R}^m \in \mathcal{F}_I$ is shared across domains.
3. [Continuous Exogenous Noise] $\epsilon \sim \mathcal{N}(0, I)$ without loss of generality.

In practice, invertibility can be relaxed using pseudo-invertible or approximately invertible functions, as seen with a VAE in Section 5.2. The autoregressive assumption ensures that the invertible function properly represents a DAG causal graph. While it assumes a fixed ordering of variables, we note that there is no such restriction on g and thus g can absorb any reordering of the variables to match the autoregressive structure of f_d . Thus, in view of the observation function g , this autoregressive assumption does not reduce expressivity of this model class. The shared observation function g will be critical for producing useful constraints on ILD but it does not inherently reduce expressivity as g could (in theory) just be the identity. As another example, suppose we have an ILD model (g, \mathcal{F}) where g is the identity. We could construct other models (g', \mathcal{F}') that produce the same distributions, where g' is an arbitrary invertible function and $f'_d = g'^{-1} \circ f_d$ (see Def. 5 for the formalization of distribution equivalence). Lastly, we assume the exogenous noise distribution is standard Gaussian, which is made mostly for convenience and can be made without loss of generality using the invertible Rosenblatt transformation [Rosenblatt, 1952, Melchers and Beck, 2018, Chapter B].

Definition 4 (Intervention of ILD). *We define the intervention set of an ILD \mathcal{F} as the union of any pairs of SCM in \mathcal{F} .*

$$\mathcal{I}(\mathcal{F}) \triangleq \bigcup_{f_d, f_{d'} \in \mathcal{F}} \mathcal{I}(f_d, f_{d'}) = \bigcup_{d \leq N_d} \mathcal{I}(f_1, f_d). \quad (5)$$

Definition 5 (Distribution Equivalence). *Two ILDs (g, \mathcal{F}) and (g', \mathcal{F}') are distributionally equivalent, denoted by $(g, \mathcal{F}) \simeq_D (g', \mathcal{F}')$, if the induced domain distributions are equal, i.e., for all $d \in [N_d]$,*

$$p_{\mathcal{N}}(f_d^{-1} \circ g^{-1}(\mathbf{x})) \Big|_{J_{f_d^{-1} \circ g^{-1}}(\mathbf{x})} = p_{\mathcal{N}}(f'_d{}^{-1} \circ g'^{-1}(\mathbf{x})) \Big|_{J_{f'_d{}^{-1} \circ g'^{-1}}(\mathbf{x})}. \quad (6)$$

The distributional equivalence defines a true equivalence relation because (6) has the properties of reflexivity, symmetry, and transitivity by the properties of the equality of measure.

Compared to prior SCM works that operate in the observed space, the non-identifiability of ILDs is extenuated because we are considering *latent* SCMs. A natural *necessary* (though certainly not sufficient) condition for estimation is that the ILD matches the observed distributions, which in practice is implemented as minimizing a distribution divergence with respect to the observed samples.

3 ILD Counterfactual

In this section, we aim to explore counterfactual equivalence ILDs consistent with the ground truth latent causal model. This section unfolds as follows:

1. We define the *domain counterfactual equivalence* and characterize the counterfactual equivalence using ILD.
2. We narrow our search to the canonical ILD, an exponentially smaller search space than original ILD, affirming the presence of equivalent canonical ILD that preserve intervention set size.
3. We show that all counterfactual equivalent model shares the same intervention set size, and misspecifying the intervention set size will lead to sub-optimal counterfactual. Our theory indicates a bias-variance tradeoff effect.

3.1 Domain Counterfactual

While distributional equivalence is a natural and common constraint for learning causal models, we now focus on our core contribution in the space of characterizing domain counterfactually equivalent models. We first provide a natural definition of this equivalence and prove that it is an equivalence relation. We proceed with briefly discussing the idea of a domain counterfactual.

The main idea of domain counterfactuals is that we can invert the causal model to retrieve the exogenous noise variables from the observed variables and domain label and then push these exogenous noise variables through the target domain SCM and the observation function. We formalize this for ILDs in the following definition.

Definition 6 (Domain Counterfactual). *Given an ILD (g, \mathcal{F}) , a counterfactual of \mathbf{x} from domain d projected into the target domain d' can be written as:*

$$\mathbf{x}_{d \rightarrow d'} \triangleq g \circ f_{d'} \circ f_d^{-1} \circ g^{-1}(\mathbf{x}), \text{ where } f_d, f_{d'} \in \mathcal{F}. \quad (7)$$

(7) can be interpreted as first projecting the sample into the latent space, i.e., $g^{-1}(\cdot)$, recovering the exogenous noise variables via $f_d^{-1}(\cdot)$, intervening by switching to the d' causal model $f_{d'}(\cdot)$ and then projecting back to the observed space via $g(\cdot)$. Given this notion of a domain counterfactual, we now provide an equivalence relation that will define which ILDs have the same domain counterfactuals.

Definition 7 (Domain Counterfactual Equivalence). *Two ILDs (g, \mathcal{F}) and (g', \mathcal{F}') are counterfactually equivalent, denoted by $(g, \mathcal{F}) \simeq_C (g', \mathcal{F}')$, if all counterfactuals are equal, i.e., for all d, d' , there holds*

$$g \circ f_{d'} \circ f_d^{-1} \circ g^{-1} = g' \circ f'_{d'} \circ f'_d{}^{-1} \circ g'^{-1}. \quad (8)$$

Lemma 1 (Equivalence relation of counterfactual equivalence). *Domain counterfactually equivalent, denoted by $(g, \mathcal{F}) \simeq_C (g', \mathcal{F}')$ is an equivalence relation, i.e., the relation satisfies reflexivity, symmetry, and transitivity.*

See Appendix B.1 for equivalence relation proof. While Definition 7 succinctly defines the equivalence classes of ILDs, it does not give much insight into the structure of the equivalence classes.

To fill this gap in characterizing these domain counterfactual equivalence classes, we now present one of our main theoretic results. Namely, we prove that an alternative property is both *necessary and sufficient* to be counterfactually equivalent.

Theorem 1 (Characterization of Counterfactual Equivalence). *Two ILDs are domain counterfactually equivalent, i.e., $(g, \mathcal{F}) \simeq_C (g', \mathcal{F}')$ if and only if:*

$$\exists h_1, h_2 \in \mathcal{F}_I \text{ s.t. } g' = g \circ h_1^{-1} \in \mathcal{F}_I \text{ and } f'_d = h_1 \circ f_d \circ h_2 \in \mathcal{F}_{IA}, \forall d. \quad (9)$$

See Appendix B.2 for proofs. Importantly, Theorem 1 can be used to *construct* counterfactually equivalent models and *verify* if two models are domain counterfactually equivalent (or determine they are not equivalent). More generally, this characterization exposes that the set of counterfactually equivalent models is actually very large. In fact, for any two invertible functions h_1 and h_2 that satisfy the implicit autoregressive constraint, i.e., for all $d, h_1 \circ f_d \circ h_2 \in \mathcal{F}_A$, we can construct a counterfactually equivalent model. In the next section, we demonstrate how to employ this novel characterization to establish a smaller set of domain counterfactual models, which we refer to as canonical domain counterfactuals.

3.2 Canonical ILD

We will now define the *canonical* ILD that allows each domain counterfactual equivalence class to be represented by a much smaller set of ILDs.

Definition 8 (Canonical Domain Counterfactual Model). *An ILD (g, \mathcal{F}) is in canonical counterfactual form, denoted by $(g, \mathcal{F}) \in \mathcal{C}$, if the following two properties hold:*

1. [Identity Domain] *The SCM corresponding to one domain is the identity, i.e., $\exists d, f_d = \text{Id}$. The identical domain w.l.o.g. can be the first domain, i.e., $f_1 = \text{Id}$.*
2. [Last Variables Intervened] *Only last variables are intervened, i.e., $j \in \mathcal{I}(\mathcal{F}), \forall j > m - |\mathcal{I}(\mathcal{F})|$.*

For canonical ILD, all the intervened nodes' descendant are also intervened nodes, and all the unintervened nodes follows standard Gaussian distribution. While this definition may seem quite restrictive, in our next key result, we show that (surprisingly) *any* ILD has an counterfactual and distributional equivalent *canonical* ILD.

Theorem 2 (Existence of Equivalent Canonical ILD). *Given an ILD (g, \mathcal{F}) , there exists a canonical ILD (g', \mathcal{F}') that is both counterfactually and distributionally equivalent to (g, \mathcal{F}) while maintaining the size of the intervention set, i.e.,*

$$\forall (g, \mathcal{F}), \exists (g', \mathcal{F}') \in \mathcal{C} \text{ s.t. } (g', \mathcal{F}') \simeq_{C,D} (g, \mathcal{F}), \text{ and } |\mathcal{I}(\mathcal{F})| = |\mathcal{I}(\mathcal{F}')|. \quad (10)$$

See Appendix B.3 for full proof. An example to elucidate Theorem 2 can be found at Example 1.

Corollary 3 (Relaxed Canonical Existence). *Given an ILD (g, \mathcal{F}) , there exists another ILD (g', \mathcal{F}') that only satisfies the last variable property (Property 2) of Definition 8 that is both counterfactually and distributionally equivalent to (g, \mathcal{F}) while maintaining the size of the intervention set.*

We omit the proof, which is done by applying the inverse of step 1 in the proof of Theorem 2.

The existence of canonical equivalent ILD indicates that we could only search for the canonical form to get a good domain counterfactual model which significantly simplifies the search space, which might help algorithm convergence.

3.3 Intervention Sparsity Analysis

In this section, assuming that the underlying ground truth intervention set size is k^* , We show that if we impose an intervention set size k where $k \neq k^*$ in our ILD, it is infeasible to find a counterfactual equivalent ILD.

Theorem 4 (Canonical ILD and Shared Intervention Sparsity). *Given an ILD (g, \mathcal{F}) , all canonical ILDs that are distributionally and counterfactually equivalent to (g, \mathcal{F}) have the same intervention set, i.e.,*

$$\mathcal{I}(\mathcal{F}) = \mathcal{I}(\mathcal{F}'), \quad \forall (g', \mathcal{F}') \in \mathcal{C} \text{ and } (g', \mathcal{F}') \simeq_{C,D} (g, \mathcal{F}). \quad (11)$$

Proof see Appendix B.4 for the proof. This theorem reveals that all counterfactually and distributionally equivalent canonical ILDs share the same intervention set size k . The significance of this proof is that we show that all ILDs counterfactually and distributionally equivalent to the underlying true model must share the same intervention set size k .

Corollary 5. *Given an ILD (g, \mathcal{F}) , all ILDs that are counterfactual and distributional equivalent to (g, \mathcal{F}) shares the same intervention set size. i.e.,*

$$\mathcal{I}(\mathcal{F}) = \mathcal{I}(\mathcal{F}'), \quad \forall (g', \mathcal{F}') \simeq_{C,D} (g, \mathcal{F}). \quad (12)$$

Corollary 5 uses canonical ILDs as bridges to connect every pairs of counterfactually and distributionally equivalent ILDs according to Lemma 1.

4 Related Work

Causal Representation Learning Causal representation learning is a rapidly developing field that aims to discover the underlying causal mechanisms that drive observed patterns in data and learn representations of data that are causally informative [Schölkopf et al., 2021]. This is in contrast to traditional representation learning, which does not consider the causal relationships between variables. An extensive review can be found in Schölkopf et al. [2021]. As this is a highly difficult task, most works make assumptions on the problem structure such as access to atomic hard interventions as well as the observation function being linear [Squires et al., 2023, Varici et al., 2023]. Other works such as [Brehmer et al., 2022, Ahuja et al., 2022, Von Kügelgen et al., 2021] assume a weakly-supervised setting where one can train on counterfactual pairs (x, \tilde{x}) during training. In our work, we aim to maximize the practicality of our assumptions while still maintaining our theoretical goal of equivalent domain counterfactuals (as seen in Table 1).

Counterfactual Generation Counterfactual examples are answers to hypothetical queries such as “What would the outcome have been if we were in setting B instead of A ?”. A line of works focus on the identifiability of counterfactual queries [Nasr-Esfahany et al., 2023, Shah et al., 2022]. For example, given knowledge of the ground-truth causal structure, Nasr-Esfahany et al. [2023] are able to recover the structural causal models up to equivalence. However, they do not consider the latent causal setting and they assume some prior knowledge of underlying causal structures such as the backdoor criterion. There is a weaker form of counterfactual generation which does not use causal reasoning but instead uses generative models to generate counterfactuals [Nemirovsky et al., 2022, Zhu et al., 2017, Choi et al., 2018, Zhou et al., 2023, Kulinski and Inouye, 2023]. These typically involve training a generative model which has a meaningful latent representation that can be intervened on to guide a counterfactual generation [Ilse et al., 2020]. As these works do not directly incorporate causal learning in their frameworks, we consider them out of scope for this paper. Another branch of works try to estimate causal effect without trying to learn the underlying causal structure, which typically assume all variables are observable [Louizos et al., 2017].

Causal Discovery Causal discovery focus on identifying the causal relationships from observational data. Peters et al. [2016], Heinze-Deml et al. [2018] achieve this via the invariant mechanism between certain variable and its direct causes. Some other works try to identify nonlinear ICA with access to auxiliary variables [Hyvarinen et al., 2019, Khemakhem et al., 2020]. Most of these works do not assume the latent SCM setting. Another branch of works aim at learning the latent causal structure [Xie et al., 2022]. However, they typically require strong assumption such as linearity.

5 Experiments

We have shown theoretically the benefit of our canonical ILD characterization and restriction of intervention sparsity. In this section, we empirically test whether our theory could guide us to design better models for producing domain counterfactuals while only having access to observational data \mathbf{x} and the corresponding domain label d . In our simulated experiment, under the scenario where all of our modeling assumptions hold, we try to answer the following questions: (1) When we know the ground truth sparsity, does sparse canonical ILD lead to better domain counterfactual generation over naïve ML approaches (dense models)? (2) What would happen if there is a mismatch of sparsity between the dataset and modeling and what is a good model design strategy in practice? After this simulated experiment, we perform experiments on image datasets to determine if sparse canonical models are still advantageous in this more realistic setting. In this case, we assume the latent causal model lies in a lower dimensional space than the observed space and thus we use autoencoders to approximate an observation function that is invertible on a lower-dimensional manifold.

5.1 Simulated Dataset

Experiment Setup To extensively address our questions against diverse causal mechanism settings, for each experiment, we generate 10 distinct ground truth ILDs. The ground truth latent SCM $f_d^* \in \mathcal{F}_{IA}$ takes the form $f_d^*(\epsilon) = F_d^* \epsilon + b_d^* \mathbb{1}_{\mathcal{I}}$ where $F_d^* = (I - L_d^*)^{-1}$, $L_d^* \in \mathbb{R}^{m \times m}$ is a domain-specific lower triangular matrix that satisfies the sparsity constraint, $b_d^* \in \mathbb{R}$ is a domain-specific bias, $\mathbb{1}_{\mathcal{I}}$ is an indicator vector where entries corresponding to the intervention set are 1, and L_d^* and b_d^* are randomly generated for each experiment. We use maximum likelihood estimation to train two ILDs (like training of a normalizing flow): *ILD-Relax-Can* which represents our relaxed canonical ILD form in Corollary 3 and a baseline model, *ILD-Dense*, which has no sparsity restrictions on its latent SCM. To evaluate the models, we compute the mean square error between the estimated counterfactual and ground truth counterfactual. More details on datasets and models, and illustrating figures of the models can be found in Appendix D.1.

Result To answer whether sparse canonical ILD provides any benefit in domain counterfactual generation, we first look at the simplest case where the latent causal structure of the dataset and our model exactly match. In Figure 1a, we notice that when the ground truth intervention set \mathcal{I}^* is $\{5, 6\}$ (i.e. the last two nodes), *ILD-Relax-Can* significantly outperforms *ILD-Dense*. Then we create a few harder and more practical tasks where the intervention set size is still 2 but not constrained to the last few nodes. Again, in Figure 1a, we observe that no matter which two nodes are intervened on, *ILD-Relax-Can* performs much better than the naïve ML approach *ILD-Dense*. This first checks that restricting model structure to the specific canonical form does not harm the optimization even though the ground truth structure is different. Furthermore, it validates the benefit of our model design for domain counterfactual generation. More results with different number of domains and latent dimensions can be found in Appendix D.2, which all show that *ILD-Relax-Can* consistently perform better than *ILD-Dense*. We also include an illustrating figure visualizing how *ILD-Relax-Can* achieves lower counterfactual error. We then transition to the more practical scenario where the

true sparsity $|\mathcal{I}^*|$ is unknown. In Figure 1b, at first glance, we observe a trend of the decrease in counterfactual error as we decrease $|\mathcal{I}|$. For the case where $|\mathcal{I}| \geq |\mathcal{I}^*|$ (i.e. when $|\mathcal{I}| = 2, 3, 4$), this aligns with our intuition that the smaller

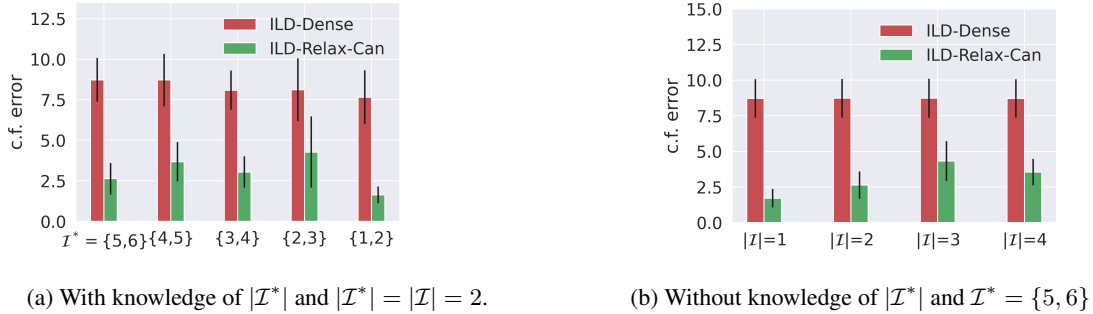


Figure 1: Simulated experiment results ($N_d = 3$) averaged over 10 runs with different ground truth SCMs and the error bar represents the standard error. (a) This shows *ILD-Relax-Can* is consistently better than *ILD-Dense* regardless of intervened nodes in the dataset. (b) Here we test varying $|\mathcal{I}|$ while holding \mathcal{I}^* fixed. The performance of *ILD-Relax-Can* approaches to that of *ILD-Dense* as we increase $|\mathcal{I}|$. An unexpected result is that *ILD-Relax-Can* performs best when $|\mathcal{I}| = 1$ and that results from a worse data fitting which is more carefully investigated in Appendix D.2.

search space of *ILD-Relax-Can* leads to a higher chance of finding model with low counterfactual error. For the case where $|\mathcal{I}| = 1$, we notice that it perform better than the canonical model that matches the true sparsity. Though this $|\mathcal{I}| = 1$ is a biased model (i.e., it cannot be distributionally equivalent to the true model per our Theorem 4), the reduction in variance seems to be enough to enable comparable or better counterfactuals on average. We further check the performance of the data fitting and see a significant drop in that of *ILD-Relax-Can*. This suggests that the performance in data fitting can be used as an indicator for whether we find the appropriate $|\mathcal{I}|$. More results about data fitting performance and experiments with different setups could be found in Appendix D.2, and they all lead to the conclusion that *ILD-Relax-Can* produces better counterfactuals than *ILD-Dense* even though we do not know $|\mathcal{I}^*|$.

5.2 Image-based Counterfactual Experiments

Here we seek to learn domain counterfactuals in the more realistic image regime. Following the manifold hypothesis [Gorban and Tyukin, 2018, Schölkopf et al., 2021], we assume that the causal interactions in this regime happen through lower-dimensional semantic latent factors as opposed to high-dimensional pixel-level interactions. To allow for learning of the lower dimensional latent space, we relax the invertibility constraint of our image-based ILD to only require pseudoinvertibility and test our models in this practical setting.

High-dim ILD Modeling We modify the ILD models from Section 5.1 to fit a VAE [Kingma and Welling, 2013] structure where the variational encoder, (g^+, \mathcal{F}^+) , first projects to the latent space via g^+ to produce the latent encoding z , which is then passed to two domain-specific latent causal models $f_{d,\mu}^+, f_{d,\sigma}^+$ which produce the parameters of posterior noise distribution. The decoder, (g, \mathcal{F}) , follows the typical ILD structure: $g \circ f_d$, where g and f_d can be viewed as pseudoinverse of $f_{d,\mu}^+$ and g^+ . A detailed description and diagram of the models can be found in Figure 16, but informally, these modified ILD models can be seen as training a VAE *per* domain with the restriction that each VAE shares parameters for its initial encoder and final decoder layers (i.e. g is shared). As an additional baseline, we compare against the naïve setup, which we call *ILD-Independent*, where each VAE has no shared parameters (i.e. a separate g is learned for each domain). These models were trained using the β -VAE framework [Higgins et al., 2017]. Further details can be found in the Appendix E.4. After training, we can perform domain counterfactuals as described in Equation (7).

Dataset We apply our methods to five image-based datasets: Rotated MNIST (RMNIST), Rotated FashionMNIST (RFMNIST)[Xiao et al., 2017], Colored Rotated MNIST (CRMNIST), 3D Shapes [Burgess and Kim, 2018] and Causal3DIdent [Von Kügelgen et al., 2021], which all have both domain information (e.g., the rotation of the MNIST digit) and class information (e.g., the digit number). For each dataset, we split the data into disjoint domains (e.g., each rotation in CRMNIST constitutes a different domain) and define class variables which are generated independently of domains (e.g., digit class in CRMNIST), to evaluate our model’s capability of generating domain counterfactuals. Specifically, for RMNIST, RFMNIST and 3D Shapes, all latent variables are independently generated, and for CRMNIST and Causal3DIdent, there is a more complicated causal graph containing the domain, class and other la-

latent variables. Further details on each dataset and (assumed) ground-truth latent causal graphs could be found in Appendix E.1 and Appendix E.3.

	CRMNIST				3D Shapes				Causal3DIdent			
	Comp.	Rev.	Eff.	Pre.	Comp.	Rev.	Eff.	Pre.	Comp.	Rev.	Eff.	Pre.
<i>ILD-Independent</i>	87.24	59.88	94.65	60.39	99.79	32.56	94.97	32.49	88.15	51.43	91.05	51.94
<i>ILD-Dense</i>	88.18	62.29	92.72	59.60	99.76	32.60	80.92	32.64	83.59	49.17	92.17	48.83
<i>ILD-Relax-Can</i>	92.10	85.74	94.48	72.95	99.85	79.84	96.72	64.99	86.00	79.73	84.15	79.73

Table 2: Quantitative result for **Composition** (Comp.), **Reversibility** (Rev.), **Preservation** (Pre.), and **Effectiveness** (Eff.), where higher is better. CRMNIST, 3D Shapes, Causal3DIdent are averaged 20, 5, 10 runs respectively. Best models are bold (within 1 standard deviation) and due to space constraint, expanded tables with additional datasets and standard deviation are in Appendix E.5.

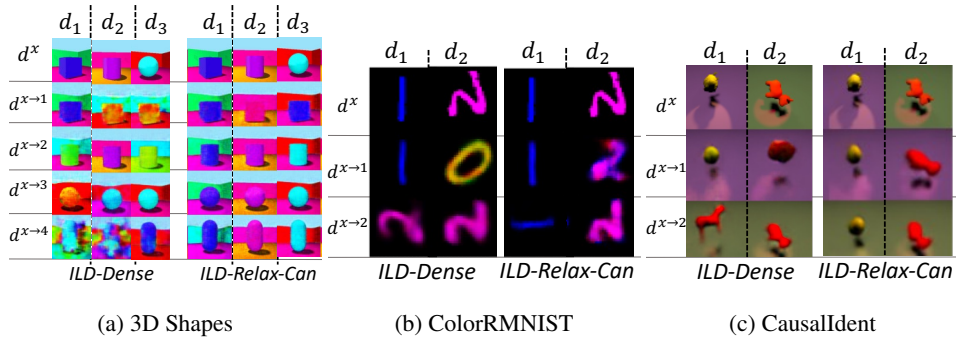


Figure 2: Domain counterfactuals with 3D Shapes, CRMNIST and CausalIdent. Expanded figures can be found in Appendix E.5 (a) For 3D Shapes, only the object shape should change with domain counterfactuals – the other latent factors such as the hue of object, floor, background, should not change. (b) For CRMNIST, as the domain changes, the rotation should change while the digit should not change. (c) For CausalIdent, as the domain changes, the color of the background should change while holding all else unchanged. *ILD-Relax-Can* clearly performs better than the baseline *ILD-Dense* in terms of preserving non-domain features while changing domains for all datasets.

Metrics Inspired by the work in Monteiro et al. [2023], we evaluate the image-based counterfactuals with latent SCMs via the following metrics: Effectiveness (whether the counterfactual truly changes the domain, $\mathbb{E}_{(x,d)}[\mathbb{1}_{h_{\text{domain}}(\hat{x}^{d \rightarrow d'})=d'}]$), Preservation (whether the domain counterfactual *only* changes domain-specific information, $\mathbb{E}_{(x,d)}[\mathbb{1}_{h_{\text{class}}(\hat{x}^{d \rightarrow d'})=y}]$), Composition (whether the counterfactual model is invertible, $\mathbb{E}_{(x,d)}[\mathbb{1}_{h_{\text{class}}(\hat{x}^{d \rightarrow d})=y}]$), and Reversibility (whether the counterfactual model is cycle-consistent, $\mathbb{E}_{(x,d)}[\mathbb{1}_{h_{\text{class}}(\hat{x}^{d \rightarrow d' \rightarrow d})=y}]$) where h_{domain} and h_{class} represents pretrained domain classifier and class classifier respectively. For example, in the case of CRMNIST, a model might be able to rotate the image but cannot preserve the digit class during rotation, which would be high in effectiveness but low in preservation score. Details on the computation of these metrics and causal interpretations can be found in Appendix E.2 and Appendix E.3 respectively.

Result Due to space constraint, we put all results with RMNIST and RFMNIST in Appendix E.5. In Figure 2 we can see examples of domain counterfactuals for both *ILD-Dense* and *ILD-Relax-Can*. We note that no latent information other than the domain label was seen during training, thus suggesting the intervention sparsity is what allowed the canonical models to preserve important non-domain-specific information such as class information when generating domain counterfactuals. In Table 2, we include quantitative results using our metrics, which shows *ILD-Relax-Can* having significantly better reversibility and preservation while maintaining similar levels of counterfactual effectiveness and composition than the non-sparse counterparts. In Appendix E.5, we further investigate our model’s sensitivity to the choice of sparsity by tracking how each metric change w.r.t. $|\mathcal{I}|$. We observe that reversibility and preservation tends to decrease while effectiveness tends to increase as we increase $|\mathcal{I}|$, which aligns with our findings here as *ILD-Dense* is equivalent to making \mathcal{I} contain all latent nodes. In summary, our results here indicate our theory-inspired model design leads to better domain counterfactual generation in the practical pseudo-invertible setting.

6 Conclusion

In this paper, we prove a necessary and sufficient characterization of domain counterfactual equivalence with more practical assumptions in comparison to existing works. Given this characterization, we show that any ILD model can be written in a canonical form, and we further prove all ILD models can be split into disjoint equivalence classes based on their sparsity. Then we empirically validate that our theory-inspired model design leads to better counterfactual estimation with extensive simulated and higher dimensional image-based experiments. We discuss limitation of our methods in Appendix F. We hope our theory could give inspiration to the design of practical algorithms and models which bridge the gap between causal reasoning and machine learning.

Acknowledgement

Z.Z., R.B., S.K., and D.I. acknowledge support from NSF (IIS-2212097), ARL (W911NF-2020-221), and ONR (N00014-23-C-1016). M.K. acknowledges support from NSF CAREER 2239375.

References

- Kartik Ahuja, Jason S Hartford, and Yoshua Bengio. Weakly supervised representation learning with sparse perturbations. *Advances in Neural Information Processing Systems*, 35:15516–15528, 2022.
- 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.
- Chris Burgess and Hyunjik Kim. 3d shapes dataset. <https://github.com/deepmind/3dshapes-dataset/>, 2018.
- Christopher P Burgess, Irina Higgins, Arka Pal, Loic Matthey, Nick Watters, Guillaume Desjardins, and Alexander Lerchner. Understanding disentangling in *beta*-vae. *arXiv preprint arXiv:1804.03599*, 2018.
- Zhengming Chen, Feng Xie, Jie Qiao, Zhifeng Hao, Kun Zhang, and Ruichu Cai. Identification of linear latent variable model with arbitrary distribution. In *Thirty-Sixth AAAI Conference on Artificial Intelligence, AAAI 2022, Thirty-Fourth Conference on Innovative Applications of Artificial Intelligence, IAAI 2022, The Twelfth Symposium on Educational Advances in Artificial Intelligence, EAAI 2022 Virtual Event, February 22 - March 1, 2022*, pages 6350–6357. AAAI Press, 2022. URL <https://ojs.aaai.org/index.php/AAAI/article/view/20585>.
- Yunjey Choi, Minje Choi, Munyoung Kim, Jung-Woo Ha, Sunghun Kim, and Jaegul Choo. Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8789–8797, 2018.
- Alexander N Gorban and Ivan Yu Tyukin. Blessing of dimensionality: mathematical foundations of the statistical physics of data. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 376(2118):20170237, 2018.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- Christina Heinze-Deml, Jonas Peters, and Nicolai Meinshausen. Invariant causal prediction for nonlinear models. *Journal of Causal Inference*, 6(2):20170016, 2018.
- Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. *beta*-vae: Learning basic visual concepts with a constrained variational framework. In *International conference on learning representations*, 2017.
- Biwei Huang, Charles Jia Han Low, Feng Xie, Clark Glymour, and Kun Zhang. Latent hierarchical causal structure discovery with rank constraints. *CoRR*, abs/2210.01798, 2022. doi: 10.48550/arXiv.2210.01798. URL <https://doi.org/10.48550/arXiv.2210.01798>.
- Aapo Hyvarinen, Hiroaki Sasaki, and Richard Turner. Nonlinear ica using auxiliary variables and generalized contrastive learning. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 859–868. PMLR, 2019.
- Maximilian Ilse, Jakub M Tomczak, Christos Louizos, and Max Welling. Diva: Domain invariant variational autoencoders. In *Medical Imaging with Deep Learning*, pages 322–348. PMLR, 2020.
- 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.

- Ilyes Khemakhem, Ricardo Monti, Robert Leech, and Aapo Hyvarinen. Causal autoregressive flows. In *International conference on artificial intelligence and statistics*, pages 3520–3528. PMLR, 2021.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- Murat Kocaoglu, Christopher Snyder, Alexandros G. Dimakis, and Sriram Vishwanath. Causalgan: Learning causal implicit generative models with adversarial training. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018. URL <https://openreview.net/forum?id=BJE-4xWOW>.
- Sean Kulinski and David I Inouye. Towards explaining distribution shifts. In *International Conference on Machine Learning*, pages 17931–17952. PMLR, 2023.
- Matt J Kusner, Joshua Loftus, Chris Russell, and Ricardo Silva. Counterfactual fairness. *Advances in neural information processing systems*, 30, 2017.
- Yuhang Liu, Zhen Zhang, Dong Gong, Mingming Gong, Biwei Huang, Anton van den Hengel, Kun Zhang, and Javen Qinfeng Shi. Identifying weight-variant latent causal models. *arXiv preprint arXiv:2208.14153*, 2022a.
- Yuhang Liu, Zhen Zhang, Dong Gong, Mingming Gong, Biwei Huang, Kun Zhang, and Javen Qinfeng Shi. Identifying latent causal content for multi-source domain adaptation. *CoRR*, abs/2208.14161, 2022b. doi: 10.48550/arXiv.2208.14161. URL <https://doi.org/10.48550/arXiv.2208.14161>.
- Christos Louizos, Uri Shalit, Joris M. Mooij, David A. Sontag, Richard S. Zemel, and Max Welling. Causal effect inference with deep latent-variable models. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 6446–6456, 2017. URL <https://proceedings.neurips.cc/paper/2017/hash/94b5bde6de88ddf9cde6748ad2523d1-Abstract.html>.
- Robert E Melchers and André T Beck. *Structural reliability analysis and prediction*. John wiley & sons, 2018.
- Miguel Aires Barros Monteiro, Fabio De Sousa Ribeiro, Nick Pawlowski, Daniel Coelho de Castro, and Ben Glocker. Measuring axiomatic soundness of counterfactual image models. *ArXiv*, abs/2303.01274, 2023. URL <https://api.semanticscholar.org/CorpusID:257280401>.
- Raha Moraffah, Mansooreh Karami, Ruocheng Guo, Adrienne Raglin, and Huan Liu. Causal interpretability for machine learning-problems, methods and evaluation. *ACM SIGKDD Explorations Newsletter*, 22(1):18–33, 2020.
- Arash Nasr-Esfahany, Mohammad Alizadeh, and Devavrat Shah. Counterfactual identifiability of bijective causal models. *arXiv preprint arXiv:2302.02228*, 2023.
- Daniel Nemirovsky, Nicolas Thiebaud, Ye Xu, and Abhishek Gupta. CounterGAN: Generating counterfactuals for real-time recourse and interpretability using residual GANs. In *Uncertainty in Artificial Intelligence*, pages 1488–1497. PMLR, 2022.
- Jonas Peters, Peter Bühlmann, and Nicolai Meinshausen. Causal inference by using invariant prediction: identification and confidence intervals. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 78(5):947–1012, 2016.
- Murray Rosenblatt. Remarks on a multivariate transformation. *The annals of mathematical statistics*, 23(3):470–472, 1952.
- Axel Sauer and Andreas Geiger. Counterfactual generative networks. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=BXewfAYMmJw>.
- 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.
- Abhin Shah, Raaz Dwivedi, Devavrat Shah, and Gregory W Wornell. On counterfactual inference with unobserved confounding. *arXiv preprint arXiv:2211.08209*, 2022.
- Chandler Squires, Anna Seigal, Salil S. Bhate, and Caroline Uhler. Linear causal disentanglement via interventions. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett, editors, *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pages 32540–32560. PMLR, 2023. URL <https://proceedings.mlr.press/v202/squires23a.html>.
- Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in neural information processing systems*, 30, 2017.

- Burak Varici, Emre Acarturk, Karthikeyan Shanmugam, Abhishek Kumar, and Ali Tajer. Score-based causal representation learning with interventions. *arXiv preprint arXiv:2301.08230*, 2023.
- 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.
- Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- Feng Xie, Biwei Huang, Zhengming Chen, Yangbo He, Zhi Geng, and Kun Zhang. Identification of linear non-gaussian latent hierarchical structure. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato, editors, *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pages 24370–24387. PMLR, 2022. URL <https://proceedings.mlr.press/v162/xie22a.html>.
- Feng Xie, Yan Zeng, Zhengming Chen, Yangbo He, Zhi Geng, and Kun Zhang. Causal discovery of 1-factor measurement models in linear latent variable models with arbitrary noise distributions. *Neurocomputing*, 526:48–61, 2023. ISSN 0925-2312. doi: <https://doi.org/10.1016/j.neucom.2023.01.034>. URL <https://www.sciencedirect.com/science/article/pii/S0925231223000449>.
- Yuin Yang, AmirEmad Ghassami, Mohamed S. Nafea, Negar Kiyavash, Kun Zhang, and Ilya Shpitser. Causal discovery in linear latent variable models subject to measurement error. *CoRR*, abs/2211.03984, 2022. doi: 10.48550/arXiv.2211.03984. URL <https://doi.org/10.48550/arXiv.2211.03984>.
- Zeyu Zhou, Sheikh Shams Azam, Christopher Brinton, and David I Inouye. Efficient federated domain translation. In *The Eleventh International Conference on Learning Representations, ICLR, 2023*.
- Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2223–2232, 2017.

Appendix

Table of Contents

A Auxiliary results and proofs of Section 2	13
A.1 Proof of Proposition 1	14
A.2 Proof of Proposition 2	14
A.3 Proof of Proposition 3	15
B Auxiliary results and proofs of Section 3	16
B.1 Proof of Lemma 1	16
B.2 Proof of Theorem 1	16
B.3 Proof of Theorem 2	18
B.4 Proof of Theorem 4	21
C Proofs of Lemmata	22
C.1 Miscellaneous Proofs	22
C.2 Proof of Invertible Composition Equivalence Lemma 3	22
C.3 Proof of swapping Lemma 4	23
D Simulated Experiment	25
D.1 Experiment Details	25
D.2 Additional Simulated Experiment Results	26
E Image Counterfactual Experiments	31
E.1 Dataset Descriptions	31
E.2 Metrics	31
E.3 Causal Interpretation of our experiments	32
E.4 Experiment Details	32
E.5 Additional Results	33
F Limitations	34

A Auxiliary results and proofs of Section 2

In this section, we prove the lemmas in Section 2 which capture important properties of the ILD Model. Before proving Proposition 1, we first introduce another lemma that is useful later in proving Proposition 2.

Lemma 2 (Invertible Upper Subfunctions). *The upper subfunctions of $f \in \mathcal{F}_{IA}$ are also invertible, i.e., $\bar{f}^{(j)}(\epsilon_{\leq j}) \triangleq [f(\epsilon_{\leq j}, \cdot)]_{\leq j}$ is an invertible function of $\epsilon_{\leq j}$.*

Proof. We will prove this by induction on k where $j = m - k$. For $k = 0$, it is trivial because $\bar{f}^{(\leq m)} \equiv f \in \mathcal{F}_I$. We will prove the inductive step by contradiction. Suppose $\bar{f}^{(\leq m-k)}$ is not invertible. This would mean it is not injective and/or not surjective.

If $\bar{f}^{(j)}$ is not injective, then $\exists \epsilon_{\leq j} \neq \epsilon'_{\leq j}$ such that $\bar{f}_{\leq j}(\epsilon_{\leq j}) = \bar{f}_{\leq j}(\epsilon'_{\leq j})$. We would then have for some $\epsilon_{> j}$ (e.g., all zeros):

$$\begin{aligned}
 & \bar{f}^{(\leq j+1)}(\epsilon_{\leq j}, \epsilon_{j+1}) \\
 &= [\bar{f}^{(\leq j)}(\epsilon_{\leq j}), [f(\epsilon_{\leq j}, \epsilon_{> j})]_{j+1}]^\top \\
 &= [\bar{f}^{(\leq j)}(\epsilon'_{\leq j}), [f(\epsilon_{\leq j}, \epsilon_{> j})]_{j+1}]^\top \\
 &= \bar{f}^{(\leq j+1)}(\epsilon'_{\leq j}, \epsilon_{j+1}),
 \end{aligned} \tag{13}$$

but this would contradict the fact that $\bar{f}_{\leq j+1}$ is invertible by the inductive hypothesis.

If $\bar{f}^{(\leq j)}$ is not surjective, then $\exists \mathbf{x}_{\leq j}$ such that $\forall \epsilon_{\leq j}, \bar{f}^{(\leq j)}(\epsilon_{\leq j}) \neq \mathbf{x}_{\leq j}$. We would then have that $\forall \epsilon_{\leq j}, \epsilon_{> j}$

$$\bar{f}^{(j+1)}(\epsilon_{\leq j}, \epsilon_{j+1}) = [\bar{f}^{(j)}(\epsilon_{\leq j}), [f(\epsilon_{> j})]_{j+1}]^\top \neq [\mathbf{x}_{\leq j}, x_{j+1}]^\top. \quad (14)$$

but this would contradict the fact inductive hypothesis that \bar{f}_{j+1} is surjective. Therefore, $\bar{f}^{(j)}$ must be invertible for all $j \in [m]$. \square

A.1 Proof of Proposition 1

The proof leverages the invertible Rosenblatt transformation [Rosenblatt, 1952, Melchers and Beck, 2018, Chapter B] that can transform any distribution to the uniform distribution or vice versa using its inverse. Given an ordering of a set of random variables, i.e., $\mathbf{x} = [x_1, x_2, \dots, x_m]^\top$, the Rosenblatt transformation is defined as follows:

$$\begin{aligned} u_1 &:= F_1(x_1) \\ u_2 &:= F_2(x_2|x_1) \\ u_3 &:= F_3(x_3|x_1, x_2) \\ &\vdots \\ u_m &:= F_m(x_m|x_1, x_2, \dots, x_{m-1}), \end{aligned} \quad (15)$$

where $F_j(x_j|\mathbf{x}_{< j})$ is the conditional CDF of x_j given \mathbf{x}_j , i.e., the CDF corresponding to the distribution $p(x_j|\mathbf{x}_{< j})$. Its inverse can be written as follows:

$$\begin{aligned} x_1 &= F_1^{-1}(u_1) \\ x_2 &= F_2^{-1}(u_2|x_1) \\ x_3 &= F_3^{-1}(u_3|x_1, x_2) \\ &\vdots \\ x_m &= F_m^{-1}(u_m|x_1, x_2, \dots, x_{m-1}), \end{aligned} \quad (16)$$

where $F_j^{-1}(u_j|\mathbf{x}_{< j})$ is the conditional inverse CDF corresponding to the conditional CDF $F_j(x_j|\mathbf{x}_{< j})$. Let $F_p(\mathbf{x})$ denote the Rosenblatt transformation for distribution p , and let $F_p^{-1}(\mathbf{u})$ denote its inverse as defined above. Assuming the random variables are continuous, the Rosenblatt transformation transforms the samples from any distribution to samples from the Uniform distribution (i.e., the push-forward of the Rosenblatt transformation is the uniform distribution and the pushforward of a uniform distribution through the inverse Rosenblatt is the distribution p).

Proof. Given any continuous target distribution p , we can construct an invertible SCM whose observed distribution is p . Specifically, if we let q denote the exogenous noise distribution, then the following invertible and autoregressive function f —which defines an invertible SCM via Proposition 2—can be used to match the SCM distribution to p :

$$f(\epsilon) = F_p \circ F_q^{-1}(\epsilon), \quad (17)$$

where F_q^{-1} maps to the uniform distribution and then F_p maps to the target distribution per the properties of the Rosenblatt transformation. The function is invertible since both functions are invertible. Additionally, both functions are autoregressive and thus the composition is autoregressive. Therefore, f represents a valid invertible SCM whose observed distribution is p . \square

A.2 Proof of Proposition 2

The unique representation is given by:

$$f(\epsilon) = \left[\tilde{f}^{(1)}(\epsilon_1), \tilde{f}^{(2)}(\epsilon_2, \underbrace{\tilde{f}^{(1)}(\epsilon_1)}_{\text{recover } z_1}), \tilde{f}^{(3)}(\epsilon_3, \underbrace{\tilde{f}^{(1)}(\epsilon_1), \tilde{f}^{(2)}(\epsilon_2, \tilde{f}^{(1)}(\epsilon_1))}_{\text{recover } \mathbf{z}_{\text{Pa}(j)}}), \dots \right]^\top, \quad (18)$$

where for all j ,

$$\tilde{f}^{(j)}(\epsilon_j, \mathbf{z}_{\text{Pa}(j)}) = [f(\underbrace{[f^{-1}(\mathbf{z}_{< j}, \cdot)]_{\text{Pa}(j)}}_{\text{recover } \epsilon_{< j} \text{ from } \mathbf{z}_{\text{Pa}(j)}}), \epsilon_j, \cdot]_j. \quad (19)$$

Proof. First, let's define a relaxed version of SCM where each z_j is a function of all $z_{<j}$ instead of $z_{\text{Pa}(j)}$, i.e. $z_j = \widehat{f}^{(j)}(\epsilon_j, z_{<j})$. Then we will prove the following statement:

$$f(\epsilon) = \left[\widehat{f}^{(1)}(\epsilon_1), \widehat{f}^{(2)}(\epsilon_2, \underbrace{\widehat{f}^{(1)}(\epsilon_1)}_{\text{recover } z_1}), \widehat{f}^{(3)}(\epsilon_3, \underbrace{\widehat{f}^{(1)}(\epsilon_1), \widehat{f}^{(2)}(\epsilon_2, \widehat{f}^{(1)}(\epsilon_1))}_{\text{recover } z_{<3}}), \dots \right]^\top, \quad (20)$$

$$\widehat{f}^{(j)}(\epsilon_j, z_{<j}) = [f(\underbrace{[f^{-1}(z_{<j}, \cdot)]_{<j}}_{\text{recover } \epsilon_{<j} \text{ from } z_{<j}}, \epsilon_j, \cdot)]_j. \quad (21)$$

We first prove one direction. Given an invertible SCM defined by its causal mechanisms $\{\widehat{f}^{(j)}(\epsilon_j, z_{<j})\}_{j=1}^m$, the observed variables are given recursively as:

$$z_j = \widehat{f}^{(j)}(\epsilon_j, z_{<j}). \quad (22)$$

We now define the corresponding f as in the lemma:

$$f(\epsilon) \triangleq \left[\widehat{f}^{(1)}(\epsilon_1), \widehat{f}^{(2)}(\epsilon_2, \underbrace{\widehat{f}^{(1)}(\epsilon_1)}_{\text{recover } z_1}), \widehat{f}^{(3)}(\epsilon_3, \underbrace{\widehat{f}^{(1)}(\epsilon_1), \widehat{f}^{(2)}(\epsilon_2, \widehat{f}^{(1)}(\epsilon_1))}_{\text{recover } z_{<3}}), \dots \right]^\top. \quad (23)$$

We need to prove that the observed variables are equivalent to the given SCM. Formally, we will prove by induction on $j \in [m]$ the hypothesis that $[f(\epsilon)]_j = \widehat{f}^{(j)}(\epsilon_j, z_{<j}) = z_j$, $\forall \epsilon \in \mathbb{R}^m$. The base case is trivial from the definition in (23), i.e., $\forall \epsilon \in \mathbb{R}^m$, $[f(\epsilon)]_1 = \widehat{f}^{(1)}(\epsilon_1) = z_1$. For the inductive step, we have the following:

$$[f(\epsilon)]_{j+1} = \widehat{f}^{(j+1)}(\epsilon_{j+1}, \underbrace{\widehat{f}^{(1)}(\epsilon_1)}_{z_1}, \underbrace{\widehat{f}^{(2)}(\epsilon_2, \widehat{f}^{(1)}(\epsilon_1))}_{z_2}, \dots) = \widehat{f}^{(j+1)}(\epsilon_{j+1}, z_{<j+1}) = z_{j+1} \quad (24)$$

where the first equals is by (23), the second is by the inductive hypothesis, and the last is by definition of the SCM.

Now we prove the other direction. Given an invertible autoregressive function $f \in \mathcal{F}_I \cap \mathcal{F}_A$, we define the following recursive set of mechanism functions:

$$\forall j, z_j \equiv \widehat{f}^{(j)}(\epsilon_j, z_{<j}) \triangleq [f([f^{-1}(z_{<j}, \cdot)]_{<j}, \epsilon_j, \cdot)]_j. \quad (25)$$

Again, we will prove that these functional forms are equivalent via induction on j for the hypothesis that $\widehat{f}^{(j)}(\epsilon_j, z_{<j}) = [f(\epsilon)]_j = z_j$. The base case is trivial based on (25):

$$\widehat{f}^{(1)}(\epsilon_1) = [f([f^{-1}(z_{<1}, \cdot)]_{<1}, \epsilon_1, \cdot)]_1 = [f(\epsilon_1, \cdot)]_1 = z_1 \quad (26)$$

For the inductive step, we use the definition of $\bar{f}_{<j}$ and its inverse from Lemma 2 and derive the final result:

$$\widehat{f}^{(j+1)}(\epsilon_{j+1}, z_{<j+1}) = [f([f^{-1}(z_{<j}, \cdot)]_{<j}, \epsilon_j, \cdot)]_j = [f(\bar{f}_{<j}^{-1}(z_{<j}), \epsilon_j, \cdot)]_j = [f(\epsilon_{<j}, \epsilon_j, \cdot)]_j = z_j. \quad (27)$$

□

A.3 Proof of Proposition 3

Proof. Step 1: Prove $\{j : [f^{-1}]_j \neq [f'^{-1}]_j\} \subset \mathcal{I}(\tilde{f}, \tilde{f}')$.

For all $j \in \{j : [f^{-1}]_j \neq [f'^{-1}]_j\}$, there exists some z , such that

$$[f^{-1}(z)]_j \neq [f'^{-1}(z)]_j, \quad (28)$$

given that f, f' are auto-regressive function, we conclude there exists some $(z_{<j}, z_j)$ such that

$$\epsilon_j = [f^{-1}(z_{<j}, z_j, \cdot)]_j \neq [f'^{-1}(z_{<j}, z_j, \cdot)]_j = \epsilon'_j. \quad (29)$$

we have, for ϵ_j, ϵ'_j and such $z_{<j}$ there holds

$$\begin{aligned} \widehat{f}^{(j)}(\epsilon_j, z_{<j}) &\stackrel{(29)}{=} z_j \\ &\stackrel{(29)}{=} \widehat{f}'^{(j)}(\epsilon'_j, z_{<j}) \\ &= [f'([f'^{-1}(z_{<j}, \cdot)]_{<j}, \epsilon'_j, \cdot)]_j \\ &\stackrel{(a)}{\neq} [f'([f'^{-1}(z_{<j}, \cdot)]_{<j}, \epsilon_j, \cdot)]_j \\ &= \widehat{f}'^{(j)}(\epsilon_j, z_{<j}). \end{aligned} \quad (30)$$

where (a) comes from the $f' \in \mathcal{F}_I$. Thus it implies $j \in \mathcal{I}(\tilde{f}, \tilde{f}')$.

Step 2: Prove $\mathcal{I}(\tilde{f}, \tilde{f}') \subset \{j : [f^{-1}]_j \neq [f'^{-1}]_j\}$.

For all $j \in \mathcal{I}(\tilde{f}, \tilde{f}')$, there exists some $(\epsilon_j, \mathbf{z}_{<j})$, such that

$$z_j \triangleq \widehat{f}^{(j)}(\epsilon_j, \mathbf{z}_{<j}) \neq \widehat{f}'^{(j)}(\epsilon_j, \mathbf{z}_{<j}) \triangleq z'_j, \quad (31)$$

Define

$$\mathbf{z}_{\leq j} \triangleq [\mathbf{z}_{<j}, z_j] \quad \text{and} \quad \mathbf{z}'_{\leq j} \triangleq [\mathbf{z}_{<j}, z'_j], \quad (32)$$

then we have

$$[f^{-1}(\mathbf{z}_{\leq j}, \cdot)]_j = \epsilon_j = [f'^{-1}(\mathbf{z}'_{\leq j}, \cdot)]_j, \quad (33)$$

given that $f, f' \in \mathcal{F}_I$, we conclude,

$$[f^{-1}(\mathbf{z}_{\leq j}, \cdot)]_j \neq [f'^{-1}(\mathbf{z}_{\leq j}, \cdot)]_j, \quad (34)$$

which implies $j \in \{j : [f^{-1}]_j \neq [f'^{-1}]_j\}$. \square

B Auxiliary results and proofs of Section 3

B.1 Proof of Lemma 1

Proof. We only need to prove that it satisfies reflexivity, symmetry, and transitivity.

1. Reflexivity - Letting $g' = g$ and $\mathcal{F}' = \mathcal{F}$ in the definition, it is trivial to see that $\forall d, d'$

$$g \circ f_{d'} \circ f_d^{-1} \circ g^{-1} = g \circ f_{d'} \circ f_d^{-1} \circ g^{-1},$$

and thus $(g, \mathcal{F}) \simeq_C (g', \mathcal{F}')$.

2. Symmetry - Similarly, it is trivial to see that $\forall d, d'$,

$$g \circ f_{d'} \circ f_d^{-1} \circ g^{-1} = g' \circ f'_{d'} \circ f'_d{}^{-1} \circ g'^{-1} \iff g' \circ f'_{d'} \circ f'_d{}^{-1} \circ g'^{-1} = g \circ f_{d'} \circ f_d^{-1} \circ g^{-1}, \quad (35)$$

and thus $(g, \mathcal{F}) \simeq_C (g', \mathcal{F}') \iff (g', \mathcal{F}') \simeq_C (g, \mathcal{F})$.

3. Transitivity - For (g, \mathcal{F}) , (g', \mathcal{F}') and (g'', \mathcal{F}'') , we can derive the transitive property by applying the property twice to the first two and the last two pairs $\forall d, d'$:

$$g \circ f_{d'} \circ f_d^{-1} \circ g^{-1} = g' \circ f'_{d'} \circ f'_d{}^{-1} \circ g'^{-1} = g'' \circ f''_{d'} \circ f''_d{}^{-1} \circ g''^{-1},$$

which means that $(g, \mathcal{F}) \simeq_C (g'', \mathcal{F}'')$. \square

B.2 Proof of Theorem 1

The proof of [Theorem 1](#) relies heavily on one the following key lemma that provides a necessary and sufficient condition for the composition of two invertible functions to be equal.

Lemma 3 (Invertible Composition Equivalence). *For two pairs of invertible functions (f_1, f_2) and (f'_1, f'_2) , the following two conditions are equivalent:*

1. *The compositions are equal:*

$$f_1 \circ f_2 = f'_1 \circ f'_2.$$

2. *There exists an intermediate invertible function h s.t.*

$$f'_1 = f_1 \circ h^{-1}, f'_2 = h \circ f_2. \quad (36)$$

See Appendix [subsection C.2](#) for the proof of this Lemma.

Proof of Theorem 1. The basic idea is to use repeated application of [Lemma 3](#) under the constraint that h_1 and h_2 must be shared across for all d and g and g^{-1} must be inverses of each other.

For one direction as in [Lemma 3](#), if (9) holds, it is nearly trivial to show (8), for all d, d' :

$$g' \circ f'_{d'} \circ f'_d{}^{-1} \circ g'^{-1} = (g \circ h_1^{-1}) \circ (h_1 \circ f_{d'} \circ h_2) \circ (h_2^{-1} \circ f_d^{-1} \circ h_1^{-1}) \circ (h_1 \circ g^{-1}) = g \circ f_{d'} \circ f_d^{-1} \circ g^{-1}.$$

To prove the other direction, let us define the following functions for a specific (d, d') (we will treat the case of all (d, d') afterwards): $f_1 \triangleq g^{-1}, f_2 \triangleq f_d^{-1}, f_3 \triangleq f_{d'}$, and $f_4 \triangleq g$ and similarly f'_1, f'_2, f'_3 , and f'_4 for the other side. Given these definitions, we can write the property as:

$$f_4 \circ f_3 \circ f_2 \circ f_1 = f'_4 \circ f'_3 \circ f'_2 \circ f'_1.$$

By recursively applying [Lemma 3](#) for each of the three function compositions, we arrive at the following fact that there exist h_1, h_2, h_3 , such that :

$$\begin{aligned} f'_1 &= h_1 \circ f_1 \text{ and } f'_4 \circ f'_3 \circ f'_2 = f_4 \circ f_3 \circ f_2 \circ h_1^{-1} \\ f'_2 &= h_2 \circ f_2 \circ h_1^{-1} \text{ and } f'_4 \circ f'_3 = f_4 \circ f_3 \circ h_2^{-1} \\ f'_3 &= h_3 \circ f_3 \circ h_2^{-1} \text{ and } f'_4 = f_4 \circ h_3^{-1}. \end{aligned}$$

By using the definitions of f_1, f_2 , etc., we can now derive the following:

$$\begin{aligned} g' &= g \circ h_3^{-1} \\ f'_{d'} &= h_3 \circ f_{d'} \circ h_2^{-1} \\ f'_d{}^{-1} &= h_2 \circ f_d^{-1} \circ h_1^{-1} \\ g'^{-1} &= h_1 \circ g^{-1}. \end{aligned}$$

We can connect the first and the last equality to derive that $h_3 = h_1$:

$$g'^{-1} = (g \circ h_3^{-1})^{-1} = h_3 \circ g^{-1} = h_1 \circ g^{-1}.$$

Thus, there are only two free functions. Specifically, for any fixed pair of (d, d') there exist $h_{1,d,d'} (\equiv h_{3,d,d'})$ and $h_{2,d,d'}$ such that

$$g' = g \circ h_{1,d,d'}^{-1}, f'_d = h_{1,d,d'} \circ f_d \circ h_{2,d,d'}^{-1}, \text{ and } f'_{d'} = h_{1,d,d'} \circ f_{d'} \circ h_{2,d,d'}^{-1}. \quad (37)$$

Finally, we tackle the case of all (d, d') by assuming that there could be unique functions $h_{1,d,d'}$ and $h_{2,d,d'}$ for all pairs of (d, d') and show that they are in fact equal. Because the condition holds for all pairs (d, d') , we know that for any particular (d, d') and (d'', d) , we have the following two set of equations based on the proof above. First, for counterfactual $d \rightarrow d'$, we have

$$g' \circ f'_{d'} \circ f'_d{}^{-1} \circ g'^{-1} = g \circ f_{d'} \circ f_d^{-1} \circ g^{-1}.$$

By applying (37), we get that there exist $h_{1,d,d'}, h_{2,d,d'}$ such that

$$\begin{aligned} g' &= g \circ h_{1,d,d'}^{-1} \\ f'_d &= h_{1,d,d'} \circ f_d \circ h_{2,d,d'}^{-1} \\ f'_{d'} &= h_{1,d,d'} \circ f_{d'} \circ h_{2,d,d'}^{-1}. \end{aligned}$$

Similarly, for (d'', d) , there exist $h_{1,d'',d}, h_{2,d'',d}$ such that

$$\begin{aligned} g' &= g \circ h_{1,d'',d}^{-1} \\ f'_{d''} &= h_{1,d'',d} \circ f_{d''} \circ h_{2,d'',d}^{-1} \\ f'_d &= h_{1,d'',d} \circ f_d \circ h_{2,d'',d}^{-1}. \end{aligned}$$

By equating the RHS for the g' equations above, we have

$$g \circ h_{1,d,d'}^{-1} = g \circ h_{1,d'',d}^{-1}.$$

Thus $h_{1,d,d'} = h_{1,d'',d}$. Using this fact and similarly by equating the RHS for the f'_d equations above, we can derive:

$$f'_d = h_{1,d,d'} \circ f_d \circ h_{2,d,d'}^{-1} = h_{1,d'',d} \circ f_d \circ h_{2,d'',d}^{-1} = h_{1,d,d'} \circ f_d \circ h_{2,d'',d}^{-1},$$

which shows that

$$\begin{aligned} h_{2,d,d'}^{-1} &= h_{2,d'',d}^{-1} \\ h_{2,d,d'} &= h_{2,d'',d}. \end{aligned}$$

By applying these facts to all possible triples of (d, d', d'') , we can conclude that for all d, d'

$$\begin{aligned} h_{1,d,d'} &= h_1 \\ h_{2,d,d'} &= h_2, \end{aligned}$$

i.e., these intermediate functions must be independent of d and d' . Finally, we can adjust notation so that for all d ,

$$\begin{aligned} f'_d &= \tilde{h}_1 \circ f_d \circ \tilde{h}_2 \triangleq h_1 \circ f_d \circ h_2 \\ g' &= g \circ \tilde{h}_1^{-1} \triangleq g \circ h_1^{-1}. \end{aligned}$$

□

B.3 Proof of Theorem 2

Lemma 4 (Swapping Lemma). *Given that the first canonical counterfactual property is satisfied, i.e., $f_1 = \text{Id}$, denote f' as SCM constructed by $f' = h_1 \circ f \circ h_2$, where $h_1 = h_2$ denote swapping the j -th feature with j' -th feature. Then there exists g' such that*

$$(g, \mathcal{F}) \simeq_C (g', \mathcal{F}'), \quad f'_1 = \text{Id}, \quad \mathcal{I}(\mathcal{F}') = (\mathcal{I}(\mathcal{F}) \setminus \{j\}) \cup \{j'\}.$$

if the following conditions hold

$$j \in \mathcal{I}(\mathcal{F}) \text{ and } \forall \tilde{j} : j < \tilde{j} \leq j', \quad \tilde{j} \notin \mathcal{I}(\mathcal{F}).$$

Built upon swapping Lemma, we move to our main result on the existence of equivalent Canonical ILD. See [subsection C.3](#) for proofs.

Proof of Theorem 2. At high level the proof is organized in the following two steps.

(Step 1) we use [Theorem 1](#) to construct an equivalent counterfactual $(g^{(0)}, \mathcal{F}^{(0)}) \simeq_C (g, \mathcal{F})$ by choosing two invertible functions $h_1 = f_1^{-1}$ and $h_2 = \text{Id}$. In this way, [Theorem 1](#) implies

$$\begin{aligned} f_1^{(0)} &= h_1 \circ f_1 \circ h_2 = f_1^{-1} \circ f_1 \circ \text{Id} = \text{Id} \\ \forall d > 1, \quad f_d^{(0)} &= h_1 \circ f_d \circ h_2 = f_1^{-1} \circ f_d \circ \text{Id} = f_1^{-1} \circ f_d, \quad \text{and} \quad g^{(0)} = g \circ h_1^{-1} = g \circ f_1. \end{aligned}$$

Equipped with $(g^{(0)}, \mathcal{F}^{(0)})$, we can show that part I of [Def. 8](#) is satisfied, i.e., $f_1^{(0)} = \text{Id}$. Choosing $h_2 = \text{Id}$, we could prove **(Step 1)** could guarantee the distribution equivalence.

(Step 2) we can further construct a series of equivalent counterfactuals iteratively to gradually satisfy part II of [Def. 8](#). Specifically, in **(Step 2)**, we recursively construct, for all iteration $k \in \{1, 2, \dots, k^{\text{last}}\}$. We use superscript $f^{(k)}$, $g^{(k)}$ and $\mathcal{F}^{(k)}$ to denote the k -th swap operation. Specifically, in swap k , we have for all $f_d \in \mathcal{F}$,

$$f_d^{(k)} \triangleq h_{j(k) \leftrightarrow j'(k)} \circ f_d^{(k-1)} \circ h_{j(k) \leftrightarrow j'(k)},$$

and

$$g^{(k)} \triangleq g^{(k-1)} \circ h_{j(k) \leftrightarrow j'(k)}^{-1} = g^{(k-1)} \circ h_{j(k) \leftrightarrow j'(k)},$$

where $h_{j(k) \leftrightarrow j'(k)}$ denotes swapping the $j(k)$ -th and $j'(k)$ -th feature values, i.e.,

$$h_{j \leftrightarrow j'}(\mathbf{x}) \triangleq [x_1, x_2, \dots, x_{j-1}, x_{j'}, x_{j+1}, \dots, x_{j'-1}, x_j, x_{j'+1}, \dots, x_m]^T, \quad (38)$$

and further define

$$j'(k) \triangleq \max \left\{ j, j \notin \mathcal{I} \left(\mathcal{F}^{(k)} \right) \right\}, \text{ and } j(k) \triangleq \max \left\{ j < j'(k), j \in \mathcal{I} \left(\mathcal{F}^{(k)} \right) \right\}. \quad (39)$$

In high level, at each iteration, we seek the largest index $j'(k)$ which does not lies in the previous intervention set $\mathcal{I} \left(\mathcal{F}^{(k)} \right)$, and swap it with the largest index $j(k)$ which is smaller than $j'(k)$. We terminate at k when $\{j < j'(k), j \in \mathcal{I} \left(\mathcal{F}^{(k)} \right)\} = \emptyset$.

By the definition of $j'(k), j(k)$ in (39), we can show that **1)** for each swap step k , there holds

$$j(k) \in \mathcal{I} \left(\mathcal{F}^{(k)} \right), \text{ and } \forall \tilde{j} : j(k) < \tilde{j} \leq j'(k), \tilde{j} \notin \mathcal{I} \left(\mathcal{F}^{(k)} \right), \quad (40)$$

which implies **Lemma 4** can be applied to ensure the counterfactual equivalence at each step.

2) When meeting the stopping criterion at step k^{last} , i.e.,

$$\left\{ j < j'(k^{\text{last}}), j \in \mathcal{I} \left(\mathcal{F}^{(k^{\text{last}-1)} \right) \right\} = \emptyset, \quad (41)$$

there holds

$$\forall j \in \mathcal{I} \left(\mathcal{F}^{(k^{\text{last}-1)} \right), \quad j > m - \left| \mathcal{I} \left(\mathcal{F}^{(k^{\text{last}-1)} \right) \right|,$$

i.e., $(g^{(k^{\text{last}-1)}}, \mathcal{F}^{(k^{\text{last}-1)})$ is in canonical form. Chaining **1)** and **2)**, we conclude

$$\exists (g', \mathcal{F}') \triangleq (g^{(k^{\text{last}-1)}, \mathcal{F}^{(k^{\text{last}-1)})} \in \mathcal{C} \text{ s.t. } (g', \mathcal{F}') \simeq_C (g, \mathcal{F}).$$

Note that $g^{(k)} \circ f_d^{(k)} = g^{(k-1)} \circ f_d^{(k-1)} \circ h_{j(k) \leftrightarrow j'(k)}$, and linear operator $h_{j(k) \leftrightarrow j'(k)}$ is orthogonal, then iteratively, we conclude $(g', \mathcal{F}') \simeq_D (g, \mathcal{F})$.

To prove **1)**, observe in (39), $j(k)$ is the largest index in the intervention set which is smaller than $j'(k)$. This simply implies (40).

To prove **2)**, suppose when meeting the stopping criterion at step k^{last} , there holds

$$\exists j \in \mathcal{I} \left(\mathcal{F}^{(k^{\text{last}-1)} \right) \text{ such that } j \leq m - \left| \mathcal{I} \left(\mathcal{F}^{(k^{\text{last}-1)} \right) \right|. \quad (42)$$

It implies that

$$\exists \hat{j} \notin \mathcal{I} \left(\mathcal{F}^{(k^{\text{last}-1)} \right) \text{ and } \hat{j} \in \left\{ m - \left| \mathcal{I} \left(\mathcal{F}^{(k^{\text{last}-1)} \right) \right| + 1, \dots, m \right\}.$$

Then we can choose $j'(k) = \hat{j}$, implying $j \in \left\{ j < j'(k), j \in \mathcal{I} \left(\mathcal{F}^{(k^{\text{last}-1)} \right) \right\} \neq \emptyset$, contradict to (41). This concludes the proof of part I in **Theorem 2**.

It remains to prove that the construction of $f^{(0)}$ in the **step 1** does not change the intervention set.

1) For any $j \notin \mathcal{I}(\mathcal{F})$, for any pairs d, d' , we have $[f_d^{-1}]_j = [f_{d'}^{-1}]_j$, based on the construction of $f^{(0)}$, we have

$$\left[f_d^{(0)-1} \right]_j = [f_d^{-1} \circ f_1]_j = [f_{d'}^{-1} \circ f_1]_j = \left[f_{d'}^{(0)-1} \right]_j \quad (43)$$

thus, $\mathcal{I}(f_d^{(0)}, f_{d'}^{(0)}) \subset \mathcal{I}(f_d, f_{d'})$.

2) For any $j \in \mathcal{I}(\mathcal{F})$, there exists d, d' and z , such that $[f_d^{-1}(z)]_j \neq [f_{d'}^{-1}(z)]_j$. Note that f_1 is a bijective function, there exists z' such that $z = f_1(z')$, we have

$$\begin{aligned} & [f_d^{-1}(z)]_j \neq [f_{d'}^{-1}(z)]_j \\ \Leftrightarrow & [f_d^{-1}(f_1(z'))]_j \neq [f_{d'}^{-1}(f_1(z'))]_j \\ \Leftrightarrow & \left[f_d^{(0)-1}(z') \right]_j \neq \left[f_{d'}^{(0)-1}(z') \right]_j \\ \Leftrightarrow & j \in \mathcal{I} \left(f_d^{(0)}, f_{d'}^{(0)} \right) \end{aligned}$$

thus $\mathcal{I} \left(f_d^{(0)}, f_{d'}^{(0)} \right) \supset \mathcal{I}(f_d, f_{d'})$. Combining **1)** and **2)**, we have $\mathcal{I} \left(f_d^{(0)}, f_{d'}^{(0)} \right) = \mathcal{I}(f_d, f_{d'})$. This show that the construction of **step 1** does not change the intervention set, combining the fact in **step 1**, we iteratively used swapping **Lemma 4**, and swapping **Lemma 4** does not change the intervention set size, i.e., $\mathcal{I}(\mathcal{F}') = (\mathcal{I}(\mathcal{F}) \setminus \{j\}) \cup \{j'\}$, we conclude that $\left| \mathcal{I} \left(f_d^{(0)}, f_{d'}^{(0)} \right) \right| = \left| \mathcal{I}(f_d, f_{d'}) \right|$. This completes the proof. \square

To help understanding, we design a simple linear ILD model to explain the theorem procedure.

Example 1. Suppose we have a 4-dimensional ILD model (g, \mathcal{F}) containing 2 domains, where

$$f_1 \triangleq \begin{bmatrix} 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 \end{bmatrix}, f_2 \triangleq \begin{bmatrix} 1 & 0 & 0 & 0 \\ 2 & 2 & 0 & 0 \\ 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 \end{bmatrix}, g \text{ invertible.}$$

Following the proof of **Theorem 2**, we have Following **Step 1** in the proof of **Theorem 2**, we have $h_1 = f_1^{-1}$,

$$f_1^{(0)} = f_1^{-1} \circ f_1, f_2^{(0)} = f_1^{-1} \circ f_2 \\ g^{(0)} = g \circ f_1,$$

$$f_1^{(0)} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}, f_2^{(0)} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 1 & 2 & 0 & 0 \\ -1 & -1 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix},$$

$$g^{(0)} = g \circ \begin{bmatrix} 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 \end{bmatrix}$$

Notice that $\mathcal{I}(f^{(0)}) = \{2, 3\}$. Following **Step 2** in the proof of **Theorem 2**, we first swap $j = 3$ and $j' = 4$,

$$h_{3 \leftrightarrow 4} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{bmatrix}, g^{(1)} = g \circ \begin{bmatrix} 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 0 \end{bmatrix}$$

We have $f^{(2)} \triangleq h_{3 \leftrightarrow 4} \circ f^{(1)} \circ h_{3 \leftrightarrow 4}$

$$f_1^{(2)} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}, f_2^{(2)} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 1 & 2 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ -1 & -1 & 0 & 1 \end{bmatrix}.$$

Notice that $\mathcal{I}(f^{(1)}) = \{2, 4\}$. Following **Step 2** in the proof of **Theorem 2**, we first swap $j = 2$ and $j' = 3$,

$$h_{2 \leftrightarrow 3} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}, g^{(2)} = g \circ \begin{bmatrix} 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 0 \end{bmatrix}$$

We have $f^{(3)} \triangleq h_{2 \leftrightarrow 3} \circ f^{(2)} \circ h_{2 \leftrightarrow 3}$

$$f_1^{(2)} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}, f_2^{(2)} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 2 & 0 \\ -1 & 0 & -1 & 1 \end{bmatrix}.$$

$$g^{(2)} = g \circ \begin{bmatrix} 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 0 \end{bmatrix}$$

Notice that $(g^{(2)}, f^{(2)})$ is the canonical form. They are counterfactually equivalent to each other by checking definition.

B.4 Proof of Theorem 4

Lemma 5. *If $f : \mathbb{R}^m \rightarrow \mathbb{R}^m \in \mathcal{F}_{IA}$, then $[f(\mathbf{x})]_k$ must be a non-constant function of x_k .*

Proof. We prove this by contradiction. Suppose k is the first index that $[f(\mathbf{x})]_k = \tilde{f}(x_1, \dots, x_{k-1})$. Since k is the smallest index, $[f(\mathbf{x})]_{<k}$ is uniquely determined by $[\mathbf{x}]_{\leq k}$. The remaining $m - k$ dimension outputs could not be bijective to $m - k + 1$ inputs. \square

Proof of Theorem 4. In the proof, we denote F as a non constant function without specifying the expression.

Step 1: Characterization of counterfactual equivalence for canonical forms. Theorem 1 states that there exists $h_1, h_2 \in \mathcal{F}_I$, such that for all d ,

$$f'_d = h_1 \circ f_d \circ h_2. \quad (44)$$

Furthermore, by the definition of canonical form (Def. 8), we have

$$f'_1 = \text{Id}, f_1 = \text{Id}. \quad (45)$$

Plugging this into (44), we have

$$\text{Id} = h_1 \circ \text{Id} \circ h_2.$$

Thus,

$$h_1^{-1} = h_2 \triangleq h.$$

Plugging this into (44), for all d , we have

$$f_d'^{-1} = h^{-1} \circ f_d^{-1} \circ h. \quad (46)$$

Step 2: Counterfactual equivalence between canonical forms maintain the intervention set. The goal of this step is to prove that h is a bridge satisfying the following property: for any $i \notin \mathcal{I}(f'_1, f'_d)$, for all x , there exists a unique j , such that $[h^{-1}(\mathbf{x})]_i$ only depends on x_j . In addition, we can prove such j satisfies $j \notin \mathcal{I}(f_1, f_d)$.

We start with writing the i -th output of $f_d'^{-1}(\mathbf{x})$ as the following

$$[f_d'^{-1}(\mathbf{x})]_i \stackrel{(46)}{=} [h^{-1}(f_d^{-1}(h(\mathbf{x})))]_i \quad (47)$$

$$= [h^{-1}([f_d^{-1}(h(\mathbf{x}))]_1, [f_d^{-1}(h(\mathbf{x}))]_2, \dots, [f_d^{-1}(h(\mathbf{x}))]_m)]_i \quad (48)$$

$$\stackrel{(a)}{=} \left[h^{-1} \left(\tilde{f}_{d,1}^{-1}([h(\mathbf{x}))]_1), \tilde{f}_{d,2}^{-1}([h(\mathbf{x}))]_1, [h(\mathbf{x}))]_2), \dots, \tilde{f}_{d,m}^{-1}([h(\mathbf{x}))]_1, \dots, [h(\mathbf{x}))]_m) \right]_i, \quad (49)$$

where in step (a), we used autoregressiveness of f_d^{-1} , and $\tilde{f}_{d,k}^{-1}$ is defined as a function from \mathbb{R}^k to \mathbb{R} . According to Lemma 5, $\tilde{f}_{d,k}^{-1}(\mathbf{x})$ is a non-constant function of x_k .

Step 2.1: We show i and j must be one-to-one mapping of h . We proof this by contradiction.

Suppose h^{-1} maps more than one index to i -th index, w.l.o.g, we could assume j_1 and j_2 . That is, $[h^{-1}(\mathbf{u})]_i$ depends on u_{j_1} and u_{j_2} . Take $\mathbf{u} = f_d^{-1}(h(\mathbf{x}))$, then we have

$$[f_d'^{-1}(\mathbf{x})]_i = F \left(\tilde{f}_{d,j_1}^{-1}([h(\mathbf{x}))]_1, \dots, [h(\mathbf{x}))]_{j_1}), \tilde{f}_{d,j_2}^{-1}([h(\mathbf{x}))]_1, \dots, [h(\mathbf{x}))]_{j_2}) \right) \quad (50)$$

Due to that $f_d^{-1} \in \mathcal{F}_{IA}$, from Lemma 5, we have

$$\tilde{f}_{d,j_1}^{-1}([h(\mathbf{x}))]_1, \dots, [h(\mathbf{x}))]_{j_1}) = F([h(\mathbf{x}))]_{j_1}, \cdot) \quad (51)$$

$$\tilde{f}_{d,j_2}^{-1}([h(\mathbf{x}))]_1, \dots, [h(\mathbf{x}))]_{j_2}) = F([h(\mathbf{x}))]_{j_2}, \cdot). \quad (52)$$

Plug (51), (52) into (50), we have

$$[f_d'^{-1}(\mathbf{x})]_i = F([h(\mathbf{x}))]_{j_1}, [h(\mathbf{x}))]_{j_2}, \cdot) \quad (53)$$

Given that $h \in \mathcal{F}_{IA}$, we conclude $([h(\mathbf{x}))]_{j_1}, [h(\mathbf{x}))]_{j_2})$ depend at least two distinct indices. That is, there exists i_1, i_2 such that

$$([h(\mathbf{x}))]_{j_1}, [h(\mathbf{x}))]_{j_2}) = F(x_{i_1}, x_{i_2}). \quad (54)$$

That implies $[f_d'^{-1}(\mathbf{x})]_i$ is a nontrivial function of (x_{i_1}, x_{i_2}) . This leads to the contradiction that $i \notin \mathcal{I}(f'_1, f'_d)$, where for all \mathbf{x} , there holds

$$[f_d'^{-1}(\mathbf{x})]_i = x_i \quad (55)$$

Step 2.2: We show such j is not in the intervention set between f_1 and f_d . We prove this by contradiction as well.

Step 2.1 implies

$$[f_d'^{-1}(\mathbf{x})]_i = F\left(\tilde{f}_{d,j}^{-1}([h(\mathbf{x})]_1, \dots, [h(\mathbf{x})]_j)\right) = F'([h(\mathbf{x})]_j, \cdot), \quad (56)$$

Suppose $j \in \mathcal{I}(f_1, f_d)$, then $f_d^{-1}([h(\mathbf{x})]_j)$ Recall that $f_d^{-1} \in \mathcal{F}_A$,

then $[f_d^{-1}(h(\mathbf{x}))]_j$ must be a non-constant function of $[h(\mathbf{x})]_j$ and $[h(\mathbf{x})]_{j'}$ for some $j' < j$, i.e.,

$$[f_d^{-1}(h(\mathbf{x}))]_j = \tilde{f}_{d,j}^{-1}([h(\mathbf{x})]_1, \dots, [h(\mathbf{x})]_j) = F([h(\mathbf{x})]_{j'}, [h(\mathbf{x})]_j, \cdot). \quad (57)$$

Similarly, we know that $[h(\mathbf{x})]_{j'}$ and $[h(\mathbf{x})]_j$ must be nontrivial functions of x_{i_3} and x_{i_4} , which $i_3 \neq i_4$. However, we know $[f_d'^{-1}(\mathbf{x})]_i$ is a function of x_i exclusively, which leads to contradiction. This shows that the number of non-intervened node in f_d' must not be greater than that in f_d , i.e.,

$$\mathcal{I}(f_1', f_d') \geq \mathcal{I}(f_1, f_d), \forall d. \quad (58)$$

We further notice that the symmetric relationship between f_d and f_d' , we could also have

$$\mathcal{I}(f_1', f_d') \leq \mathcal{I}(f_1, f_d), \forall d. \quad (59)$$

Union among on d , we have

$$\mathcal{I}(f') = \mathcal{I}(f). \quad (60)$$

□

C Proofs of Lemmata

C.1 Miscellaneous Proofs

Lemma 6 (Invertible function rewrite). *Given any two invertible functions $f : \mathcal{X} \rightarrow \mathcal{Y}$ and $f' : \mathcal{X} \rightarrow \mathcal{Y}$, f' can be decomposed into the composition of f and another invertible function. Specifically, f' can be decomposed in the following two ways:*

$$f' \equiv f \circ h_{\mathcal{X}} \quad (61)$$

$$f' \equiv h_{\mathcal{Y}} \circ f, \quad (62)$$

where $h_{\mathcal{X}} \triangleq f^{-1} \circ f' : \mathcal{X} \rightarrow \mathcal{X}$ and $h_{\mathcal{Y}} \triangleq f' \circ f^{-1} : \mathcal{Y} \rightarrow \mathcal{Y}$ are both invertible functions.

Proof of Lemma 6. The proof is straightforward. We first note that $h_{\mathcal{X}}$ and $h_{\mathcal{Y}}$ are invertible because they are compositions of invertible functions. Then, we have that:

$$f \circ h_{\mathcal{X}} = f \circ f^{-1} \circ f' = f' \quad (63)$$

$$h_{\mathcal{Y}} \circ f = f' \circ f^{-1} \circ f = f'. \quad (64)$$

□

C.2 Proof of Invertible Composition Equivalence Lemma 3

Proof of Lemma 3. For notational simplicity in this proof, we will let $g \triangleq f_1$, $f \triangleq f_2$, $g' \triangleq f_1'$ and $f' \triangleq f_2'$ —note that g and f are just arbitrary invertible functions in this proof. Furthermore, without loss of generality, we will prove for the property that there exists h such that $g' = g \circ h$, $f' = h^{-1} \circ f$, which is equivalent to that there exists h , such that $h : g' = g \circ h^{-1}$, $f' = h \circ f$. Thus, in the new notation, we are seeking to prove that $g \circ f = g' \circ f'$ iff there exists h such that

$$g' = g \circ h, f' = h^{-1} \circ f. \quad (65)$$

For the first direction, if there exists h satisfying $g' = g \circ h$, $f' = h^{-1} \circ f$, then it is easy to show that $g \circ f = g' \circ f'$:

$$g' \circ f' = g \circ h \circ h^{-1} \circ f = g \circ f. \quad (66)$$

For the other direction, we will prove by contradiction. First, using [Lemma 6](#), we can first rewrite g' and f' using the two uniquely determined invertible functions h_1 and h_2 :

$$g' = g \circ h_1 \tag{67}$$

$$f' = h_2 \circ f. \tag{68}$$

Now, suppose that $g \circ f = g' \circ f'$ but there does not exist h such that $g' = g \circ h, f' = h^{-1} \circ f$. By the first assumption and the facts above, we can derive the following:

$$g \circ f = g' \circ f' = g \circ h_1 \circ h_2 \circ f. \tag{69}$$

Thus $f = h_1 \circ h_2 \circ f$, which gives

$$h_1^{-1} = h_2. \tag{70}$$

From the second assumption, i.e., there does not exist h such that $g' = g \circ h, f' = h^{-1} \circ f$. Thus for all h such that $g' = g \circ h$, it holds

$$f' \neq h^{-1} \circ f. \tag{71}$$

Thus for specific $h_1^{-1} = h_2$, it also holds

$$f' \neq h_1^{-1} \circ f. \tag{72}$$

Then from [\(68\)](#), we have

$$h_2 \circ f \neq h_1^{-1} \circ f, \tag{73}$$

which gives us

$$h_1^{-1} \neq h_2. \tag{74}$$

But this leads to a direct contradiction of [\(70\)](#). Therefore, $g \circ f = g' \circ f'$, iff there exists h such that

$$g' = g \circ h, f' = h^{-1} \circ f.$$

□

C.3 Proof of swapping [Lemma 4](#)

Before proving the swapping Lemma, we introduce a Lemma which is useful for proving [Lemma 4](#).

Lemma 7. For an ILD with $f_1 = \text{Id}$, $\mathcal{I}(f_d, f_1) = \left\{ j : [f_d]_j \neq [f_1]_j \right\}$.

Proof. First, we show that for any $j \in \mathcal{I}(f_1, f_d)$, there exists \mathbf{x} , such that

$$[f_1(\mathbf{x})]_j \neq [f_d(\mathbf{x})]_j. \tag{75}$$

From [Proposition 3](#), we know that for any j , there exists \mathbf{x} , such that $[f_1^{-1}(\mathbf{x})]_j \neq [f_d^{-1}(\mathbf{x})]_j$. We define $\mathbf{x}' \triangleq f_d^{-1}(\mathbf{x})$. Thus we have $[\mathbf{x}]_j \neq [\mathbf{x}']_j$, Then we have

$$\begin{aligned} [f_1(\mathbf{x}')]_j &= [\mathbf{x}']_j = x'_j \\ [f_d(\mathbf{x}')]_j &= [\mathbf{x}]_j = x_j. \end{aligned}$$

Thus there exists \mathbf{x}' , $[f_1(\mathbf{x}')]_j \neq [f_d(\mathbf{x}')]_j$. On the other direction, we show that for all $j \notin \mathcal{I}(f_1, f_d)$, $[f_1]_j = [f_d]_j$. It is equal to show that if there exist \mathbf{x} such that $[f_1(\mathbf{x})]_j \neq [f_d(\mathbf{x})]_j$, then $j \in \mathcal{I}(f_1, f_d)$. We define $\mathbf{x}' \triangleq f_d(\mathbf{x})$, thus $[\mathbf{x}]_j \neq [\mathbf{x}']_j$, we have

$$\begin{aligned} [f_1^{-1}(\mathbf{x}')]_j &= [\mathbf{x}']_j = x'_j \\ [f_d^{-1}(\mathbf{x}')]_j &= [\mathbf{x}]_j = x_j. \end{aligned}$$

Thus there exists \mathbf{x}' , $[f_1^{-1}(\mathbf{x}')]_j \neq [f_d^{-1}(\mathbf{x}')]_j$.

□

Proof of Lemma 4. For function $f_d \in \mathcal{F}_A$, we denote $f_{d,j}(\mathbf{x}_{\leq j}) \triangleq [f_d(\mathbf{x})]_j$. First, note that because j' is not intervened, i.e., $[f_d]_{j'} = [f_1]_{j'}$, then we can derive that it's corresponding conditional function is independent of all but the j' -th value:

$$[f_d(\mathbf{x}_{\leq j'}, \cdot)]_{j'} = [f_1(\mathbf{x}_{\leq j'}, \cdot)]_{j'} = x_{j'}. \quad (76)$$

For the new model, we choose the invertible functions as swapping the j -th and j' -th feature values, i.e.,

$$h_1(\mathbf{x}) \triangleq [x_1, x_2, \dots, x_{j-1}, x_{j'}, x_{j+1}, \dots, x_{j'-1}, x_j, x_{j'+1}, \dots, x_m]^T \quad (77)$$

and similarly for h_2 , i.e., $h_2 \triangleq h_1$. Because h_1 and h_2 are invertible, we know that the new model will be in the same counterfactual equivalence class by **Theorem 1**. Construct $g' \triangleq g \circ h_1^{-1}$, and then for all d ,

$$\begin{aligned} f'_d(x) &= h_1 \circ f_d \circ h_2(x) \\ &= h_1 \circ f_d([x_1, x_2, \dots, x_{j-1}, x_{j'}, x_{j+1}, \dots, x_{j'-1}, x_j, x_{j'+1}, \dots, x_m]^T) \\ &= h_1 \circ f_d([y_1, y_2, \dots, y_{j-1}, y_j, y_{j+1}, \dots, y_{j'-1}, y_{j'}, y_{j'+1}, \dots, y_m]^T) \\ &= h_1 \circ [f_{d,i}(\mathbf{y}_{\leq i})]_{i=1}^m \\ &= \left[f_{d,1}(\mathbf{y}_1), \dots, f_{d,j-1}(\mathbf{y}_{\leq j-1}), f_{d,j'}(\mathbf{y}_{\leq j'}), f_{d,j+1}(\mathbf{y}_{\leq j+1}), \dots, \right. \\ &\quad \left. f_{d,j'-1}(\mathbf{y}_{\leq j'-1}), f_{d,j}(\mathbf{y}_{\leq j}), f_{d,j'+1}(\mathbf{y}_{\leq j'+1}), \dots, f_{d,m}(\mathbf{y}_{\leq m}) \right], \end{aligned}$$

where we define $\mathbf{y} \triangleq h_2^{-1}(\mathbf{x})$.

We now need to check that the first canonical counterfactual property still holds.

$$f'_1 = h_1 \circ f_1 \circ h_2 = h_1 \circ \text{Id} \circ h_2 = h_1 \circ h_2 = \text{Id}, \quad (78)$$

where the last equals is because swap operations are self-invertible.

We move to check that the autoregressive property still holds for other domain SCMs.

1) For the j -th feature, we have that:

$$[f'_d(\mathbf{x})]_j = f_{d,j'}(\mathbf{y}_{\leq j'}) = f_{d,j'}(x_1, \dots, x_{j-1}, x_{j'}, x_{j+1}, \dots, x_{j'-1}, x_j) = x_j$$

where the last equals is because the $f_{d,j'}(\mathbf{y}_{\leq j'}) = y_{j'} = x_j$. This clearly satisfies the autoregressive property as $[f'_d]_j$ only depends on x_j .

2) For the j' -th feature, we have that:

$$[f'_d(\mathbf{x})]_{j'} = f_{d,j}(\mathbf{y}_{\leq j}) = f_{d,j}(x_1, \dots, x_{j-1}, x_{j'})$$

where again this satisfies the autoregressive property because all input indices are less than j' because $j < j'$. Now we handle the cases for other variables. If $\tilde{j} < j$, then we have the following:

$$[f'_d]_{\tilde{j}} = [h_1 \circ f_d \circ h_2]_{\tilde{j}} = [f_d \circ h_2]_{\tilde{j}} = f_{d,\tilde{j}}([h_2(\mathbf{x})]_{\leq \tilde{j}}) = f_{d,\tilde{j}}(x_1, \dots, x_{\tilde{j}}) \quad (79)$$

3) Similarly if $j < \tilde{j} < j'$:

$$[f'_d]_{\tilde{j}} = f_{d,\tilde{j}}(x_1, \dots, x_{j-1}, x_{j'}, x_{j+1}, \dots, x_{\tilde{j}}) = x_{\tilde{j}}, \quad (80)$$

where we use the fact that there are no intervening nodes in between j and j' .

4) Finally, for $\tilde{j} > j'$, we have:

$$[f'_d]_{\tilde{j}} = f_{d,\tilde{j}}(x_1, \dots, x_{j-1}, x_{j'}, x_{j+1}, \dots, x_{j'-1}, x_j, x_{j'+1}, \dots, x_{\tilde{j}}), \quad (81)$$

which is still autoregressive because $\tilde{j} > j'$ and $\tilde{j} > j$. Thus, the new f'_d is autoregressive and is thus a valid model.

It remains to prove that $\mathcal{I}(\mathcal{F}') = (\mathcal{I}(\mathcal{F}) \setminus \{j\}) \cup \{j'\}$.

1) When $k < j$, we have for all d ,

$$[f'_d]_k = f_{d,k}(\mathbf{y}_{\leq k}) = f_{d,k}(\mathbf{x}_{\leq k}) = [f_d]_k,$$

then for all $k \in \mathcal{I}(\mathcal{F})$, there exists d_0 , such that

$$[f'_{d_0}]_k = [f'_{d_0}]_k \neq [f_1]_k = [f_1^{-1}]_k.$$

Thus, $k \in \mathcal{I}(\mathcal{F}')$.

If $k \notin \mathcal{I}(\mathcal{F})$, we have for all d ,

$$[f'_d{}^{-1}]_k = [f'_d]_k = [f_1]_k = [f_1{}^{-1}]_k.$$

Thus $k \notin \mathcal{I}(\mathcal{F}')$.

2) When $j \leq k < j'$, we have $\forall d, [f'_d(\mathbf{x})]_k = x_k \Rightarrow [f'_d{}^{-1}(\mathbf{x})]_k = x_k$. Thus we have $\forall d, [f'_d{}^{-1}]_k = [f_1{}^{-1}]_k$, which means for all $j \leq k < j'$, $k \notin \mathcal{I}(\mathcal{F}')$.

3) When $k = j'$, we have $\forall d, [f'_d]_{j'} = f_{d,j}(x_1, \dots, x_{j-1}, x_{j'})$. Furthermore, since, $j \in \mathcal{I}(\mathcal{F})$, we have $\exists d_0, [f_{d_0}]_j \neq [f_1]_j$ by Lemma 7. Thus $[f'_{d_0}]_{j'} = [f_{d_0}]_j \neq [f_1]_j = [f_1]_{j'} \Rightarrow j' \in \mathcal{I}(\mathcal{F}')$ also by Lemma 7.

4) When $k > j'$, if $k \in \mathcal{I}(\mathcal{F})$, $\exists d_1, d_2, [f_{d_1}]_k \neq [f_{d_2}]_k$, Chaining with (81), we have $[f'_{d_1}]_k \neq [f'_{d_2}]_k$. Thus, $k \in \mathcal{I}(\mathcal{F}')$ by Lemma 7. Similarly, if $k \notin \mathcal{I}(\mathcal{F})$, then $k \notin \mathcal{I}(\mathcal{F}')$. To summarize, $\mathcal{I}(\mathcal{F}') = (\mathcal{I}(\mathcal{F}) \setminus \{j\}) \cup \{j'\}$. \square

D Simulated Experiment

D.1 Experiment Details

Dataset The ground truth latent SCM $f_d^* \in \mathcal{F}_{IA}$ takes the form $f_d^*(\epsilon) = F_d^* \epsilon + b_d^* \mathbb{1}_{\mathcal{I}}$ where $F_d^* = (I - L_d^*)^{-1}$, $L_d^* \in \mathbb{R}^{m \times m}$ is domain-specific lower triangular matrix that satisfies sparsity constraint, $b_d^* \in \mathbb{R}$ is a domain-specific bias and $\mathbb{1}_{\mathcal{I}}$ is an indicator vector where any entries corresponding to the intervention set are 1. To be specific, $[L_d^*]_{i,j} \sim \mathcal{N}(0, 1)$ and $b_d^* \sim \text{Uniform}(-2\sqrt{m/|\mathcal{I}|}, 2\sqrt{m/|\mathcal{I}|})$. The observation function takes the form $g^*(\mathbf{x}) = G^* \text{LeakyReLU}(\mathbf{x})$ where $G^* \in \mathbb{R}^{m \times m}$ and the slope of LeakyReLU is 0.5. To allow for similar scaling across problem settings, we set the determinant of G^* to be 1 and standardize the intermediate output of the LeakyReLU. The generated F_d^*, b_d^*, G^* all vary with random seeds and all experiments are repeated for 10 different seeds. We generate 100,000 samples from each domain for the training set and 1,000 samples from each domain in the validation and test set.

Model We test with two ILD models: *ILD-Relax-Can* which represents the relaxed canonical ILD form from Cor. 3 and a baseline model, *ILD-Dense* which has no sparsity restrictions on its latent SCM. To be specific, the latent SCM of *ILD-Dense* could be any model in \mathcal{F}_{IA} . We use \mathcal{I} and \mathcal{I}^* to represent the intervention set of the model and dataset, respectively. We note that for *ILD-Dense*, \mathcal{I} contains all nodes and for *ILD-Relax-Can*, \mathcal{I} contains only the last few nodes. Both models follow a similar structure as the ground truth. To be specific, the latent SCM takes the form $f_d(\epsilon) = F_d \epsilon + \mathbf{b}_d$ where $F_d = (I - L_d)^{-1} S_d$, $L_d \in \mathbb{R}^{m \times m}$, $S_d \in \mathbb{R}^{m \times m}$, and $\mathbf{b}_d \in \mathbb{R}^m$. The observation takes the form $g(\mathbf{x}) = G \text{LeakyReLU}(\mathbf{x}) + \mathbf{b}$ where $G \in \mathbb{R}^{m \times m}$, $\mathbf{b} \in \mathbb{R}^m$, and the slope of LeakyReLU is 0.5. In Figure 3a and Figure 3b, we add an illustration of the latent SCM for *ILD-Dense* and *ILD-Relax-Can* respectively. We emphasize a few main differences between the dataset and models here: (1) *ILD-Relax-Can*, \mathcal{I} only contains the last few nodes while for the dataset, \mathcal{I}^* could contain any node we specify. We note that *ILD-Dense* is equivalent to a *ILD-Relax-Can* with all nodes in its intervention set. (2) There is no constraint on the determinant of G and standardization in $g(\mathbf{x})$. (3) The bias added to all dimensions in the ground truth model is the same scalar value, but the bias in the model is allowed to vary for each axis. (4) In the model, g is allowed a learnable bias.

Algorithm In the simulated experiments, our algorithm only tries to fit the observed distribution for all models. As all models are strictly invertible, we fit the distribution via maximum likelihood estimation (MLE). To be specific, the objective is as below

$$\max_{g, f_1, \dots, f_{N_d}} \mathbb{E}_d[p(\mathbf{x}|d)] \quad (82)$$

where $p(\mathbf{x}|d) = p_{\mathcal{N}}(f_d^{-1} \circ g^{-1}(\mathbf{x})) |J_{f_d^{-1} \circ g^{-1}}(\mathbf{x})|$.

Metric To evaluate the models, we compute the mean square error between the estimated counterfactual and ground truth counterfactual, i.e. $\text{Error} = \frac{2}{N_d(N_d-1)} \sum_{d' \neq d} \sum_d \|g^* \circ f_{d'}^* \circ (f_d^*)^{-1} \circ (g^*)^{-1}(\mathbf{x}_d) - g \circ f_{d'} \circ f_d^{-1} \circ g^{-1}(\mathbf{x}_d)\|^2$. As in practice, we can only check data fitting instead of counterfactual estimation, and we report the counterfactual error computed with the test dataset when the likelihood computed with the validation set is highest.

Training details We use Adam optimizer for both f and g with a lr = 0.001, $\beta_1 = 0.5$, $\beta_2 = 0.999$, and a batch size of is 500. We run all experiments for 50,000 iterations and compute validation likelihood and test counterfactual error every 100 steps. f is randomly initialized. Regarding g , G is initialized as an identity matrix and \mathbf{b} is initialized as 0.

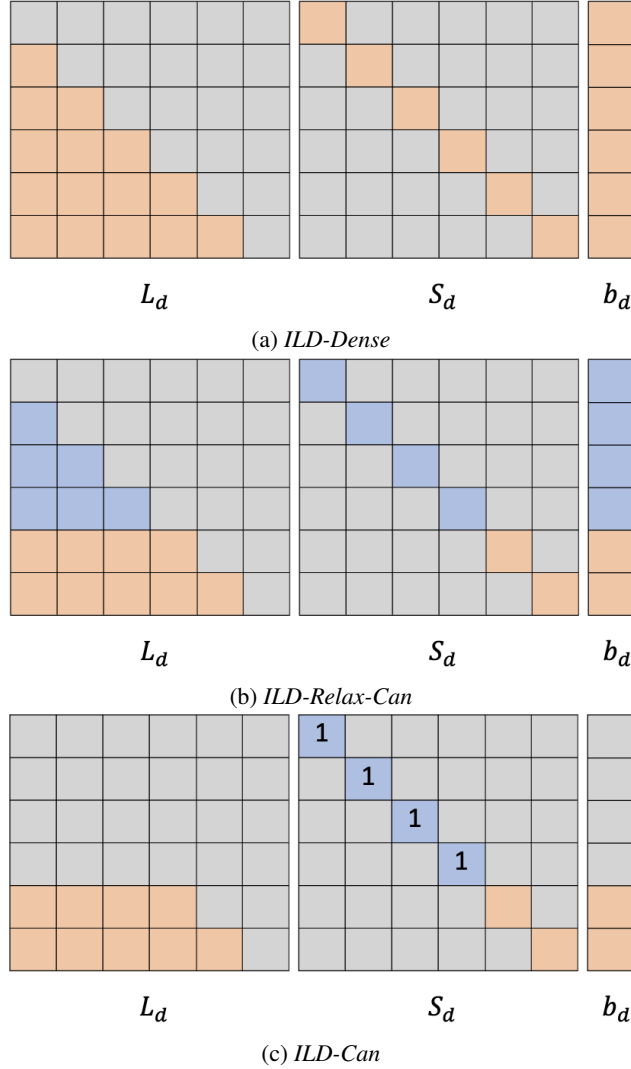


Figure 3: An illustration of the matrices/vector used to create f_d across the three ILD models when $m = 6$ and $|\mathcal{I}| = 2$. These are used such that $f_d(\epsilon) = F_d \epsilon + b_d$ where $F_d = (I - L_d)^{-1} S_d$. The grey elements are 0, the orange elements are parameters that are different for different domains, and the blue elements are parameters shared across domains. We specify the value if it is a fixed number other than 0. Note that we don't implement *ILD-Can* in our experiments. We include it here only for illustration of our theory.

D.2 Additional Simulated Experiment Results

For better organization here, we split our experiment into three cases as introduced below. The first two cases point to the question: given the fact that we use the correct sparsity, does sparse canonical form model designing provide benefits in generating domain counterfactuals? The third case investigates the more practical scenario where we don't have any knowledge of the ground truth sparsity and we explore what would be a better model design practice in this case.

Case 0: Exact match between dataset and models In this section, we investigate the performance of *ILD-Dense* and *ILD-Relax-Can* while assuming that the ground truth intervention set only contains the last few nodes and we choose the correct size of the intervention set.

To understand how the true intervention set affects the gap between *ILD-Dense* and *ILD-Relax-Can*, we varied the size of the ground truth intervention. In Figure 4, we observe that the performance gap tends to be largest when the true intervention set is the most sparse and the performance of *ILD-Relax-Can* approaches to the performance

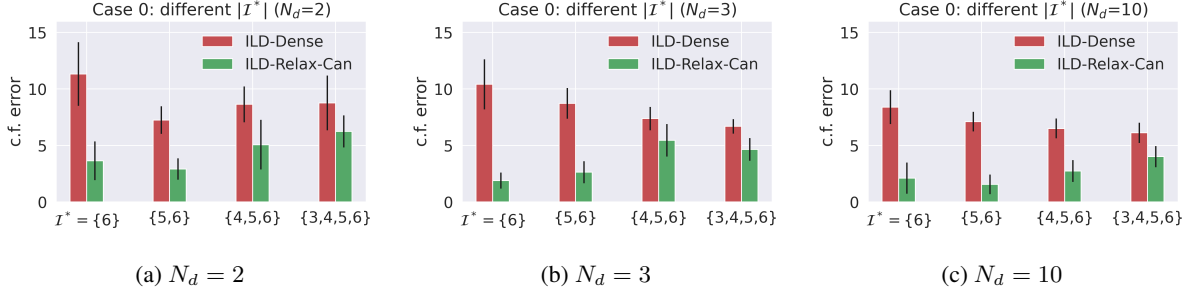


Figure 4: Case 0: Test counterfactual error with different \mathcal{I}^* . To understand how the true intervention set affects the gap between *ILD-Dense* and *ILD-Relax-Can*, we varied the size of the ground truth intervention. It can be observed that the performance gap tends to be largest when the true intervention set is the sparsest and the performance of *ILD-Relax-Can* approaches to the performance of *ILD-Dense* as we increase the size.

Table 3: Case 0: Test counterfactual error and validation log likelihood for each seed when $m = 10, N_d = 2$. We observe that the log likelihood of *ILD-Dense* tends to be much lower when it has a larger counterfactual error than that of *ILD-Relax-Can*.

	Seed	0	1	2	3	4	5	6	7	8	9
Counterfactual error	<i>ILD-Relax-Can</i>	4.625	0.111	0.120	0.072	4.572	10.617	4.360	6.809	0.099	0.479
	<i>ILD-Dense</i>	23.821	0.611	2.178	5.823	4.779	0.694	0.487	1.653	3.170	6.365
Log likelihood	<i>ILD-Relax-Can</i>	-6.873	-7.066	-5.672	-4.637	-0.572	-3.261	-6.062	-4.552	-1.367	-5.170
	<i>ILD-Dense</i>	-4.034	-6.434	-5.679	-4.197	0.711	-1.908	-4.180	-2.413	-1.483	-4.796

of *ILD-Dense* as we increase the size. This makes sense as *ILD-Relax-Can* is a subset of *ILD-Dense* and they are equivalent when $\mathcal{I} = \{1, 2, 3, 4, 5, 6\}$. Additionally, even when the ground truth model is relatively dense (when $|\mathcal{I}^*|$ is close to m), *ILD-Relax-Can* is still better than *ILD-Dense*. We then investigate how models perform under tasks with different numbers of domains. In Figure 5, we change the number of domains in the datasets, and we observe that the performance gap does not seem to be sensitive to the number of domains though the absolute error seems to slightly decrease with more domains. Finally, we test how our algorithm scales with dimension when the number of domains is different. In Figure 6, we notice that *ILD-Relax-Can* is significantly better than *ILD-Dense* in 9 out of 12 cases. In the next paragraphs, we further investigate the 3 cases that do not outperform *ILD-Dense* to understand if it seems to be a theoretic or algorithmic/optimization problem.

We take a further investigation on the three cases where *ILD-Relax-Can* is close to or worse than *ILD-Dense*. As shown in Figure 7, when the latent dimension is 10 and the number of domains is 2, i.e. $m = 10$ and $N_d = 2$, the validation likelihood of *ILD-Relax-Can* is much lower than *ILD-Dense* especially in comparison to that with $m = 4, 6$. We conjecture that the performance drop in terms of counterfactual error could be a result of the worse data fitting, i.e., the model does not fit the data well in terms of log-likelihood. As further evidence, we show the counterfactual error and corresponding validation log-likelihood in Table 3. We observe that the log-likelihood of *ILD-Dense* tends to be much lower when it has a larger counterfactual error than that of *ILD-Relax-Can*. As for the relatively worse performance of *ILD-Relax-Can* when $m = 4, N_d = 2$ and $m = 4, N_d = 3$, we report the counterfactual error corresponding to each seed in Table 4 and Table 5 respectively. When the latent dimension is 4 and the number of domains is 2, i.e., $m = 4, N_d = 2$, *ILD-Relax-Can* is better than *ILD-Dense* with 9 out of 10 seeds. However, it fails significantly with seed 0 and thus leads to a larger average of counterfactual error. When $m = 4, N_d = 3$, *ILD-Relax-Can* is better than *ILD-Dense* with 7 out of 10 seeds but *ILD-Relax-Can* is not significantly better than *ILD-Dense* in terms of average error. We think this is more likely an optimization issue with lower dimensions, which is not explored by our theory. We conjecture that larger models with smoother optimization landscapes will perform better as we see in the imaged-based experiments. We also note that these models are not significantly overparameterized and thus may not benefit from the traditional overparameterization that aids the performance of deep learning in many cases. Further investigation into overparameterized models may alleviate this algorithmic issue.

Despite some corner cases in which the optimization landscape may be difficult for these simple models, all the results point to the same trend that the sparse constraint and canonical form motivated by our theoretic derivation indeed aids in counterfactual performance—despite not explicitly training for counterfactual performance.

Case 1: Correct $|\mathcal{I}|$ but mismatched intervention indices In this section, we include more results in the more practical scenario where we choose the correct number of the intervened nodes but they are not necessarily the last

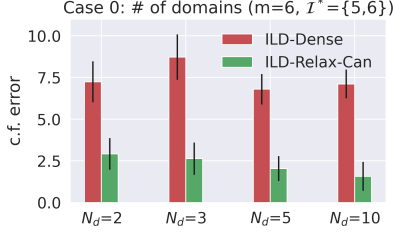


Figure 5: Case 0: Test counterfactual error with different number of domains. Here we investigate how the number of domains affects the performance gap between *ILD-Dense* and *ILD-Relax-Can*. We observe that the gap is not sensitive to the number of domains.

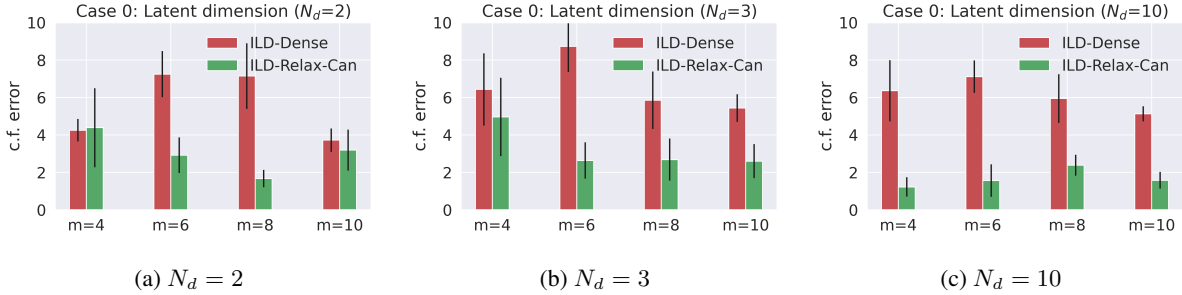


Figure 6: Case 0: Test counterfactual error with different dimension. We investigate how our algorithm scales with dimension. We observe that *ILD-Relax-Can* is significantly better than *ILD-Dense* in 9 out of 12 cases, and we also notice that there 3 cases where their performance is close to that of each other. Here the intervention set contains the last two nodes. For example, when $m = 4$, $\mathcal{I} = \{3, 4\}$, and when $m = 10$, $\mathcal{I} = \{9, 10\}$.

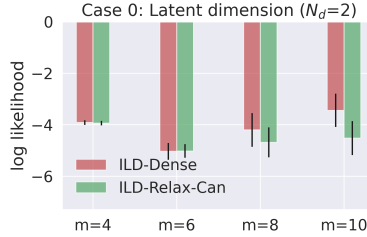


Figure 7: Case 1: Lowest validation log likelihood (same as when we report the test counterfactual error) when testing different dimension with $N_d = 2$. We observe that the likelihood gap between *ILD-Relax-Can* and *ILD-Dense* is largest when $m = 10$.

Table 4: Case 0: Test counterfactual error for each seed when $m = 4$, $N_d = 2$. *ILD-Relax-Can* is better than *ILD-Dense* except when seed is 0. However, there is a significant failure for *ILD-Relax-Can* with seed 0.

Seed	0	1	2	3	4	5	6	7	8	9
<i>ILD-Relax-Can</i>	23.790	2.309	1.747	3.180	1.265	0.864	0.779	0.227	3.325	6.362
<i>ILD-Dense</i>	3.321	3.435	2.838	4.209	5.356	6.456	1.615	2.165	5.195	7.937

Table 5: Case 0: Test counterfactual error for each seed when $m = 4$, $N_d = 3$. *ILD-Relax-Can* is better than *ILD-Dense* with seed 1, 2, 3, 5, 6, 7, 8.

Seed	0	1	2	3	4	5	6	7	8	9
<i>ILD-Relax-Can</i>	23.821	0.611	2.178	5.823	4.779	0.694	0.487	1.653	3.170	6.365
<i>ILD-Dense</i>	24.472	3.658	2.925	5.785	3.260	5.795	3.878	4.560	4.009	5.965

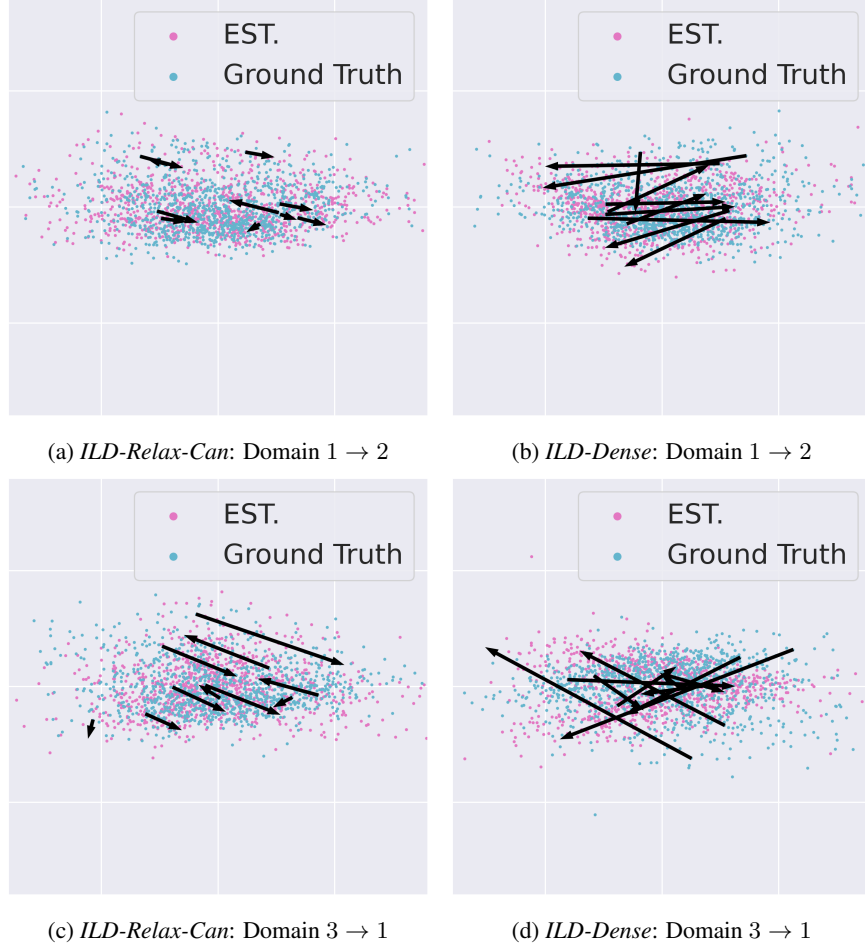


Figure 8: Visualization of counterfactual error when $m = 6, N_d = 3, |\mathcal{I}| = 2, \mathcal{I}^* = \{1, 2\}$. In each plot, we find the first two principle components and project the data along that direction. We select 10 points, then find the corresponding ground truth counterfactual and estimated counterfactual. The black arrow points from ground truth to estimated counterfactual.

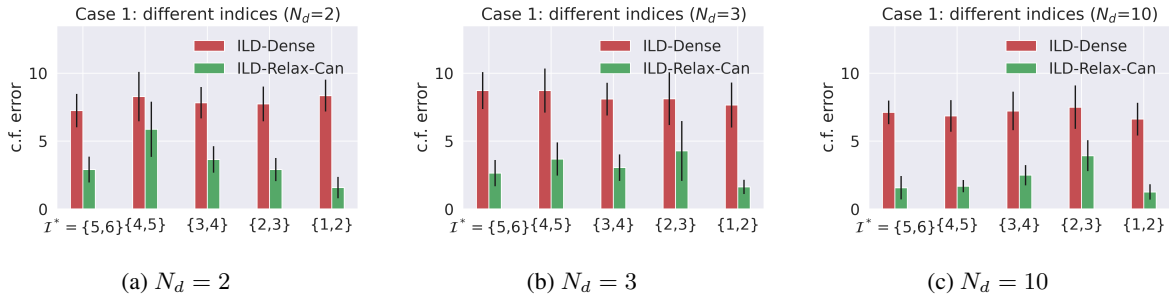


Figure 9: Case 1: Test counterfactual error with different indices. Here we observe that *ILD-Relax-Can* performs consistently better than *ILD-Dense*. When $N_d = 2$ and $\mathcal{I} = \{4, 5\}$, the performance of *ILD-Relax-Can* gets relatively higher because it fails significantly in one case as shown in Table 6.

few nodes in the latent SCM. This experiment is related to our canonical ILD theory, i.e., that there exists a canonical counterfactual model (where the intervened nodes are the last ones) corresponding to any true non-canonical ILD that has the same sparsity. As a starting point, we first illustrate the existence of a canonical model we try to find in Figure 11.

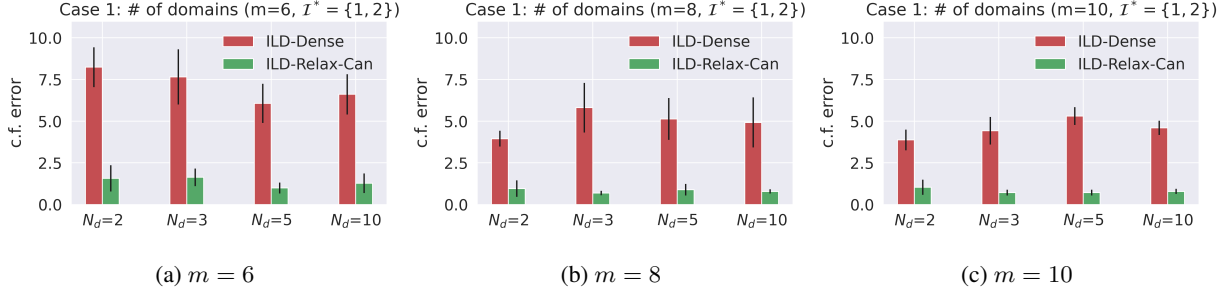


Figure 10: Case 1: Test counterfactual error with different number of domains when $\mathcal{I} = \{1, 2\}$. *ILD-Relax-Can* performs consistently well with different number of domains and latent dimension.

Table 6: Case 1: Test counterfactual error and validation log likelihood for each seed when $N_d = 2$ and $\mathcal{I} = \{4, 5\}$. When seed is 5, the error of *ILD-Relax-Can* is much larger than that of *ILD-Dense*. In the meanwhile, we notice that the log likelihood of *ILD-Relax-Can* is much lower than that of *ILD-Dense* which indicates *ILD-Relax-Can* fails to fit the observed distribution well. When seed is 6, there is also a gap in log likelihood. But both models perform very badly in terms of counterfactual error in this case, and we conjecture this results from a very hard dataset.

	Seed	0	1	2	3	4	5	6	7	8	9
Counterfactual error	<i>ILD-Relax-Can</i>	1.395	0.862	1.338	0.193	7.557	12.422	21.762	3.879	2.352	0.479
	<i>ILD-Dense</i>	8.610	5.979	4.134	2.983	9.795	4.719	24.232	5.327	8.497	8.500
Log likelihood	<i>ILD-Relax-Can</i>	-4.441	-5.737	-4.448	-5.504	-4.393	-3.376	-5.187	-5.073	-4.033	-4.102
	<i>ILD-Dense</i>	-4.170	-5.632	-4.316	-5.458	-4.174	-2.181	-4.052	-5.010	-5.270	-4.302

To investigate the effect of different indices of the intervened nodes, in Figure 9, we change the true intervention set \mathcal{I}^* while keeping the number of intervened nodes $|\mathcal{I}^*|$ the same. We observe that *ILD-Relax-Can* is consistently better than *ILD-Dense* regardless of which nodes are intervened except for one case. When the number of domains is 2 and $\mathcal{I}^* = \{4, 5\}$, we find the gap is much smaller mainly because *ILD-Relax-Can* fails to fit the observed distribution in one case as shown in Table 6. We then test the effect of the number of domains with different latent dimensions in Figure 10. We observe that our model performs consistently well with different numbers of domains and latent dimensions. In Figure 8, we visualize how *ILD-Relax-Can* leads to a lower counterfactual error in comparison to *ILD-Dense*. As shown in Figure 8a and Figure 8b, *ILD-Relax-Can* clearly does better in counterfactual estimation. In Figure 8c and Figure 8d, both of them have a relatively larger error. However, *ILD-Relax-Can* tends to find a closer solution while *ILD-Dense* matches distribution more randomly. This could result from the large search space of *ILD-Dense* and it can easily encodes a transformation such as rotation which will not hurt distribution fitting but will lead to a significant counterfactual error.

Even though we do not know the specific nodes being intervened on, similar to Case 1, we show that sparse constraint leads to better counterfactual estimation.

Case 3: Intervention set size mismatch In this section, we include more results in the most difficult cases where we have no knowledge of the dataset. To investigate what will happen if there is a mismatch of the number of intervened nodes between the true model and the approximation, i.e., $|\mathcal{I}| \neq |\mathcal{I}^*|$, we first change \mathcal{I}^* while keeping the model unchanged, i.e., \mathcal{I} is fixed. As shown in Figure 12, the performance gap between *ILD-Relax-Can* and *ILD-Dense* become smaller as the dataset becomes less sparse while *ILD-Relax-Can* outperforms *ILD-Dense* in all cases. We then change \mathcal{I} while keeping \mathcal{I}^* unchanged. As shown in Figure 13, the performance of *ILD-Relax-Can* approaches to that of *ILD-Dense* as we increase $|\mathcal{I}|$. A somewhat surprising result is that *ILD-Relax-Can* has the lowest counterfactual error when $|\mathcal{I}| = 1$. However, as we check data fitting in Figure 14, we can tell *ILD-Relax-Can* fails to fit the observed distribution in this case. We conjecture the main reason for this is that our theory does not guarantee the existence of a distributionally and counterfactually equivalent canonical model in those cases as we are using a model that is more sparse than the ground truth dataset. Hence, we cannot rely on the counterfactual estimation when the observed distribution is not fitted.

In summary, we observe that *ILD-Relax-Can* always tends to get a lower counterfactual error even though we choose a wrong size of intervention set, i.e. $|\mathcal{I}| \neq |\mathcal{I}^*|$. However, we also observe that in the cases where our model is more sparse than ground truth, the data fitting performance of *ILD-Relax-Can* would drop more significantly. We believe this could also be a good indicator of whether we find a reasonable $|\mathcal{I}|$.

E Image Counterfactual Experiments

E.1 Dataset Descriptions

Rotated MNIST and FashionMNIST We split the MNIST trainset into 90% training data, 10% validation, and for testing we use the MNIST test set. Within each dataset, we create the domain-specific data by replicating all samples and applying a fixed θ_d counterclockwise rotation to within that domain. Specifically we generate data from 5 domains by applying rotation of $0^\circ, 15^\circ, 30^\circ, 45^\circ, 60^\circ$. For Rotated FashionMNIST, we use the same setup as the RMNIST setup, except we used the Fashion MNIST [Xiao et al., 2017] dataset. This dataset is structured similar to the MNIST dataset but is designed to require more complex modeling [Xiao et al., 2017].

3D Shapes This is a dataset of 3D shapes that are procedurally generated from 6 independent latent factors: floor hue, wall hue, object hue, scale, shape, and orientation [Burgess and Kim, 2018]. In our experiment, we only choose samples with one fixed scale. We then split the four object shapes into four separate domains and set the 10 object colors as the class label. The causal graph for this dataset can be seen in Fig. 15c, and following this, we should expect to see only the object shape change when the domain is changed. Similar to the RMNIST experiment, we use 90% of the samples for training and 10% of the samples for validation.

Color Rotated MNIST (CRMNIST) This is an extension of the RMNIST dataset which introduces a latent color variable whose parents are the latent domain-specific variable and latent class variable. Similar to RMNIST, the latent domain variable corresponds to the rotation of the given digit, except here $d_1 = 0^\circ$ rotation, $d_2 = 90^\circ$ rotation, and the class labels are restricted to digits $y \in \{0, 1, 2\}$. For each sample, there is a 50% chance the color is determined by the combination of class and digit label and a 50% chance the color is randomly chosen. For example if $\epsilon \sim \mathcal{N}(0, 1)$,

$$f_{z_c}(y, d, \epsilon) = \begin{cases} \text{red,} & \text{if } y = 0, d = 1, \epsilon > 0 \\ \text{green,} & \text{if } y = 0, d = 2, \epsilon > 0 \\ \text{blue,} & \text{if } y = 1, d = 1, \epsilon > 0 \\ \dots & \\ \text{Random(red, green, blue, yellow, cyan, pink),} & \text{if } \epsilon < 0 \end{cases}$$

The causal graph for this dataset can be seen in Fig. 15b. Similar to the RMNIST experiment, we use 90% of the samples for training and 10% of the samples for validation.

Causal3DIdent Dataset This is a benchmark dataset from [Von Kügelgen et al., 2021] which contains rendered images of realistic 3d objects on a colored background that contain hallmarks of natural environments (e.g. shadows, different lighting conditions, etc.) which are generated via a causal graph imposed over latent variables (the causal graph can be seen in Figure 15d). Similar to [Von Kügelgen et al., 2021], we chose the shape of the 3D object to be the class label, and we defined the background color as the domain label. In the original dataset, the range of the background hue was $[\frac{-\pi}{2}, \frac{\pi}{2}]$ and to convert it to a binary domain variable, we binned the background hue values into bins $[\frac{-\pi}{2}, -0.8]$ and $[\frac{\pi}{2}, 0.8]$. These ranges were chosen to be distinct enough that we can easily distinguish between the domains yet large enough to keep the majority of original samples in this altered dataset. We split the 18k binned samples into 90% training data and 10% validation data for our experiment.

E.2 Metrics

Inspired by the work in Monteiro et al. [2023], we define four metrics (Effectiveness, Preservation, Composition, and Preservation) specifically for the image-based counterfactuals with latent SCMs. The key idea is to check if the correct latent information is changed when generating domain counterfactuals (e.g., domain-specific information is changed, while all else is preserved). Since we don't have direct access to the ground truth value of latent variables (nor their counterfactual values), we use a domain classifier h_{domain} and class classifier h_{class} to measure if the intended change has taken place.

Effectiveness: The idea is to check if the domain-specific variables change as wish in the counterfactual samples.

$$\mathbb{E}_{(x,d)} \left[\mathbb{1}_{h_{\text{domain}}(\hat{x}^{d \rightarrow d'})=d'} \right]$$

Preservation: This checks if the semantically meaningful content (i.e. the class information) that is independent of the domain is left unchanged while the domain is changed.

$$\mathbb{E}_{(x,d)} \left[\mathbb{1}_{h_{\text{class}}(\hat{x}^{d \rightarrow d'})=y} \right]$$

Composition: We check if our model is invertible on the image manifold, thus satisfying the pseudoinvertibility criteria.

$$\mathbb{E}_{(x,d)} [\mathbb{1}_{h_{\text{class}}(\hat{x}^{d \rightarrow d})=y}]$$

Reversibility: This metric checks if our model is cycle-consistent, or in other words, checking if the mapping between the observation and the counterfactual is deterministic.

$$\mathbb{E}_{(x,d)} [\mathbb{1}_{h_{\text{class}}(\hat{x}^{d \rightarrow d' \rightarrow d})=y}]$$

For the domain classifier h_{domain} and class classifier h_{class} , we used pretrained ResNet18 models [He et al., 2016] that were fine-tuned by classifying *clean* samples (i.e. not counterfactuals) for 25 epochs with the Adam optimizer, a learning rate of 1e-3, and a random data augmentation with probabilities: 50%: no augmentation, 17%: sharpness adjustment (factor=2), 17%: gaussian blur (kernel size=3), 17%: gaussian blur (kernel size=5). A reminder that for MNIST/FMNIST/ColorRotated MNIST, the domain is rotation and the label is the original label of images (digits/type of clothes), for 3D shapes, the domain is object shape and the label is object color, and for Causal3DIdent, the domain is hue of the background and the label is the object shape.

E.3 Causal Interpretation of our experiments

In this section, we introduce the causal interpretation of our experiments. To evaluate the model’s capability of generating good domain counterfactuals, for each dataset, we have one domain latent variable and choose one class latent variable that are generated independently of the domain latent variable. As an example, for RMNIST, we choose rotation as the domain latent variable and digit class as the class latent variable. As indicated in Figure 15a, for the counterfactual query “Given we observe image in this domain, what would have happened if it is in another domain?”, we should expect the image to be rotated while the class remain unchanged. Specifically, we want to check $\forall d, d', d'' \in \mathcal{D}, \mathbb{P}(Z_{\text{rot}}(D = d')|X = x, D = d) = \mathbb{P}(Z_{\text{rot}}(D = d'')|X = x, D = d)$ iff $d' = d''$ and $\forall d, d', d'' \in \mathcal{D}, \mathbb{P}(Z_y(D = d')|X = x, D = d) = \mathbb{P}(Z_y(D = d'')|X = x, D = d)$ where \mathcal{D} is the set of all domains. However, in practice we cannot directly get the value of those latent variables. This motivates our choice of evaluation metric of training a domain classifier and class classifier to detect if the domain latent variable is changed and class latent variable (which we call class) is preserved in the counterfactuals.

For RMNIST/RFMNIST, we choose rotation as the domain variable and digit/clothes class as the class variable. For 3D Shapes, we choose object shape as the domain variable and hue of objects as the class variable. For CRMNIST, we choose rotation as the domain variable and digit class as the class variable. For CausalIdent, we choose the hue of the background as the domain variable and object class as the class variable. In the case of 3D Shapes, we can technically choose anything other than object shape as the class variable. However, for simplicity, we choose one of them. In the case of CRMNIST, we cannot choose Z_{color} because it will change after we change the domain. In the case of Causal3DIdent, we can choose anything but the hue of the object, though we figure Z_y is easier to check and can reduce error caused by classifier proxies.

We also want to note that other than observational image, access to domain information is also important for answering this query. For example, in the case of RMNSIT, given an image that looks like digit “9”, for the question “what would have happened if it is in domain 90°”, the fact that the current digit is in domain 0° (which means it is indeed digit “9”) or the current digit is in domain 180° domain 0° (which means it is digit “6”) would lead to different answer.

E.4 Experiment Details

Model setup The relaxation to pseudo invertibility allows us to modify the ILD models to fit a VAE [Kingma and Welling, 2013] structure. The overall VAE structure can be seen in Fig. 16, where the variational encoder first projects to the latent space via g^+ to produce the latent encoding z , which is then passed to two domain-specific autoregressive models $f_{d,\mu}^+, f_{d,\sigma}^+$ which produce the mean and variance parameters (respectively) of the Gaussian posterior distribution. The decoder of the VAE follows the structure typical ILD structure: $g \circ f_d$. Here, g^+ can be viewed as the pseudoinverse of the observation function g and f_d can be viewed as a pseudoinverse of $f_{d,\mu}^+$. During training, the exogenous noise variable ϵ is then found via sampling from the posterior distribution ($\epsilon \sim \mathcal{N}(\mu_d, \sigma_d)$) which can be viewed as a stochastic SCM, however, to reduce noise when producing counterfactuals, when performing inference the exogenous variable is set to the mean of the latent posterior distribution (i.e. $\epsilon = \mu_d$). In all experiments, g and g^+ follow the β -VAE architecture seen in Higgins et al. [2017] (with the exception that in the Causal3DIdent experiment, g and g^+ follow the base VQ-VAE architecture [Van Den Oord et al., 2017] without the quantizer), and the structure of the f models is determined by the type of ILD model used (e.g., independent, dense, or relaxed canonical) and matches that seen in the simulated experiments and visualized in Fig. 3. For the f models which enforce sparsity

	RMNIST				RFMNIST			
	Comp.	Rev.	Eff.	Pre.	Comp.	Rev.	Eff.	Pre.
<i>ILD-Independent</i>	99.79 ± 0.44	32.56 ± 0.20	94.97 ± 4.71	32.49 ± 0.22	69.75 ± 1.86	22.36 ± 0.76	99.62 ± 0.37	22.54 ± 1.19
<i>ILD-Dense</i>	99.76 ± 0.28	32.60 ± 0.21	80.92 ± 2.21	32.64 ± 0.23	71.20 ± 3.39	24.23 ± 2.51	98.51 ± 0.93	23.98 ± 2.18
<i>ILD-Relax-Can</i>	99.85 ± 0.27	79.84 ± 17.54	96.72 ± 1.89	64.99 ± 9.83	71.79 ± 4.55	70.44 ± 3.54	98.82 ± 0.73	62.15 ± 6.65

Table 7: Quantitative result with RMNIST and RFMNIST, where higher is better. They are both averaged over 20 runs.

	CRMNIST				3D Shapes			
	Comp.	Rev.	Eff.	Pre.	Comp.	Rev.	Eff.	Pre.
<i>ILD-Independent</i>	87.24 ± 11.98	59.88 ± 6.46	94.65 ± 15.34	60.39 ± 6.95	99.79 ± 0.44	32.56 ± 0.20	94.97 ± 4.71	32.49 ± 0.22
<i>ILD-Dense</i>	88.18 ± 17.84	62.29 ± 10.51	92.72 ± 15.52	59.60 ± 8.92	99.76 ± 0.28	32.60 ± 0.21	80.92 ± 2.21	32.64 ± 0.23
<i>ILD-Relax-Can</i>	92.10 ± 13.24	85.74 ± 13.33	94.48 ± 10.71	72.95 ± 12.42	99.85 ± 0.27	79.84 ± 17.54	96.72 ± 1.89	64.99 ± 9.83

Table 8: Quantitative result with CRMNIST and 3D Shapes, where higher is better. CRMNIST are averaged over 20 runs and 3D Shapes are averaged over 5 runs.

(i.e. *ILD-Relax-Can*), we use a sparsity level, $|\mathcal{I}|$, of 5. We also introduce an additional baseline, *ILD-Independent*, which has an architecture similar to the *ILD-Dense* baseline, with the exception that the g and g^+ functions are no longer shared across domains. The *ILD-Independent* baseline can be seen as training an independent β -VAE for each domain, where each β -VAE an autoregressive f_{dense} model as its last (first) layer for the encoder (decoder), respectively. For experiment with RMNIST, RFMNIST, 3D Shapes and Causal3DIdent, we choose $m = 20$ and for CRMNIST, we choose $m = 10$.

Training We train each ILD model for 300K,300K,300K,500K,200K for RMNIST, RFMNIST, CRMNIST, 3D Shapes and Causal3DIdent respectively using the Adam optimizer [Kingma and Ba, 2014] with $\beta_1 = 0.5, \beta_2 = 0.999$, and a batch size of 1024. The learning rate for g and g^+ is 10^{-4} , and all f models use 10^{-3} . During training, we calculate two loss terms: a reconstruction loss $\ell_{recon} = \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2$ where $\hat{\mathbf{x}}$ is the reconstructed image of \mathbf{x} and the ℓ_{align} alignment loss which measures the KL-divergence between the posterior distribution $Q_d(\epsilon|\mathbf{x})$ and the prior $P(\epsilon)$. Following the β -VAE loss calculation in Higgins et al. [2017], we apply a β_{KLD} upscaling to the alignment loss such that $\ell_{total} = \ell_{recon} + \beta_{KLD} * \ell_{align}$. For all MNIST-like experiments, we use $\beta_{KLD} = 1000$, which we found leads to the lowest counterfactual error on the validation datasets across all models; this also matches the β_{KLD} used in Burgess et al. [2018], and for 3DShape and Causal3DIdent we found $\beta_{KLD} = 10$ leads to the lowest counterfactual error.

E.5 Additional Results

The quantitative results in Table 7, Table 8, and Table 9 match the visual result seen in Figure 17, Figure 18, Figure 19, where almost across all datasets the *ILD-Relax-Can* model seems to find a proper latent causal structure that can disentangle the domain information from the class information – unlike the baseline models which seem to commonly change the class during counterfactual. We again note that the training process for all of the models only include the typical VAE invertibility loss (i.e. reconstruction loss) and latent alignment loss (i.e. the KL-divergence between the latent prior and posterior distributions) and do not specifically include any counterfactual training. Thus, we conjecture the enforcing of sparsity in the canonical models correctly biased these models in a manner that preserved important non-domain-specific information when performing counterfactuals. In Figure 22, Figure 23, Figure 24 and Figure 25, we track the change of our metrics w.r.t $|\mathcal{I}|$ (we did not do this investigation for Causal3DIdent because that the training of that model takes much longer time). We observe that as we increase $|\mathcal{I}|$, the reversibility and preservation tends to decrease while the effectiveness tends to increase. We conjecture that this is because as $|\mathcal{I}|$ increases, there is less constraint on the original optimization problem (fitting the observational distribution) which could potentially increase the performance. However, it leads to lower chance in finding a proper latent causal structure for domain counterfactual generation, which results in the decrease in preservation. *ILD-Dense* can be regarded as an extreme case of this. In

	Causal3DIdent			
	Comp.	Rev.	Eff.	Pre.
<i>ILD-Independent</i>	88.15 ± 5.0	51.43±2.7	91.05 ± 17.7	51.94±3.0
<i>ILD-Dense</i>	83.59 ± 5.4	49.17±2.5	92.17 ± 13.6	48.83±3.0
<i>ILD-Relax-Can</i>	86.00 ± 5.6	79.73 ± 6.6	84.15 ± 23.5	79.73 ± 8.6

Table 9: Quantitative result with Causal3DIdent, where higher is better. Causal3DIdent are averaged over 10 runs.

summary, we validate the practicality of our model design in the pseudoinvertible setting with extensive study on 5 image-based datasets.

F Limitations

In this paper, we first prove the existence of distributionally and counterfactually equivalent models. Then we investigate how hard it would be to learn such models in practice when the only objective in the algorithm is to fit the observed distribution. In our extensive simulated experiments and image-based experiments, we find that the sparsity constraint inspired by our theory helps the model achieve more accurate counterfactual estimation. From a theoretic side, while our theory proves the existence of canonical ILD models, however, we have not proven identifiability of the latent causal model or observation function. Indeed, we conjecture that complete identifiability of latent causal models is likely infeasible in our setup except under very strong constraints. A deeper investigation into the conditions for identifiability or proof of non-identifiability would be interesting future directions.

A practical problem we noticed in our simulated experiments is that sometimes the sparse model is harder to fit, i.e., its log-likelihood is worse than the dense model, even if we only consider the cases where the true model is in the model class being optimized (e.g., the sparsity of the model is at least as large as the sparsity of the ground truth model). We conjecture that this results from a harder loss landscape as we add more constraints to the model. We believe a more careful investigation of the model and algorithm could be an interesting and important future work. For example, if we use a more significantly overparameterized model, there are chances that the training of *ILD-Relax-Can* would become easier. Additionally, the addition of further loss terms could aid in the training of these models, such as, assuming access to some ground truth domain counterfactuals (e.g., the same patient received imaging at multiple hospitals) could be used to penalize our model when it changes latent variables which do not change under the ground truth counterfactuals.

In our experiments, we aimed to test the effects of breaking some of our assumptions (e.g., “what if our model is not strictly invertible”), and while our models still performed better in these cases, there are likely cases where the breaking of our assumptions can cause our models to fail to produce faithful counterfactuals. For example, in a case where there is a very large difference between domains and there is no sparsity in the domain shifts, then it is likely that the constraints constituted by our sparsity assumption will make the sparse models struggle to fit the observed distributions.

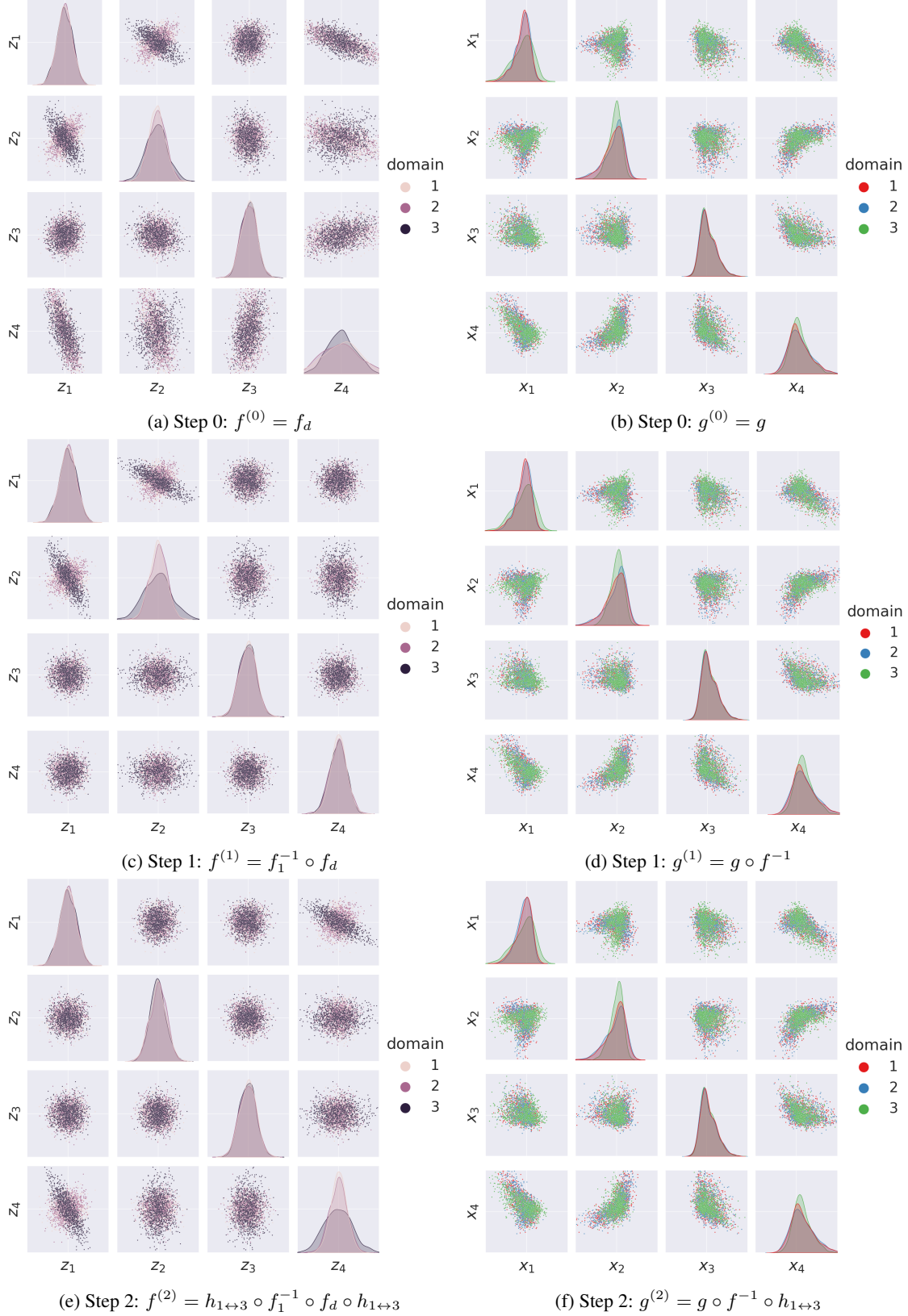


Figure 11: An illustration of the existence of a distributionally and counterfactually equivalent model in canonical form when $m = 4$ and $\mathcal{I} = \{2\}$. $h_{1 \leftrightarrow 3}$ represents a swapping matrix. $g^{(2)} \circ f^{(2)}$ is one of the canonical models we try to find. Note that the observed distributions in the right column are always the same while the latent distributions on the left change. In particular, the canonical ILD model on the bottom left has independent distributions for the first three variables and is only the non-identity on the last node.

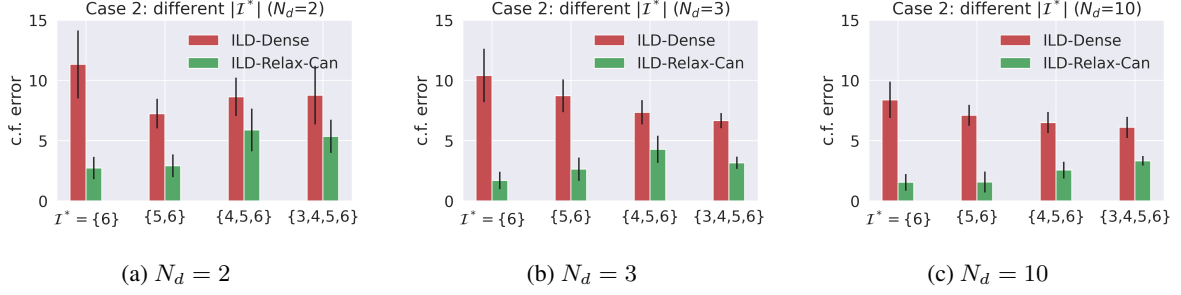


Figure 12: Case 2: Test counterfactual error with different $|\mathcal{I}^*|$ and fixed $|\mathcal{I}| = 2$. The performance of *ILD-Relax-Can* gets worse as the dataset becomes less sparse. But it is still better than *ILD-Dense*. Note that when $|\mathcal{I}| = 2$ and $\mathcal{I}^* = \{6\}$, the ground truth canonical model is still a subset of the models we search over.

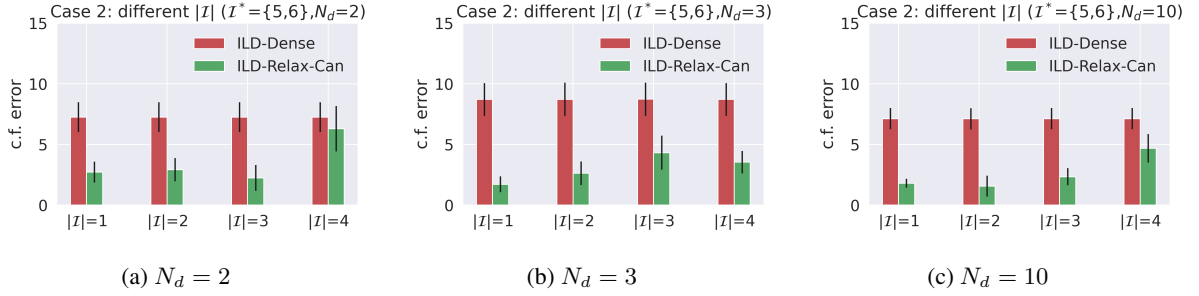


Figure 13: Case 2: Test counterfactual error with different $|\mathcal{I}|$ and fixed \mathcal{I}^* . The performance of *ILD-Relax-Can* approaches to that of *ILD-Dense* as we increase $|\mathcal{I}|$.

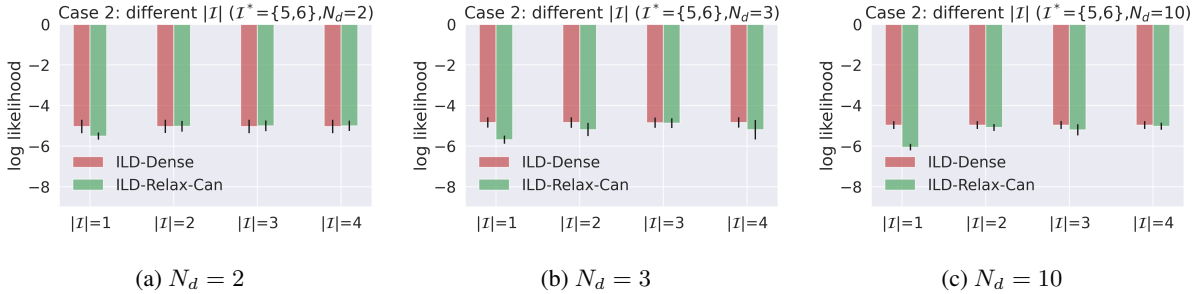


Figure 14: Case 2: Lowest validation log likelihood with different $|\mathcal{I}|$ and fixed \mathcal{I}^* . When $|\mathcal{I}| = 1$, there is a more significant gap between *ILD-Relax-Can* and *ILD-Dense* with all N_d which indicates *ILD-Relax-Can* might fail to fit the observed distribution.

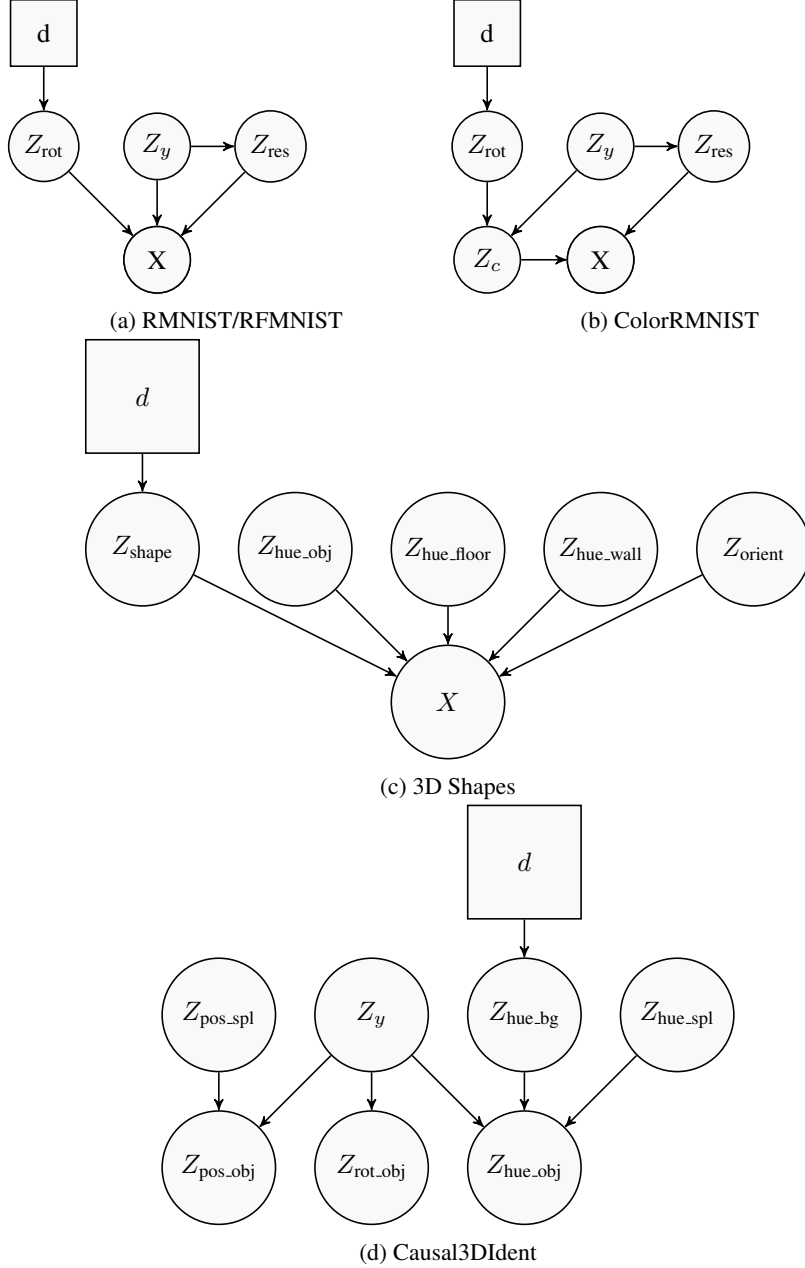


Figure 15: (a) RMNIST/RFMNIST. Here Z_{rot} represents the rotation of the image, Z_y represents the original RMNIST/RFMNIST class, Z_{res} contains other detail information such as writing style, which is controlled by how MNIST dataset was originally created. (b) Z_c represents the color of the digit while others are the same as (a). (c) 3D Shapes. Z_{shape} represents the object shape. $Z_{hue.obj}$, $Z_{hue.floor}$, $Z_{hue.wall}$ represent the hue of the object, floor and wall respectively. Z_{orient} represents the orientation of the object. (d) Causal3DIdent. Z_y represents the object class. $Z_{hue.obj}$, $Z_{hue.bg}$, $Z_{hue.spl}$ represent the hue of the object, background and spotlight respectively. $Z_{pos.obj}$, $Z_{pos.spl}$ represent the position of the object and spotlight respectively. $Z_{rot.obj}$ represents the rotation of the object. X is not shown in the graph but all nodes should point to it.

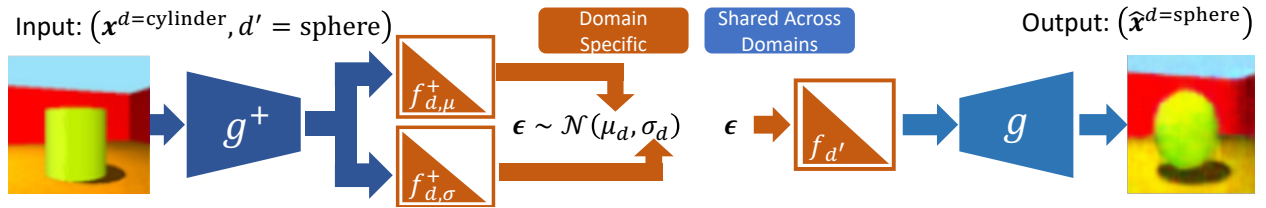


Figure 16: The model structure for the pseudo-invertible ILD model used in the high-dimensional experiments. The overall structure matches that of a VAE where the encoder (left) first projects to the latent space via g^+ (the pseudo-inverse of the observation function g). This latent encoding is then passed to two domain-specific autoregressive models $f_{d,\mu}^+, f_{d,\sigma}^+$ which produce the mean and variance parameters (respectively) of the Gaussian posterior distribution. During training, the exogenous noise variable ϵ is then found via sampling from the posterior distribution ($\epsilon \sim \mathcal{N}(\mu_d, \sigma_d)$) which can be viewed as a stochastic SCM, however, during inference the exogenous variable is set to the mean of the latent posterior distribution (i.e. $\epsilon := \mu_d$) to reduce noise when producing counterfactuals. The decoder (right) follows the usual VAE decoder structure, with the exception that the initial linear layer is an autoregressive function of the ϵ input. The structure of all the f models is determined by the type of ILD model used (e.g., dense, canonical, or relaxed canonical) and matches that seen in Fig. 3.

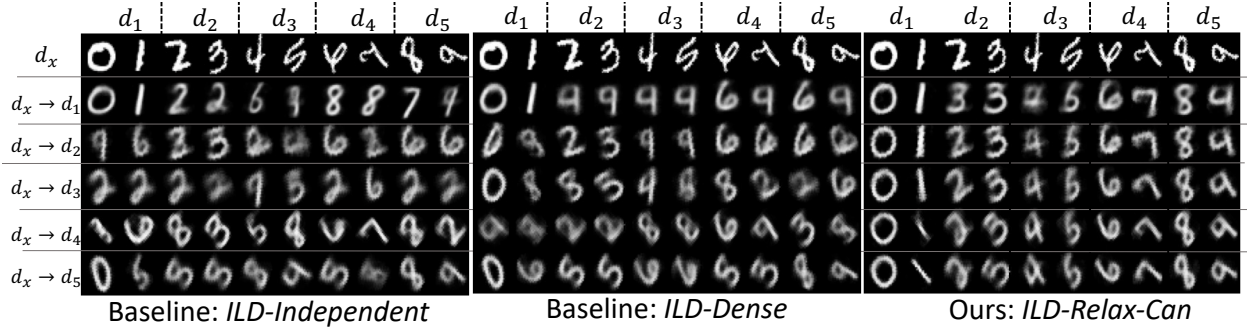


Figure 17: Counterfactual plots for the three relaxed ILD models, where across the columns we show examples of two clothing classes (e.g., “handbag” or “boot”) from each domain and each row corresponds to the counterfactual to a different domain. It can be seen that while all models correctly recover the rotation for each domain counterfactual, the baseline models usually change the class label during counterfactual, while *ILD-Relax-Can* tends to preserve the clothing label, despite not being privy to any label information during training.

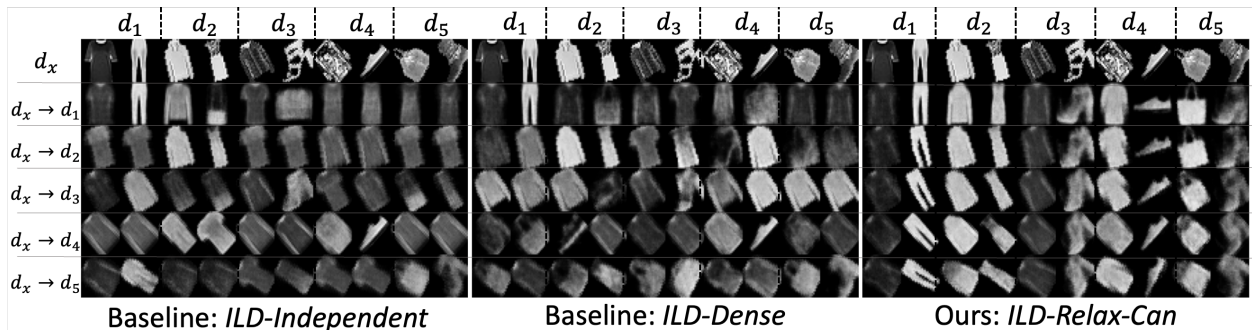


Figure 18: Counterfactual plots for the three ILD models, where across the columns we show examples of two classes from each domain and each row corresponds to the counterfactual to a different RMNIST domain. It can be seen that while all four models correctly recover the rotation for each domain counterfactual, the baseline models usually change the digit label during counterfactual, while *ILD-Relax-Can* tends to preserve the digit label, despite not being privy to any label information during training.

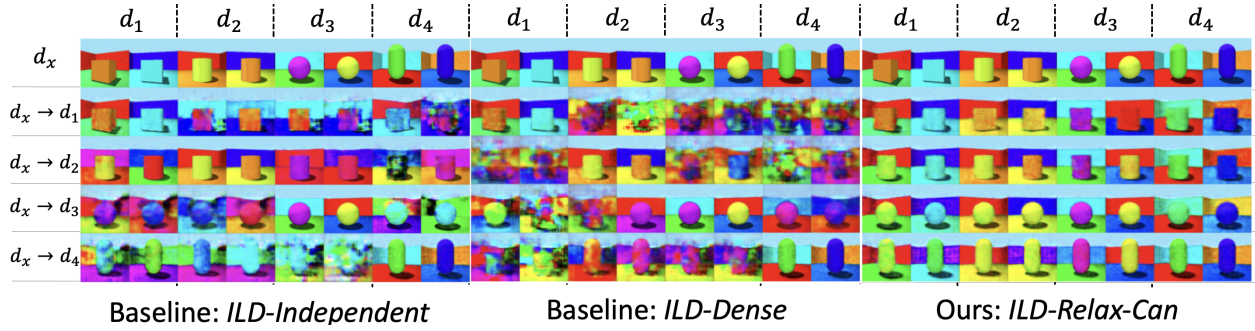


Figure 19: Counterfactual plots for the three ILD models, where across the columns we show examples of two classes from each domain and each row corresponds to the counterfactual to a different object shape domain.

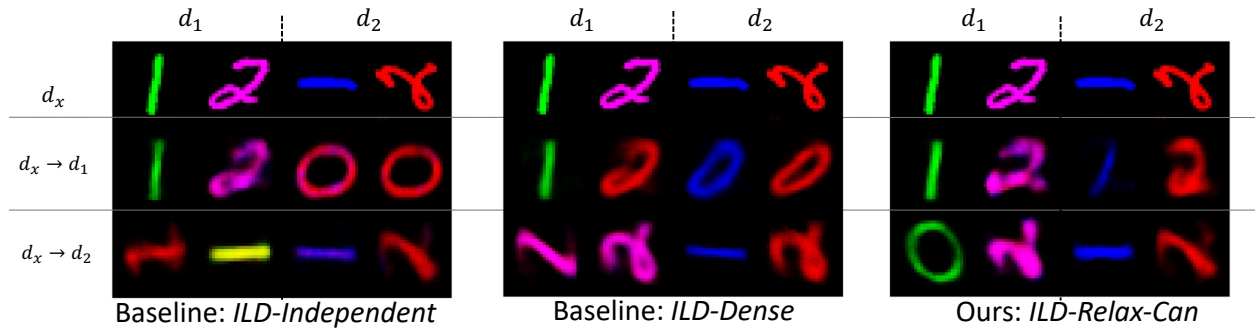


Figure 20: Counterfactual plots for the three ILD models, where across the columns we show examples of two classes from each domain and each row corresponds to the counterfactual to a different rotation domain.

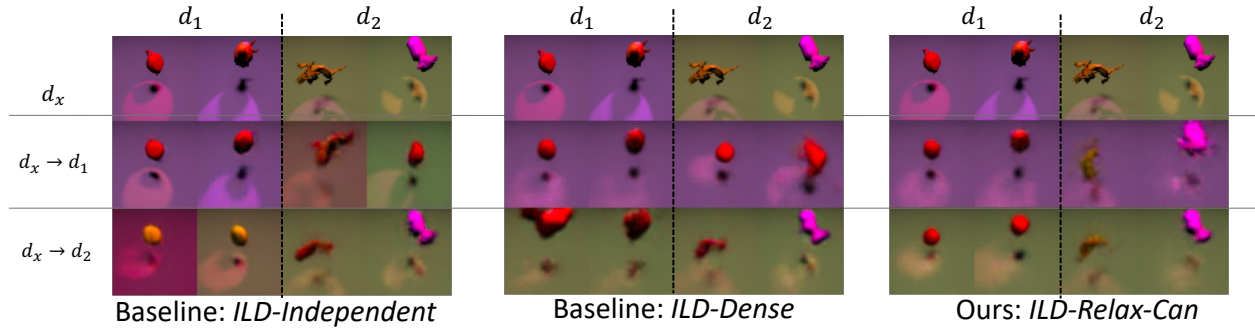


Figure 21: Counterfactual plots for the three ILD models, where across the columns we show examples of two classes from each domain and each row corresponds to the counterfactual to a different background hue domain.

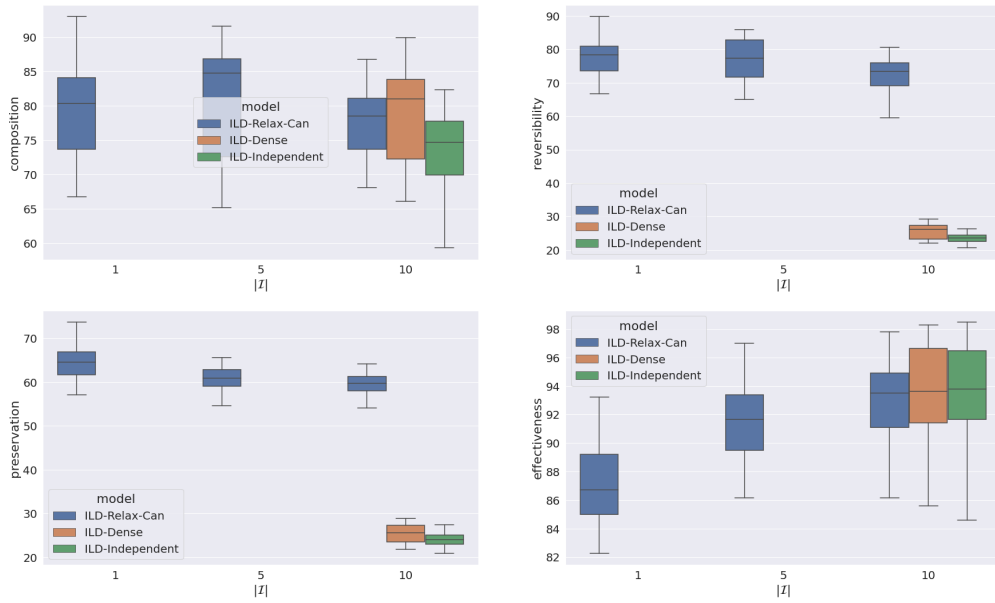


Figure 22: Change of metrics w.r.t $|Z|$ for RMNIST. Results are with 20 runs and we remove outliers when plotting.

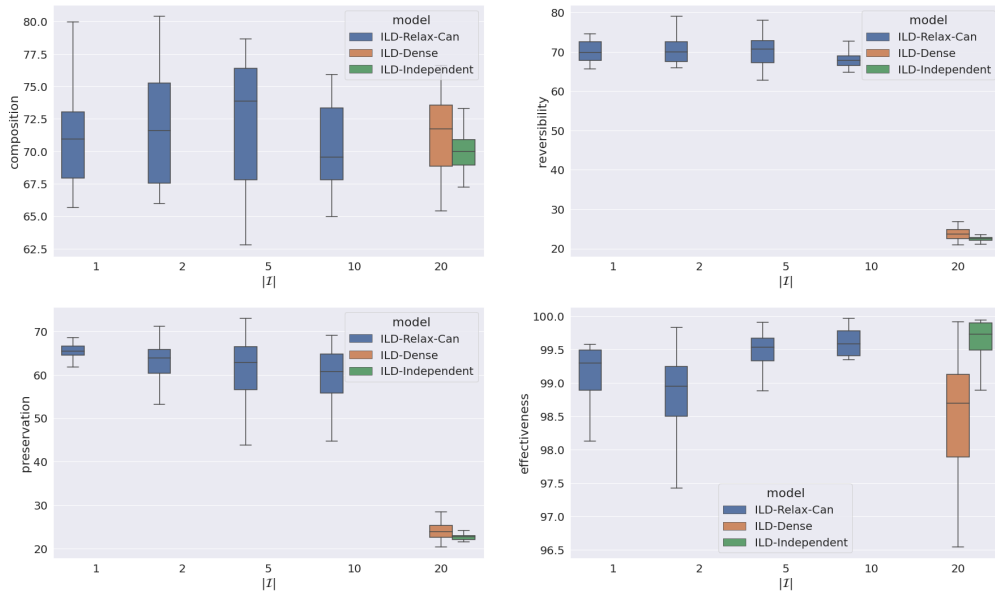


Figure 23: Change of metrics w.r.t $|Z|$ for RFMNIST. Results are with 20 runs and we remove outliers when plotting.

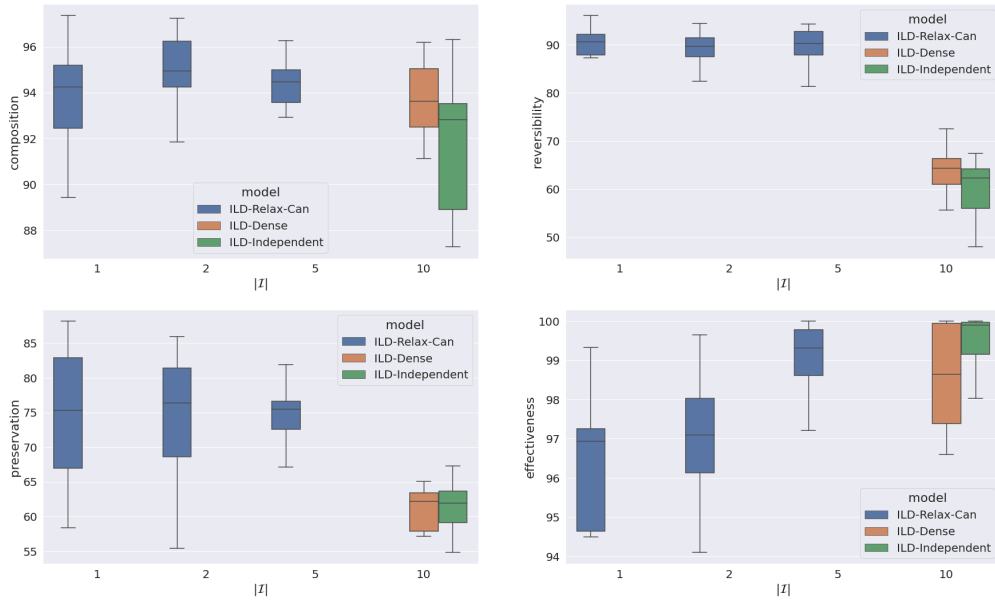


Figure 24: Change of metrics w.r.t $|Z|$ for CRMNIST. Results are with 20 runs and we remove outliers when plotting.

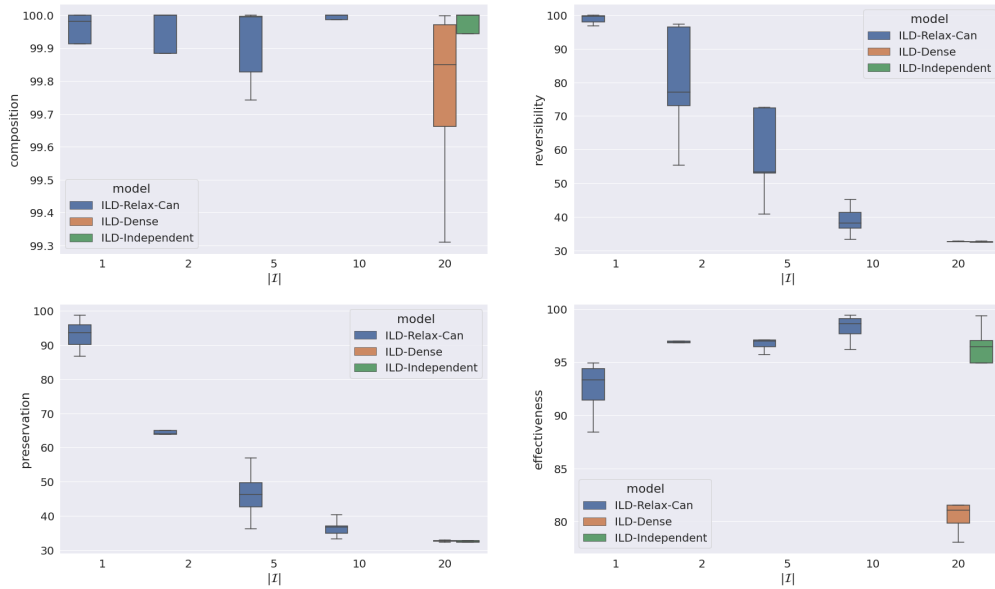


Figure 25: Change of metrics w.r.t $|Z|$ for 3D Shapes. Results are with 5 runs and we remove outliers when plotting.