

# INVARIANT CAUSAL REPRESENTATION LEARNING FOR OUT-OF-DISTRIBUTION GENERALIZATION

Chaochao Lu<sup>1,2</sup>, Yuhuai Wu<sup>3,4†</sup>, José Miguel Hernández-Lobato<sup>1,5\*</sup>, Bernhard Schölkopf<sup>2\*</sup>

## ABSTRACT

Due to spurious correlations, machine learning systems often fail to generalize to environments whose distributions differ from the ones used at training time. Prior work addressing this, either explicitly or implicitly, attempted to find a data representation that has an invariant relationship with the target. This is done by leveraging a diverse set of training environments to reduce the effect of spurious features and build an invariant predictor. However, these methods have generalization guarantees only when both data representation and classifiers come from a linear model class. We propose invariant Causal Representation Learning (iCaRL), an approach that enables out-of-distribution (OOD) generalization in the nonlinear setting (i.e., nonlinear representations and nonlinear classifiers). It builds upon a practical and general assumption: the prior over the data representation (i.e., a set of latent variables encoding the data) given the target and the environment belongs to general exponential family distributions, i.e., a more flexible conditionally non-factorized prior that can actually capture complicated dependences between the latent variables. Based on this, we show that it is possible to identify the data representation up to simple transformations. We also show that all direct causes of the target can be fully discovered, which further enables us to obtain generalization guarantees in the nonlinear setting. Experiments on both synthetic and real-world datasets demonstrate that our approach outperforms a variety of baseline methods.

## 1 INTRODUCTION

Modern machine learning algorithms still lack robustness, and may fail to generalize outside of a specific training distribution because they learn easy-to-fit spurious correlations which are prone to change between training and testing environments. We recall the widely used example of classifying images of camels and cows (Beery et al., 2018). Here, the training dataset has a selection bias, i.e., many pictures of cows are taken on green pastures, while most pictures of camels happen to be in deserts. After training, it is found that the model builds on spurious correlations, i.e., it relates green pastures with cows and deserts with camels, and fails to recognize images of cows on the beach.

To address this problem, a natural idea is to identify which features of the training data present domain-varying spurious correlations with labels and which features describe true correlations of interest that are stable across domains. In the example above, the former are the features describing the context (e.g., pastures and deserts), whilst the latter are the features describing the animals (e.g., animal shape). By exploiting the varying degrees of spurious correlation naturally present in training data collected from multiple environments, one can try to identify stable features and build invariant predictors. Invariant risk minimization (IRM) seeks to find data representations (Arjovsky et al., 2019) or features (Rojas-Carulla et al., 2018) for which the optimal predictor is invariant across all environments. The general formulation of IRM is a challenging bi-leveled optimization problem, and theoretical guarantees require constraining both data representations and classifiers to be linear (Arjovsky et al., 2019, Theorem 9), or considering the special case of feature selection (Rojas-Carulla et al., 2018, Theorem 4). Ahuja et al. (2020a) study the problem from the perspective of game theory, with an approach termed invariant risk minimization games (IRMG). They show that the set of Nash equilibria for a proposed game is equivalent to the set of invariant predictors for any finite number of environments, even with nonlinear data representations and nonlinear classifiers. However, these

<sup>1</sup>University of Cambridge, <sup>2</sup>MPI for Intelligent Systems, <sup>3</sup>Stanford University, <sup>4</sup>Google Research, <sup>5</sup>The Alan Turing Institute, <sup>†</sup>Work done at University of Toronto, <sup>\*</sup>Equal Supervision, Correspondence at c1641@cam.ac.uk.

theoretical results in the nonlinear setting only guarantee that one can learn invariant predictors from training environments, but do not guarantee that the learned invariant predictors can generalize well across all environments including unseen testing environments.

We propose invariant Causal Representation Learning (iCaRL), a novel approach that enables out-of-distribution (OOD) generalization in the nonlinear setting (i.e., nonlinear representations and nonlinear classifiers<sup>1</sup>). We achieve this by extending and using methods from representation learning and graphical causal discovery. In more detail, we first introduce our main *general assumption*: when conditioning on the target (e.g., labels) and the environment (represented as an index), the prior over the data representation (i.e., a set of latent variables encoding the data) belongs to a *general exponential family*. Unlike the conditionally factorized prior assumed in recent identifiable variational autoencoders (iVAE) (Khemakhem et al., 2020a), this is a more flexible conditionally non-factorized prior, which can actually capture complicated dependences between the latent variables. We then extend iVAE to the case in which the latent variable prior belongs to such a *general exponential family*. The combination of this result and the previous general assumption allows us to guarantee that the data representation can be identified up to simple transformations. We then show that the direct causes of the target can be fully discovered by analyzing all possible graphs in a structural equation model setting. Once they are discovered, the challenging bi-leveled optimization problem in IRM and IRMG can be reduced to two simpler independent optimization problems, that is, learning the data representation and learning the optimal classifier can be performed separately. This leads to a practical algorithm and enables us to obtain generalization guarantees in the nonlinear setting.

Overall, we make a number of key contributions: (1) We propose a general framework for out-of-distribution generalization in the nonlinear setting with the theoretical guarantees on both identifiability and generalizability; (2) We propose a general assumption on the underlying causal diagram for prediction (Assumption 1 and Fig. 1c), which covers many real-world scenarios (Section 3.2); (3) We propose a general assumption on the prior over the latent variables (Assumption 2), i.e., a more flexible conditionally non-factorized prior; (4) We prove that an extended iVAE with this conditionally non-factorized prior is also identifiable (Theorems 1, 2&3); (5) We prove that our framework has the theoretical guarantees for OOD generalization in the nonlinear setting (Proposition 1).

## 2 PRELIMINARIES

### 2.1 IDENTIFIABLE VARIATIONAL AUTOENCODERS

Variational autoencoders (VAEs, see Appendix B) (Kingma & Welling, 2013; Rezende et al., 2014) lack identifiability guarantees. Consider a VAE model where  $\mathbf{X} \in \mathbb{R}^d$  stands for the observed variables (data) and  $\mathbf{Z} \in \mathbb{R}^n$  for the latent variables. Khemakhem et al. (2020a) show that a VAE with an unconditional prior distribution  $p_\theta(\mathbf{Z})$  over the latent variables is unidentifiable. However, they also show that it is possible to obtain an identifiable model if one posits a conditionally factorized prior distribution over the latent variables,  $p_\theta(\mathbf{Z}|\mathbf{U})$ , where  $\mathbf{U} \in \mathbb{R}^m$  is an additional observed variable (Hyvärinen et al., 2019). Specifically, let  $\theta = (\mathbf{f}, \mathbf{T}, \lambda) \in \Theta$  be the parameters of the conditional generative model

$$p_\theta(\mathbf{X}, \mathbf{Z}|\mathbf{U}) = p_f(\mathbf{X}|\mathbf{Z})p_{\mathbf{T},\lambda}(\mathbf{Z}|\mathbf{U}), \quad (1)$$

where  $p_f(\mathbf{X}|\mathbf{Z}) = p_\epsilon(\mathbf{X} - \mathbf{f}(\mathbf{Z}))$  in which  $\epsilon$  is an independent noise variable with probability density function  $p_\epsilon(\epsilon)$ . Importantly, the prior  $p_{\mathbf{T},\lambda}(\mathbf{Z}|\mathbf{U})$  is assumed to be conditionally factorial, where each element of  $Z_i \in \mathbf{Z}$  has a univariate exponential family distribution given  $\mathbf{U}$ . The conditioning on  $\mathbf{U}$  is through an arbitrary function  $\lambda(\mathbf{U})$  (e.g., a neural net) that outputs the individual exponential family parameters  $\lambda_i(\mathbf{U})$  for each  $Z_i$ . The prior probability density thus takes the form

$$p_{\mathbf{T},\lambda}(\mathbf{Z}|\mathbf{U}) = \prod_i Q_i(Z_i)/C_i(\mathbf{U}) \exp \left[ \sum_{j=1}^k T_{i,j}(Z_i)\lambda_{i,j}(\mathbf{U}) \right], \quad (2)$$

where  $Q_i$  is the base measure,  $Z_i$  the  $i$ -th dimension of  $\mathbf{Z}$ ,  $C_i(\mathbf{U})$  the normalizing constant,  $\mathbf{T}_i = (T_{i,1}, \dots, T_{i,k})$  the sufficient statistics,  $\lambda_i(\mathbf{U}) = (\lambda_{i,1}(\mathbf{U}), \dots, \lambda_{i,k}(\mathbf{U}))$  the corresponding natural parameters depending on  $\mathbf{U}$ , and  $k$  the dimension of each sufficient statistic that is fixed in advance. It is worth noting that this prior is restrictive as it is factorial and therefore cannot capture dependencies. As in VAEs, the model parameters are estimated by maximizing the corresponding evidence lower bound (ELBO),

$$\mathcal{L}_{\text{iVAE}}(\theta, \phi) := \mathbb{E}_{p_D} \left[ \mathbb{E}_{q_\phi(\mathbf{Z}|\mathbf{X}, \mathbf{U})} [\log p_f(\mathbf{X}|\mathbf{Z}) + \log p_{\mathbf{T},\lambda}(\mathbf{Z}|\mathbf{U}) - \log q_\phi(\mathbf{Z}|\mathbf{X}, \mathbf{U})] \right], \quad (3)$$

<sup>1</sup>In fact, we are not restricted to the classification case and allow the target to be either continuous or categorical, which will be formally defined in Section 2.2.

where we denote by  $p_D$  the empirical data distribution given by the dataset  $\mathcal{D} = \{(\mathbf{X}^{(i)}, \mathbf{U}^{(i)})\}_{i=1}^N$  and  $q_\phi(\mathbf{Z}|\mathbf{X}, \mathbf{U})$  denotes an approximate conditional distribution for  $\mathbf{Z}$  given by a recognition network with parameters  $\phi$ . This approach is called identifiable VAE (iVAE). Most importantly, it can be proved that under the conditions stated in Theorem 2 of Khemakhem et al. (2020a), iVAE can identify the latent variables  $\mathbf{Z}$  up to a permutation and a simple componentwise transformation, see Appendix F.

## 2.2 INVARIANT RISK MINIMIZATION

Arjovsky et al. (2019) introduced invariant risk minimization (IRM), whose goal is to construct an *invariant predictor*  $f$  that performs well across all environments  $\mathcal{E}_{all}$  by exploiting data collected from multiple environments  $\mathcal{E}_{tr}$ , where  $\mathcal{E}_{tr} \subseteq \mathcal{E}_{all}$ . Technically, they consider datasets  $D_e := \{(\mathbf{x}_i^e, \mathbf{y}_i^e)\}_{i=1}^{n_e}$  from multiple training environments  $e \in \mathcal{E}_{tr}$ , where  $\mathbf{x}_i^e \in \mathcal{X} \subseteq \mathbb{R}^d$  is the input observation and its corresponding label is  $\mathbf{y}_i^e \in \mathcal{Y} \subseteq \mathbb{R}^s$ .<sup>2</sup> The dataset  $D_e$ , collected from environment  $e$ , consists of examples identically and independently distributed according to some probability distribution  $P(\mathbf{X}^e, \mathbf{Y}^e)$ . The goal of IRM is to use these multiple datasets to learn a predictor  $\mathbf{Y} = f(\mathbf{X})$  that performs well for all the environments. Here we define the risk reached by  $f$  in environment  $e$  as  $R^e(f) = \mathbb{E}_{\mathbf{X}^e, \mathbf{Y}^e} [\ell(f(\mathbf{X}^e), \mathbf{Y}^e)]$ , where  $\ell(\cdot)$  is a loss function. Then, the invariant predictor can be formally defined as follows:

**Definition 1** (Invariant Predictor (Arjovsky et al., 2019)). *We say that a data representation  $\Phi \in \mathcal{H}_\Phi : \mathcal{X} \rightarrow \mathcal{F}$  elicits an invariant predictor  $w \circ \Phi$  across environments  $\mathcal{E}$  if there is a classifier  $w \in \mathcal{H}_w : \mathcal{F} \rightarrow \mathcal{Y}$  simultaneously optimal for all environments, that is,  $w \in \arg \min_{\bar{w} \in \mathcal{H}_w} R^e(\bar{w} \circ \Phi)$  for all  $e \in \mathcal{E}$ , where  $\circ$  means function composition.*

Mathematically, IRM can be phrased as the following constrained optimization problem:

$$\min_{\Phi \in \mathcal{H}_\Phi, w \in \mathcal{H}_w} \sum_{e \in \mathcal{E}_{tr}} R^e(w \circ \Phi) \quad \text{s.t. } w \in \arg \min_{\bar{w} \in \mathcal{H}_w} R^e(\bar{w} \circ \Phi), \forall e \in \mathcal{E}_{tr}. \quad (4)$$

Since this is a generally infeasible bi-leveled optimization problem, Arjovsky et al. (2019) rephrased it as a tractable penalized optimization problem by transferring the inner optimization routine to a penalty term. The main generalization result (Theorem 9 in Arjovsky et al. (2019)) states that if both  $\Phi$  and  $w$  come from the class of linear models (i.e.,  $\mathcal{H}_\Phi = \mathbb{R}^{n \times n}$  and  $\mathcal{H}_w = \mathbb{R}^{n \times 1}$ ), under certain conditions on the diversity of training environments (Assumption 8 in Arjovsky et al. (2019)) and the data generation, the invariant predictor  $w \circ \Phi$  obtained by solving Eq. (4) remains invariant in  $\mathcal{E}_{all}$ .

## 3 PROBLEM SETUP

### 3.1 A MOTIVATING EXAMPLE

In this section, we extend the example which was introduced by Wright (1921) and discussed by Arjovsky et al. (2019), and provide a further in-depth analysis.

**Model 1.** *Consider a structural equation model (SEM) with a discrete environment variable  $E$  that modulates the noises in the structural assignments connecting the other variables (cf. Fig. 1a below):  $Z_1 \leftarrow \text{Gaussian}(0, \sigma_1(E))$ ,  $Y \leftarrow Z_1 + \text{Gaussian}(0, \sigma_2(E))$ ,  $Z_2 \leftarrow Y + \text{Gaussian}(0, \sigma_3(E))$ , where  $\text{Gaussian}(0, \sigma)$  denotes a Gaussian random variable with zero mean and standard deviation  $\sigma$ , and  $\sigma_1, \dots, \sigma_3$  are functions of the value  $e \in \mathcal{E}_{all}$  taken by the environment variable  $E$ .*

To ease exposition, here we consider the simple scenario in which  $\mathcal{E}_{all}$  only contains all modifications varying the noises of  $Z_1$ ,  $Z_2$  and  $Y$  within a finite range, i.e.,  $\sigma_i(e) \in [0, \sigma_{\max}^2]$ . Then, to predict  $Y$  from  $(Z_1, Z_2)$  using a least-square predictor  $\hat{Y}^e = \hat{\alpha}_1 Z_1^e + \hat{\alpha}_2 Z_2^e$  for environment  $e$ , we can

- Case 1: regress from  $Z_1^e$ , to obtain  $\hat{\alpha}_1 = 1$  and  $\hat{\alpha}_2 = 0$ ,
- Case 2: regress from  $Z_2^e$ , to obtain  $\hat{\alpha}_1 = 0$  and  $\hat{\alpha}_2 = \frac{\sigma_1(e) + \sigma_2(e)}{\sigma_1(e) + \sigma_2(e) + \sigma_3(e)}$ ,
- Case 3: regress from  $(Z_1^e, Z_2^e)$ , to obtain  $\hat{\alpha}_1 = \frac{\sigma_3(e)}{\sigma_2(e) + \sigma_3(e)}$  and  $\hat{\alpha}_2 = \frac{\sigma_2(e)}{\sigma_2(e) + \sigma_3(e)}$ .

In the generic scenario (i.e.,  $\sigma_1(e) \neq 0$ ,  $\sigma_2(e) \neq 0$ , and  $\sigma_3(e) \neq 0$ ), the regression using  $Z_1$  in Case 1 is an invariant correlation: it is the only regression whose coefficients do not vary with  $e$ . By contrast,

<sup>2</sup>The setup applies to both continuous and categorical data. If any observation or label is categorical, we could one-hot encode it.

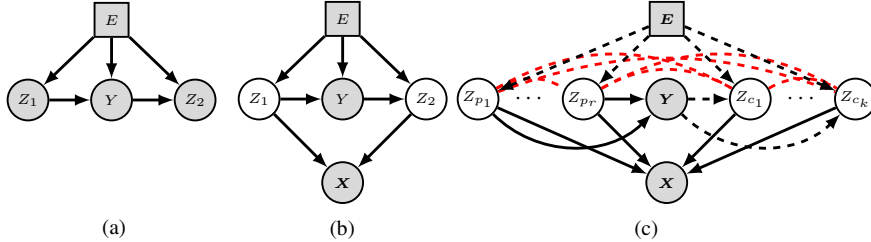


Figure 1: (a) Causal structure of Model 1. (b) A more practical extension of Model 1, where  $Z_1$  and  $Z_2$  are not directly observed and  $X$  is their observation. (c) A general version of (b), where we assume there exist multiple unobserved variables. Each of them could be either a parent, a child of  $Y$ , or has no direct connection with  $Y$ . We allow for arbitrary connections between the latent variables (red dashed lines) as long as the resulting causal diagram including  $Y$  is a directed acyclic graph (DAG). Grey nodes denote observed variables and white nodes represent unobserved variables. Dashed lines denote the edges which might vary across environments and even be absent in some scenarios, whilst solid lines indicate that they are invariant across all the environments.

the regressions in both Case 2 and Case 3 have coefficients that depend on  $e$ . Therefore, only the invariant correlation in Case 1 will generalize well to new test environments.

Another way to understand Model 1 is through its graphical representation<sup>3</sup>, as shown in Fig. 1a. We treat the environment as a random variable  $E$ , where  $E$  could be any information specific to the environment (Storkey, 2009; Peters et al., 2015; Zhang et al., 2017; Huang et al., 2020). Unless stated otherwise, for simplicity, we let  $E$  be the environment index, i.e.,  $E \in \{1, \dots, N\}$ , where  $N$  is the number of training environments. A more realistic version appearing in many settings is shown in Fig. 1b, where the true variables  $\{Z_1, Z_2\}$  are unobserved and we can only observe their transformation  $X$ . In this case, Invariant Causal Prediction (ICP) (Peters et al., 2015) will fail when applied to  $X$ , even when  $Y$  is not affected by  $E$  (i.e., the edge  $E \rightarrow Y$  is removed). The reason is that each variable (i.e., each dimension) of  $X$  is jointly influenced by both  $Z_1$  and  $Z_2$  so that ICP is unable to find the variables containing the information only about  $Z_1$  by searching for a subset of variables  $X$ . By contrast, both IRM and IRMG work, as long as the transformation is linear. These findings are also empirically illustrated in Appendix K.1. We now go even further and consider a more general causal graph in which  $Y$  can have more than one parent or child.

### 3.2 ASSUMPTIONS ON THE CAUSAL GRAPH

We extend the causal graph in Fig. 1b to a more general setting<sup>4</sup>, as encapsulated in Fig. 1c. In particular, we now have  $X \in \mathcal{X} \subseteq \mathbb{R}^d$ ,  $Y \in \mathcal{Y} \subseteq \mathbb{R}^s$ ,  $Z = (Z_{p_1}, \dots, Z_{p_r}, Z_{c_1}, \dots, Z_{c_k}) \in \mathcal{Z} \subseteq \mathbb{R}^n$ , where  $n = r + k$ , and  $\{Z_i\}_{i \in I_p = \{p_1, \dots, p_r\}}$  and  $\{Z_j\}_{j \in I_c = \{c_1, \dots, c_k\}}$  are multiple scalar causal factors and non-causal factors<sup>5</sup> of  $Y$ , respectively. We denote  $Z_p \doteq (Z_{p_1}, \dots, Z_{p_r})$  and  $Z_c \doteq (Z_{c_1}, \dots, Z_{c_k})$  for the ease of clarification. We also assume that  $Z$  is of lower dimension than  $X$ , that is,  $n \leq d$ . We allow for arbitrary connections between the latent variables  $Z$  as long as the resulting causal diagram including  $Y$  is a directed acyclic graph (DAG). We use dashed lines to indicate the causal mechanisms which might vary across environments and even be absent in some scenarios, whilst solid lines indicate that they are invariant across all the environments. To sum up, we assume that the underlying causal graph encapsulated in Fig. 1c satisfies the following assumption<sup>6</sup>:

**Assumption 1.** (a)  $Z_i$  depends on one or both of  $Y$  and  $E$  for any  $i$ ; (b) The causal graph containing  $Z$  and  $Y$  is a DAG; (c)  $X \perp\!\!\!\perp Y, E | Z$ , implying that  $p(X|Z)$  is invariant across all the environments; (d)  $Y \perp\!\!\!\perp E | Z_p$ , implying that  $p(Y|Z_p)$  is invariant across all the environments.

One may be wondering how practical Assumption 1 is in real world applications. Let us explore this in more detail. Assumption 1a rules out all the useless  $Z_i$  in the task of predicting  $Y$ . This is because if Assumption 1a is violated, meaning that  $Z_i$  is independent of  $Y$  and  $E$  and has no influence in predicting  $Y$ , then such  $Z_i$  should be viewed as noise and thus eliminated during learning. Assumption 1b is a common assumption in causal discovery (Spirtes et al., 2000; Pearl, 2009; Peters et al., 2015). It also makes sense in Assumption 1c that the generative mechanism  $p(X|Z)$  is

<sup>3</sup>The relation between SEM and its graphical representation is formally defined in Appendix D.

<sup>4</sup>For simplicity, we do not explicitly consider unobserved confounders in this paper. In particular, we assume that there are no unobserved confounders between  $Z, Y, X$ , and  $E$ .

<sup>5</sup>This means that  $Z_j \in I_c$  could be either an effect of  $Y$ , independent of  $Y$ , or spuriously correlated with  $Y$  via a third set of confounders (i.e., both  $Z_j$  and  $Y$  are affected by a subset of  $\{Z_i\}_{i \neq j}$  and  $E$ ).

<sup>6</sup>For generality, we replace  $E$  with  $\mathbf{E}$  to additionally include the case of multi-dimensional variables.

**Algorithm 1:** Invariant Causal Representation Learning (iCaRL)

**Phase 1:** We first learn a NF-iVAE model, including the decoder and its corresponding encoder, by optimizing the objective function in (10) on the data  $\{\mathbf{X}, \mathbf{Y}, \mathbf{E}\}$ . Then, we use the mean of the NF-iVAE encoder to infer the latent variables  $\mathbf{Z}$  from observations  $\{\mathbf{X}, \mathbf{Y}, \mathbf{E}\}$ . The latent variables are guaranteed to be identified up to a permutation and simple transformation.

**Phase 2:** After inferring  $\mathbf{Z}$ , we first conduct the PC algorithm to learn a Markov equivalence class of DAGs, and then discover direct causes (parents) of  $\mathbf{Y}$  among its neighbors by testing all pairs of latent variables with (conditional) independence testing, i.e., finding a set of latent variables in which each pair of  $Z_i$  and  $Z_j$  satisfies that the dependency between them increases after additionally conditioning on  $\mathbf{Y}$ .

**Phase 3:** Having obtained  $\text{Pa}(\mathbf{Y})$ , we can solve (11) to learn the invariant classifier  $w$ . When in a new environment, we first infer  $\text{Pa}(\mathbf{Y})$  from  $\mathbf{X}$  by solving (12) and then leverage the learned  $w$  for prediction.

invariant across all the environments. Otherwise, it is impossible to infer  $\mathbf{Z}$  from  $\mathbf{X}$  in any unseen environment. Assumption 1d is a widely-used default assumption in OOD generalization (Peters et al., 2015; Arjovsky et al., 2019). In fact, Assumption 1d can be further relaxed to the more practical one that  $\mathbb{E}[\mathbf{Y}|\mathbf{Z}_p]$  is invariant across all the environments. That is, given  $\mathbf{Z}_p$ , we allow  $\mathbf{E}$  to only affect the amount of noise in the distribution of  $\mathbf{Y}$ , because that would not change the expected value of  $\mathbf{Y}$  since the mean of the noise would be zero. Apparently, Assumption 1, together with the causal graph in Fig. 1c, covers most scenarios (e.g., the ones of Zhang et al. (2013); von Kügelgen et al. (2020); Sun et al. (2020); Ahuja et al. (2021); von Kügelgen et al. (2021)) and is a very flexible model for causal analysis when predicting  $\mathbf{Y}$  from  $\mathbf{X}$ .

### 3.3 ASSUMPTIONS ON THE PRIOR

When the underlying causal graph satisfies Assumption 1, our primary assumption leading to identifiability in this general setting is that the conditional prior  $p(\mathbf{Z}|\mathbf{Y}, \mathbf{E})$  belongs to a general exponential family. This is formalized as follows:

**Assumption 2.**  $p(\mathbf{Z}|\mathbf{Y}, \mathbf{E})$  belongs to a general exponential family with parameter vector given by an arbitrary function  $\lambda(\mathbf{Y}, \mathbf{E})$  and sufficient statistics  $\mathbf{T}(\mathbf{Z}) = [\mathbf{T}_f(\mathbf{Z})^T, \mathbf{T}_{NN}(\mathbf{Z})^T]^T$  given by the concatenation of a) the sufficient statistics  $\mathbf{T}_f(\mathbf{Z}) = [\mathbf{T}_1(Z_1)^T, \dots, \mathbf{T}_n(Z_n)^T]^T$  of a factorized exponential family, where all the  $\mathbf{T}_i(Z_i)$  have dimension larger or equal to 2, and b) the output  $\mathbf{T}_{NN}(\mathbf{Z})$  of a neural network with ReLU activations. The resulting density function is thus given by

$$p_{\mathbf{T}, \lambda}(\mathbf{Z}|\mathbf{Y}, \mathbf{E}) = \mathcal{Q}(\mathbf{Z})/\mathcal{C}(\mathbf{Y}, \mathbf{E}) \exp[\mathbf{T}(\mathbf{Z})^T \lambda(\mathbf{Y}, \mathbf{E})], \quad (5)$$

where  $\mathcal{Q}$  is the base measure and  $\mathcal{C}$  the normalizing constant.

A neural network with ReLU activation has universal approximation power. Therefore, the term  $\mathbf{T}_{NN}(\mathbf{Z})$  in the above prior distribution will allow us to capture arbitrary dependencies between the latent variables. The distribution in Eq. (5) is more flexible than the conditionally factorized prior assumed by iVAEs. However, surprisingly, the identifiability results of iVAEs also hold when using the more flexible prior in Eq. (5), as we will show in Section 4.1. This motivates using an extended iVAE model with the above prior to model data generated by the ground truth model in Fig. 1c. However, in this case, the data generating model and the learned model might have different priors. For example, in the ground truth model, the prior for each  $Z_{i \in I_p}$  might be  $p(Z_i|\mathbf{E})$ , when  $Z_i$  is only caused by  $\mathbf{E}$ . By contrast, in the extended iVAE model the prior is  $p(\mathbf{Z}|\mathbf{Y}, \mathbf{E})$ . Is this going to affect the identifiability of the latent variables? Well, in practice not because the posterior distribution for  $\mathbf{Z}$  given  $\mathbf{X}, \mathbf{Y}$  and  $\mathbf{E}$  would be equivalent in both models (up to identifiability guarantees).

## 4 INVARIANT CAUSAL REPRESENTATION LEARNING

We now introduce our algorithm, invariant Causal Representation Learning (iCaRL), which consists of 3 phases as summarized in Algorithm 1. The idea is that we first identify  $\mathbf{Z}$  by using an extended iVAE model under Assumptions 1&2 (Phase 1), then discover direct causes of  $\mathbf{Y}$  among the identified  $\mathbf{Z}$  (Phase 2), and finally learn an invariant predictor for  $\mathbf{Y}$  from the discovered causes (Phase 3).

### 4.1 PHASE 1: IDENTIFYING LATENT VARIABLES USING NF-iVAE

In this section, we describe our identifiable model, namely NF-iVAE, which is an extended iVAE with a general non-factorized prior that is able to capture complex dependences between the latent variables. Technically, in the general setting under Assumption 1, it is straightforward to obtain a

corresponding generative model by directly substituting  $U$  with  $(Y, E)$  in Eq. (1):

$$p_{\theta}(\mathbf{X}, \mathbf{Z}|\mathbf{Y}, \mathbf{E}) = p_f(\mathbf{X}|\mathbf{Z})p_{T,\lambda}(\mathbf{Z}|\mathbf{Y}, \mathbf{E}), \quad (6)$$

$$p_f(\mathbf{X}|\mathbf{Z}) = p_{\epsilon}(\mathbf{X} - \mathbf{f}(\mathbf{Z})). \quad (7)$$

The corresponding ELBO is

$$\mathcal{L}_{\text{phase1}}^{\text{ELBO}}(\theta, \phi) := \mathbb{E}_{p_D} \left[ \mathbb{E}_{q_{\phi}(\mathbf{Z}|\mathbf{X}, \mathbf{Y}, \mathbf{E})} [\log p_f(\mathbf{X}|\mathbf{Z}) + \log p_{T,\lambda}(\mathbf{Z}|\mathbf{Y}, \mathbf{E}) - \log q_{\phi}(\mathbf{Z}|\mathbf{X}, \mathbf{Y}, \mathbf{E})] \right]. \quad (8)$$

To obtain an identifiability result, we assume that the prior  $p_{T,\lambda}(\mathbf{Z}|\mathbf{Y}, \mathbf{E})$  satisfies Assumption 2 (i.e., Eq. (5)). Since the prior is a general multivariate exponential family distribution with unknown normalization constant, we cannot learn its parameters  $(T, \lambda)$  by directly maximizing Eq. (8). Instead, we use score matching, a well-known method for training unnormalized probabilistic models (Hyvärinen, 2005; Vincent, 2011), and learn  $(T, \lambda)$  by minimizing

$$\mathcal{L}_{\text{phase1}}^{\text{SM}}(T, \lambda) := \mathbb{E}_{p_D} \left[ \mathbb{E}_{q_{\phi}(\mathbf{Z}|\mathbf{X}, \mathbf{Y}, \mathbf{E})} [|\nabla_{\mathbf{Z}} \log q_{\phi}(\mathbf{Z}|\mathbf{X}, \mathbf{Y}, \mathbf{E}) - \nabla_{\mathbf{Z}} \log p_{T,\lambda}(\mathbf{Z}|\mathbf{Y}, \mathbf{E})|^2] \right]. \quad (9)$$

In practice, we can use a simple trick of partial integration to simplify the evaluation of Eq. (9), see Appendix C. Furthermore, we can jointly learn  $(\theta, \phi)$  by combining Eq. (8) and Eq. (9) in the following objective:

$$\mathcal{L}_{\text{phase1}}(\theta, \phi) = \mathcal{L}_{\text{phase1}}^{\text{ELBO}}(\mathbf{f}, \hat{T}, \hat{\lambda}, \phi) - \mathcal{L}_{\text{phase1}}^{\text{SM}}(\hat{\mathbf{f}}, T, \lambda, \hat{\phi}), \quad (10)$$

where  $\hat{\mathbf{f}}, \hat{T}, \hat{\lambda}, \hat{\phi}$  are copies of  $\mathbf{f}, T, \lambda, \phi$  that are treated as constants and whose gradient is not calculated during learning. More details can be found in Appendix M.

We now state our main theoretical results:

**Theorem 1.** *Assume that we observe data sampled from a generative model defined according to Eqs. (5-7), with parameters  $\theta := (\mathbf{f}, T, \lambda)$ , where  $p_{T,\lambda}(\mathbf{Z}|\mathbf{Y}, \mathbf{E})$  satisfies Assumption 2. Furthermore, assume the following holds: (i) The set  $\{\mathbf{X} \in \mathcal{O} | \varphi_{\epsilon}(\mathbf{X}) = 0\}$  has measure zero, where  $\varphi_{\epsilon}$  is the characteristic function of the density  $p_{\epsilon}$  defined in Eq. (7). (ii) Function  $\mathbf{f}$  in Eq. (7) is injective, and has all second-order cross derivatives. (iii) The sufficient statistics in  $T_f$  are all twice differentiable. (iv) There exist  $k + 1$  distinct points  $(\mathbf{Y}, \mathbf{E})^0, \dots, (\mathbf{Y}, \mathbf{E})^k$  such that the matrix  $L = (\lambda((\mathbf{Y}, \mathbf{E})^1) - \lambda((\mathbf{Y}, \mathbf{E})^0), \dots, \lambda((\mathbf{Y}, \mathbf{E})^k) - \lambda((\mathbf{Y}, \mathbf{E})^0))$  of size  $k \times k$  is invertible, where  $k$  is the dimension of  $T$ . Then the parameters  $\theta$  are identifiable up to a permutation and a “simple transformation” of the latent variables  $Z$ , defined as a componentwise nonlinearity making each recovered  $T_i(Z_i)$  in  $T_f(Z)$  equal to the original up to a linear operation.*

Note that, this theorem is inspired by but beyond the main results of iVAEs in that the former is predicated on Assumption 2 which is more flexible than the conditionally factorized prior assumed in iVAEs. It results in several key changes in the proof, clarified in Appendix H. Interestingly, from (iv) we can further see that  $E$  is unnecessary when there exist  $k + 1$  distinct points  $\mathbf{Y}^0, \dots, \mathbf{Y}^k$  such that the matrix  $L = (\lambda(\mathbf{Y}^1) - \lambda(\mathbf{Y}^0), \dots, \lambda(\mathbf{Y}^k) - \lambda(\mathbf{Y}^0))$  of size  $k \times k$  is invertible. Not requiring  $E$  would make our approach even more applicable.

We further have the following consistency result for the estimation.

**Theorem 2.** *Assume that the following holds: (i) The family of distributions  $q_{\phi}(\mathbf{Z}|\mathbf{X}, \mathbf{Y}, \mathbf{E})$  contains  $p_{\theta}(\mathbf{Z}|\mathbf{X}, \mathbf{Y}, \mathbf{E})$ , and  $q_{\phi}(\mathbf{Z}|\mathbf{X}, \mathbf{Y}, \mathbf{E}) > 0$  everywhere. (ii) We maximize  $\mathcal{L}_{\text{phase1}}(\theta, \phi)$  with respect to both  $\theta$  and  $\phi$ . Then in the limit of infinite data, we learn the true parameters  $\theta^*$  up to a permutation and simple transformation of the latent variables  $Z$ .*

As a consequence of Theorems 1&2, we have:

**Theorem 3.** *Assume the hypotheses of Theorem 1 and Theorem 2 hold, then in the limit of infinite data, we identify the true latent variables  $Z^*$  up to a permutation and simple transformation.*

Theorem 3 states that we can use NF-iVAE to infer the true  $Z^*$  up to a permutation and simple transformation. We use the mean of  $q_{\phi}(\mathbf{Z}|\mathbf{X}, \mathbf{Y}, \mathbf{E})$  for this task. Note that the noise  $\epsilon$  may introduce uncertainty in the estimation of  $Z$ . However, when  $X$  is high dimensional and  $Z$  is low dimensional (as common in real world applications),  $q_{\phi}(\mathbf{Z}|\mathbf{X}, \mathbf{Y}, \mathbf{E})$  will be highly concentrated and we will still be able to estimate  $Z$  with high accuracy. The good results obtained by our method in various experiments seem to corroborate this. All three theorems are proven in Appendix H.

## 4.2 PHASE 2: DISCOVERING DIRECT CAUSES

After estimating  $Z$  for each data point, the next step is to determine which components of  $Z$  are direct causes of  $Y$ . We denote these components by  $\text{Pa}(Y)$ . We first conduct the PC algorithm (Spirtes et al., 2000) to learn a Markov equivalence class of DAGs, which gives us the direct neighbors of  $Y$ ,

denoted by  $\text{Ne}(\mathbf{Y})$ . Then, from Assumption 1, one observation is that in the generic case, for any two latent variables  $Z_i$  and  $Z_j$  from  $\text{Ne}(\mathbf{Y})$ , only when both are causes of  $\mathbf{Y}$  do we have that the dependency between them increases after additionally conditioning on  $\mathbf{Y}$ . Thus, when there exist at least two causal latent variables in  $\text{Ne}(\mathbf{Y})$ , we can test all pairs of latent variables<sup>7</sup> with conditional independence testing<sup>8</sup> (Zhang et al., 2012) to discover  $\text{Pa}(\mathbf{Y})$  by comparing  $p$ -values from the two tests:  $\text{IndTest}(Z_i, Z_j | \mathbf{E})$  and  $\text{IndTest}(Z_i, Z_j | \mathbf{Y}, \mathbf{E})$ , where  $\text{IndTest}$  denotes (conditional) independence test. Conversely, if no such a pair is found, it implies that there is at most one causal latent variable. This is a highly unlikely case in real world applications, which is left to Appendix G.

### 4.3 PHASE 3: LEARNING AN INVARIANT PREDICTOR

After having obtained the causal latent variables  $\text{Pa}(\mathbf{Y})$  for  $\mathbf{Y}$  across training environments, we can learn  $w$  by solving the following optimization problem:

$$\min_{w \in \mathcal{H}_w} \sum_{e \in \mathcal{E}_{tr}} R^e(w) = \min_{w \in \mathcal{H}_w} \sum_{e \in \mathcal{E}_{tr}} \mathbb{E}_{\text{Pa}(\mathbf{Y}^e), \mathbf{Y}^e} [\ell_1(w(\text{Pa}(\mathbf{Y}^e)), \mathbf{Y}^e)], \quad (11)$$

where  $\ell_1(\cdot)$  could be any loss. Since we assume that  $\mathbb{E}(\mathbf{Y} | \text{Pa}(\mathbf{Y}))$  is invariant across  $\mathcal{E}_{all}$  (the relaxed version of Assumption 1d), the learned  $w$  is guaranteed to perform well across  $\mathcal{E}_{all}$ .

The remaining question is how to infer  $\text{Pa}(\mathbf{Y})$  (i.e.,  $\mathbf{Z}_p$ ) from  $\mathbf{X}$  in a new environment. This can be implemented by leveraging the learned  $p(\mathbf{X} | \mathbf{Z})$ . The rationale behind is that  $p(\mathbf{X} | \mathbf{Z})$  is assumed to be invariant across  $\mathcal{E}_{all}$  (Assumption 1c). In light of this idea, we follow Sun et al. (2020) and infer  $\mathbf{Z}_p$  from  $\mathbf{X}$  in any new testing environment by solving the following optimization problem:

$$\max_{\mathbf{Z}_p, \mathbf{Z}_c} \log p_f(\mathbf{X} | \mathbf{Z}_p, \mathbf{Z}_c) + \lambda_1 \|\mathbf{Z}_p\|_2^2 + \lambda_2 \|\mathbf{Z}_c\|_2^2, \quad (12)$$

where the hyperparameters  $\lambda_1 > 0$  and  $\lambda_2 > 0$  control the learned  $\mathbf{Z}_p$  and  $\mathbf{Z}_c$  in a reasonable scale, and both are selected on training/validation data. For optimization, we follow Schott et al. (2018) to first use values of  $\mathbf{Z}$  sampled from the training set as initial points and then use Adam to optimize for several iterations.<sup>9</sup> Note that the noise  $\epsilon$  will introduce uncertainty in the estimation of  $\mathbf{Z}_p$  and  $\mathbf{Z}_c$  from  $\mathbf{X}$ . However, as we mentioned before (below Theorem 3), this noise is not going to affect the estimation much because the likelihood will be highly concentrated around the ground truth values, as corroborated by our good empirical results.

A key question is if iCaRL performs well across  $\mathcal{E}_{all}$  even though it uses only data from  $\mathcal{E}_{tr}$ . That is, does iCaRL enable OOD generalization, as defined by Arjovsky et al. (2019)? The answer is positive since Theorem A.1 in Arjovsky (2021) indicates that i) any predictor  $w \circ \Phi$  with optimal OOD generalization uses only  $\text{Pa}(\mathbf{Y})$  to compute  $\Phi$  and ii) the classifier  $w$  in this optimal predictor can be estimated using data from any environment  $e$  for which the distribution of  $\text{Pa}(\mathbf{Y}^e)$  has full support, which will always be the case since the conditional prior in Eq. (5) has full support. Finally, Theorem A.1 in Arjovsky (2021) also indicates that iii) the optimal predictor will be invariant across  $\mathcal{E}_{all}$ . Key to these results is that  $\text{Pa}(\mathbf{Y}^e)$  are available when solving (12). This requires first to identify the latent variables  $\mathbf{Z}$  from  $\mathbf{X}$ ,  $\mathbf{Y}$  and  $\mathbf{E}$  and second, to discover the direct causes of  $\mathbf{Y}$ . The hypotheses of Theorems 1 and 2 and Assumption 1 provide this guarantee. We therefore have the following result whose proof is in Appendix H.

**Proposition 1.** *Under Assumption 1 and the assumptions of Theorems 1 and 2, the predictor learned by iCaRL across  $\mathcal{E}_{tr}$  in the limit of infinite data has optimal OOD generalization across  $\mathcal{E}_{all}$ .*

## 5 EXPERIMENTS

We compare our approach with a variety of methods on both synthetic and real-world datasets. In all comparisons, unless stated otherwise, we average performance over ten runs. **Due to space constraints, we only highlight some key results while pointing to the extensive appendices for more information.** The supplement contains all the details of the experiments, e.g., datasets (Appendix J), implementation (Appendix M), hyperparameters and architectures (Appendix N), etc.

### 5.1 SYNTHETIC DATA

To verify the identifiability of NF-iVAE, we conduct a series of experiments on synthetic data generated according to the causal graph shown in Fig. 2e. Details of the ground truth data generating

<sup>7</sup>We only need to consider those variables in  $\text{Ne}(\mathbf{Y})$  whose edges connecting to  $\mathbf{Y}$  are not oriented by PC.

<sup>8</sup>These conditional independence tests can be performed in parallel to largely accelerate the testing procedure. See more in Appendix G.

<sup>9</sup>Note that, (12) can be optimized either independently for each data point or for a number of data points.

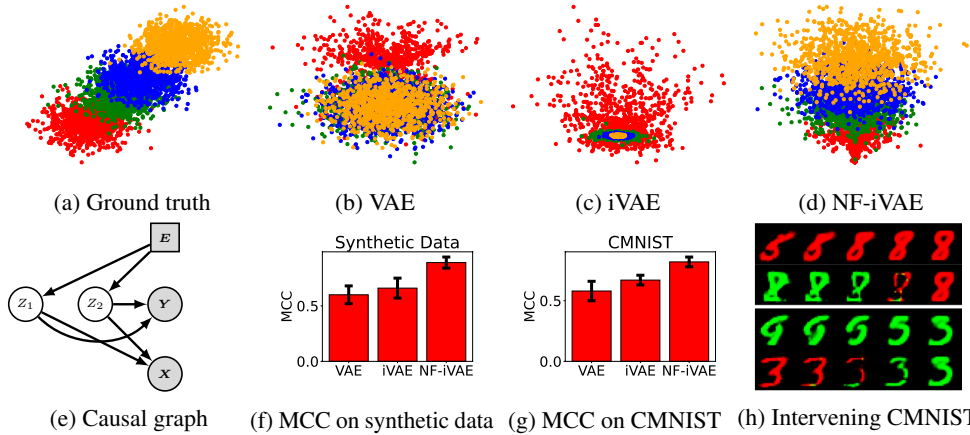


Figure 2: (a-d) Visualization of the samples (i.e.,  $\hat{\mathbf{Z}} = (\hat{Z}_1, \hat{Z}_2)$ ) in latent space recovered through different algorithms: (a) Samples from the true distribution; (b-c) Samples from the posterior inferred using VAE and iVAE, respectively. Apparently, our method (d) can recover the original data up to a permutation and a simple componentwise transformation. (e) The causal structure with  $\mathbf{Y}$  having two causes describes the data generating process of the synthetic dataset. (f) Mean correlation coefficient (MCC) scores for VAE, iVAE, and NF-iVAE on synthetic data. (g) MCC scores for VAE, iVAE, and NF-iVAE on CMNIST. (h) The effects on the CMNIST images of digit 8 (top two rows) and digit 3 (bottom two rows) when intervening on a causal factor  $Z_{i \in I_p}$  and on a non-causal factor (effect)  $Z_{j \in I_e}$ , respectively.

process are given in Appendix K. The reason we choose this setting is that it is the simplest case satisfying our requirements: a) For ease of visualization, the latent space had better be 2-dimensional; b) To introduce the non-factorized prior given  $\mathbf{Y}$  and  $\mathbf{E}$  (i.e.,  $Z_i \not\perp\!\!\!\perp Z_j | \mathbf{Y}, \mathbf{E}$ ),  $\mathbf{Y}$  has at least two causes. We draw 1000 samples from each of the four environments  $\mathbf{E} = \{0.2, 2, 3, 5\}$ , and thus the whole synthetic dataset consists of 4000 samples. The task is to recover the true latent variable  $\mathbf{Z} = (Z_1, Z_2)$  using the samples of  $\mathbf{X}$ ,  $\mathbf{E}$ , and  $\mathbf{Y}$ . We compare with two widely-used baselines: VAE (Kingma & Welling, 2013) (without identifiability guarantees) and iVAE (Khemakhem et al., 2020a) (with a conditionally factorized prior for identifiability). Through the aforementioned theoretical analysis, it is evident that our method has a more general assumption on the prior leading to identifiability. This is demonstrated empirically in Figs. 2b-2d. Our method NF-iVAE can recover the original data  $\mathbf{Z}$  up to a permutation and a simple componentwise transformation, whereas all the other methods fail because they are unable to handle the non-factorized case in which  $Z_i \not\perp\!\!\!\perp Z_j | \mathbf{Y}, \mathbf{E}$ . We also compute the mean correlation coefficient (MCC) used in Khemakhem et al. (2020a), which can be obtained by calculating the correlation coefficient between all pairs of true and recovered latent factors and then solving a linear sum assignment problem by assigning each recovered latent factor to the true latent factor with which it best correlates. By definition, higher MCC scores indicate stronger identifiability. From Fig. 2f, we can see that the MCC score for NF-iVAE is significantly greater than those of VAE and iVAE, indicating much stronger identifiability. Note that, in Appendix K we additionally compare with more methods, whose differences are further summarized in a table.

## 5.2 COLORED MNIST, COLORED FASHION MNIST, AND VLCS

In this section, we first report experiments on two datasets used in IRM and IRMG: Colored MNIST (CMNIST) and Colored Fashion MNIST (CFMNIST). We follow the same setting of Ahuja et al. (2020a) to create these two datasets (see the details in Appendix J). The task is to predict a binary label assigned to each image which is originally grayscale but artificially colored in a way that the color is correlated strongly but spuriously with the class label. For all the experiments on these two datasets, we set the number of the latent variables to  $n = 10$ .

Likewise, we investigate the identifiability of NF-iVAE on CMNIST by computing the MCC score between samples of the true latent variable and of the recovered latent variable. Since the true latent variable on CMNIST is inaccessible to us, we follow Khemakhem et al. (2020b) and compute an average MCC score between samples of latent variables recovered by different models trained with different random initialization. As shown in Fig. 2g, it is evident that the MCC score for NF-iVAE greatly outperforms the others, showing that the latent variables recovered by NF-iVAE have much better identifiability.

Furthermore, we demonstrate the ability of iCaRL to discover the causal latent variables (Phase 2) by visualizing the generated images through performing intervention upon a causal latent variable and a non-causal latent variable, respectively. Fig. 2h shows how intervening upon each of them affects the image. Obviously, intervening on a causal latent variable affects the shape of the digit but not its color



Table 1: Colored Fashion MNIST. Comparisons in terms of accuracy (%) (mean  $\pm$  std deviation).

METHOD	TRAIN	TEST
ERM	83.17 $\pm$ 1.01	22.46 $\pm$ 0.68
ERM 1	81.33 $\pm$ 1.35	33.34 $\pm$ 8.85
ERM 2	84.39 $\pm$ 1.89	13.16 $\pm$ 0.82
ROBUST MIN MAX	82.81 $\pm$ 0.11	29.22 $\pm$ 8.56
F-IRM GAME	62.31 $\pm$ 2.35	69.25 $\pm$ 5.82
V-IRM GAME	68.96 $\pm$ 0.95	70.19 $\pm$ 1.47
IRM	75.01 $\pm$ 0.25	55.25 $\pm$ 12.42
<b>iCaRL (ours)</b>	<b>74.96 <math>\pm</math> 0.37</b>	<b>73.61 <math>\pm</math> 0.63</b>
ERM GRAYSCALE	74.79 $\pm$ 0.37	74.67 $\pm$ 0.48
OPTIMAL	75	75

Table 2: VLCS. Comparisons in terms of accuracy (%) (mean  $\pm$  std deviation).

METHOD	TEST
ERM	77.4 $\pm$ 0.3
IRM	78.1 $\pm$ 0.0
DRO (Sagawa et al., 2019)	77.2 $\pm$ 0.6
Mixup (Yan et al., 2020)	77.7 $\pm$ 0.4
CORAL (Sun & Saenko, 2016)	77.7 $\pm$ 0.5
MMD (Li et al., 2018b)	76.7 $\pm$ 0.9
DANN (Ganin et al., 2016)	78.7 $\pm$ 0.3
C-DANN (Li et al., 2018c)	78.2 $\pm$ 0.4
LaCIM (Sun et al., 2020)	78.4 $\pm$ 0.5
<b>iCaRL (ours)</b>	<b>81.8 <math>\pm</math> 0.6</b>

(top plots), whilst intervening on a non-causal latent variable, which is an effect in Fig. 2h, affects the color of the digit only (bottom plots). This visually verifies the results of iCaRL in Phase 2.

In terms of the OOD generalization performance, we compare iCaRL with 1) IRM, 2) two variants of IRMG: F-IRM Game (with  $\Phi$  fixed to the identity) and V-IRM Game (with a variable  $\Phi$ ), 3) three variants of ERM: ERM (on entire training data), ERM  $e$  (on each environment  $e$ ), and ERM GRAYSCALE (on data with no spurious correlations), and 4) ROBUST MIN MAZ (minimizing the maximum loss across the multiple environments). Table 1 shows that iCaRL outperforms all other baselines on CFMNIST. It is worth emphasising that the train and test accuracies of iCaRL closely approach the ones of ERM GRAYSCALE and OPTIMAL, implying that iCaRL approximately learns the true invariant causal representations with almost no correlation with the spurious color feature. We can draw similar conclusions from the results on CMNIST (Appendix L).

We also report the results on one of the widely used realistic datasets for OOD generalization: VLCS (Fang et al., 2013). This dataset consists of 10, 729 photographic images of dimension (3, 224, 224) and 5 classes from four domains: Caltech101, LabelMe, SUN09, and VOC2007. We used the exact experimental setting that is described in Gulrajani & Lopez-Paz (2020). We provide results averaged over all possible train and test environment combination for one of the commonly used hyper-parameter tuning procedure: train domain validation. As shown in Table 2, iCaRL achieves state-of-the-art performance when compared to those most popular domain generalization alternatives. We further include experimental evidence on another popular dataset: PACS (Li et al., 2017a), all the details of which are placed in Appendix L.

## 6 RELATED WORK

Invariant Causal Prediction (ICP), aims to find the *causal feature set* (i.e., all direct causes of a target variable of interest) (Peters et al., 2015) by exploiting the invariance property in causality which has been discussed under the term “autonomy”, “modularity”, and “stability” (Haavelmo, 1944; Aldrich, 1989; Hoover, 1990; Pearl, 2009; Dawid et al., 2010; Schölkopf et al., 2012). This invariance property assumed in ICP and its nonlinear extension (Heinze-Deml et al., 2018) is limited, because no intervention is allowed on the target variable  $Y$ . Besides, ICP methods implicitly assume that the variables of interest  $Z$  are given. The works of Magliacane et al. (2018) and Subbaswamy et al. (2019) attempt to find invariant predictors that are maximally predictive using conditional independence tests and other graph-theoretic tools, both of which also assume that the  $Z$  are given and further assume that additional information about the structure over  $Z$  is known. Mitrovic et al. (2020) analyze data augmentations in self-supervised learning from the perspective of invariant causal mechanisms. Arjovsky et al. (2019) reformulate this invariance as an optimization-based problem, allowing us to learn an invariant data representation from  $X$  constrained to be a linear transformation of  $Z$ . The risks of this approach have been discussed in Rosenfeld et al. (2020); Kamath et al. (2021); Nagarajan et al. (2020) and its sample complexity is analyzed in Ahuja et al. (2020b).

Another line of related work is in the field of domain generalization, which we discuss in Appendix A.

## 7 CONCLUSION

We have proposed a novel framework to learn invariant predictors from a diverse set of training environments. It is based on a practical and general assumption: the prior over the data representation belongs to a general exponential family when conditioning on the target and the environment. This assumption leads to guarantees that the components in the representation can be identified up to a permutation and simple transformation. This allows us to discover all the direct causes of the target, which enables generalization guarantees in the nonlinear setting. We hope our framework will inspire new ways to address the OOD generalization problem through a causal lens.

## ACKNOWLEDGEMENTS

We are thankful to Wenlin Chen for contributing to Step III in the proof of Theorem 4 and for generalizing the proof in Theorem 5 to the case in which each  $T_i(Z_i)$  in  $T_f(\mathcal{Z})$  contains arbitrary sufficient statistics instead of just  $Z_i$  and  $Z_i^2$ . We also thank Ilyes Khemakhem, Aapo Hyvärinen, and Arthur Gretton for their helpful discussions, and the anonymous reviewers for their constructive comments on an earlier version of this paper.

## REFERENCES

- Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dandelion Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. URL <https://www.tensorflow.org/>. Software available from tensorflow.org.
- Kartik Ahuja, Karthikeyan Shanmugam, Kush Varshney, and Amit Dhurandhar. Invariant risk minimization games. *arXiv preprint arXiv:2002.04692*, 2020a.
- Kartik Ahuja, Jun Wang, Amit Dhurandhar, Karthikeyan Shanmugam, and Kush R Varshney. Empirical or invariant risk minimization? a sample complexity perspective. *arXiv preprint arXiv:2010.16412*, 2020b.
- Kartik Ahuja, Ethan Caballero, Dinghuai Zhang, Yoshua Bengio, Ioannis Mitliagkas, and Irina Rish. Invariance principle meets information bottleneck for out-of-distribution generalization. *arXiv preprint arXiv:2106.06607*, 2021.
- John Aldrich. Autonomy. *Oxford Economic Papers*, 41(1):15–34, 1989.
- Martín Arjovsky. Out of distribution generalization in machine learning. *CoRR*, abs/2103.02667, 2021.
- Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization. *arXiv preprint arXiv:1907.02893*, 2019.
- Yogesh Balaji, Swami Sankaranarayanan, and Rama Chellappa. Metareg: Towards domain generalization using meta-regularization. *Advances in Neural Information Processing Systems*, 31: 998–1008, 2018.
- Sara Beery, Grant Van Horn, and Pietro Perona. Recognition in terra incognita. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 456–473, 2018.
- Shai Ben-David, John Blitzer, Koby Crammer, Fernando Pereira, et al. Analysis of representations for domain adaptation. *Advances in neural information processing systems*, 19:137, 2007.
- Shai Ben-David, John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wortman Vaughan. A theory of learning from different domains. *Machine learning*, 79(1):151–175, 2010.
- Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1798–1828, 2013.
- Pierre Comon. Independent component analysis, a new concept? *Signal processing*, 36(3):287–314, 1994.
- Koby Crammer, Michael Kearns, and Jennifer Wortman. Learning from multiple sources. *Journal of Machine Learning Research*, 9(8), 2008.
- Povilas Daniusis, Dominik Janzing, Joris Mooij, Jakob Zscheischler, Bastian Steudel, Kun Zhang, and Bernhard Schölkopf. Inferring deterministic causal relations. *arXiv preprint arXiv:1203.3475*, 2012.

- A Philip Dawid, Vanessa Didelez, et al. Identifying the consequences of dynamic treatment strategies: A decision-theoretic overview. *Statistics Surveys*, 4:184–231, 2010.
- Zhengming Ding and Yun Fu. Deep domain generalization with structured low-rank constraint. *IEEE Transactions on Image Processing*, 27(1):304–313, 2017.
- Qi Dou, Daniel C Castro, Konstantinos Kamnitsas, and Ben Glocker. Domain generalization via model-agnostic learning of semantic features. *arXiv preprint arXiv:1910.13580*, 2019.
- Sarah Erfani, Mahsa Baktashmotlagh, Masud Moshtaghi, Xuan Nguyen, Christopher Leckie, James Bailey, and Rao Kotagiri. Robust domain generalisation by enforcing distribution invariance. In *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence (IJCAI-16)*, pp. 1455–1461. AAAI Press, 2016.
- Chen Fang, Ye Xu, and Daniel N Rockmore. Unbiased metric learning: On the utilization of multiple datasets and web images for softening bias. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 1657–1664, 2013.
- José AR Fonollosa. Conditional distribution variability measures for causality detection. In *Cause Effect Pairs in Machine Learning*, pp. 339–347. Springer, 2019.
- Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks. *The Journal of Machine Learning Research*, 17(1):2096–2030, 2016.
- Muhammad Ghifary, W. Bastiaan Kleijn, Mengjie Zhang, and David Balduzzi. Domain generalization for object recognition with multi-task autoencoders. In *ICCV*, pp. 2551–2559. IEEE Computer Society, 2015.
- Ishaan Gulrajani and David Lopez-Paz. In search of lost domain generalization. *arXiv preprint arXiv:2007.01434*, 2020.
- Trygve Haavelmo. The probability approach in econometrics. *Econometrica: Journal of the Econometric Society*, pp. iii–115, 1944.
- Christina Heinze-Deml, Jonas Peters, and Nicolai Meinshausen. Invariant causal prediction for nonlinear models. *Journal of Causal Inference*, 6(2), 2018.
- Kevin D Hoover. The logic of causal inference: Econometrics and the conditional analysis of causation. *Economics & Philosophy*, 6(2):207–234, 1990.
- Patrik O Hoyer, Dominik Janzing, Joris M Mooij, Jonas Peters, and Bernhard Schölkopf. Nonlinear causal discovery with additive noise models. In *Advances in neural information processing systems*, pp. 689–696, 2009.
- Biwei Huang, Kun Zhang, Jiji Zhang, Joseph Ramsey, Ruben Sanchez-Romero, Clark Glymour, and Bernhard Schölkopf. Causal discovery from heterogeneous/nonstationary data. *Journal of Machine Learning Research*, 21(89):1–53, 2020.
- Aapo Hyvärinen. Estimation of non-normalized statistical models by score matching. *Journal of Machine Learning Research*, 6(4), 2005.
- Aapo Hyvärinen and Erkki Oja. Independent component analysis: algorithms and applications. *Neural networks*, 13(4-5):411–430, 2000.
- Aapo Hyvärinen, Hiroaki Sasaki, and Richard Turner. Nonlinear ICA using auxiliary variables and generalized contrastive learning. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 859–868, 2019.
- Dominik Janzing, Joris Mooij, Kun Zhang, Jan Lemeire, Jakob Zscheischler, Povilas Daniušis, Bastian Steudel, and Bernhard Schölkopf. Information-geometric approach to inferring causal directions. *Artificial Intelligence*, 182:1–31, 2012.
- Prithish Kamath, Akilesh Tangella, Danica J Sutherland, and Nathan Srebro. Does invariant risk minimization capture invariance? *arXiv preprint arXiv:2101.01134*, 2021.

- Ilyes Khemakhem, Diederik Kingma, Ricardo Monti, and Aapo Hyvärinen. Variational autoencoders and nonlinear ICA: A unifying framework. In *International Conference on Artificial Intelligence and Statistics*, pp. 2207–2217, 2020a.
- Ilyes Khemakhem, Ricardo Monti, Diederik Kingma, and Aapo Hyvärinen. Ice-beem: Identifiable conditional energy-based deep models based on nonlinear ica. *Advances in Neural Information Processing Systems*, 33, 2020b.
- Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy M Hospedales. Deeper, broader and artier domain generalization. In *Proceedings of the IEEE international conference on computer vision*, pp. 5542–5550, 2017a.
- Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy Hospedales. Learning to generalize: Meta-learning for domain generalization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018a.
- Da Li, Jianshu Zhang, Yongxin Yang, Cong Liu, Yi-Zhe Song, and Timothy M Hospedales. Episodic training for domain generalization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 1446–1455, 2019a.
- Haoliang Li, Sinno Jialin Pan, Shiqi Wang, and Alex C Kot. Domain generalization with adversarial feature learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 5400–5409, 2018b.
- Wen Li, Zheng Xu, Dong Xu, Dengxin Dai, and Luc Van Gool. Domain generalization and adaptation using low rank exemplar svms. *IEEE transactions on pattern analysis and machine intelligence*, 40(5):1114–1127, 2017b.
- Ya Li, Xinmei Tian, Mingming Gong, Yajing Liu, Tongliang Liu, Kun Zhang, and Dacheng Tao. Deep domain generalization via conditional invariant adversarial networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 624–639, 2018c.
- Yiying Li, Yongxin Yang, Wei Zhou, and Timothy Hospedales. Feature-critic networks for heterogeneous domain generalization. In *International Conference on Machine Learning*, pp. 3915–3924. PMLR, 2019b.
- Francesco Locatello, Stefan Bauer, Mario Lucic, Gunnar Raetsch, Sylvain Gelly, Bernhard Schölkopf, and Olivier Bachem. Challenging common assumptions in the unsupervised learning of disentangled representations. In *international conference on machine learning*, pp. 4114–4124, 2019.
- Sara Magliacane, Thijs van Ommen, Tom Claassen, Stephan Bongers, Philip Versteeg, and Joris M Mooij. Domain adaptation by using causal inference to predict invariant conditional distributions. In *Advances in Neural Information Processing Systems*, pp. 10846–10856, 2018.
- Massimiliano Mancini, Samuel Rota Bulò, Barbara Caputo, and Elisa Ricci. Best sources forward: domain generalization through source-specific nets. In *2018 25th IEEE international conference on image processing (ICIP)*, pp. 1353–1357. IEEE, 2018.
- Jovana Mitrovic, Brian McWilliams, Jacob Walker, Lars Buesing, and Charles Blundell. Representation learning via invariant causal mechanisms. *arXiv preprint arXiv:2010.07922*, 2020.
- Saeid Motiian, Marco Piccirilli, Donald A. Adjeroh, and Gianfranco Doretto. Unified deep supervised domain adaptation and generalization. In *ICCV*, pp. 5716–5726. IEEE Computer Society, 2017.
- Krikamol Muandet, David Balduzzi, and Bernhard Schölkopf. Domain generalization via invariant feature representation. In *ICML*, volume 28 of *JMLR Workshop and Conference Proceedings*, pp. 10–18, 2013.
- Vaishnavh Nagarajan, Anders Andreassen, and Behnam Neyshabur. Understanding the failure modes of out-of-distribution generalization. *arXiv preprint arXiv:2010.15775*, 2020.

- Li Niu, Wen Li, and Dong Xu. Multi-view domain generalization for visual recognition. In *Proceedings of the IEEE international conference on computer vision*, pp. 4193–4201, 2015.
- Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *IEEE Transactions on knowledge and data engineering*, 22(10):1345–1359, 2009.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems 32*, pp. 8024–8035. Curran Associates, Inc., 2019.
- Vishal M Patel, Raghuraman Gopalan, Ruonan Li, and Rama Chellappa. Visual domain adaptation: A survey of recent advances. *IEEE signal processing magazine*, 32(3):53–69, 2015.
- Judea Pearl. *Causality*. Cambridge university press, 2009.
- Jonas Peters, Peter Bühlmann, and Nicolai Meinshausen. Causal inference using invariant prediction: identification and confidence intervals. *arXiv preprint arXiv:1501.01332*, 2015.
- Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. *Elements of causal inference*. The MIT Press, 2017.
- Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. *arXiv preprint arXiv:1401.4082*, 2014.
- M. Rojas-Carulla, B. Schölkopf, R. Turner, and J. Peters. Invariant models for causal transfer learning. *Journal of Machine Learning Research*, 19(36):1–34, 2018.
- Elan Rosenfeld, Pradeep Ravikumar, and Andrej Risteski. The risks of invariant risk minimization. *arXiv preprint arXiv:2010.05761*, 2020.
- Shiori Sagawa, Pang Wei Koh, Tatsunori B Hashimoto, and Percy Liang. Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization. *arXiv preprint arXiv:1911.08731*, 2019.
- B. Schölkopf, D. Janzing, J. Peters, E. Sgouritsa, K. Zhang, and J. M. Mooij. On causal and anticausal learning. In *Proceedings of the 29th International Conference on Machine Learning (ICML)*, pp. 1255–1262, 2012.
- Lukas Schott, Jonas Rauber, Matthias Bethge, and Wieland Brendel. Towards the first adversarially robust neural network model on mnist. *arXiv preprint arXiv:1805.09190*, 2018.
- Peter Spirtes, Clark N Glymour, Richard Scheines, and David Heckerman. *Causation, prediction, and search*. MIT press, 2000.
- Amos Storkey. *When Training and Test Sets Are Different: Characterizing Learning Transfer*. 2009.
- Adarsh Subbaswamy, Bryant Chen, and Suchi Saria. A universal hierarchy of shift-stable distributions and the tradeoff between stability and performance. *arXiv preprint arXiv:1905.11374*, 2019.
- Baochen Sun and Kate Saenko. Deep coral: Correlation alignment for deep domain adaptation. In *European conference on computer vision*, pp. 443–450. Springer, 2016.
- Xinwei Sun, Botong Wu, Chang Liu, Xiangyu Zheng, Wei Chen, Tao Qin, and Tie-yan Liu. Latent causal invariant model. *arXiv preprint arXiv:2011.02203*, 2020.
- Hemanth Venkateswara, Jose Eusebio, Shayok Chakraborty, and Sethuraman Panchanathan. Deep hashing network for unsupervised domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 5018–5027, 2017.
- Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural computation*, 23(7):1661–1674, 2011.

- 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. *arXiv preprint arXiv:2106.04619*, 2021.
- J. von Kügelgen, A. Mey, M. Loog, and B. Schölkopf. Semi-supervised learning, causality and the conditional cluster assumption. *Conference on Uncertainty in Artificial Intelligence (UAI)*, 2020.
- Changzhang Wang, You Zhou, Qiang Zhao, and Zhi Geng. Discovering and orienting the edges connected to a target variable in a dag via a sequential local learning approach. *Computational Statistics & Data Analysis*, 77:252–266, 2014.
- Haohan Wang, Aaksha Meghawat, Louis-Philippe Morency, and Eric P Xing. Select-additive learning: Improving generalization in multimodal sentiment analysis. In *2017 IEEE International Conference on Multimedia and Expo (ICME)*, pp. 949–954. IEEE, 2017.
- Haohan Wang, Zexue He, Zachary C Lipton, and Eric P Xing. Learning robust representations by projecting superficial statistics out. *arXiv preprint arXiv:1903.06256*, 2019.
- Garrett Wilson and Diane J Cook. A survey of unsupervised deep domain adaptation. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 11(5):1–46, 2020.
- Sewall Wright. Correlation and causation. *J. agric. Res.*, 20:557–580, 1921.
- Shen Yan, Huan Song, Nanxiang Li, Lincan Zou, and Liu Ren. Improve unsupervised domain adaptation with mixup training. *arXiv preprint arXiv:2001.00677*, 2020.
- Kun Zhang, Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. Kernel-based conditional independence test and application in causal discovery. *arXiv preprint arXiv:1202.3775*, 2012.
- Kun Zhang, Bernhard Schölkopf, Krikamol Muandet, and Zhikun Wang. Domain adaptation under target and conditional shift. In *International Conference on Machine Learning*, pp. 819–827. PMLR, 2013.
- Kun Zhang, Mingming Gong, and Bernhard Schölkopf. Multi-source domain adaptation: A causal view. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 29, 2015.
- Kun Zhang, Biwei Huang, Jiji Zhang, Clark Glymour, and Bernhard Schölkopf. Causal discovery from nonstationary/heterogeneous data: Skeleton estimation and orientation determination. In *IJCAI: Proceedings of the Conference*, volume 2017, pp. 1347. NIH Public Access, 2017.
- Han Zhao, Remi Tachet Des Combes, Kun Zhang, and Geoffrey Gordon. On learning invariant representations for domain adaptation. In *International Conference on Machine Learning*, pp. 7523–7532. PMLR, 2019.

## A DOMAIN GENERALIZATION

The goal of domain generalization (DG) (Muandet et al., 2013) is OOD generalization: learning a predictor that performs well at unseen test domains. Unlike domain adaptation (Pan & Yang, 2009; Ben-David et al., 2007; 2010; Crammer et al., 2008; Patel et al., 2015; Zhao et al., 2019; Wilson & Cook, 2020; Zhang et al., 2015), DG assumes that the test domain data are not available during training. One thread of DG is to explore techniques from kernel methods (Muandet et al., 2013; Niu et al., 2015; Erfani et al., 2016; Li et al., 2017b). Muandet et al. (2013) propose a kernel-based optimization algorithm that learns an invariant transformation by minimizing the discrepancy among domains and preventing the loss of relationship between input and output features. Another line of DG work is using end-to-end methods from deep learning: (a) reducing the differences of representations across domains through adversarial or similar techniques (Ghifary et al., 2015; Wang et al., 2017; Motiian et al., 2017; Li et al., 2018b;c); (b) projecting out superficial domain-specific statistics to reduce sensitivity to the domain (Wang et al., 2019); (c) fusing representations from an ensemble of models across domains (Ding & Fu, 2017; Mancini et al., 2018). Meta-learning can also be applied to domain generalization, by dividing source domains into meta-training and meta-test sets, and aiming for a low generalization error on meta-test sets after training on meta-training sets (Balaji et al., 2018; Dou et al., 2019; Li et al., 2018a; 2019a;b). Recently, an extensive empirical survey of many DG algorithm (Gulrajani & Lopez-Paz, 2020) suggested that with current models and data augmentation techniques, plain ERM may be competitive with the state-of-the-art. It is worth noting that Sun et al. (2020) also propose an approach to learning latent causal factors for prediction. However, their assumptions over the underlying causal graph are restricted due to two reasons: 1) they only consider the scenarios when  $\mathbf{Z}$  and  $\mathbf{Y}$  are generated concurrently, which excludes the cases in which some part of  $\mathbf{Z}$  could also be affected by  $\mathbf{Y}$  in some manner; 2) They assume that the causal latent factors  $\mathbf{Z}_p$  and the non-causal latent factors  $\mathbf{Z}_c$  are independent when conditioning on  $\mathbf{E}$ , that is,  $\mathbf{Z}_p \perp\!\!\!\perp \mathbf{Z}_c | \mathbf{E}$ . In practice, during model learning, they actually further assume that  $Z_i \perp\!\!\!\perp Z_j | \mathbf{E}$  for any  $i \neq j$  so that VAE could be leveraged to learn the model. In this sense, their approach has the same issue as the one in iVAE, i.e., unable to deal with the non-factorized cases. This point is also verified in Table 2, where our approach greatly outperforms theirs.

## B VARIATIONAL AUTOENCODERS

We briefly describe the framework of variational autoencoders (VAEs), which allows us to efficiently learn deep latent-variable models and their corresponding inference models (Kingma & Welling, 2013; Rezende et al., 2014). Consider a simple latent variable model where  $\mathbf{X} \in \mathbb{R}^d$  stands for an observed variable and  $\mathbf{Z} \in \mathbb{R}^n$  for a latent variable. A VAE method learns a full generative model  $p_\theta(\mathbf{X}, \mathbf{Z}) = p_\theta(\mathbf{X}|\mathbf{Z})p_\theta(\mathbf{Z})$  and an inference model  $q_\phi(\mathbf{Z}|\mathbf{X})$ , typically a factorized Gaussian distribution whose mean and variance parameters are given by the output of a neural network with input  $\mathbf{X}$ . This inference model approximates the posterior  $p_\theta(\mathbf{Z}|\mathbf{X})$ , where  $\theta$  is a vector of parameters of the generative model,  $\phi$  a vector of parameters of the inference model, and  $p_\theta(\mathbf{Z})$  is a prior distribution over the latent variables. Instead of maximizing the data log-likelihood, we maximize its lower bound  $\mathcal{L}_{\text{VAE}}(\theta, \phi)$ :

$$\log p_\theta(\mathbf{X}) \geq \mathcal{L}_{\text{VAE}}(\theta, \phi) := \mathbb{E}_{q_\phi(\mathbf{Z}|\mathbf{X})} [\log p_\theta(\mathbf{X}|\mathbf{Z})] - \text{KL}(q_\phi(\mathbf{Z}|\mathbf{X})||p_\theta(\mathbf{Z})),$$

where we have used Jensen’s inequality, and  $\text{KL}(\cdot||\cdot)$  denotes the Kullback-Leibler divergence between two distributions.

## C DERIVATION

### C.1 iVAE

In Section 2.1, the evidence lower bound of iVAE is defined by

$$\begin{aligned} \mathcal{L}_{\text{iVAE}}(\theta, \phi) &:= \mathbb{E}_{p_D} [\mathbb{E}_{q_\phi(\mathbf{Z}|\mathbf{X}, U)} [\log p_f(\mathbf{X}|\mathbf{Z})] - \text{KL}(q_\phi(\mathbf{Z}|\mathbf{X}, U)||p_{T, \lambda}(\mathbf{Z}|U))] \\ &= \mathbb{E}_{p_D} [\mathbb{E}_{q_\phi(\mathbf{Z}|\mathbf{X}, U)} [\log p_f(\mathbf{X}|\mathbf{Z}) + \log p_{T, \lambda}(\mathbf{Z}|U) - \log q_\phi(\mathbf{Z}|\mathbf{X}, U)]] . \end{aligned}$$

## C.2 SCORE MATCHING

In Section 4.1, we follow Hyvärinen (2005) and use a simple trick of partial integration to simplify the evaluation of the score matching objective  $\mathcal{L}_{\text{phase1}}^{\text{SM}}$  in Eq. (10) of the main text:

$$\begin{aligned} \mathcal{L}_{\text{phase1}}^{\text{SM}}(\mathbf{T}, \boldsymbol{\lambda}) &:= \mathbb{E}_{p_D} \left[ \mathbb{E}_{q_\phi(\mathbf{Z}|\mathbf{X}, \mathbf{Y}, \mathbf{E})} \left[ \|\nabla_{\mathbf{Z}} \log q_\phi(\mathbf{Z}|\mathbf{X}, \mathbf{Y}, \mathbf{E}) - \nabla_{\mathbf{Z}} \log p_{\mathbf{T}, \boldsymbol{\lambda}}(\mathbf{Z}|\mathbf{Y}, \mathbf{E})\|^2 \right] \right] \\ &= \mathbb{E}_{p_D} \left[ \mathbb{E}_{q_\phi(\mathbf{Z}|\mathbf{X}, \mathbf{Y}, \mathbf{E})} \left[ \sum_{j=1}^n \left[ \frac{\partial^2 p_{\mathbf{T}, \boldsymbol{\lambda}}(\mathbf{Z}|\mathbf{Y}, \mathbf{E})}{\partial Z_j^2} + \frac{1}{2} \left( \frac{\partial p_{\mathbf{T}, \boldsymbol{\lambda}}(\mathbf{Z}|\mathbf{Y}, \mathbf{E})}{\partial Z_j} \right)^2 \right] \right] \right] + \text{const.} \end{aligned}$$

where the last equality is due to the Theorem 1 in Hyvärinen (2005).

## D DEFINITIONS FOR SEM AND IRM

**Definition 2.** A structural equation model (SEM)  $\mathcal{M} := (\mathcal{S}, N)$  governing the random vector  $\mathbf{Z} = (Z_1, \dots, Z_d)$  is a set of structural equations:

$$\mathcal{S}_i : Z_i \leftarrow f_i(\text{Pa}(Z_i), N_i),$$

where  $\text{Pa}(Z_i) \subseteq \{Z_1, \dots, Z_d\} \setminus \{Z_i\}$  are called the parents of  $Z_i$ , and the  $N_i$  are independent noise random variables. We say that “ $Z_i$  causes  $Z_j$ ” if  $Z_i \in \text{Pa}(Z_j)$ . We call causal graph of  $\mathbf{Z}$  to the graph obtained by drawing i) one node for each  $Z_i$ , and ii) one edge from  $Z_i$  to  $Z_j$  if  $Z_i \in \text{Pa}(Z_j)$ . We assume acyclic causal graphs.

**Definition 3.** Consider a SEM  $\mathcal{M} := (\mathcal{S}, N)$ . An intervention  $e$  on  $\mathcal{M}$  consists of replacing one or several of its structural equations to obtain an intervened SEM  $\mathcal{M}^e := (\mathcal{S}^e, N^e)$ , with structural equations:

$$\mathcal{S}_i^e : Z_i^e \leftarrow f_i^e(\text{Pa}^e(Z_i^e), N_i^e),$$

The variable  $\mathbf{Z}^e$  is intervened if  $\mathcal{S}_i \neq \mathcal{S}_i^e$  or  $N_i \neq N_i^e$ .

**Definition 4.** Consider a structural equation model (SEM)  $\mathcal{S}$  governing the random vector  $(Z_1, \dots, Z_n, \mathbf{Y})$ , and the learning goal of predicting  $\mathbf{Y}$  from  $\mathbf{Z}$ . Then, the set of all environments  $\mathcal{E}_{\text{all}}(\mathcal{S})$  indexes all the interventional distributions  $P(\mathbf{Z}^e, \mathbf{Y}^e)$  obtainable by valid interventions  $e$ . An intervention  $e \in \mathcal{E}_{\text{all}}(\mathcal{S})$  is valid as long as (i) the causal graph remains acyclic, (ii)  $\mathbb{E}[\mathbf{Y}^e | \text{Pa}(\mathbf{Y})] = \mathbb{E}[\mathbf{Y} | \text{Pa}(\mathbf{Y})]$ , and (iii)  $\forall [\mathbf{Y}^e | \text{Pa}(\mathbf{Y})]$  remains within a finite range.

## E DEFINITIONS AND LEMMAS FOR THE EXPONENTIAL FAMILIES

**Definition 5. (Exponential family)** A multivariate exponential family is a set of distributions whose probability density function can be written as

$$p(\mathbf{Z}) = \frac{\mathcal{Q}(\mathbf{Z})}{\mathcal{C}(\boldsymbol{\theta})} \exp(\langle \mathbf{T}(\mathbf{Z}), \boldsymbol{\theta} \rangle), \quad (13)$$

where  $\mathcal{Q} : \mathcal{Z} \rightarrow \mathbb{R}$  is the base measure,  $\mathcal{C}(\boldsymbol{\theta})$  is the normalizing constant,  $\mathbf{T} : \mathcal{Z} \rightarrow \mathbb{R}^k$  is the sufficient statistics, and  $\boldsymbol{\theta} \in \mathbb{R}^k$  is the natural parameter. The size  $k \geq n$  is the dimension of the sufficient statistics  $\mathbf{T}$  and depends on the latent space dimension  $n$ . Note that  $k$  is fixed given  $n$ .

**Definition 6. (Strongly exponential distributions)** A multivariate exponential family distribution

$$p(\mathbf{Z}) = \frac{\mathcal{Q}(\mathbf{Z})}{\mathcal{C}(\boldsymbol{\theta})} \exp(\langle \mathbf{T}(\mathbf{Z}), \boldsymbol{\theta} \rangle) \quad (14)$$

is strongly exponential, if

$$(\exists \boldsymbol{\theta} \in \mathbb{R}^k \text{ s.t. } \langle \mathbf{T}(\mathbf{Z}), \boldsymbol{\theta} \rangle = \text{const.}, \forall \mathbf{Z} \in \mathcal{Z}) \implies (l(\mathcal{Z}) = 0 \text{ or } \boldsymbol{\theta} = \mathbf{0}), \quad \forall \mathcal{Z} \subset \mathbb{R}^n, \quad (15)$$

where  $l$  is the Lebesgue measure.

The density of a strongly exponential distribution has almost surely the exponential component and can only be reduced to the base measure on a set of measure zero. Note that all common multivariate exponential family distributions (e.g. multivariate Gaussian) are strongly exponential.

## F DEFINITIONS FOR IDENTIFIABILITY

**Definition 7.** Let  $\Theta$  be the domain of the parameters  $\boldsymbol{\theta} = \{\mathbf{f}, \mathbf{T}, \boldsymbol{\lambda}\}$ . Let  $\sim$  be an equivalence relation on  $\Theta$ . A deep generative model is said to be  $\sim$ -identifiable if

$$p_{\boldsymbol{\theta}}(\mathbf{X}) = p_{\tilde{\boldsymbol{\theta}}}(\mathbf{X}) \implies \boldsymbol{\theta} \sim \tilde{\boldsymbol{\theta}}. \quad (16)$$



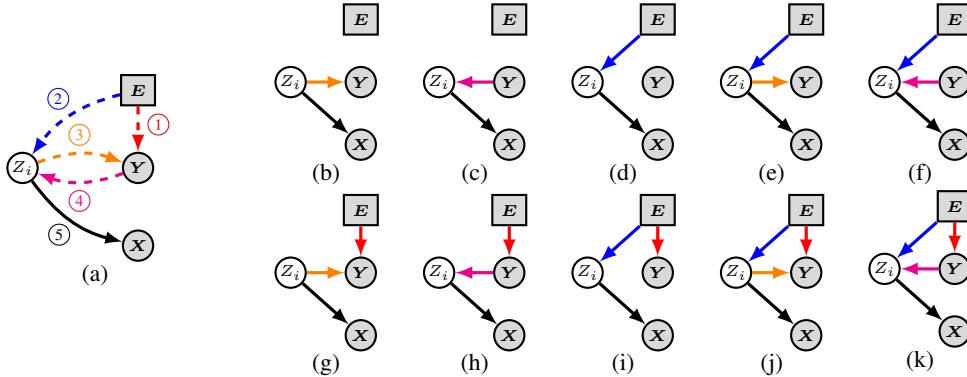


Figure 3: (a) General causal structure over  $\{X_i, Y, X, E\}$ , where the arrow from  $Z_i$  to  $X$  is a must-have connection and the other four might not be necessary. (b) Ten possible causal structures from (a) under Assumptions 1&2.

The elements in the quotient space  $\Theta \setminus \sim$  are called the identifiability classes.

**Definition 8.** Let  $\sim_A$  be an equivalence relation on  $\Theta$  defined by:

$$(\mathbf{f}, \mathbf{T}, \boldsymbol{\lambda}) \sim_A (\tilde{\mathbf{f}}, \tilde{\mathbf{T}}, \tilde{\boldsymbol{\lambda}}) \iff \exists A, \mathbf{c} \text{ s.t. } \mathbf{T}(\mathbf{f}^{-1}(\mathbf{X})) = A\tilde{\mathbf{T}}(\tilde{\mathbf{f}}^{-1}(\mathbf{X})) + \mathbf{c}, \forall \mathbf{X} \in \mathcal{O}, \quad (17)$$

where  $A \in \mathbb{R}^{k \times k}$  is an invertible matrix and  $\mathbf{c} \in \mathbb{R}^k$  is a vector.

**Definition 9.** Let  $\sim_P$  be an equivalence relation on  $\Theta$  defined by:

$$(\mathbf{f}, \mathbf{T}, \boldsymbol{\lambda}) \sim_P (\tilde{\mathbf{f}}, \tilde{\mathbf{T}}, \tilde{\boldsymbol{\lambda}}) \iff \exists P, \mathbf{c} \text{ s.t. } \mathbf{T}(\mathbf{f}^{-1}(\mathbf{X})) = P\tilde{\mathbf{T}}(\tilde{\mathbf{f}}^{-1}(\mathbf{X})) + \mathbf{c}, \forall \mathbf{X} \in \mathcal{O}, \quad (18)$$

where  $P \in \mathbb{R}^{k \times k}$  is a block permutation matrix and  $\mathbf{c} \in \mathbb{R}^k$  is a vector.

## G PHASE 2: DISCOVERING SINGLE CAUSE

### G.1 SOME SPECIAL CASES IN MULTI-CAUSE SETTINGS

These conditional independence tests can be performed in parallel to largely accelerate the testing procedure. Note that, in practice, it might occur that there exist some  $Z_i$  which is independent of any other  $Z_j$  when conditioning on  $Y$  and  $E$ . It is probably because such  $Z_i$  is a deterministic transformation of  $Y$ . In this special case, we can use IGCI (Danusis et al., 2012; Janzing et al., 2012) to determine whether or not  $Z_i$  is a cause of  $Y$ . Also note that, in some scenarios in which the dependence signals between a pair of causal latent variables might be weak due to the data issue, we can test the conditional independence of such a latent variable with all the other causal latent variables. If it is conditionally dependent on more than half of them or its average p-value is larger than the pre-specified threshold, we will select it as one cause of  $Y$ .

### G.2 DISCOVERING SINGLE CAUSE

In the single-cause case, by following Wang et al. (2014), we leverage the *MB-by-MB* algorithm to first construct a *local* structure around  $Y$  and then discover the single parent of  $Y$  according to the constructed *local* graph. One obvious advantage of this approach is in efficiency, because there is no need to construct the whole causal graph containing all the latent variables and  $Y$ .

We could have an even more efficient solution to the single-cause case in some special scenarios where we assume that  $Z_i \perp\!\!\!\perp Z_j | Y, E$  for any  $i \neq j$ . In fact, this assumption covers more scenarios than the common assumption that  $Z_i \perp\!\!\!\perp Z_j$  for  $i \neq j$  in latent variable models, e.g., disentanglement (Bengio et al., 2013; Locatello et al., 2019), autoencoders (Kingma & Welling, 2013; Rezende et al., 2014), ICA (Comon, 1994; Hyvärinen & Oja, 2000), etc. If  $Y$  is caused by at most one of  $Z_i$  and  $Z_j$ , and no matter whether  $Z_i$  and  $Z_j$  are caused by  $E$  or not, then  $Z_i \perp\!\!\!\perp Z_j | Y, E$  holds, but we may well have  $Z_i \not\perp\!\!\!\perp Z_j$  (e.g., if  $Y$  causes both  $Z_i$  and  $Z_j$ , or if there is a chain  $Z_i \rightarrow Y \rightarrow Z_j$ ). If  $Y$  causes or is caused by at most one of  $Z_i$  and  $Z_j$ , and at most one of  $Z_i$  and  $Z_j$  is caused by  $E$ , then both  $Z_i \perp\!\!\!\perp Z_j$  and  $Z_i \perp\!\!\!\perp Z_j | Y, E$  hold. If  $\{Y, Z_i, Z_j\}$  form a collider  $Z_i \rightarrow Y \leftarrow Z_j$ , and no matter whether  $Z_i$  and  $Z_j$  are caused by  $E$  or not, then  $Z_i \perp\!\!\!\perp Z_j | E$  hold, but we may have  $Z_i \not\perp\!\!\!\perp Z_j$  (e.g., when both  $Z_i$  and  $Z_j$  are caused by  $E$ ).

Under this assumption, we are able to separately look into each  $Z_i$  given  $\mathbf{Y}$  and  $\mathbf{E}$ , without considering any other  $Z_j$ . Fig. 5a shows that there exist only five possible connections between  $Z_i$ ,  $\mathbf{Y}$ ,  $\mathbf{E}$ , and  $\mathbf{X}$ . Among them, only the arrow from  $Z_i$  to  $\mathbf{X}$  must exist, because  $\mathbf{X}$  is generated from  $\mathbf{Z}$ . The other four arrows might not be present, with the exception that there must be at least one connection between  $Z_i$  and  $\mathbf{Y}$  or  $\mathbf{E}$  (Assumption 1a). This leaves ten possible structures, shown in Figs. 3b-3k.

Given data  $\{\hat{Z}_i, \mathbf{Y}, \mathbf{E}, \mathbf{X}\}$  in which the  $\hat{Z}_i$  are obtained using  $q_\phi(\mathbf{Z}|\mathbf{X}, \mathbf{Y}, \mathbf{E})$  (for example, as given by the mean of this distribution), we are able to distinguish all ten structures in Figs. 3b-3k by using causal discovery algorithms (Peters et al., 2017; Zhang et al., 2017; Huang et al., 2020) and performing conditional independence tests (Spirtes et al., 2000; Zhang et al., 2012). This is summarized in Proposition 2 below. Its proof can be found in Appendix H, which also describes the specific assumptions made. In practice, we can assess in parallel whether or not each  $Z_i$  is a direct cause of  $\mathbf{Y}$ , which accelerates this phase significantly.

**Proposition 2.** *Under the assumptions stated in Appendix H, the ten structures shown in Figs. 3b-3k can be identified using causal discovery methods consistent in the infinite sample limit.*

Note that there are only four cases in which  $Z_i$  is a parent of  $\mathbf{Y}$  (i.e., Figs. 3b, 3e, 3g, and 3j). We can identify these by applying rules 1.2, 1.6, 2.1, and 3.1 from Appendix H.5.

## H PROOFS

### H.1 PROOF OF THEOREM 1

The proof of this theorem consists of three parts.

In Part I, we prove that the parameters  $\theta$  are  $\sim_A$  identifiable (Definition 8) by using assumption (i), the first half of assumption (ii), and assumption (iv) of Theorem 1.

In Part II, based on the result in Part I, we further prove that the parameters  $\theta$  are  $\sim_P$  identifiable (Definition 9) by additionally using Assumption 2, the second half of assumption (ii) and assumption (iii) of Theorem 1.

In Part III, we combine the results (Theorems 4&5) in both Part I and Part II into one theorem (Theorem 1), which completes the proof.

It is worth noting that, compared to the proof in iVAE, the main changes in our proof consist of

- Part I, In step III.
  - It has been updated to account for vectors of sufficient statistics whose entries can be arbitrary functions of all entries in the random variable vector, while in the previous proof the sufficient statistics contained entries that are functions of individual entries in the random variable vector.
  - The assumption of “The sufficient statistics in  $\mathbf{T}$  are all linearly independent.” is not required in our proof, but it is in the proof of iVAE.
- Part II.
  - It has been updated to account for the part of the sufficient statistics which is the output of a deep neural network with ReLU activation functions.

#### H.1.1 PART I

For notational simplicity, in the proof of this part we denote  $(\mathbf{Y}, \mathbf{E})$  by  $\mathbf{U}$ . Hence, our generative model defined according to Eqs. (6-8) in the main text now becomes:

$$p_\theta(\mathbf{X}, \mathbf{Z}|\mathbf{U}) = p_f(\mathbf{X}|\mathbf{Z})p_{\mathbf{T}, \lambda}(\mathbf{Z}|\mathbf{U}), \quad (19)$$

$$p_f(\mathbf{X}|\mathbf{Z}) = p_\epsilon(\mathbf{X} - \mathbf{f}(\mathbf{Z})), \quad (20)$$

$$p_{\mathbf{T}, \lambda}(\mathbf{Z}|\mathbf{U}) = \mathcal{Q}(\mathbf{Z})/\mathcal{C}(\mathbf{U}) \exp[\mathbf{T}(\mathbf{Z})^\top \lambda(\mathbf{U})]. \quad (21)$$

**Theorem 4.** *Suppose that we observe data sampled from a deep generative model defined according to Eqs. (19-21) with parameters  $(\mathbf{f}, \mathbf{T}, \lambda)$ . Assume that*

- (i) *The set  $\{\mathbf{X} \in \mathbb{O} | \varphi_\epsilon(\mathbf{X}) = 0\}$  has measure zero, where  $\varphi_\epsilon$  is the characteristic function of the density  $p_\epsilon$  defined in Eq. (20);*

(ii) The mixing function  $\mathbf{f}$  in Eq. (20) is injective;

(iii) There exist  $k + 1$  points  $\mathbf{U}^0, \mathbf{U}^1, \dots, \mathbf{U}^k \in \mathcal{U}$  such that the matrix

$$L = [\boldsymbol{\lambda}(\mathbf{U}^1) - \boldsymbol{\lambda}(\mathbf{U}^0), \dots, \boldsymbol{\lambda}(\mathbf{U}^k) - \boldsymbol{\lambda}(\mathbf{U}^0)] \in \mathbb{R}^{k \times k} \quad (22)$$

is invertible.

Then the parameters  $\{\mathbf{f}, \mathbf{T}, \boldsymbol{\lambda}\}$  are  $\sim_A$  identifiable.

*Proof.* Define  $\text{vol}(B) = \sqrt{\det(B^T B)}$  for any full column rank matrix  $B$ . Suppose that we have two sets of parameters  $\boldsymbol{\theta} = (\mathbf{f}, \mathbf{T}, \boldsymbol{\lambda})$  and  $\tilde{\boldsymbol{\theta}} = (\tilde{\mathbf{f}}, \tilde{\mathbf{T}}, \tilde{\boldsymbol{\lambda}})$  such that  $p_{\boldsymbol{\theta}}(\mathbf{X}|\mathbf{U}) = p_{\tilde{\boldsymbol{\theta}}}(\mathbf{X}|\mathbf{U})$ ,  $\forall (\mathbf{X}, \mathbf{U}) \in \mathcal{O} \times \mathcal{U}$ . We want to show  $\boldsymbol{\theta} \sim_A \tilde{\boldsymbol{\theta}}$ .

**Step I.** The proof of this step is similar to Step I in the proof of Theorem 1 in [Khemakhem et al. \(2020a\)](#). We transform the equality of the marginal distributions over observed data into the equality of noise-free distributions. For all pairs  $(\mathbf{X}, \mathbf{U}) \in \mathcal{O} \times \mathcal{U}$ , we have

$$p_{\boldsymbol{\theta}}(\mathbf{X}|\mathbf{U}) = p_{\tilde{\boldsymbol{\theta}}}(\mathbf{X}|\mathbf{U}) \quad (23)$$

$$\implies \int_{\mathcal{Z}} p_{\mathbf{f}}(\mathbf{X}|\mathbf{Z}) p_{\mathbf{T}, \boldsymbol{\lambda}}(\mathbf{Z}|\mathbf{U}) d\mathbf{Z} = \int_{\mathcal{Z}} p_{\tilde{\mathbf{f}}}(\mathbf{X}|\mathbf{Z}) p_{\tilde{\mathbf{T}}, \tilde{\boldsymbol{\lambda}}}(\mathbf{Z}|\mathbf{U}) d\mathbf{Z} \quad (24)$$

$$\implies \int_{\mathcal{Z}} p_{\varepsilon}(\mathbf{X} - \mathbf{f}(\mathbf{Z})) p_{\mathbf{T}, \boldsymbol{\lambda}}(\mathbf{Z}|\mathbf{U}) d\mathbf{Z} = \int_{\mathcal{Z}} p_{\varepsilon}(\mathbf{X} - \tilde{\mathbf{f}}(\mathbf{Z})) p_{\tilde{\mathbf{T}}, \tilde{\boldsymbol{\lambda}}}(\mathbf{Z}|\mathbf{U}) d\mathbf{Z} \quad (25)$$

$$\implies \int_{\mathcal{O}} p_{\varepsilon}(\mathbf{X} - \bar{\mathbf{X}}) p_{\mathbf{T}, \boldsymbol{\lambda}}(\mathbf{f}^{-1}(\bar{\mathbf{X}})|\mathbf{U}) \text{vol}(J_{\mathbf{f}^{-1}}(\bar{\mathbf{X}})) d\bar{\mathbf{X}} = \int_{\mathcal{O}} p_{\varepsilon}(\mathbf{X} - \bar{\mathbf{X}}) p_{\tilde{\mathbf{T}}, \tilde{\boldsymbol{\lambda}}}(\tilde{\mathbf{f}}^{-1}(\bar{\mathbf{X}})|\mathbf{U}) \text{vol}(J_{\tilde{\mathbf{f}}^{-1}}(\bar{\mathbf{X}})) d\bar{\mathbf{X}} \quad (26)$$

$$\implies \int_{\mathbb{R}^d} p_{\varepsilon}(\mathbf{X} - \bar{\mathbf{X}}) \tilde{p}_{\mathbf{f}, \mathbf{T}, \boldsymbol{\lambda}, \mathbf{U}}(\bar{\mathbf{X}}) d\bar{\mathbf{X}} = \int_{\mathbb{R}^d} p_{\varepsilon}(\mathbf{X} - \bar{\mathbf{X}}) \tilde{p}_{\tilde{\mathbf{f}}, \tilde{\mathbf{T}}, \tilde{\boldsymbol{\lambda}}, \tilde{\mathbf{U}}}(\bar{\mathbf{X}}) d\bar{\mathbf{X}} \quad (27)$$

$$\implies (\tilde{p}_{\mathbf{f}, \mathbf{T}, \boldsymbol{\lambda}, \mathbf{U}} * p_{\varepsilon})(\mathbf{X}) = (\tilde{p}_{\tilde{\mathbf{f}}, \tilde{\mathbf{T}}, \tilde{\boldsymbol{\lambda}}, \tilde{\mathbf{U}}} * p_{\varepsilon})(\mathbf{X}) \quad (28)$$

$$\implies F[\tilde{p}_{\mathbf{f}, \mathbf{T}, \boldsymbol{\lambda}, \mathbf{U}}](\boldsymbol{\omega}) \varphi_{\varepsilon}(\boldsymbol{\omega}) = F[\tilde{p}_{\tilde{\mathbf{f}}, \tilde{\mathbf{T}}, \tilde{\boldsymbol{\lambda}}, \tilde{\mathbf{U}}}] (\boldsymbol{\omega}) \varphi_{\varepsilon}(\boldsymbol{\omega}) \quad (29)$$

$$\implies F[\tilde{p}_{\mathbf{f}, \mathbf{T}, \boldsymbol{\lambda}, \mathbf{U}}](\boldsymbol{\omega}) = F[\tilde{p}_{\tilde{\mathbf{f}}, \tilde{\mathbf{T}}, \tilde{\boldsymbol{\lambda}}, \tilde{\mathbf{U}}}] (\boldsymbol{\omega}) \quad (30)$$

$$\implies \tilde{p}_{\mathbf{f}, \mathbf{T}, \boldsymbol{\lambda}, \mathbf{U}}(\mathbf{X}) = \tilde{p}_{\tilde{\mathbf{f}}, \tilde{\mathbf{T}}, \tilde{\boldsymbol{\lambda}}, \tilde{\mathbf{U}}}(\mathbf{X}) \quad (31)$$

where

- in Eq. (26),  $J$  denotes the Jacobian, and we made the change of variable  $\bar{\mathbf{X}} = \mathbf{f}(\mathbf{Z})$  on the left hand side, and  $\bar{\mathbf{X}} = \tilde{\mathbf{f}}(\mathbf{Z})$  on the right hand side.
- in Eq. (27), we introduced
 
$$\tilde{p}_{\mathbf{f}, \mathbf{T}, \boldsymbol{\lambda}, \mathbf{U}}(\mathbf{X}) \triangleq p_{\mathbf{T}, \boldsymbol{\lambda}}(\mathbf{f}^{-1}(\mathbf{X})|\mathbf{U}) \text{vol}(J_{\mathbf{f}^{-1}}(\mathbf{X})) \mathbb{I}_{\mathcal{O}}(\mathbf{X}) \quad (32)$$
 on the left hand side, and similarly on the right hand side.
- in Eq. (28), we used  $*$  for the convolution operator.
- in Eq. (29), we used  $F[\cdot]$  to designate the Fourier transform, and where  $\varphi_{\varepsilon} = F[p_{\varepsilon}]$  (by definition of the characteristic function).
- in Eq. (30), we dropped  $\varphi_{\varepsilon}(\boldsymbol{\omega})$  from both sides as it is non-zero almost everywhere (by assumption (i)).

**Step II.** In this step, we remove all terms that are either a function of  $\mathbf{X}$  or  $\mathbf{U}$ . First, by replacing both sides of Eq. (31) by their corresponding expressions from Eq. (32), we have

$$p_{\mathbf{T}, \boldsymbol{\lambda}}(\mathbf{f}^{-1}(\mathbf{X})|\mathbf{U}) \text{vol}(J_{\mathbf{f}^{-1}}(\mathbf{X})) = p_{\tilde{\mathbf{T}}, \tilde{\boldsymbol{\lambda}}}(\tilde{\mathbf{f}}^{-1}(\mathbf{X})|\mathbf{U}) \text{vol}(J_{\tilde{\mathbf{f}}^{-1}}(\mathbf{X})). \quad (33)$$

Then, by taking logarithm on both sides of Eq. (33) and replacing  $p_{\mathbf{T}, \boldsymbol{\lambda}}$  by its expression from Eq. (21), we obtain

$$\begin{aligned} & \log \text{vol}(J_{\mathbf{f}^{-1}}(\mathbf{X})) + \log Q(\mathbf{f}^{-1}(\mathbf{X})) - \log Z(\mathbf{U}) + \langle \mathbf{T}(\mathbf{f}^{-1}(\mathbf{X})), \boldsymbol{\lambda}(\mathbf{U}) \rangle \\ &= \log \text{vol}(J_{\tilde{\mathbf{f}}^{-1}}(\mathbf{X})) + \log \tilde{Q}(\tilde{\mathbf{f}}^{-1}(\mathbf{X})) - \log \tilde{Z}(\mathbf{U}) + \langle \tilde{\mathbf{T}}(\tilde{\mathbf{f}}^{-1}(\mathbf{X})), \tilde{\boldsymbol{\lambda}}(\mathbf{U}) \rangle. \end{aligned} \quad (34)$$

Let  $\mathbf{U}^0, \mathbf{U}^1, \dots, \mathbf{U}^k \in \mathcal{U}$  be the  $k + 1$  points defined in assumption (iii). We evaluate the above equation at these points to obtain  $k + 1$  equations, and subtract the first equation from the remaining

$k$  equations to obtain:

$$\begin{aligned} \langle \mathbf{T}(\mathbf{f}^{-1}(\mathbf{X})), \boldsymbol{\lambda}(\mathbf{U}^l) - \boldsymbol{\lambda}(\mathbf{U}^0) \rangle + \log \frac{Z(\mathbf{U}^0)}{Z(\mathbf{U}^l)} \\ = \langle \tilde{\mathbf{T}}(\tilde{\mathbf{f}}^{-1}(\mathbf{X})), \tilde{\boldsymbol{\lambda}}(\mathbf{U}^l) - \tilde{\boldsymbol{\lambda}}(\mathbf{U}^0) \rangle + \log \frac{\tilde{Z}(\mathbf{U}^0)}{\tilde{Z}(\mathbf{U}^l)}, \quad l = 1, \dots, k. \end{aligned} \quad (35)$$

Let  $L$  be defined as in assumption (iii) and  $\tilde{L}$  defined similarly for  $\tilde{\boldsymbol{\lambda}}$ . Note that  $L$  is invertible by assumption, but  $\tilde{L}$  is not necessarily invertible. Letting  $\mathbf{b} \in \mathbb{R}^k$  in which  $b_l = \log \frac{\tilde{Z}(\mathbf{U}^0)Z(\mathbf{U}^l)}{\tilde{Z}(\mathbf{U}^l)Z(\mathbf{U}^0)}$ , we have

$$L^T \mathbf{T}(\mathbf{f}^{-1}(\mathbf{X})) = \tilde{L}^T \tilde{\mathbf{T}}(\tilde{\mathbf{f}}^{-1}(\mathbf{X})) + \mathbf{b}. \quad (36)$$

Left multiplying both sides of the above equation by  $L^{-T}$  gives

$$\mathbf{T}(\mathbf{f}^{-1}(\mathbf{X})) = A \tilde{\mathbf{T}}(\tilde{\mathbf{f}}^{-1}(\mathbf{X})) + \mathbf{c}, \quad (37)$$

where  $A = L^{-T} \tilde{L} \in \mathbb{R}^{k \times k}$  and  $\mathbf{c} = L^{-T} \mathbf{b} \in \mathbb{R}^k$ .

**Step III.** To complete the proof, we need to show that  $A$  is invertible. Let  $\mathbf{Z}_l \in \mathcal{Z}$ ,  $\mathbf{X}_l = \mathbf{f}(\mathbf{Z}_l)$ ,  $l = 0, \dots, k$ . We evaluate Eq. (37) at these  $k + 1$  points to obtain  $k + 1$  equations and subtract the first equation from the remaining  $k$  equations to obtain

$$\begin{aligned} \underbrace{[\mathbf{T}(\mathbf{Z}_1) - \mathbf{T}(\mathbf{Z}_0), \dots, \mathbf{T}(\mathbf{Z}_k) - \mathbf{T}(\mathbf{Z}_0)]}_{\triangleq R \in \mathbb{R}^{k \times k}} \\ = A \underbrace{[\tilde{\mathbf{T}}(\tilde{\mathbf{f}}^{-1}(\mathbf{X}_1)) - \tilde{\mathbf{T}}(\tilde{\mathbf{f}}^{-1}(\mathbf{X}_0)), \dots, \tilde{\mathbf{T}}(\tilde{\mathbf{f}}^{-1}(\mathbf{X}_k)) - \tilde{\mathbf{T}}(\tilde{\mathbf{f}}^{-1}(\mathbf{X}_0))]}_{\triangleq \tilde{R} \in \mathbb{R}^{k \times k}}. \end{aligned} \quad (38)$$

We need to show that for a given  $\mathbf{Z}_0 \in \mathcal{Z}$ , there exist  $k$  points  $\mathbf{Z}_1, \dots, \mathbf{Z}_k \in \mathcal{Z}$  such that the columns of  $R$  are linearly independent. Suppose, for contradiction, that the columns of  $R$  would never be linearly independent for any choice of  $\mathbf{Z}_1, \dots, \mathbf{Z}_k \in \mathcal{Z}$ . Then the function  $\mathbf{g}(\mathbf{Z}) \triangleq \mathbf{T}(\mathbf{Z}) - \mathbf{T}(\mathbf{Z}_0)$  would live in a  $k - 1$  or lower dimensional subspace, and thus we could find a non-zero vector  $\boldsymbol{\lambda} \in \mathbb{R}^k$  orthogonal to that subspace. This would imply that  $\langle \mathbf{T}(\mathbf{Z}) - \mathbf{T}(\mathbf{Z}_0), \boldsymbol{\lambda} \rangle = 0$  and thus  $\langle \mathbf{T}(\mathbf{Z}), \boldsymbol{\lambda} \rangle = \langle \mathbf{T}(\mathbf{Z}_0), \boldsymbol{\lambda} \rangle = \text{const}$ ,  $\forall \mathbf{Z} \in \mathcal{Z}$ , which contradicts the assumption that the prior is strongly exponential (Definition 6). Therefore, we have shown that there exist  $k + 1$  points  $\mathbf{Z}_0, \mathbf{Z}_1, \dots, \mathbf{Z}_k \in \mathcal{Z}$  such that  $R$  is invertible. Since  $R = A \tilde{R}$  and  $A$  is not a function of  $\mathbf{Z}$ ,  $A$  must be invertible. This completes the proof.  $\square$

## H.1.2 PART II

**Theorem 5.** *Suppose that all assumptions in Theorem 4 hold. Let the sufficient statistics  $\mathbf{T}(\mathbf{Z}) = [\mathbf{T}_f(\mathbf{Z})^T, \mathbf{T}_{NN}(\mathbf{Z})^T]^T$  given by the concatenation of a) the sufficient statistics  $\mathbf{T}_f(\mathbf{Z}) = [\mathbf{T}_1(\mathbf{Z}_1)^T, \dots, \mathbf{T}_n(\mathbf{Z}_n)^T]^T$  of a factorized exponential family, where all the  $\mathbf{T}_i(\mathbf{Z}_i)$  have dimension larger or equal to 2, and b) the output  $\mathbf{T}_{NN}(\mathbf{Z})$  of a neural network with ReLU activations. (note that a neural network with ReLU activation has universal approximation power and should be able to capture any dependencies of interest). Let  $k'$  be the dimension of  $\mathbf{T}_f$  and thus that  $k' \geq 2n$ . Assume that*

- (i) *the sufficient statistics  $\mathbf{T}_f$  have all second-order own derivatives;*
- (ii) *the mixing function  $\mathbf{f}$  has all second-order cross derivatives.*

*Then the parameters  $\{\mathbf{f}, \mathbf{T}, \boldsymbol{\lambda}\}$  are  $\sim_P$  identifiable.*

*Proof.* Let  $\mathbf{v} = \tilde{\mathbf{f}}^{-1} \circ \mathbf{f} : \mathcal{Z} \rightarrow \mathcal{Z}$ . Since all assumptions in Theorem 4 hold, we have

$$\mathbf{T}(\mathbf{Z}) = A \tilde{\mathbf{T}}(\mathbf{v}(\mathbf{Z})) + \mathbf{c}, \quad (39)$$

where  $A \in \mathbb{R}^{k \times k}$  is invertible. We want to show that  $A$  is a block permutation matrix.

**Step I.** In this step, we show that  $\mathbf{v}$  is a componentwise function. First we differentiate both sides of Eq. (39) with respect to  $Z_s$  and  $Z_t$  ( $s \neq t$ ) to obtain

$$\frac{\partial \mathbf{T}(\mathbf{Z})}{\partial Z_s} = A \sum_{i=1}^n \frac{\partial \tilde{\mathbf{T}}(\mathbf{v}(\mathbf{Z}))}{\partial v_i(\mathbf{Z})} \cdot \frac{\partial v_i(\mathbf{Z})}{\partial Z_s} \quad (40)$$

$$\frac{\partial^2 \mathbf{T}(\mathbf{Z})}{\partial Z_s \partial Z_t} = A \sum_{i=1}^n \sum_{j=1}^n \frac{\partial^2 \tilde{\mathbf{T}}(\mathbf{v}(\mathbf{Z}))}{\partial v_i(\mathbf{Z}) \partial v_j(\mathbf{Z})} \cdot \frac{\partial v_j(\mathbf{Z})}{\partial Z_t} \cdot \frac{\partial v_i(\mathbf{Z})}{\partial Z_s} + A \sum_{i=1}^n \frac{\partial \tilde{\mathbf{T}}(\mathbf{v}(\mathbf{Z}))}{\partial v_i(\mathbf{Z})} \cdot \frac{\partial^2 v_i(\mathbf{Z})}{\partial Z_s \partial Z_t}. \quad (41)$$

By construction, the second-order cross derivatives of  $\mathbf{T}$  and  $\tilde{\mathbf{T}}$  are all zero. Therefore, we have

$$\mathbf{0} = A \sum_{i=1}^n \frac{\partial^2 \tilde{\mathbf{T}}(\mathbf{v}(\mathbf{Z}))}{\partial v_i(\mathbf{Z})^2} \cdot \frac{\partial v_i(\mathbf{Z})}{\partial Z_t} \cdot \frac{\partial v_i(\mathbf{Z})}{\partial Z_s} + A \sum_{i=1}^n \frac{\partial \tilde{\mathbf{T}}(\mathbf{v}(\mathbf{Z}))}{\partial v_i(\mathbf{Z})} \cdot \frac{\partial^2 v_i(\mathbf{Z})}{\partial Z_s \partial Z_t}. \quad (42)$$

The above equation can be written in the matrix-vector form:

$$\mathbf{0} = A \tilde{\mathbf{T}}''(\mathbf{Z}) \mathbf{v}'_{s,t}(\mathbf{Z}) + A \tilde{\mathbf{T}}'(\mathbf{Z}) \mathbf{v}''_{s,t}(\mathbf{Z}), \quad (43)$$

where

$$\tilde{\mathbf{T}}''(\mathbf{Z}) = \left[ \frac{\partial^2 \tilde{\mathbf{T}}(\mathbf{v}(\mathbf{Z}))}{\partial v_1(\mathbf{Z})^2}, \dots, \frac{\partial^2 \tilde{\mathbf{T}}(\mathbf{v}(\mathbf{Z}))}{\partial v_n(\mathbf{Z})^2} \right] \in \mathbb{R}^{k \times n} \quad (44)$$

$$\mathbf{v}'_{s,t}(\mathbf{Z}) = \left[ \frac{\partial v_1(\mathbf{Z})}{\partial Z_t} \cdot \frac{\partial v_1(\mathbf{Z})}{\partial Z_s}, \dots, \frac{\partial v_n(\mathbf{Z})}{\partial Z_t} \cdot \frac{\partial v_n(\mathbf{Z})}{\partial Z_s} \right]^T \in \mathbb{R}^n, \quad (45)$$

and

$$\tilde{\mathbf{T}}'(\mathbf{Z}) = \left[ \frac{\partial \tilde{\mathbf{T}}(\mathbf{v}(\mathbf{Z}))}{\partial v_1(\mathbf{Z})}, \dots, \frac{\partial \tilde{\mathbf{T}}(\mathbf{v}(\mathbf{Z}))}{\partial v_n(\mathbf{Z})} \right] \in \mathbb{R}^{k \times n} \quad (46)$$

$$\mathbf{v}''_{s,t}(\mathbf{Z}) = \left[ \frac{\partial^2 v_1(\mathbf{Z})}{\partial Z_s \partial Z_t}, \dots, \frac{\partial^2 v_n(\mathbf{Z})}{\partial Z_s \partial Z_t} \right]^T \in \mathbb{R}^n. \quad (47)$$

Now by concatenating

$$\tilde{\mathbf{T}}'''(\mathbf{Z}) = [\tilde{\mathbf{T}}''(\mathbf{Z}), \tilde{\mathbf{T}}'(\mathbf{Z})] \in \mathbb{R}^{k \times 2n} \quad \text{and} \quad \mathbf{v}''_{s,t}(\mathbf{Z}) = [\mathbf{v}'_{s,t}(\mathbf{Z})^T, \mathbf{v}''_{s,t}(\mathbf{Z})^T]^T \in \mathbb{R}^{2n}, \quad (48)$$

we obtain

$$\mathbf{0} = A \tilde{\mathbf{T}}'''(\mathbf{Z}) \mathbf{v}''_{s,t}(\mathbf{Z}). \quad (49)$$

Finally, we take the rows of  $\tilde{\mathbf{T}}'''(\mathbf{Z})$  that corresponds to the factorized strongly exponential family distribution part and denote them by  $\tilde{\mathbf{T}}_f'''(\mathbf{Z}) \in \mathbb{R}^{k' \times 2n}$ . By Lemma 5 in the iVAE paper (Khemakhem et al., 2020a) and the assumption that  $k' \geq 2n$ , we have that the rank of  $\tilde{\mathbf{T}}_f'''(\mathbf{Z})$  is  $2n$ . Since  $k' \geq 2n$ , the rank of  $\tilde{\mathbf{T}}'''(\mathbf{Z})$  is also  $2n$ . Since the rank of  $A$  is  $k$ , the rank of  $A \tilde{\mathbf{T}}'''(\mathbf{Z}) \in \mathbb{R}^{k \times 2n}$  is  $2n$ . This implies that  $\mathbf{v}''_{s,t}(\mathbf{Z})$  must be a zero vector. In particular, we have that  $\mathbf{v}'_{s,t}(\mathbf{Z}) = \mathbf{0}$ ,  $\forall s \neq t$ . Therefore, we have shown that  $\mathbf{v}$  is a componentwise function.

**Step II.** To complete the proof, we need to show that  $A$  is a block permutation matrix. Without loss of generality, we assume that the permutation in  $\mathbf{v}$  is the identity. That is  $\mathbf{v}(\mathbf{Z}) = [v_1(Z_1), \dots, v_n(Z_n)]^T$  for some nonlinear univariate scalar functions  $v_1, \dots, v_n$ . Since  $\mathbf{f}$  and  $\tilde{\mathbf{f}}$  are bijective, we have that  $\mathbf{v}$  is also bijective and  $\mathbf{v}^{-1}(\mathbf{Z}) = [v_1^{-1}(Z_1), \dots, v_n^{-1}(Z_n)]^T$ . We denote  $\tilde{\mathbf{T}}(\mathbf{v}(\mathbf{Z})) = \tilde{\mathbf{T}}(\mathbf{v}(\mathbf{Z})) + A^{-1}\mathbf{c}$  and plug it into Eq. (39) to obtain  $\mathbf{T}(\mathbf{Z}) = A \tilde{\mathbf{T}}(\mathbf{v}(\mathbf{Z}))$ . Applying  $\mathbf{v}^{-1}$  to the variables  $\mathbf{Z}$  at both sides gives

$$\mathbf{T}(\mathbf{v}^{-1}(\mathbf{Z})) = A \tilde{\mathbf{T}}(\mathbf{Z}). \quad (50)$$

Let  $t$  be the index of an entry in the sufficient statistics  $\mathbf{T}$  that corresponds to the the factorized strongly exponential family distribution part  $\mathbf{T}_f$ . For all  $s \neq t$ , we have

$$0 = \frac{\partial \mathbf{T}(\mathbf{v}^{-1}(\mathbf{Z}))_t}{\partial Z_s} = \sum_{j=1}^k a_{tj} \frac{\partial \tilde{\mathbf{T}}(\mathbf{Z})_j}{\partial Z_s}. \quad (51)$$

Since the entries of  $\tilde{\mathbf{T}}$  are linearly independent (if they were not linearly independent, then  $\tilde{\mathbf{T}}$  can be compressed into a smaller vector by removing the redundant entries), we have that  $a_{tj}$  is zero for any  $j$  such that  $\frac{\partial \tilde{\mathbf{T}}(\mathbf{Z})_j}{\partial Z_s} \neq 0$ . This includes the entries  $j$  in the sufficient statistics  $\tilde{\mathbf{T}}$  that corresponds to 1) the factorized strongly exponential family distribution part which do not depend on  $Z_t$ ; and 2) the neural network part.

Therefore, when  $t$  is the index of an entry in the sufficient statistics  $\mathbf{T}$  that corresponds to factor  $i$  in the factorized strongly exponential family distribution part  $\mathbf{T}_f$ , the only non-zero  $a_{tj}$  are the ones that map between  $\mathbf{T}_i(Z_i)$  and  $\tilde{\mathbf{T}}_i(v_i(Z_i))$ , where  $\mathbf{T}_i$  are the factors in  $\mathbf{T}_f$  that only depend on  $Z_i$  and  $\tilde{\mathbf{T}}_i$  is defined similarly. Therefore, we can construct an invertible submatrix  $A'_i$  with all non-zero elements  $a_{tj}$  for all  $t$  that corresponds to factor  $i$ , such that

$$\mathbf{T}_i(Z_i) = A'_i \tilde{\mathbf{T}}_i(v_i(Z_i)) = A'_i \tilde{\mathbf{T}}_i(v_i(Z_i)) + \mathbf{c}_i, \quad i = 1, \dots, n, \quad (52)$$

where  $\tilde{\mathbf{T}}_i$  are the factors in  $\tilde{\mathbf{T}}_f$  that only depends on  $Z_i$ , and  $\mathbf{c}_i$  are the corresponding elements of  $\mathbf{c}$ . This means that the matrix  $A$  is a block permutation matrix. For each  $i = 1, \dots, n$ , the block  $A'_i$  of  $A$  affinely transforms  $\mathbf{T}_i(Z_i)$  into  $\tilde{\mathbf{T}}_i(v_i(Z_i))$ . There is also an additional block  $A'$  which affinely transforms  $\mathbf{T}_{NN}(\mathbf{Z})$  into  $\tilde{\mathbf{T}}_{NN}(\mathbf{v}(\mathbf{Z}))$ . This completes the proof.  $\square$

### H.1.3 PART III

Now we combine Theorem 4 in Part I and Theorem 5 in Part II into one theorem, which completes the proof of Theorem 1.

## H.2 PROOF OF THEOREM 2

We recall that the loss function in Phase 1 is as follows:

$$\mathcal{L}_{\text{phase1}}(\boldsymbol{\theta}, \phi) = \mathcal{L}_{\text{phase1}}^{\text{ELBO}}(\hat{\mathbf{f}}, \hat{\mathbf{T}}, \hat{\boldsymbol{\lambda}}, \phi) - \mathcal{L}_{\text{phase1}}^{\text{SM}}(\hat{\mathbf{f}}, \mathbf{T}, \boldsymbol{\lambda}, \hat{\phi}), \quad (53)$$

where

$$\mathcal{L}_{\text{phase1}}^{\text{ELBO}}(\boldsymbol{\theta}, \phi) := \mathbb{E}_{p_D} \left[ \mathbb{E}_{q_\phi(\mathbf{Z}|\mathbf{X}, \mathbf{Y}, \mathbf{E})} [\log p_f(\mathbf{X}|\mathbf{Z}) + \log p_{\mathbf{T}, \boldsymbol{\lambda}}(\mathbf{Z}|\mathbf{Y}, \mathbf{E}) - \log q_\phi(\mathbf{Z}|\mathbf{X}, \mathbf{Y}, \mathbf{E})] \right], \quad (54)$$

$$\mathcal{L}_{\text{phase1}}^{\text{SM}}(\mathbf{T}, \boldsymbol{\lambda}) := \mathbb{E}_{p_D} \left[ \mathbb{E}_{q_\phi(\mathbf{Z}|\mathbf{X}, \mathbf{Y}, \mathbf{E})} [|\nabla_{\mathbf{Z}} \log q_\phi(\mathbf{Z}|\mathbf{X}, \mathbf{Y}, \mathbf{E}) - \nabla_{\mathbf{Z}} \log p_{\mathbf{T}, \boldsymbol{\lambda}}(\mathbf{Z}|\mathbf{Y}, \mathbf{E})|^2] \right]. \quad (55)$$

*Proof.* If the family of  $q_\phi(\mathbf{Z}|\mathbf{X}, \mathbf{Y}, \mathbf{E})$  is flexible enough to contain  $p_\theta(\mathbf{Z}|\mathbf{X}, \mathbf{Y}, \mathbf{E})$ , then by optimizing the loss over its parameter  $\phi$ , we will minimize the score matching term  $\mathcal{L}_{\text{phase1}}^{\text{SM}}$  in Eq. (55), which will eventually reach zero. If we assume that the model is not degenerate and that  $q_\phi > 0$  everywhere, then having that  $\mathcal{L}_{\text{phase1}}^{\text{SM}} = 0$  implies that  $\nabla_{\mathbf{Z}} \log q_\phi(\mathbf{Z}|\mathbf{X}, \mathbf{Y}, \mathbf{E})$  and  $\nabla_{\mathbf{Z}} \log p_{\mathbf{T}, \boldsymbol{\lambda}}(\mathbf{Z}|\mathbf{Y}, \mathbf{E})$  are equal. This implies that  $\log q_\phi(\mathbf{Z}|\mathbf{X}, \mathbf{Y}, \mathbf{E}) = \log p_{\mathbf{T}, \boldsymbol{\lambda}}(\mathbf{Z}|\mathbf{Y}, \mathbf{E}) + c$  for some constant  $c$ . But  $c$  is necessarily 0 because both  $q_\phi(\mathbf{Z}|\mathbf{X}, \mathbf{Y}, \mathbf{E})$  and  $p_{\mathbf{T}, \boldsymbol{\lambda}}(\mathbf{Z}|\mathbf{Y}, \mathbf{E})$  are pdf's. Therefore, the ELBO term  $\mathcal{L}_{\text{phase1}}^{\text{ELBO}}$  in Eq. (54) will be equal to the log-likelihood, meaning that the loss  $\mathcal{L}_{\text{phase1}}$  in Eq. (53) will be equal to the log-likelihood. Under this circumstance, the estimation in Eq. (53) inherits all the properties of maximum likelihood estimation (MLE). In this particular case, since our identifiability is guaranteed up to a permutation and componentwise transformation, the consistency of MLE means that we converge to the true parameter  $\boldsymbol{\theta}^*$  up to a permutation and componentwise transformation in the limit of infinite data. Because true identifiability is one of the assumptions for MLE consistency, replacing it by identifiability up to a permutation and componentwise transformation does not change the proof but only the conclusion.  $\square$

## H.3 PROOF OF THEOREM 3

*Proof.* Theorem 1 and Theorem 2 guarantee that in the limit of infinite data, NF-iVAE can learn the true parameters  $\boldsymbol{\theta}^* := (\mathbf{f}^*, \mathbf{T}^*, \boldsymbol{\lambda}^*)$  up to a permutation and componentwise transformation of the latent variables. Let  $(\hat{\mathbf{f}}, \hat{\mathbf{T}}, \hat{\boldsymbol{\lambda}})$  be the parameters obtained by NF-iVAE. We, therefore, have  $(\hat{\mathbf{f}}, \hat{\mathbf{T}}, \hat{\boldsymbol{\lambda}}) \sim_P (\mathbf{f}^*, \mathbf{T}^*, \boldsymbol{\lambda}^*)$ , where  $\sim_P$  denotes the equivalence up to a permutation and componentwise transformation. If there were no noise, this would mean that the learned  $\hat{\mathbf{f}}$  transforms  $\mathbf{X}$  into  $\hat{\mathbf{Z}} = \hat{\mathbf{f}}^{-1}(\mathbf{X})$  that are equal to  $\mathbf{Z}^* = (\mathbf{f}^*)^{-1}(\mathbf{X})$  up to a permutation and componentwise transformation (Definition 9). If with noise, we obtain the posterior distribution of the latent variables up to an analogous indeterminacy.  $\square$

## H.4 PROOF OF PROPOSITION 1

*Proof.* Theorem A.1 in (Arjovsky, 2021) has showed that i) any predictor  $w \circ \Phi$  with optimal OOD generalization uses only  $\text{Pa}(\mathbf{Y})$  to compute  $\Phi$ ; ii) the classifier  $w$  in this optimal predictor can be estimated using data from any environment  $e$  for which the distribution of  $\text{Pa}(\mathbf{Y}^e)$  has full support;

iii) the optimal predictor will be invariant across  $\mathcal{E}_{all}$ . In iCaRL, the hypotheses of Theorems 1 and 2 and Assumption 1 guarantee that  $\text{Pa}(\mathbf{Y})$  can be recovered by first identifying the latent variables  $\mathbf{Z}$  from  $\mathbf{X}$ ,  $\mathbf{Y}$  and  $\mathbf{E}$  and then discovering the direct causes of  $\mathbf{Y}$  through solving Eq. (12). Furthermore, since the conditional prior in Eq. (5) of the main text has full support, the distribution of  $\text{Pa}(\mathbf{Y}^e)$  always has full support. Also, under Assumption 1 we have that  $p(\mathbf{Y}|\text{Pa}(\mathbf{Y}))$  is invariant across  $\mathcal{E}_{all}$ . Hence, the classifier  $w$  in this optimal predictor can be estimated using data from any environment  $e$ . We therefore have that the resulting optimal predictor will be invariant across  $\mathcal{E}_{all}$ . This completes the proof.  $\square$

## H.5 PROOF OF PROPOSITION 2

*Proof.* The following rules can be independently performed to distinguish all the 10 structures shown in Figs. 3b-3k. For clarity, we divide them into three groups. Note that, since these rules rely on different algorithms of causal discovery and conditional independence tests, unless stated otherwise, we assume that the assumptions of these algorithms are satisfied during the proof process.

**Group 1** All the six structures in this group can be discovered only by performing conditional independence tests.

- **Rule 1.1** If  $Z_i \perp\!\!\!\perp \mathbf{Y}$ ,  $Z_i \not\perp\!\!\!\perp \mathbf{E}$ , and  $\mathbf{E} \perp\!\!\!\perp \mathbf{Y}$ , then Fig. 3d is discovered.
- **Rule 1.2** If  $Z_i \not\perp\!\!\!\perp \mathbf{Y}$ ,  $Z_i \perp\!\!\!\perp \mathbf{E}$ , and  $\mathbf{E} \not\perp\!\!\!\perp \mathbf{Y}$ , then Fig. 3g is discovered.
- **Rule 1.3** If  $Z_i \not\perp\!\!\!\perp \mathbf{Y}$ ,  $Z_i \not\perp\!\!\!\perp \mathbf{E}$ , and  $\mathbf{E} \perp\!\!\!\perp \mathbf{Y}$ , then Fig. 3f is discovered.
- **Rule 1.4** If  $Z_i \not\perp\!\!\!\perp \mathbf{Y}$ ,  $Z_i \not\perp\!\!\!\perp \mathbf{E}$ ,  $\mathbf{E} \not\perp\!\!\!\perp \mathbf{Y}$ , and  $Z_i \perp\!\!\!\perp \mathbf{Y}|\mathbf{E}$ , then Fig. 3i is discovered.
- **Rule 1.5** If  $Z_i \not\perp\!\!\!\perp \mathbf{Y}$ ,  $Z_i \not\perp\!\!\!\perp \mathbf{E}$ ,  $\mathbf{E} \not\perp\!\!\!\perp \mathbf{Y}$ , and  $Z_i \perp\!\!\!\perp \mathbf{E}|\mathbf{Y}$ , then Fig. 3h is discovered.
- **Rule 1.6** If  $Z_i \not\perp\!\!\!\perp \mathbf{Y}$ ,  $Z_i \not\perp\!\!\!\perp \mathbf{E}$ ,  $\mathbf{E} \not\perp\!\!\!\perp \mathbf{Y}$ , and  $\mathbf{Y} \perp\!\!\!\perp \mathbf{E}|Z_i$ , then Fig. 3e is discovered.

**Group 2** If  $Z_i \not\perp\!\!\!\perp \mathbf{Y}$ ,  $Z_i \perp\!\!\!\perp \mathbf{E}$ , and  $\mathbf{E} \perp\!\!\!\perp \mathbf{Y}$ , then we can discover both Fig. 3b and Fig. 3c. These two structures cannot be further distinguished only by conditional independence tests, because they come from the same Markov equivalence class. Fortunately, we can further distinguish them by running binary causal discovery algorithms (Peters et al., 2017), e.g., ANM (Hoyer et al., 2009) for continuous data and CDS (Fonollosa, 2019) for continuous and discrete data.

- **Rule 2.1** If  $Z_i \not\perp\!\!\!\perp \mathbf{Y}$ ,  $Z_i \perp\!\!\!\perp \mathbf{E}$ , and  $\mathbf{E} \perp\!\!\!\perp \mathbf{Y}$ , and a chosen binary causal discovery algorithm prefers  $Z_i \rightarrow \mathbf{Y}$  to  $Z_i \leftarrow \mathbf{Y}$ , then Fig. 3b is discovered.
- **Rule 2.2** If  $Z_i \not\perp\!\!\!\perp \mathbf{Y}$ ,  $Z_i \perp\!\!\!\perp \mathbf{E}$ , and  $\mathbf{E} \perp\!\!\!\perp \mathbf{Y}$ , and a chosen binary causal discovery algorithm prefers  $Z_i \leftarrow \mathbf{Y}$  to  $Z_i \rightarrow \mathbf{Y}$ , then Fig. 3c is discovered.

**Group 3** If  $Z_i \not\perp\!\!\!\perp \mathbf{Y}$ ,  $Z_i \not\perp\!\!\!\perp \mathbf{E}$ ,  $\mathbf{E} \not\perp\!\!\!\perp \mathbf{Y}$ ,  $Z_i \not\perp\!\!\!\perp \mathbf{Y}|\mathbf{E}$ ,  $Z_i \not\perp\!\!\!\perp \mathbf{E}|\mathbf{Y}$ , and  $\mathbf{Y} \not\perp\!\!\!\perp \mathbf{E}|Z_i$ , then we can discover both Fig. 3j and Fig. 3k. These two structures cannot be further distinguished only by conditional independence tests, because they come from the same Markov equivalence class. They also cannot be distinguished by any binary causal discovery algorithm, since both  $Z_i$  and  $\mathbf{Y}$  are affected by  $\mathbf{E}$ . Fortunately, Zhang et al. (2017) provided a heuristic solution to this case based on the invariance of causal mechanisms, i.e.,  $P(\text{cause})$  and  $P(\text{effect}|\text{cause})$  change independently. The detailed description of their method is given in Section 4.2 of Zhang et al. (2017). For convenience, here we directly borrow their final result. Zhang et al. (2017) states that determining the causal direction between  $Z_i$  and  $\mathbf{Y}$  in Fig. 3j and Fig. 3k is finally reduced to calculating the following term:

$$\Delta_{Z_i \rightarrow \mathbf{Y}} = \left\langle \log \frac{\bar{P}(\mathbf{Y}|Z_i)}{\langle \hat{P}(\mathbf{Y}|Z_i) \rangle} \right\rangle, \quad (56)$$

where  $\langle \cdot \rangle$  denotes the sample average,  $\bar{P}(\mathbf{Y}|Z_i)$  is the empirical estimate of  $P(\mathbf{Y}|Z_i)$  on all data points, and  $\langle \hat{P}(\mathbf{Y}|Z_i) \rangle$  denotes the sample average of  $\hat{P}(\mathbf{Y}|Z_i)$ , which is the estimate of  $P(\mathbf{Y}|Z_i)$  in each environment. We take the direction for which  $\Delta$  is smaller to be the causal direction.

- **Rule 3.1** If  $Z_i \not\perp\!\!\!\perp Y$ ,  $Z_i \not\perp\!\!\!\perp E$ ,  $E \not\perp\!\!\!\perp Y$ ,  $Z_i \not\perp\!\!\!\perp Y|E$ ,  $Z_i \not\perp\!\!\!\perp E|Y$ ,  $Y \not\perp\!\!\!\perp E|Z_i$ , and  $\Delta_{Z_i \rightarrow Y}$  is smaller than  $\Delta_{Y \rightarrow Z_i}$ , then Fig. 3j is discovered.
- **Rule 3.2** If  $Z_i \not\perp\!\!\!\perp Y$ ,  $Z_i \not\perp\!\!\!\perp E$ ,  $E \not\perp\!\!\!\perp Y$ ,  $Z_i \not\perp\!\!\!\perp Y|E$ ,  $Z_i \not\perp\!\!\!\perp E|Y$ ,  $Y \not\perp\!\!\!\perp E|Z_i$ , and  $\Delta_{Y \rightarrow Z_i}$  is smaller than  $\Delta_{Z_i \rightarrow Y}$ , then Fig. 3k is discovered.

□

## I ILLUSTRATIONS FOR MODEL LEARNING

As described in Section 3.3, in the ground truth model, the prior for each  $Z_{i \in I_p}$  is either  $p(Z_i|E)$  or  $p(Z_i)$ , depending on whether  $Z_i$  is caused by  $E$  or not. By contrast, in the NF-iVAE model the prior is  $p(Z|Y, E)$ . Is this going to affect the identifiability of the latent variables? Well, in practice not because the posterior distribution for  $Z$  given  $X, Y$  and  $E$  would be equivalent in both models (up to identifiability guarantees).

## J DATASETS

For convenience and completeness, we provide descriptions of Colored MNIST Digits and Colored Fashion MNIST here. Please refer to the original papers (Arjovsky et al., 2019; Ahuja et al., 2020a; Gulrajani & Lopez-Paz, 2020; Venkateswara et al., 2017) for more details.

### J.1 SYNTHETIC DATA

For the nonlinear transformation, we use the MLP:

- Input layer: Input batch (*batch size, input dimension*)
- Layer 1: Fully connected layer, output size = 6, activation = ReLU
- Output layer: Fully connected layer, output size = 10

### J.2 COLORED MNIST DIGITS

We use the exact same environment as in Arjovsky et al. (2019). Arjovsky et al. (2019) propose to create an environment for training to classify digits in MNIST data<sup>10</sup>, where the images in MNIST are now colored in such a way that the colors spuriously correlate with the labels. The task is to classify whether the digit is less than 5 (not including 5). There are three environments (two training containing 30,000 points each, one test containing 10,000 points) We add noise to the preliminary label ( $\tilde{y} = 0$  if the digit is between 0-4 and  $\tilde{y} = 1$  if the digit is between 5-9) by flipping it with 25 percent probability to construct the final labels. We sample the color id  $z$  by flipping the final labels with probability  $p_e$ , where  $p_e$  is 0.2 in the first environment, 0.1 in the second environment, and 0.9 in the third environment. The third environment is the testing environment. We color the digit red if  $z = 1$  or green if  $z = 0$ .

### J.3 COLORED FASHION MNIST

We modify the fashion MNIST dataset<sup>11</sup> in a manner similar to the MNIST digits dataset. Fashion MNIST data has images from different categories: “t-shirt”, “trouser”, “pullover”, “dress”, “coat”, “sandal”, “shirt”, “sneaker”, “bag”, “ankle boots”. We add colors to the images in such a way that the colors correlate with the labels. The task is to classify whether the image is that of foot wear or a clothing item. There are three environments (two training, one test) We add noise to the preliminary label ( $\tilde{y} = 0$ : “t-shirt”, “trouser”, “pullover”, “dress”, “coat”, “shirt” and  $\tilde{y} = 1$ : “sandal”, “sneaker”, “ankle boots”) by flipping it with 25 percent probability to construct the final label. We sample the color id  $z$  by flipping the noisy label with probability  $p_e$ , where  $p_e$  is 0.2 in the first environment, 0.1

<sup>10</sup>[https://www.tensorflow.org/api\\_docs/python/tf/keras/datasets/mnist/load\\_data](https://www.tensorflow.org/api_docs/python/tf/keras/datasets/mnist/load_data)

<sup>11</sup>[https://www.tensorflow.org/api\\_docs/python/tf/keras/datasets/fashion\\_mnist/load\\_data](https://www.tensorflow.org/api_docs/python/tf/keras/datasets/fashion_mnist/load_data)



Table 3: Results on synthetic data. Comparisons are in terms of MSE (mean  $\pm$  std deviation).

$g(\cdot)$	METHOD	TRAIN ( $\sigma_3 = \{0.2, 2\}$ )	TEST ( $\sigma_3 = 100$ )
Identity	ERM	0.00 $\pm$ 0.00	<b>0.00 <math>\pm</math> 0.00</b>
	IRM	0.00 $\pm$ 0.00	<b>0.00 <math>\pm</math> 0.00</b>
	F-IRM GAME	0.98 $\pm$ 0.23	1.03 $\pm$ 0.04
	V-IRM GAME	0.99 $\pm$ 2.74	1.07 $\pm$ 2.26
	<b>iCaRL (ours)</b>	0.01 $\pm$ 0.03	1.00 $\pm$ 0.01
Linear	ERM	0.00 $\pm$ 0.00	<b>0.00 <math>\pm</math> 0.00</b>
	IRM	0.00 $\pm$ 0.00	<b>0.00 <math>\pm</math> 0.00</b>
	F-IRM GAME	0.99 $\pm$ 0.01	1.08 $\pm$ 0.06
	V-IRM GAME	1.00 $\pm$ 5.98	1.05 $\pm$ 0.04
	<b>iCaRL (ours)</b>	0.01 $\pm$ 0.03	1.01 $\pm$ 0.04
Nonlinear	ERM	0.06 $\pm$ 0.01	220.79 $\pm$ 229.97
	IRM	0.08 $\pm$ 0.01	149.60 $\pm$ 104.85
	F-IRM GAME	1.06 $\pm$ 0.09	196.59 $\pm$ 150.71
	V-IRM GAME	1.00 $\pm$ 0.01	170.46 $\pm$ 125.62
	<b>iCaRL (ours)</b>	0.29 $\pm$ 0.04	<b>28.16 <math>\pm</math> 2.54</b>

in the second environment, and 0.9 in the third environment, which is the test environment. We color the object red if  $z = 1$  or green if  $z = 0$ .

## K IN-DEPTH ANALYSIS ON SYNTHETIC DATA

### K.1 COMPARISONS WITH STATE-OF-THE-ART

We first conduct a series of experiments on synthetic data generated according to an extension of the SEM in Model 1. The extension is to map the variables  $\mathbf{Z} := (Z_1, Z_2)$  into a 10 dimensional observation  $\mathbf{X}$  through a linear or nonlinear transformation. Our goal is to predict  $Y$  from  $\mathbf{X}$ , where  $\mathbf{X} = g(\mathbf{Z})$ . We consider three transformations: (a) *Identity*:  $g(\cdot)$  is the identity matrix  $\mathbf{I} \in \mathbb{R}^{2 \times 2}$ , i.e.,  $\mathbf{X} = g(\mathbf{Z}) = \mathbf{Z}$ . (b) *Linear*:  $g(\cdot)$  is a random matrix  $\mathbf{S} \in \mathbb{R}^{2 \times 10}$ , i.e.,  $\mathbf{X} = g(\mathbf{Z}) = \mathbf{Z} \cdot \mathbf{S}$ . (c) *Nonlinear*:  $g(\cdot)$  is given by a neural network with 2-dimensional input and 10-dimensional output, whose parameters are randomly set in advance. Since this is a regression task, we use the mean squared error (MSE) as a metric of performance. Note that, in this problem, there is only one causal latent variable  $Z_1$ , meaning that the conditional prior in Eq. (5) will not exhibit dependencies. Because of this, in this case we do not include a  $T_{NN}(\mathbf{Z})$  term in our NF-iVAE conditional prior. In the following section we do consider settings with many potential causal latent variables and, in that case, we do include  $T_{NN}(\mathbf{Z})$  in the NF-iVAE prior.

We consider a simple scenario in which we fix  $\sigma_1 = 1$  and  $\sigma_2 = 0$  for all environments and only allow  $\sigma_3$  to vary across environments. In this case,  $\sigma_3$  controls the spurious correlations between  $Z_2$  and  $Y$ . Each experiment draws 1000 samples from each of the three environments  $\sigma_3 = \{0.2, 2, 100\}$ , where the first two are for training and the third for testing. We compare with the following baselines:<sup>12</sup> ERM, and two variants of IRMG: F-IRM Game (with  $\Phi$  fixed to the identity) and V-IRM Game (with variable  $\Phi$ ).

As shown in Table 3, in the cases of *Identity* and *Linear*, our approach is better than IRMG but only comparable with ERM and IRM. This might be because the identifiability result up to a simple nonlinear transformation renders the problem more difficult by converting the original identity or linear problem into a nonlinear problem. In the *Nonlinear* case, the gains of iCaRL are very clear.

Table 4: Comparisons of assumptions on the prior leading to identifiability in different algorithms.

METHOD	Assumption on the Prior for Identifiability
VAE (Kingma & Welling, 2013)	Non-identifiability with $p_{T,\lambda}(\mathbf{Z}) = \prod_i p(Z_i) \stackrel{e.g.}{=} \mathcal{N}(\mathbf{0}, \mathbf{I})$
iVAE (Khemakhem et al., 2020a)	$p_{T,\lambda}(\mathbf{Z} \mathbf{Y}, \mathbf{E}) = \prod_i Q_i(Z_i)/C_i(\mathbf{Y}, \mathbf{E}) \exp[\sum_{j=1}^k T_{i,j}(Z_i)\lambda_{i,j}(\mathbf{Y}, \mathbf{E})]$
ICE-BeeM (Khemakhem et al., 2020b)	$p_{T,\lambda}(\mathbf{Z} \mathbf{Y}, \mathbf{E}) = Q(\mathbf{Z})/C(\mathbf{Y}, \mathbf{E}) \prod_i \exp[\sum_{j=1}^k T_{i,j}(Z_i)\lambda_{i,j}(\mathbf{Y}, \mathbf{E})]$
NF-iVAE	$p_{T,\lambda}(\mathbf{Z} \mathbf{Y}, \mathbf{E}) = Q(\mathbf{Z})/C(\mathbf{Y}, \mathbf{E}) \exp[\mathbf{T}(\mathbf{Z})^T \boldsymbol{\lambda}(\mathbf{Y}, \mathbf{E})]$

<sup>12</sup>We also tried ICP, but ICP was unable to find any parent of  $Y$  even in the identity case.

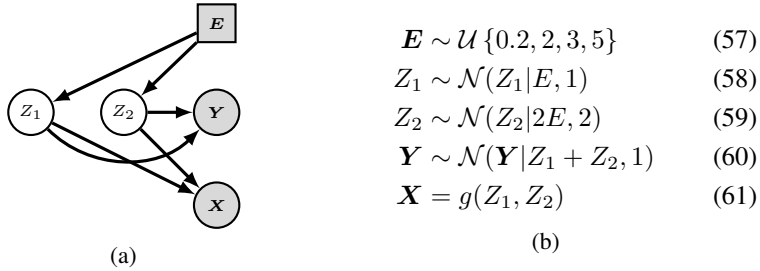


Figure 4: (a) Causal structure with  $Y$  having two causes. (b) Data generating process corresponding to (a), where  $\mathcal{U}\{\cdot\}$  denotes the discrete uniform distribution,  $\mathcal{N}(\cdot)$  the Gaussian distribution, and  $g(\cdot)$  is given by a neural network with 2-dimensional input and 10-dimensional output, whose parameters are randomly set in advance.

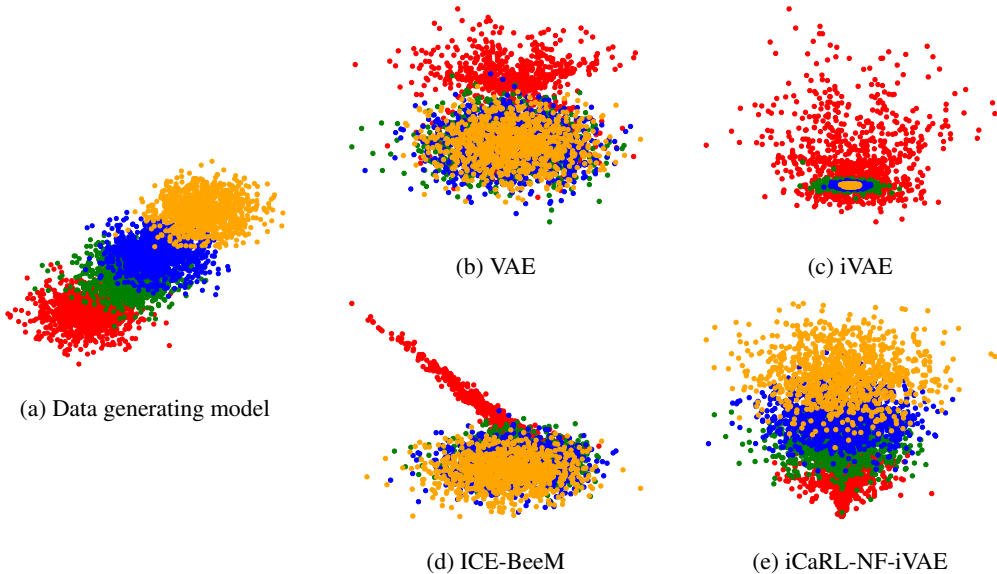


Figure 5: Visualization of the samples (i.e.,  $\hat{Z} = (\hat{Z}_1, \hat{Z}_2)$ ) in latent space generated through different algorithms. (a) Samples from the true distribution. (b-e) Samples from the posterior of different algorithms. Apparently, our method (e) can recover the original data up to a permutation and a simple componentwise transformation.

### K.2 VISUALIZATION OF IDENTIFIABILITY OF NF-iVAE

To further verify identifiability of NF-iVAE, we conduct a series of experiments on synthetic data generated according to the data generating process (Fig. 4b) corresponding to the causal graph shown in Fig. 4a. The reason we choose this setting is that it is the simplest case satisfying our requirements: a) For ease of visualization, the latent space had better be 2-dimensional; b) To introduce the non-factorized prior given  $Y$  and  $E$  (i.e.,  $Z_i \not\perp Z_j | Y, E$ ),  $Y$  has at least two causes.

We draw 1000 samples from each of the four environments  $E = \{0.2, 2, 3, 5\}$ , and thus the whole synthetic dataset consists of 4000 samples. We compare with the following baselines: VAE Kingma & Welling (2013) (without identifiability guarantees), iVAE Khemakhem et al. (2020a) (with a conditionally factorized prior for identifiability), and ICE-BeeM Khemakhem et al. (2020b). It is worth noting that in ICE-BeeM the primary assumption leading to identifiability is similar to that in iVAE, where the base measure  $Q(Z)$  could be arbitrary to capture the dependences between latent variables but the exponential term still has to factorize across components (dimensions). All these are summarized in Table 4. Clearly, from the table we can see that our method has a more general assumption on the prior leading to identifiability. This is also demonstrated empirically in Fig. 5. Our method iCaRL can recover the original data  $Z$  up to a permutation and a simple componentwise

Table 5: Colored MNIST. Comparisons in terms of accuracy (%) (mean  $\pm$  std deviation).

METHOD	TRAIN	TEST
ERM	84.88 $\pm$ 0.16	10.45 $\pm$ 0.66
ERM 1	84.84 $\pm$ 0.21	10.86 $\pm$ 0.52
ERM 2	84.95 $\pm$ 0.20	10.05 $\pm$ 0.23
ROBUST MIN MAX	84.25 $\pm$ 0.43	15.24 $\pm$ 2.45
F-IRM GAME	63.37 $\pm$ 1.14	59.91 $\pm$ 2.69
V-IRM GAME	63.97 $\pm$ 1.03	49.06 $\pm$ 3.43
IRM	59.27 $\pm$ 4.39	62.75 $\pm$ 9.59
<b>iCaRL (ours)</b>	<b>70.56 <math>\pm</math> 0.81</b>	<b>68.75 <math>\pm</math> 1.45</b>
ERM GRAYSCALE	71.81 $\pm$ 0.47	71.36 $\pm$ 0.65
OPTIMAL	75	75

Table 6: PACS. Comparisons in terms of accuracy (%) (mean  $\pm$  std deviation).

METHOD	TEST
ERM	85.7 $\pm$ 0.5
IRM	84.4 $\pm$ 1.1
DRO (Sagawa et al., 2019)	84.1 $\pm$ 0.4
Mixup (Yan et al., 2020)	84.3 $\pm$ 0.5
CORAL (Sun & Saenko, 2016)	86.0 $\pm$ 0.2
MMD (Li et al., 2018b)	85.0 $\pm$ 0.2
DANN (Ganin et al., 2016)	84.6 $\pm$ 1.1
C-DANN (Li et al., 2018c)	82.8 $\pm$ 1.5
LaCIM (Sun et al., 2020)	83.5 $\pm$ 1.2
<b>iCaRL (ours)</b>	<b>88.7 <math>\pm</math> 0.6</b>

transformation, whereas all the other methods fail because they are unable to handle the case in which  $Z_i \not\perp Z_j | \mathbf{Y}, \mathbf{E}$ .

## L IN-DEPTH ANALYSIS ON MORE REALISTIC DATA

### L.1 COLORED MNIST

We compare iCaRL with 1) IRM, 2) two variants of IRMG: F-IRM Game (with  $\Phi$  fixed to the identity) and V-IRM Game (with a variable  $\Phi$ ), 3) three variants of ERM: ERM (on entire training data), ERM  $e$  (on each environment  $e$ ), and ERM GRAYSCALE (on data with no spurious correlations), and 4) ROBUST MIN MAZ (minimizing the maximum loss across the multiple environments). Table 1 shows that iCaRL outperforms all other baselines on Colored MNIST. However, this dataset seems more difficult because even ERM GRAYSCALE, where the spurious correlation with color is removed, falls well short of the optimum.

### L.2 PACS

We report the results on another one of the widely used realistic datasets for OOD generalization: PACS (Li et al., 2017a). This dataset consists of 9,991 images of dimension (3, 224, 224) and 7 classes from four domains: art, cartoons, photos, and sketches. We used the exact experimental setting that is described in Gulrajani & Lopez-Paz (2020). We provide results averaged over all possible train and test environment combination for one of the commonly used hyper-parameter tuning procedure: train domain validation. As shown in Table 6, iCaRL achieves state-of-the-art performance when compared to those most popular domain generalization alternatives.

## M IMPLEMENTATION DETAILS

### M.1 JOINT TRAINING

As described in Section 4.1, we can jointly learn  $(\theta, \phi)$  by optimizing the following objective:

$$\begin{aligned}
 \mathcal{L}_{\text{phase1}}(\theta, \phi) &= \mathcal{L}_{\text{phase1}}^{\text{ELBO}}(\mathbf{f}, \hat{\mathbf{T}}, \hat{\lambda}, \phi) - \mathcal{L}_{\text{phase1}}^{\text{SM}}(\hat{\mathbf{f}}, \mathbf{T}, \lambda, \hat{\phi}) & (62) \\
 &= \mathbb{E}_{p_D} \left[ \mathbb{E}_{q_{\phi}(\mathbf{Z}|\mathbf{X}, \mathbf{Y}, \mathbf{E})} \left[ \log p_{\mathbf{f}}(\mathbf{X}|\mathbf{Z}) + \log p_{\hat{\mathbf{T}}, \hat{\lambda}}(\mathbf{Z}|\mathbf{Y}, \mathbf{E}) - \log q_{\phi}(\mathbf{Z}|\mathbf{X}, \mathbf{Y}, \mathbf{E}) \right] \right] \\
 &\quad - \mathbb{E}_{p_D} \left[ \mathbb{E}_{q_{\hat{\phi}}(\mathbf{Z}|\mathbf{X}, \mathbf{Y}, \mathbf{E})} \left[ \left\| \nabla_{\mathbf{Z}} \log q_{\hat{\phi}}(\mathbf{Z}|\mathbf{X}, \mathbf{Y}, \mathbf{E}) - \nabla_{\mathbf{Z}} \log p_{\mathbf{T}, \lambda}(\mathbf{Z}|\mathbf{Y}, \mathbf{E}) \right\|^2 \right] \right] & (63) \\
 &= \mathbb{E}_{p_D} \left[ \mathbb{E}_{q_{\phi}(\mathbf{Z}|\mathbf{X}, \mathbf{Y}, \mathbf{E})} \left[ \log p_{\mathbf{f}}(\mathbf{X}|\mathbf{Z}) + \log p_{\hat{\mathbf{T}}, \hat{\lambda}}(\mathbf{Z}|\mathbf{Y}, \mathbf{E}) - \log q_{\phi}(\mathbf{Z}|\mathbf{X}, \mathbf{Y}, \mathbf{E}) \right] \right] \\
 &\quad - \mathbb{E}_{p_D} \left[ \mathbb{E}_{q_{\hat{\phi}}(\mathbf{Z}|\mathbf{X}, \mathbf{Y}, \mathbf{E})} \left[ \sum_{j=1}^n \left[ \frac{\partial^2 p_{\mathbf{T}, \lambda}(\mathbf{Z}|\mathbf{Y}, \mathbf{E})}{\partial Z_j^2} + \frac{1}{2} \left( \frac{\partial p_{\mathbf{T}, \lambda}(\mathbf{Z}|\mathbf{Y}, \mathbf{E})}{\partial Z_j} \right)^2 \right] \right] \right] \\
 &\quad + \text{const.} & (64)
 \end{aligned}$$

where the last equality is due to the equation in Appendix C.2, and  $\hat{\mathbf{f}}, \hat{\mathbf{T}}, \hat{\lambda}, \hat{\phi}$  are copies of  $\mathbf{f}, \mathbf{T}, \lambda, \phi$  that are treated as constants and whose gradient is not calculated during learning. In practice,  $\hat{\mathbf{f}}, \hat{\mathbf{T}}, \hat{\lambda}, \hat{\phi}$  can be easily implemented through either “detach” in PyTorch [Paszke et al. \(2019\)](#) or “stop\_gradient” in TensorFlow [Abadi et al. \(2015\)](#).

### M.2 THE GENERAL NON-FACTORIZED PRIOR

In the experiments, the general non-factorized prior in Assumption 2 is implemented as follows:

$$p_{\mathbf{T}, \lambda}(\mathbf{Z}|\mathbf{Y}, \mathbf{E}) = \left\langle \underbrace{\text{NN}(\mathbf{Z}; \text{param1})}_{\mathbf{T}_{\text{NN}}(\mathbf{Z})}, \underbrace{\text{NN}(\mathbf{Y}, \mathbf{E}; \text{param2})}_{\lambda_{\text{NN}}(\mathbf{Y}, \mathbf{E})} \right\rangle + \left\langle \underbrace{\text{concat}(\mathbf{Z}, \mathbf{Z}^2)}_{\mathbf{T}_f(\mathbf{Z})}, \underbrace{\text{NN}(\mathbf{Y}, \mathbf{E}; \text{param3})}_{\lambda_f(\mathbf{Y}, \mathbf{E})} \right\rangle,$$

where  $\langle \cdot, \cdot \rangle$  is the dot product of two vectors, and  $\text{concat}(\cdot, \cdot)$  means the concatenation of two vectors. Now let us explain each term in details.

Firstly,  $\text{concat}(\mathbf{Z}, \mathbf{Z}^2)$  is a vector of the latent variables and their squared values, and  $\text{NN}(\mathbf{Y}, \mathbf{E}; \text{param3})$  is a deep neural network parameterized by `param3` that computes a vector of natural parameters as a function of  $\mathbf{Y}$  and  $\mathbf{E}$ . Hence, the term  $\langle \text{concat}(\mathbf{Z}, \mathbf{Z}^2), \text{NN}(\mathbf{Y}, \mathbf{E}; \text{param3}) \rangle$  is equivalent to the factorized exponential family, which also satisfies that each  $\mathbf{T}_i(Z_i)$  has dimension larger or equal to 2.

Secondly,  $\text{NN}(\mathbf{Z}; \text{param1})$  is a neural network that receives as input a vector of latent variables and outputs another vector representing complicated nonlinear transformations of those variables.  $\text{NN}(\mathbf{Y}, \mathbf{E}; \text{param2})$  is another neural network that generates a corresponding vector of natural parameters. Hence, the term  $\langle \text{NN}(\mathbf{Z}; \text{param1}), \text{NN}(\mathbf{Y}, \mathbf{E}; \text{param2}) \rangle$  will allow this prior to capture the dependencies between the latent variables  $\mathbf{Z}$ .

## N HYPERPARAMETERS AND ARCHITECTURES

In this section, we describe the hyperparameters and architectures of different models used in different experiments. Unless stated otherwise, we have  $\lambda_1 = 1$  and  $\lambda_2 = 1$ , both of which are selected on training/validation data.

### N.1 SYNTHETIC DATA

We used Adam optimizer for training with learning rate set to 1e-3 and batch size set to 128.

#### N.1.1 ERM

##### Linear ERM

- Input layer: Input batch (*batch size*, *input dimension*)
- Output layer: Fully connected layer, output size = 1

### **Nonlinear ERM**

- Input layer: Input batch (*batch size, input dimension*)
- Layer 1: Fully connected layer, output size = 6, activation = ReLU
- Output layer: Fully connected layer, output size = 1

#### N.1.2 IRM

##### **Linear Data Representation $\Phi$**

- Input layer: Input batch (*batch size, input dimension*)
- Output layer: Fully connected layer, output size = 1

##### **Nonlinear Data Representation $\Phi$**

- Input layer: Input batch (*batch size, input dimension*)
- Layer 1: Fully connected layer, output size = 6, activation = ReLU
- Output layer: Fully connected layer, output size = 1

#### N.1.3 F-IRM GAME

##### **Linear Classifier $w$**

- Input layer: Input batch (*batch size, input dimension*)
- Output layer: Fully connected layer, output size = 1

##### **Nonlinear Classifier $w$**

- Input layer: Input batch (*batch size, input dimension*)
- Layer 1: Fully connected layer, output size = 6, activation = ReLU
- Output layer: Fully connected layer, output size = 1

#### N.1.4 V-IRM GAME

##### **Linear Data Representation $\Phi$**

- Input layer: Input batch (*batch size, input dimension*)
- Output layer: Fully connected layer, output size = 2

##### **Nonlinear Data Representation $\Phi$**

- Input layer: Input batch (*batch size, input dimension*)
- Layer 1: Fully connected layer, output size = 6, activation = ReLU
- Output layer: Fully connected layer, output size = 2

##### **Linear Classifier $w$**

- Input layer: Input batch (*batch size, 2*)
- Output layer: Fully connected layer, output size = 1

##### **Nonlinear Classifier $w$**

- Input layer: Input batch (*batch size, 2*)
- Layer 1: Fully connected layer, output size = 6, activation = ReLU
- Output layer: Fully connected layer, output size = 1

## N.1.5 iCARL

**NF-iVAE  $\lambda_f$ -Linear Prior**

- Input layer: Input batch (*batch size, input dimension*)
- Output layer: Fully connected layer, output size = 4

**NF-iVAE  $\lambda_f$ -Nonlinear Prior**

- Input layer: Input batch (*batch size, input dimension*)
- Layer 1: Fully connected layer, output size = 6, activation = ReLU
- Output layer: Fully connected layer, output size = 4

**NF-iVAE  $T_{NN}$ -Nonlinear Prior**

- Input layer: Input batch (*batch size, input dimension*)
- Layer 1: Fully connected layer, output size = 6, activation = ReLU
- Output layer: Fully connected layer, output size = 1

**NF-iVAE  $\lambda_{NN}$ -Nonlinear Prior**

- Input layer: Input batch (*batch size, input dimension*)
- Layer 1: Fully connected layer, output size = 6, activation = ReLU
- Output layer: Fully connected layer, output size = 1

**NF-iVAE Linear Encoder**

- Input layer: Input batch (*batch size, input dimension*)
- Mean Output layer: Fully connected layer, output size = 2
- Log Variance Output layer: Fully connected layer, output size = 2

**NF-iVAE Nonlinear Encoder**

- Input layer: Input batch (*batch size, input dimension*)
- Layer 1: Fully connected layer, output size = 6, activation = ReLU
- Mean Output layer: Fully connected layer, output size = 2
- Log Variance Output layer: Fully connected layer, output size = 2

**NF-iVAE Linear Decoder**

- Input layer: Input batch (*batch size, 2*)
- Mean Output layer: Fully connected layer, output size = output dimension
- Variance Output layer:  $0.01 \times \mathbf{1}$ , where  $\mathbf{1}$  is a vector full of 1 with the length of output dimension

**NF-iVAE Nonlinear Decoder**

- Input layer: Input batch (*batch size, 2*)
- Layer 1: Fully connected layer, output size = 6, activation = ReLU
- Mean Output layer: Fully connected layer, output size = output dimension
- Variance Output layer:  $0.01 \times \mathbf{1}$ , where  $\mathbf{1}$  is a vector full of 1 with the length of output dimension

**Linear Classifier  $w$** 

- Input layer: Input batch (*batch size, 1*)
- Output layer: Fully connected layer, output size = 1

**Nonlinear Classifier  $w$** 

- Input layer: Input batch (*batch size, 1*)
- Layer 1: Fully connected layer, output size = 6, activation = ReLU
- Output layer: Fully connected layer, output size = 1

## N.2 CMNIST AND CFMNIST

Considering that the results of most baselines come from IRMG (Ahuja et al., 2020a), for a fair comparison, we follow the same setting of IRMG in terms of hyper-parameters and validation considerations. For example, the batch size is set to 256, and the learning rate is  $10^{-4}$ . We also did not use the test environment data for validation. Please find more details in Ahuja et al. (2020a).

**NF-iVAE  $T_{NN}$ -Prior**

- Input layer: Input batch (*batch size, input dimension*)
- Layer 1: Fully connected layer, output size = 50, activation = ReLU
- Output layer: Fully connected layer, output size = 45

**NF-iVAE  $\lambda_{NN}$ -Prior**

- Input layer: Input batch (*batch size, input dimension*)
- Layer 1: Fully connected layer, output size = 50, activation = ReLU
- Output layer: Fully connected layer, output size = 45

**NF-iVAE  $\lambda_f$ -Prior**

- Input layer: Input batch (*batch size, input dimension*)
- Layer 1: Fully connected layer, output size = 50, activation = ReLU
- Output layer: Fully connected layer, output size = 20

**NF-iVAE  $X$ -Encoder**

- Input layer: Input batch (*batch size, 2, 28, 28*)
- Layer 1: Convolutional layer, output channels = 32, kernel size = 3, stride = 2, padding = 1, activation = ReLU
- Layer 2: Convolutional layer, output channels = 32, kernel size = 3, stride = 2, padding = 1, activation = ReLU
- Layer 3: Convolutional layer, output channels = 32, kernel size = 3, stride = 2, padding = 1, activation = ReLU
- Output layer: Flatten

**NF-iVAE  $(Y, E)$ -Encoder**

- Input layer: Input batch (*batch size, input dimension*)
- Output layer: Fully connected layer, output size = 100, activation = ReLU

**NF-iVAE ( $X, Y, E$ )-Merger/Encoder**

- Input layer: Input batch (*batch size, input dimension*)
- Layer 1: Fully connected layer, output size = 100, activation = ReLU
- Mean Output layer: Fully connected layer, output size = 10
- Log Variance Output layer: Fully connected layer, output size = 10

**NF-iVAE Decoder**

- Input layer: Input batch (*batch size, 10*)
- Layer 1: Fully connected layer, output size =  $32 \times 4 \times 4$ , activation = ReLU
- Layer 2: Reshape to (*batch size, 32, 4, 4*)
- Layer 3: Deconvolutional layer, output channels = 32, kernel size = 3, stride = 2, padding = 1, outpadding = 0, activation = ReLU
- Layer 4: Deconvolutional layer, output channels = 32, kernel size = 3, stride = 2, padding = 1, outpadding = 1, activation = ReLU
- Layer 5: Deconvolutional layer, output channels = 2, kernel size = 3, stride = 2, padding = 1, outpadding = 1
- Mean Output layer: activation = Sigmoid
- Variance Output layer:  $0.01 \times \mathbf{1}$ , where  $\mathbf{1}$  is a matrix full of 1 with the size of  $2 \times 28 \times 28$ .

**Classifier  $w$** 

- Input layer: Input batch (*batch size, 50*)
- Layer 1: Fully connected layer, output size = 100, activation = ReLU
- Output layer: Fully connected layer, output size = 1, activation = Sigmoid

**N.3 VLCS AND PACS**

We used the exact experimental setting that is described in [Gulrajani & Lopez-Paz \(2020\)](#). Specifically, we trained our model over all possible train and test environment combination for one of the commonly used hyper-parameter tuning procedure: train domain validation. We use ResNet-50 as an encoder and reverse the architecture of ResNet-50 as a decoder. We set the number of the latent variables to  $n = 50$ . We do the hyperparameter search by exactly following the guides given in [Gulrajani & Lopez-Paz \(2020\)](#).