

# CausalVAE: Disentangled Representation Learning via Neural Structural Causal Models

Mengyue Yang<sup>1,2</sup>, Furui Liu<sup>1,\*</sup>, Zhitang Chen<sup>1</sup>, Xinwei Shen<sup>3</sup>, Jianye Hao<sup>1</sup>, Jun Wang<sup>2</sup>

<sup>1</sup> Noah’s Ark Lab, Huawei, Shenzhen, China

<sup>2</sup> University College London, London, United Kingdom

<sup>3</sup> The Hong Kong University of Science and Technology, Hong Kong, China

{yangmengyue2, liufurui2, chenzhitang2, haojianye}@huawei.com

xshenal@connect.ust.hk

jun.wang@cs.ucl.ac.uk

## Abstract

*Learning disentanglement aims at finding a low dimensional representation which consists of multiple explanatory and generative factors of the observational data. The framework of variational autoencoder (VAE) is commonly used to disentangle independent factors from observations. However, in real scenarios, factors with semantics are not necessarily independent. Instead, there might be an underlying causal structure which renders these factors dependent. We thus propose a new VAE based framework named CausalVAE, which includes a Causal Layer to transform independent exogenous factors into causal endogenous ones that correspond to causally related concepts in data. We further analyze the model identifiability, showing that the proposed model learned from observations recovers the true one up to a certain degree. Experiments are conducted on various datasets, including synthetic and real word benchmark CelebA. Results show that the causal representations learned by CausalVAE are semantically interpretable, and their causal relationship as a Directed Acyclic Graph (DAG) is identified with good accuracy. Furthermore, we demonstrate that the proposed CausalVAE model is able to generate counterfactual data through “do-operation” to the causal factors.*

## 1. Introduction

Disentangled representation learning is of great importance in various applications such as computer vision, speech and natural language processing, and recommender systems [9, 21, 8]. The reason is that it might help enhance the performance of models, i.e. improving the generalizability,

robustness against adversarial attacks as well as the explainability, by learning data’s latent disentangled representation. One of the most common frameworks for disentangled representation learning is Variational Autoencoders (VAE), a deep generative model trained to disentangle the underlying explanatory factors. Disentanglement via VAE can be achieved by a regularization term of the Kullback-Leibler (KL) divergence between the posterior of the latent factors and a standard Multivariate Gaussian prior, which enforces the learned latent factors to be as independent as possible. It is expected to recover the latent variables if the observation in real world is generated by countable independent factors. To further enhance the independence, various extensions of VAE consider minimizing the mutual information among latent factors. For example, Higgins *et al.* [6] and Burgess *et al.* [3] increased the weight of the KL divergence term to enforce independence. Kim *et al.* [13, 4] further encourage the independence by reducing total correlation among factors.

Most existing works of disentangled representation learning make a common assumption that the real world observations are generated by countable independent factors. Nevertheless we argue that in many real world applications, latent factors with semantics of interest are causally related and thus we need a new framework that supports causal disentanglement.

Consider a toy example of a swinging pendulum in Fig. 1. The position of the illumination source and the angle of the pendulum are causes of the position and the length of the shadow. Through causal disentangled representation learning, we aim at learning representations that correspond to the above four concepts. Obviously, these concepts are not independent and existing methods may fail to extract those factors. Furthermore, causal disentanglement allow us to manipulate the causal system to generate counterfactual data. For example, we can manipulate the latent code of shadow to

\*Corresponding author.

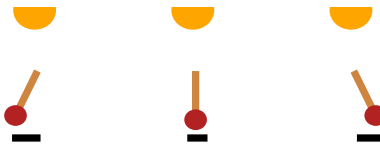


Figure 1. A swinging pendulum: an illustrative example

create new pictures without shadow even there are pendulum and light. This corresponds to the “do-operation” [25] in causality, where the system operates under the condition that certain variables are controlled by external forces. A deep generative model that supports “do-operation” is of tremendous value as it allows us to ask “what-if” questions when making decisions.

In this paper, we propose a VAE-based causal disentangled representation learning framework by introducing a novel Structural Causal Model layer (Mask Layer), which allows us to recover the latent factors with semantics and structure via a causal DAG. The input signal passes through an encoder to obtain independent exogenous factors and then a Causal Layer to generate causal representation which is taken by the decoder to reconstruct the original input. We call the whole process Causal Disentangled Representation Learning. Unlike unsupervised disentangled representation learning of which the feasibility is questionable [19], additional information is required as weak supervision signals to achieve causal representation learning. By “weak supervision”, we emphasize that in our work, the causal structure of the latent factors is automatically learned, instead of being given as a prior in [15]. To train our model, we propose a new loss function which includes the VAE evidence lower bound loss and an acyclicity constraint imposed on the learned causal graph to guarantee its “DAGness”. In addition, we analyze the identifiability of the proposed model, showing that the learned parameters of the disentangled model recover the true one up to certain degree. The contribution of our paper is three-fold. (1) We propose a new framework named CausalVAE that supports causal disentanglement and “do-operation”; (2) Theoretical justification on model identifiability is provided; (3) We conduct comprehensive experiments with synthetic and real world face images to demonstrate that the learned factors are with causal semantics and can be intervened to generate counterfactual images that do not appear in training data.

## 2. Related Works

In this section, we review state-of-the-art disentangled representation learning methods, including some recent advances on combining causality and disentangled representation learning. We also present preliminaries of causal structure learning from pure observations which is a key ingredient of our proposed CausalVAE framework.

### 2.1. Disentangled Representation Learning

Conventional disentangled representation learning methods learn mutually independent latent factors by an encoder-decoder framework. In this process, a standard normal distribution is used as a prior of the latent code. A variational posterior  $q(\mathbf{z}|\mathbf{x})$  is then used to approximate the unknown true posterior  $p(\mathbf{z}|\mathbf{x})$ . This framework was further extended by adding new independence regularization terms to the original loss function, leading to various algorithms.  $\beta$ -VAE [6] proposes an adaptation framework which adjusts the weight of KL term to balance between independence of disentangled factors and the reconstruction performance. While factor VAE [4] proposes a new framework which focuses solely on the independence of factors. Ladder VAE [17] on the other hand, leverages the structure of ladder neural network to train a structured VAE for hierarchical disentanglement. Nevertheless the aforementioned unsupervised disentangled representation learning algorithms do not perform well in some situations where there is complex causal relationship among factors. Furthermore, they are challenged for lacking inductive bias and thus the model identifiability cannot be guaranteed [19]. The identifiability problem of VAE is defined as follows: if the parameters  $\theta$  learned from data lead to a marginal distribution equal to the true one parameterized by  $\theta$ , i.e.,  $p_{\theta}(\mathbf{x}) = p_{\theta^*}(\mathbf{x})$ , then the joint distributions also match, i.e.  $p_{\theta}(\mathbf{x}, \mathbf{z}) = p_{\theta^*}(\mathbf{x}, \mathbf{z})$ . Therefore, the rotation invariance of prior  $p(\mathbf{z})$  (standard Multivariate Gaussian distribution) will lead the unidentifiability of  $p(\mathbf{z})$ . Khemakhem *et al.* [11] prove that there is infinite number of distinct models entailing the same joint distributions, which means that the underlying generative model is not identifiable through unsupervised learning. On the contrary, by leveraging a few labels, one is able to recover the true model [22, 19]. Kulkarni *et al.* [16] and Locatello *et al.* [20] use additional labels to reduce the model ambiguity. Khemakhem *et al.* [11] gives an identifiability of VAE with additional inputs, by leveraging the theory of nonlinear Independent Component Analysis (nonlinear ICA) [2].

### 2.2. Causal Discovery & Causal Disentangled Representation Learning

We refer to causal representation as ones structured by a causal graph. Discovering the causal graph from pure observations has attracted large amounts of attention in the past decades [7, 35, 29]. Methods for causal discovery use

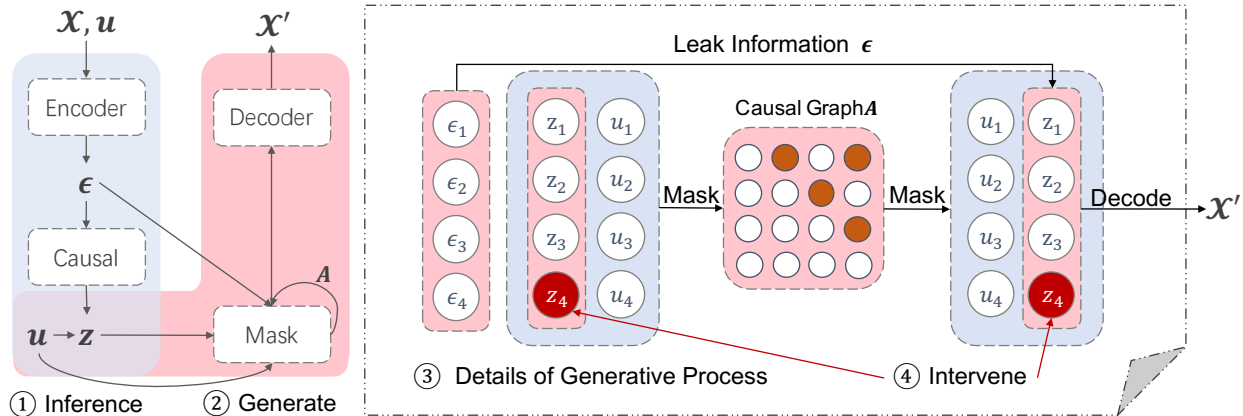


Figure 2. Model structure of CausalVAE. The encoder takes observation  $x$  as inputs to generate independent exogenous variable  $\epsilon$ , whose prior distribution is assumed to be standard Multivariate Gaussian. Then it is transformed by the Causal Layer into causal representations  $z$  (Eq. 1) with a conditional prior distribution  $p(z|u)$ . A Mask Layer is then applied to  $z$  to resemble the SCM in Eq. 2. After that,  $z$  is taken as the input of the decoder to reconstruct the observation  $x$ .

either observational data or a combination of observational and interventional data. We first introduce a set of methods based on observational data. Pearl *et al.* [25] introduced a Probabilistic Graphical Models (PGMs) based language to describe causality among variables. Shimizu *et al.* [29] proposed an effective method called LiNGAM to learn the causal graph and they prove the model identifiability under the linearity and non-Gaussianity assumption. Zheng *et al.* [36] proposed NOTEARS with a fully differentiable DAG constraint for causal structure learning, which drastically reduces a very complicated combinatorial optimization problem to a continuous optimization problem. Zhu *et al.* [38] proposed a flexible and efficient Reinforcement Learning (RL) based method to search over a DAG space for a best graph with a highest score. When interventions are doable, that is, one can manipulate the causal system and collect data under interventions, methods are proposed for causal discovery. Tillman *et al.* [33, 5] show the identifiability of learned causal structure from interventional data. Peters *et al.* [10, 26, 27] explores the structure invariance across multiple domains under interventions to identify causal edges.

Recently, the community has raised interest of combining causality and disentangled representation. Suter *et al.* [32] used causality to explain disentangled latent representations. Kocaoglu *et al.* [15] proposed a method called CausalGAN which supports “do-operation” on images but it requires the causal graph given as a prior. Instead of assuming independent latent factors, Besserve *et al.* [1] adopts dependent latent factors in the model. It relies on the principle of “independence mechanism” or modularity for disentanglement, and design a layer containing a few non-structured nodes, representing outputs of mutually independent causal mecha-

nisms [27], which contribute together to the final predictions to achieve disentanglement. In our model, we disentangle factors by causally structured layers (masking layer), and the model structure is different from theirs. Schölkopf *et al.* [28] claims the importance and necessity of causal disentangled representation learning but it still remains conceptual. To the best of our knowledge, our work is the first one that successfully implements the idea of causal disentanglement.

### 3. Causal Disentanglement in Variational Autoencoder

We start with the definition of causal representation, and then propose a new framework to achieve causal disentanglement by leveraging additional inputs, e.g. labels of concepts. Firstly, we give an overview of our proposed CausalVAE model structure in Fig. 2. A Causal Layer, which essentially describes a Structural Causal Model (SCM) [29], is introduced to a conventional VAE network. The Causal Layer transforms the independent exogenous factors to causal endogenous factors corresponding to causally related concepts of interest. A mask mechanism [23] is then used to propagate the effect of parental variables to their children, mimicking the assignment operation of SCMs. Such a Causal Layer is the key to supporting intervention or “do-operation” to the system.

#### 3.1. Transforming Independent Exogenous Factors into Causal Representations

Our model is within the framework of VAE-based disentanglement. In addition to the encoder and the decoder structures, we introduce a Structural Causal Model (SCM) layer to learn causal representations. To formalize causal

representation, we consider  $n$  concepts of interest in data. The concepts in observations are causally structured by a Directed Acyclic Graph (DAG) with an adjacency matrix  $\mathbf{A}$ . Though a general nonlinear SCM is preferred, for simplicity, in this work, the Causal Layer exactly implements a Linear SCM as described in Eq. 1 (shown in Fig. 2 ①),

$$\mathbf{z} = \mathbf{A}^T \mathbf{z} + \boldsymbol{\epsilon} = (\mathbf{I} - \mathbf{A}^T)^{-1} \boldsymbol{\epsilon}, \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad (1)$$

where  $\mathbf{A}$  is the parameters to be learnt in this layer.  $\boldsymbol{\epsilon}$  are independent Gaussian exogenous factors and  $\mathbf{z} \in \mathbb{R}^n$  is structured causal representation of  $n$  concepts that is generated by a DAG and thus  $\mathbf{A}$  can be permuted into a strictly upper triangular matrix.

Unsupervised learning of the model might be infeasible due to the identifiability issue as discussed in [19]. To address this problem, similar to iVAE [11], we adopt additional information  $\mathbf{u}$  associated with the true causal concepts as supervising signals. In our work, we use the labels of the concepts. The additional information  $\mathbf{u}$  is utilized in two ways. Firstly, we propose a conditional prior  $p(\mathbf{z}|\mathbf{u})$  to regularize the learned posterior of  $\mathbf{z}$ . This guarantees that the learned model belongs to an identifiable family. Secondly, we also leverage  $\mathbf{u}$  to learn the causal structure  $\mathbf{A}$ . Besides learning the causal representations, we further enable the model to support intervention to the causal system to generate counterfactual data which does not exist in the training data.

### 3.2. Structural Causal Model Layer

Once the causal representation  $\mathbf{z}$  is obtained, it passes through a Mask Layer [23] to reconstruct itself. Note that this step resembles a SCM which depicts how children are generated by their corresponding parental variables. We will show why such a layer is necessary to achieve intervention. Let  $z_i$  be the  $i$ th variable in the vector  $\mathbf{z}$ . The adjacency matrix associated with the causal graph is  $\mathbf{A} = [\mathbf{A}_1 | \dots | \mathbf{A}_n]$  where  $\mathbf{A}_i \in \mathbb{R}^n$  is the weight vector such that  $A_{ji}$  encodes the causal strength from  $z_j$  to  $z_i$ . We have a set of mild nonlinear and invertible functions  $[g_1, g_2, \dots, g_n]$  that map parental variables to the child variable. Then we write

$$z_i = g_i(\mathbf{A}_i \circ \mathbf{z}; \boldsymbol{\eta}_i) + \epsilon_i, \quad (2)$$

where  $\circ$  is the element-wise multiplication and  $\boldsymbol{\eta}_i$  is the parameter of  $g_i(\cdot)$  (as shown in Fig. 2 ③). Note that according to Eq. 1, we can simply write  $z_i = \mathbf{A}_i^T \mathbf{z} + \epsilon_i$ . However, we find that adding a mild nonlinear function  $g_i$  results in more stable performances. To show how this masking works, consider a variable  $z_i$  and  $\mathbf{A}_i \circ \mathbf{z}$  equals a vector that only contains its parental information as it masks out all  $z_i$ 's non-parent variables. By minimizing the reconstruction error, the adjacency matrix  $\mathbf{A}$  and the parameter  $\boldsymbol{\eta}_i$  of the mild nonlinear function  $g_i$  are trained.

This layer makes intervention or "do-operation" possible. Intervention [25] in causality refers to modifying a certain part of a system by external forces and one is interested in the outcome of such manipulation. To intervene  $z_i$ , we set  $z_i$  on the RHS of Eq. 2 (corresponding to the  $i$ -th node of  $\mathbf{z}$  in the first layer in Fig. 2) to a fixed value, and then its effect is delivered to all its children as well as itself on the LHS of Eq. 2 (corresponding to some nodes of  $\mathbf{z}$  in the second layer). Note that intervening the cause will change the effect, whereas intervening the effect, on the other hand, does not change the cause because information can only flow into the next layer from the previous one in our model, which is aligned with the definition of causal effects.

### 3.3. A Probabilistic Generative Model for Causal-VAE

We give a probabilistic formulation of the proposed generative model (shown in Fig. 2 ②). Denote by  $\mathbf{x} \in \mathbb{R}^d$  the observed variables and  $\mathbf{u} \in \mathbb{R}^n$  the additional information.  $u_i$  is the label of the  $i$ -th concept of interest in data. Let  $\boldsymbol{\epsilon} \in \mathbb{R}^n$  be the latent exogenous independent variables and  $\mathbf{z} \in \mathbb{R}^n$  be the latent endogenous variables with semantics where  $\mathbf{z} = \mathbf{A}^T \mathbf{z} + \boldsymbol{\epsilon} = (\mathbf{I} - \mathbf{A}^T)^{-1} \boldsymbol{\epsilon}$ . For simplicity, we denote  $\mathbf{C} = (\mathbf{I} - \mathbf{A}^T)^{-1}$ .

We treat both  $\mathbf{z}$  and  $\boldsymbol{\epsilon}$  as latent variables. Consider the following conditional generative model parameterized by  $\boldsymbol{\theta} = (\mathbf{f}, \mathbf{h}, \mathbf{C}, \mathbf{T}, \boldsymbol{\lambda})$ :

$$p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z}, \boldsymbol{\epsilon} | \mathbf{u}) = p_{\boldsymbol{\theta}}(\mathbf{x} | \mathbf{z}, \boldsymbol{\epsilon}, \mathbf{u}) p_{\boldsymbol{\theta}}(\boldsymbol{\epsilon}, \mathbf{z} | \mathbf{u}). \quad (3)$$

Let  $\mathbf{f}(\mathbf{z})$  denote the decoder which is assumed to be an invertible function and  $\mathbf{h}(\mathbf{x}, \mathbf{u})$  denotes the encoder. We define the generative and inference models as follows:

$$\begin{aligned} p_{\boldsymbol{\theta}}(\mathbf{x} | \mathbf{z}, \boldsymbol{\epsilon}, \mathbf{u}) &= p_{\boldsymbol{\theta}}(\mathbf{x} | \mathbf{z}) \equiv p_{\boldsymbol{\xi}}(\mathbf{x} - \mathbf{f}(\mathbf{z})), \\ q_{\boldsymbol{\phi}}(\mathbf{z}, \boldsymbol{\epsilon} | \mathbf{x}, \mathbf{u}) &\equiv q(\mathbf{z} | \boldsymbol{\epsilon}) q_{\boldsymbol{\zeta}}(\boldsymbol{\epsilon} - \mathbf{h}(\mathbf{x}, \mathbf{u})), \end{aligned} \quad (4)$$

which is obtained by assuming the following decoding and encoding processes:

$$\mathbf{x} = \mathbf{f}(\mathbf{z}) + \boldsymbol{\xi}, \quad \boldsymbol{\epsilon} = \mathbf{h}(\mathbf{x}, \mathbf{u}) + \boldsymbol{\zeta}, \quad (5)$$

where  $\boldsymbol{\xi}$  and  $\boldsymbol{\zeta}$  are the vectors of independent noise with probability densities  $p_{\boldsymbol{\xi}}$  and  $q_{\boldsymbol{\zeta}}$ . When  $\boldsymbol{\xi}$  and  $\boldsymbol{\zeta}$  are infinitesimal, the encoder and decoder can be regarded as deterministic ones. We define the joint prior  $p_{\boldsymbol{\theta}}(\boldsymbol{\epsilon}, \mathbf{z} | \mathbf{u})$  for latent variables  $\mathbf{z}$  and  $\boldsymbol{\epsilon}$  as

$$p_{\boldsymbol{\theta}}(\boldsymbol{\epsilon}, \mathbf{z} | \mathbf{u}) = p_{\boldsymbol{\epsilon}}(\boldsymbol{\epsilon}) p_{\boldsymbol{\theta}}(\mathbf{z} | \mathbf{u}), \quad (6)$$

where  $p_{\boldsymbol{\epsilon}}(\boldsymbol{\epsilon}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$  and the prior of latent endogenous variables  $p_{\boldsymbol{\theta}}(\mathbf{z} | \mathbf{u})$  is a factorized Gaussian distribution conditioning on the additional observation  $\mathbf{u}$ , i.e.

$$p_{\boldsymbol{\theta}}(\mathbf{z} | \mathbf{u}) = \prod_i^n p_{\boldsymbol{\theta}}(z_i | u_i), p_{\boldsymbol{\theta}}(z_i | u_i) = \mathcal{N}(\lambda_1(u_i), \lambda_2^2(u_i)), \quad (7)$$

where  $\lambda_1$  and  $\lambda_2$  are an arbitrary functions. In this paper, we let  $\lambda_1(\mathbf{u}) = \mathbf{u}$  and  $\lambda_2(\mathbf{u}) \equiv 1$ . The distribution has two sufficient statistics, the mean and variance of  $\mathbf{z}$ , which are denoted by sufficient statistics  $\mathbf{T}(\mathbf{z}) = (\boldsymbol{\mu}(\mathbf{z}), \boldsymbol{\sigma}(\mathbf{z})) = (T_{1,1}(z_1), \dots, T_{n,2}(z_n))$ . We use these notations for model identifiability analysis in Section 5.

## 4. Learning Strategy

In this section, we discuss how to train the CausalVAE model in order to learn the causal representation as well as the causal graph simultaneously.

### 4.1. Evidence Lower Bound of CausalVAE

We apply variational Bayes to learn a tractable distribution  $q_\phi(\boldsymbol{\epsilon}, \mathbf{z}|\mathbf{x}, \mathbf{u})$  to approximate the true posterior  $p_\theta(\boldsymbol{\epsilon}, \mathbf{z}|\mathbf{x}, \mathbf{u})$ . Given data set  $\mathcal{X}$  with the empirical data distribution  $q_{\mathcal{X}}(\mathbf{x}, \mathbf{u})$ , the parameters  $\theta$  and  $\phi$  are learned by optimizing the following evidence lower bound (ELBO):

$$\mathbb{E}_{q_{\mathcal{X}}}[\log p_\theta(\mathbf{x}|\mathbf{u})] \geq \text{ELBO} = \mathbb{E}_{q_{\mathcal{X}}}[\mathbb{E}_{\boldsymbol{\epsilon}, \mathbf{z} \sim q_\phi}[\log p_\theta(\mathbf{x}|\mathbf{z}, \boldsymbol{\epsilon}, \mathbf{u})] - \mathcal{D}(q_\phi(\boldsymbol{\epsilon}, \mathbf{z}|\mathbf{x}, \mathbf{u})||p_\theta(\boldsymbol{\epsilon}, \mathbf{z}|\mathbf{u}))], \quad (8)$$

where  $\mathcal{D}(\cdot||\cdot)$  denotes KL divergence. Eq. 8 is intractable in general. However, thanks to the one-to-one correspondence between  $\boldsymbol{\epsilon}$  and  $\mathbf{z}$ , we simplify the variational posterior as follows:

$$\begin{aligned} q_\phi(\boldsymbol{\epsilon}, \mathbf{z}|\mathbf{x}, \mathbf{u}) &= q_\phi(\boldsymbol{\epsilon}|\mathbf{x}, \mathbf{u})\delta(\mathbf{z} = \mathbf{C}\boldsymbol{\epsilon}) \\ &= q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{u})\delta(\boldsymbol{\epsilon} = \mathbf{C}^{-1}\mathbf{z}), \end{aligned} \quad (9)$$

where  $\delta(\cdot)$  is the Dirac delta function. According to the model assumptions introduced in Section 3.3, i.e., generation process (Eq. 4) and prior (Eq. 6), we attain a neat form of ELBO loss as follows:

**Proposition 1** *ELBO defined in Eq. 8 can be written as:*

$$\begin{aligned} \text{ELBO} &= \mathbb{E}_{q_{\mathcal{X}}}[\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{u})}[\log p_\theta(\mathbf{x}|\mathbf{z})] \\ &\quad - \mathcal{D}(q_\phi(\boldsymbol{\epsilon}|\mathbf{x}, \mathbf{u})||p_\epsilon(\boldsymbol{\epsilon})) \\ &\quad - \mathcal{D}(q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{u})||p_\theta(\mathbf{z}|\mathbf{u}))]. \end{aligned} \quad (10)$$

Details of the proof are given in the Appendix A. With this form, we can easily implement a loss function to train the CausalVAE model.

### 4.2. Learning the Causal Structure of Latent Codes

In addition to the encoder and decoder, our CausalVAE model involves a Causal Layer with a DAG structure to be learned. Note that both  $\mathbf{z}$  and  $\mathbf{A}$  are unknown, to ease the training task and guarantee the identifiability of causal graph  $\mathbf{A}$ , we leverage the additional labels  $\mathbf{u}$  to construct the following constraint:

$$l_u = \mathbb{E}_{q_{\mathcal{X}}}\|\mathbf{u} - \sigma(\mathbf{A}^T\mathbf{u})\|_2^2 \leq \kappa_1, \quad (11)$$

where  $\sigma$  is a logistic function as our labels are binary and  $\kappa_1$  is the small positive constant value. This follows the idea that  $\mathbf{A}$  should also describe the causal relations among labels well. Similarly we apply the same constraint to the learned latent code  $\mathbf{z}$  as follows:

$$l_m = \mathbb{E}_{\mathbf{z} \sim q_\phi} \sum_{i=1}^n \|z_i - g_i(\mathbf{A}_i \circ \mathbf{z}; \boldsymbol{\eta}_i)\|^2 \leq \kappa_2, \quad (12)$$

where  $\kappa_2$  is the small positive constant value. Lastly, the causal adjacency matrix  $\mathbf{A}$  is constrained to be a DAG. Instead of using traditional DAG constraint that is combinatorial, we adopt a continuous differentiable constraint function [36, 37, 24, 34]. The function attains 0 if and only if the adjacency matrix  $\mathbf{A}$  corresponds to a DAG [34], i.e.

$$H(\mathbf{A}) \equiv \text{tr}((\mathbf{I} + \frac{c}{m}\mathbf{A} \circ \mathbf{A})^n) - n = 0, \quad (13)$$

where  $c$  is an arbitrary positive number. The training procedure of our CausalVAE model reduces to the following constrained optimization:

$$\begin{aligned} &\text{maximize ELBO,} \\ &\text{s.t. (11)(12)(13).} \end{aligned}$$

By lagrangian multiplier method, we have the new loss function

$$\mathcal{L} = -\text{ELBO} + \alpha H(\mathbf{A}) + \beta l_u + \gamma l_m, \quad (14)$$

where  $\alpha, \beta, \gamma$  denote regularization hyperparameters.

## 5. Identifiability Analysis

In this section, we present the identifiability of our proposed model. We adopt the  $\sim$ -identifiable [11] as follows:

**Definition 1** *Let  $\sim$  be the binary relation on  $\Theta$  defined as follows:*

$$\begin{aligned} (\mathbf{f}, \mathbf{h}, \mathbf{C}, \mathbf{T}, \boldsymbol{\lambda}) &\sim (\tilde{\mathbf{f}}, \tilde{\mathbf{h}}, \tilde{\mathbf{C}}, \tilde{\mathbf{T}}, \tilde{\boldsymbol{\lambda}}) \\ &\Leftrightarrow \exists \mathbf{B}_1, \mathbf{B}_2, \mathbf{b}_1, \mathbf{b}_2 \\ \mathbf{T}(\mathbf{h}(\mathbf{x}, \mathbf{u})) &= \mathbf{B}_1 \tilde{\mathbf{T}}(\tilde{\mathbf{h}}(\mathbf{x}, \mathbf{u})) + \mathbf{b}_1, \mathbf{T}(\mathbf{f}^{-1}(\mathbf{x})) \\ &= \mathbf{B}_2 \tilde{\mathbf{T}}(\tilde{\mathbf{f}}^{-1}(\mathbf{x})) + \mathbf{b}_2, \forall \mathbf{x} \in \mathcal{X}, \end{aligned} \quad (15)$$

where  $\mathbf{C} = (\mathbf{I} - \mathbf{A}^T)^{-1}$ . If  $\mathbf{B}_1$  is an invertible matrix and  $\mathbf{B}_2$  is an invertible diagonal matrix with diagonal elements associated to  $u_i$ . We say that the model parameter is  $\sim$ -identifiable.

Following [11], we obtain the identifiability of our causal generative model as follows.

**Theorem 1** *Assume that the data we observed are generated according Eq. 3-4 and the following assumptions hold,*

1. The set  $\{x \in \mathcal{X} | \phi_\xi(\mathbf{x}) = 0\}$  has measure zero, where  $\phi_\xi$  is the characteristic function of the density  $p_\xi$  defined in Eq. 5.
2. The decoder function  $\mathbf{f}$  is differentiable and the Jacobian matrix of  $\mathbf{f}$  is of full rank <sup>1</sup>.
3. The sufficient statistics  $T_{i,s}(z_i) \neq 0$  almost everywhere for all  $1 \leq i \leq n$  and  $1 \leq s \leq 2$ , where  $T_{i,s}(z_i)$  is the  $s$ th statistic of variable  $z_i$ .
4. The additional observations  $u_i \neq 0$ .

Then the parameters  $(\mathbf{f}, \mathbf{h}, \mathbf{C}, \mathbf{T}, \boldsymbol{\lambda})$  are  $\sim$ -identifiable.

Although the parameters  $\theta$  of true generative model are unknown during the learning process, the identifiability of generative model given by Theorem 1 guarantees the parameters  $\tilde{\theta}$  learned by hypothetical functions are in an identifiable family. This shows that the learned parameters of the generative model recover the true one up to certain degree.

In addition, all  $z_i$  in  $\mathbf{z}$  align to the additional observation of concept  $i$  and they are expected to inherit the causal relationship of causal system. That is why that it could guarantee that the  $\mathbf{z}$  are causal representation.

The identifiability of the model under supervision of additional information is obtained thanks to the conditional prior  $p_\theta(\mathbf{z}|\mathbf{u})$ . The conditional prior guarantees that sufficient statistics of  $p_\theta(\mathbf{z}|\mathbf{u})$  are related to the value of  $\mathbf{u}$ . A complete proof of **Theorem 1** is available in Appendix B.

## 6. Experiments

In this section, we conduct experiments using both synthetic dataset and real human face image dataset and we compare our CausalVAE model against existing state of the art methods on disentangled representation learning. We focus on examining whether a certain algorithm is able to learn interpretable representations and whether outcomes of intervention on learned latent code is consistent to our understanding of the causal system.

### 6.1. Dataset, Baselines & Metrics

#### 6.1.1 Datasets:

We conduct experiments on a synthetic datasets and a benchmark face dataset CelebA.

**Synthetic:** We build two synthetic datasets which include images of causally related objects. The first one is named Pendulum. Each image contains 3 entities (PENDULUM, LIGHT, SHADOW), and 4 concepts ((PENDULUM ANGLE, LIGHT ANGLE)  $\rightarrow$  (SHADOW LOCATION, SHADOW LENGTH)). The second one is named Flow. Each image contains 4 concepts (BALL SIZE  $\rightarrow$  WATER SIZE, (WATER

<sup>1</sup>(rank equals to its smaller dimension)

SIZE, HOLE) $\rightarrow$  WATER FLOW). Due to page limitation, main text only shows the results on Pendulum, and experiments on Flow and more details of two datasets are given in Appendix C.1.

**Real world benchmark:** We also use a real world dataset CelebA<sup>2</sup>[18], a widely used dataset in the computer vision community. In this dataset, there are in total 200k human face images with labels on different concepts, and we choose two subsets of causally related attributes. The first set is CelebA(SMILE), which consists of GENDER, SMILE, EYES OPEN, MOUTH OPEN. The second one is CelebA(BEARD), which consists of AGE, GENDER, BALD, BEARD. Main text only shows results on CelebA(SMILE), and more experimental results on other concepts are provided in the Appendix D.

**Baselines:** We compare our method with some state of the arts and show the results of ablation study. Baselines are categorized into supervised and unsupervised methods.

CausalVAE-unsup, LadderVAE [17] and  $\beta$ -VAE [6] are unsupervised methods. CausalVAE-unsup is a reduced version of our model whose structure is the same as CausalVAE except that the Mask Layer and the supervision conditional prior  $p(\mathbf{z}|\mathbf{u})$  are removed.

Supervised methods include disentangled representation learning method ConditionVAE [30], which does not include causal layers in the model structure and causal generative model CausalGAN [15], which needs the true causal graph to be given as a prior.

As CausalGAN does not focus on representation learning, we only compare our CausalVAE with CausalGAN on intervention experiment (results given in Appendix D.3). For these methods, the prior conditioning on the labels are given, and the dimensionality of the latent representation is the same as CausalVAE.

**Metrics:** We use Maximal Information Coefficient (MIC) and Total Information Coefficient (TIC) [14] as our evaluation metrics. Both of them indicate the degree of information relevance between the learned representation and the ground truth labels of concepts.

### 6.2. Intervention experiments

Intervention experiments aim at testing if a certain dimension of the latent representation has interpretable semantics. The value of a latent code is manipulated by "do-operation" as introduced in previous sections, and we observe how the generated image appears. Intervention is conducted by the following steps:

- A generative model is trained.
- An arbitrary image from the training set is fed to the encoder to generate a latent code  $\mathbf{z}$ .

<sup>2</sup><http://mmlab.ie.cuhk.edu.hk/projects/CelebA.html>

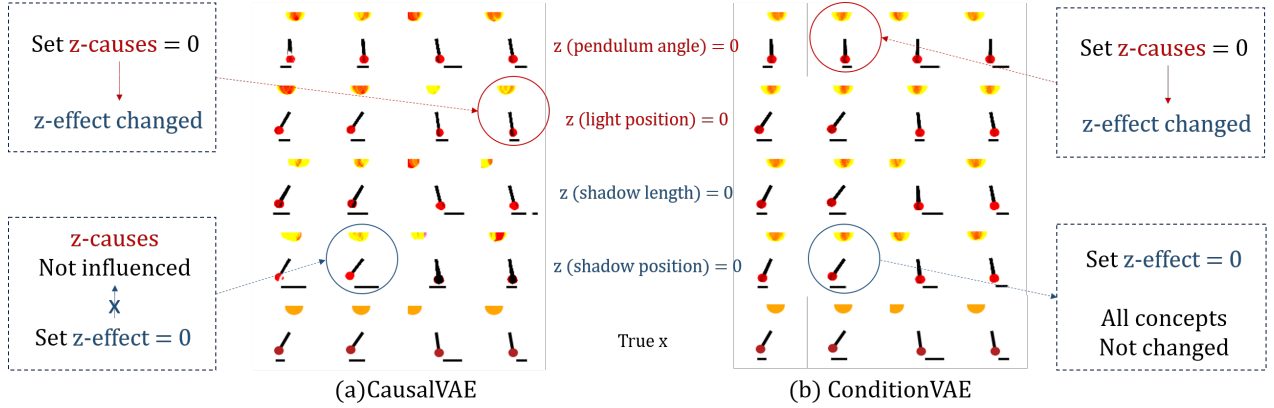


Figure 3. The results of Intervention experiments on the pendulum dataset. Each row shows the result of controlling the PENDULUM ANGLE, LIGHT ANGLE, SHADOW LENGTH, and SHADOW LOCATION respectively. The bottom row is the original input image. More intervention results on other synthetic dataset are shown in Appendix D.3.

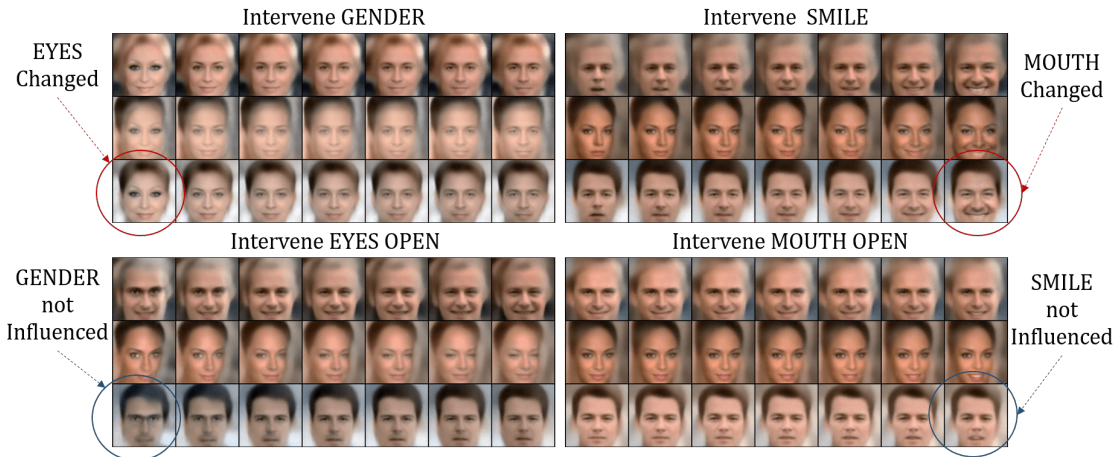


Figure 4. Results of CausalVAE model on CelebA(SMILE). The controlled factors are GENDER, SMILE, EYES OPEN and MOUTH OPEN respectively. More intervention results are shown in Appendix D.3.

- We manipulate the value of  $z_i$  corresponding to a concept of interest. For CausalVAE, as Fig. 2 ④ and Fig. 6 show, we need to manipulate both the input and output nodes of the SCM layer. Note that the effect of manipulation to a parental node will be propagated to its children.
- The intervened latent code  $\tilde{z}$  passes through the decoder to generate a new image. In the experiments, all images in the dataset are used to train our proposed model CausalVAE and other baselines.

Hyperparameters  $(\alpha, \beta, \gamma) = (1, 1, 1)$  for all experiments unless specified.

We first conduct intervention experiments on the Pendulum dataset, with 4 latent concepts and results are given in Fig. 3. We intervene a certain concept by setting the

corresponding latent code value to 0. We expect that the pattern of the manipulated concept will be fixed across all images under the same intervention. For example, when we intervene the pendulum ANGLE as shown in the first line of Fig. 3 (a), the ANGLE of pendulum of different images are almost the same. Meanwhile, we also observe that the SHADOW LOCATION and SHADOW LENGTH change in a correct way that aligns with the physics law. Note that this is also related to the concept of modularity, meaning that intervening a certain part of the generative system usually does not affect the other parts of the system. Similar phenomenon is observed in other intervention experiments, demonstrating that our model correctly implements the underlying causal system. The results of ConditionVAE, a supervised method without considering the causal structure, are given in Fig. 3 (b). There exists a problem that manipulating the latent

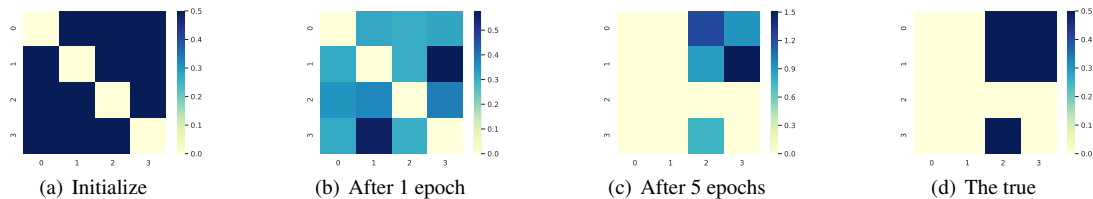


Figure 5. The learning process of causal matrix  $\mathbf{A}$ . The concepts include: GENDER, SMILE, EYES OPEN, MOUTH OPEN (top-to-bottom and left-to-right order); (c) converged  $\mathbf{A}$ , (d) ground truth .

Table 1. The MIC and TIC between learned representation  $\mathbf{z}$  and the label  $\mathbf{u}$ . The results show that among all compared methods, the learned factors of our proposed CausalVAE achieve best alignment to the concepts of interest. (Note: the metrics include mean  $\pm$  standard errors in table.)

Metrics(%)	CausalVAE		ConditionVAE		$\beta$ -VAE		CausalVAE-unsup		LadderVAE	
	MIC	TIC	MIC	TIC	MIC	TIC	MIC	TIC	MIC	TIC
Pendulum	<b>95.1 <math>\pm</math> 2.4</b>	<b>81.6 <math>\pm</math> 1.9</b>	93.8 $\pm$ 3.3	80.5 $\pm$ 1.4	22.6 $\pm$ 4.6	12.5 $\pm$ 2.2	21.2 $\pm$ 1.4	12.0 $\pm$ 1.0	22.4 $\pm$ 3.1	12.8 $\pm$ 1.2
Flow	72.1 $\pm$ 1.3	56.4 $\pm$ 1.6	<b>75.5 <math>\pm</math> 2.3</b>	<b>56.5 <math>\pm</math> 1.8</b>	23.6 $\pm$ 3.2	12.5 $\pm$ 0.6	22.8 $\pm$ 2.7	12.4 $\pm$ 1.4	34.3 $\pm$ 4.3	24.4 $\pm$ 1.5
CelebA(SMILE)	<b>83.7 <math>\pm</math> 6.2</b>	<b>71.6 <math>\pm</math> 7.2</b>	78.8 $\pm$ 10.9	66.1 $\pm$ 12.1	22.5 $\pm$ 1.2	9.92 $\pm$ 1.2	27.2 $\pm$ 5.3	14.6 $\pm$ 4.2	23.5 $\pm$ 3.0	10.3 $\pm$ 1.6
CelebA(BEARD)	<b>92.3 <math>\pm</math> 5.6</b>	<b>83.3 <math>\pm</math> 8.6</b>	89.8 $\pm$ 6.2	78.7 $\pm$ 7.7	22.4 $\pm$ 1.9	9.82 $\pm$ 2.2	11.4 $\pm$ 1.5	20.0 $\pm$ 2.2	23.5 $\pm$ 3.0	8.1 $\pm$ 1.2

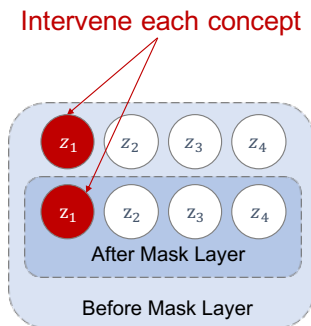


Figure 6. Intervention method

codes of effects sometimes has no influence to the whole image. This is probably because they do not explicitly consider causal disentanglement.

We also design another synthetic dataset Flow and do the same comparative experiments on that and the results support our claim. Because of page limitation, we show the results in Appendix D.

Fig. 4 demonstrates the good result of CausalVAE on real world benchmark dataset CelebA, with subfigures showing the experiments on intervening concepts GENDER, SMILE, EYES OPEN and MOUTH OPEN respectively. We observe that when we intervene the cause concept SMILE, the status of MOUTH OPEN also changes. In contrast, intervening effect concept MOUTH OPEN does not cause the cause concept SMILE to change. Table 1 records the mutual information (MIC/TIC) between the learned representation and the ground truth concept labels of all compared methods. Our model achieves best alignment with the concept labels, justifying the effectiveness of our proposed method. On the

contrary, factors learned by those compared methods have low correlation with the ground truth labels, indicating that those factors are at least not corresponding to the causal concepts of interest.

In addition, we show in Fig. 5 the learned adjacency matrix  $\mathbf{A}$ . To learn a precise causal graph, we design a pre-train process by optimizing augmented Lagrangian method [34] on Eq. 11, details are shown in Appendix C.3. As the training epoch increases, we see that the graph learned by our model quickly converges to the true one, which shows that our method is able to correctly learn the causal relationship among the factors.

## 7. Conclusion

In this paper, we investigate an important task of learning disentangled representations of causally related concepts in data, and propose a new framework called CausalVAE which includes a SCM layer to model the causal generation mechanism of data. We prove that the proposed model is fully identifiable given additional supervision signal. Experimental results with synthetic and real data show that CausalVAE successfully learns representations of causally related concepts and allows intervention to generate counterfactual outputs as expected according to our understanding of the causal system. To the best of our knowledge, our work is the first one that successfully implement causal disentanglement and is expected to bring new insights into the domain of disentangled representation learning.



## References

- [1] Michel Besserve, Rémy Sun, and Bernhard Schölkopf. Counterfactuals uncover the modular structure of deep generative models. *arXiv preprint arXiv:1812.03253*, 2018. [3](#)
- [2] Philemon Brakel and Yoshua Bengio. Learning independent features with adversarial nets for non-linear ica. *arXiv preprint arXiv:1710.05050*, 2017. [2](#)
- [3] Christopher P Burgess, Irina Higgins, Arka Pal, Loic Matthey, Nick Watters, Guillaume Desjardins, and Alexander Lerchner. Understanding disentangling in *beta*-vae. *arXiv preprint arXiv:1804.03599*, 2018. [1](#)
- [4] Tian Qi Chen, Xuechen Li, Roger B Grosse, and David K Duvenaud. Isolating sources of disentanglement in variational autoencoders. In *Advances in Neural Information Processing Systems*, pages 2610–2620, 2018. [1](#), [2](#)
- [5] David Heckerman, Dan Geiger, and David Maxwell Chickering. Learning bayesian networks: The combination of knowledge and statistical data. *CoRR*, abs/1302.6815, 2013. [3](#)
- [6] Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. beta-vae: Learning basic visual concepts with a constrained variational framework. *Iclr*, 2(5):6, 2017. [1](#), [2](#), [6](#)
- [7] 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*, pages 689–696, 2009. [2](#)
- [8] Jun-Ting Hsieh, Bingbin Liu, De-An Huang, Li F Fei-Fei, and Juan Carlos Nieves. Learning to decompose and disentangle representations for video prediction. In *Advances in Neural Information Processing Systems*, pages 517–526, 2018. [1](#)
- [9] Wei-Ning Hsu, Yu Zhang, and James Glass. Unsupervised learning of disentangled and interpretable representations from sequential data. In *Advances in neural information processing systems*, pages 1878–1889, 2017. [1](#)
- [10] Nan Rosemary Ke, Olexa Bilaniuk, Anirudh Goyal, Stefan Bauer, Hugo Larochelle, Chris Pal, and Yoshua Bengio. Learning neural causal models from unknown interventions. *CoRR*, abs/1910.01075, 2019. [3](#)
- [11] Ilyes Khemakhem, Diederik P. Kingma, and Aapo Hyvärinen. Variational autoencoders and nonlinear ICA: A unifying framework. *CoRR*, abs/1907.04809, 2019. [2](#), [4](#), [5](#), [10](#), [11](#), [12](#)
- [12] Ilyes Khemakhem, Ricardo Pio Monti, Diederik P Kingma, and Aapo Hyvärinen. Ice-beem: Identifiable conditional energy-based deep models. *arXiv preprint arXiv:2002.11537*, 2020. [12](#), [13](#)
- [13] Hyunjik Kim and Andriy Mnih. Disentangling by factorising. *arXiv preprint arXiv:1802.05983*, 2018. [1](#)
- [14] Justin B Kinney and Gurinder S Atwal. Equitability, mutual information, and the maximal information coefficient. *Proceedings of the National Academy of Sciences*, 111(9):3354–3359, 2014. [6](#)
- [15] Murat Kocaoglu, Christopher Snyder, Alexandros G. Dimakis, and Sriram Vishwanath. Causalgan: Learning causal implicit generative models with adversarial training. *CoRR*, abs/1709.02023, 2017. [2](#), [3](#), [6](#), [19](#)
- [16] Tejas D Kulkarni, William F Whitney, Pushmeet Kohli, and Josh Tenenbaum. Deep convolutional inverse graphics network. In *Advances in neural information processing systems*, pages 2539–2547, 2015. [2](#)
- [17] Daniel D. Lee, Masashi Sugiyama, Ulrike von Luxburg, Isabelle Guyon, and Roman Garnett, editors. *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, 2016. [2](#), [6](#)
- [18] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *2015 IEEE International Conference on Computer Vision, ICCV 2015, Santiago, Chile, December 7-13, 2015* [[18](#)], pages 3730–3738. [6](#), [9](#)
- [19] Francesco Locatello, Stefan Bauer, Mario Lucic, Gunnar Rätsch, Sylvain Gelly, Bernhard Schölkopf, and Olivier Bachem. Challenging common assumptions in the unsupervised learning of disentangled representations. *arXiv preprint arXiv:1811.12359*, 2018. [2](#), [4](#)
- [20] Francesco Locatello, Michael Tschannen, Stefan Bauer, Gunnar Rätsch, Bernhard Schölkopf, and Olivier Bachem. Disentangling factors of variation using few labels. *arXiv preprint arXiv:1905.01258*, 2019. [2](#)
- [21] Jianxin Ma, Chang Zhou, Peng Cui, Hongxia Yang, and Wenwu Zhu. Learning disentangled representations for recommendation. In *Advances in Neural Information Processing Systems*, pages 5712–5723, 2019. [1](#)
- [22] Emile Mathieu, Tom Rainforth, N Siddharth, and Yee Whye Teh. Disentangling disentanglement in variational autoencoders. *arXiv preprint arXiv:1812.02833*, 2018. [2](#)
- [23] Ignavier Ng, Zhuangyan Fang, Shengyu Zhu, and Zhitang Chen. Masked gradient-based causal structure learning. *arXiv preprint arXiv:1910.08527*, 2019. [3](#), [4](#)
- [24] Ignavier Ng, Shengyu Zhu, Zhitang Chen, and Zhuangyan Fang. A graph autoencoder approach to causal structure learning. *CoRR*, abs/1911.07420, 2019. [5](#)
- [25] Judea Pearl. *Causality*. Cambridge university press, 2009. [2](#), [3](#), [4](#)
- [26] Jonas Peters, Peter Bühlmann, and Nicolai Meinshausen. Causal inference using invariant prediction: identification and confidence intervals. *arXiv preprint arXiv:1501.01332*, 2015. [3](#)
- [27] Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. *Elements of causal inference*. The MIT Press, 2017. [3](#)
- [28] Bernhard Schölkopf. Causality for machine learning. *arXiv preprint arXiv:1911.10500*, 2019. [3](#)
- [29] Shohei Shimizu, Patrik O Hoyer, Aapo Hyvärinen, and Antti Kerminen. A linear non-gaussian acyclic model for causal discovery. *Journal of Machine Learning Research*, 7(Oct):2003–2030, 2006. [2](#), [3](#), [13](#)
- [30] Kihyuk Sohn, Honglak Lee, and Xinchen Yan. Learning structured output representation using deep conditional generative models. In *Advances in neural information processing systems*, pages 3483–3491, 2015. [6](#)
- [31] Peter Sorrenson, Carsten Rother, and Ullrich Köthe. Disentanglement by nonlinear ica with general incompressible-flow networks (gin). *arXiv preprint arXiv:2001.04872*, 2020. [11](#)

- [32] Raphael Suter, Dorde Miladinović, Bernhard Schölkopf, and Stefan Bauer. Robustly disentangled causal mechanisms: Validating deep representations for interventional robustness. *arXiv preprint arXiv:1811.00007*, 2018. 3
- [33] Robert E. Tillman and Peter Spirtes. Learning equivalence classes of acyclic models with latent and selection variables from multiple datasets with overlapping variables. In Geoffrey J. Gordon, David B. Dunson, and Miroslav Dudík, editors, *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics, AISTATS 2011, Fort Lauderdale, USA, April 11-13, 2011*, volume 15 of *JMLR Proceedings*, pages 3–15. JMLR.org, 2011. 3
- [34] Yue Yu, Jie Chen, Tian Gao, and Mo Yu. Dag-structure learning with graph neural networks. *arXiv preprint arXiv:1904.10098*, 2019. 5, 8, 12, 13
- [35] Kun Zhang and Aapo Hyvarinen. On the identifiability of the post-nonlinear causal model. *arXiv preprint arXiv:1205.2599*, 2012. 2
- [36] Xun Zheng, Bryon Aragam, Pradeep K Ravikumar, and Eric P Xing. Dags with no tears: Continuous optimization for structure learning. In *Advances in Neural Information Processing Systems*, pages 9472–9483, 2018. 3, 5
- [37] Shengyu Zhu and Zhitang Chen. Causal discovery with reinforcement learning. *CoRR*, abs/1906.04477, 2019. 5
- [38] Shengyu Zhu, Ignavier Ng, and Zhitang Chen. Causal discovery with reinforcement learning. In *International Conference on Learning Representations (ICLR)*, 2020. 3

## A. Proof of Proposition 1

Write the KL term in ELBO defined in Eq. 8 in the main text as

$$\begin{aligned}
& \mathcal{D}[q_\phi(\epsilon, \mathbf{z}|\mathbf{x}, \mathbf{u})\|p_\theta(\epsilon, \mathbf{z}|\mathbf{u})] \\
&= \iint q_\phi(\epsilon, \mathbf{z}|\mathbf{x}, \mathbf{u}) \log \frac{q_\phi(\epsilon, \mathbf{z}|\mathbf{x}, \mathbf{u})}{p_\epsilon(\epsilon)p_\theta(\mathbf{z}|\mathbf{u})} d\epsilon d\mathbf{z} \\
&= \iint q_\phi(\epsilon, \mathbf{z}|\mathbf{x}, \mathbf{u}) \log \frac{q_\phi(\epsilon, \mathbf{z}|\mathbf{x}, \mathbf{u})}{p_\epsilon(\epsilon)} d\epsilon d\mathbf{z} \\
&\quad + \iint q_\phi(\epsilon, \mathbf{z}|\mathbf{x}, \mathbf{u}) \log \frac{q_\phi(\epsilon, \mathbf{z}|\mathbf{x}, \mathbf{u})}{p_\theta(\mathbf{z}|\mathbf{u})} d\epsilon d\mathbf{z} \\
&\quad - \iint q_\phi(\epsilon, \mathbf{z}|\mathbf{x}, \mathbf{u}) \log q_\phi(\epsilon, \mathbf{z}|\mathbf{x}, \mathbf{u}) d\epsilon d\mathbf{z},
\end{aligned}$$

The third term in above equation could be rewritten as a constant. Details are shown as below.

$$\begin{aligned}
& - \iint q_\phi(\epsilon, \mathbf{z}|\mathbf{x}, \mathbf{u}) \log q_\phi(\epsilon, \mathbf{z}|\mathbf{x}, \mathbf{u}) d\epsilon d\mathbf{z} \\
&= - \iint q(\epsilon|x, u) \delta(z = C\epsilon) \log q(\epsilon|x, u) d\epsilon d\mathbf{z} \\
&\quad - \iint q(\epsilon|x, u) \delta(z = C\epsilon) \log \delta(z = C\epsilon) d\epsilon d\mathbf{z} \\
&= H(q_\phi(\epsilon|x, u)) - 0 = H(\mathcal{N}(\mu_\phi(x, u), s\mathbf{I})) \\
&= \text{const}, \tag{16}
\end{aligned}$$

In our method, we ignore this term in ELBO expression. Then, based on Eq. 9 in the main text, we have

$$\begin{aligned}
& \iint q_\phi(\epsilon, \mathbf{z}|\mathbf{x}, \mathbf{u}) \log \frac{q_\phi(\epsilon, \mathbf{z}|\mathbf{x}, \mathbf{u})}{p_\epsilon(\epsilon)} d\epsilon d\mathbf{z} \\
&= \int q_\phi(\epsilon|\mathbf{x}, \mathbf{u}) \log \frac{q_\phi(\epsilon|\mathbf{x}, \mathbf{u})}{p_\epsilon(\epsilon)} \int \delta(\mathbf{z} = \mathbf{C}\epsilon) d\mathbf{z} d\epsilon \\
&\quad + \int q_\phi(\epsilon|\mathbf{x}, \mathbf{u}) \int \delta(\mathbf{z} = \mathbf{C}\epsilon) \log \delta(\mathbf{z} = \mathbf{C}\epsilon) d\mathbf{z} d\epsilon \\
&= \mathcal{D}[q_\phi(\epsilon|\mathbf{x}, \mathbf{u})\|p_\epsilon(\epsilon)] + 0 \\
&= \mathcal{D}[q_\phi(\epsilon|\mathbf{x}, \mathbf{u})\|p_\epsilon(\epsilon)],
\end{aligned}$$

and

$$\begin{aligned}
& \iint q_\phi(\epsilon, \mathbf{z}|\mathbf{x}, \mathbf{u}) \log \frac{q_\phi(\epsilon, \mathbf{z}|\mathbf{x}, \mathbf{u})}{p_\theta(\mathbf{z}|\mathbf{u})} d\epsilon d\mathbf{z} \\
&= \int q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{u}) \log \frac{q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{u})}{p_\theta(\mathbf{z}|\mathbf{u})} \int \delta(\epsilon = \mathbf{C}^{-1}\mathbf{z}) d\epsilon d\mathbf{z} \\
&\quad + \int q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{u}) \int \delta(\epsilon = \mathbf{C}\mathbf{z}) \log \delta(\epsilon = \mathbf{C}^{-1}\mathbf{z}) d\epsilon d\mathbf{z} \\
&= \mathcal{D}[q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{u})\|p_\theta(\mathbf{z}|\mathbf{u})] + 0 \\
&= \mathcal{D}[q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{u})\|p_\theta(\mathbf{z}|\mathbf{u})].
\end{aligned}$$

Adding up the above two terms leads to the desired form of Proposition 1.

## B. Identifiability

### B.1. Proof of Theorem 1

The general logic of the proofing follows [11], but we focus on both encoder and decoder. In our setting, we has joint latent variables  $\epsilon, \mathbf{z}$ , and we prove identifiability of both of them.

Another different setting from iVAE is that we consider a slighter transformation matrix, since our additional observations  $\mathbf{u}$  of each concepts align to each causal representations  $\mathbf{z}$ .

#### Sketch of proof:

We analyze the identifiability of  $\epsilon$  starting with  $p_\theta(\mathbf{x}|\mathbf{u}) = p_{\tilde{\theta}}(\mathbf{x}|\mathbf{u})$ . Then we define a new invertible matrix  $\mathbf{L}$  which contains additional observation  $u_i$  in causal system, and use it to prove that the learned  $\tilde{\mathbf{T}}$  is the transformation of  $\mathbf{T}$ . Step 2: We take the inference model into consideration and analyze the identifiability of the inference model by relating the inference model to the generative model.

#### Details:

At the begining of proof, we consider a simple condition that the dimension of observation data  $d$  equals to the dimension of latent variables  $n$ .

The distribution has two sufficient statistics, the mean and variance of  $\mathbf{z}$ , which are denoted by sufficient statistics  $\mathbf{T}(\mathbf{z}) = (\boldsymbol{\mu}(\mathbf{z}), \boldsymbol{\sigma}(\mathbf{z})) = (T_{1,1}(z_1), \dots, T_{n,2}(z_n))$ . We use

these notations for model identifiability analysis in Section 5. To simplify proof process, we absorb the injective functions  $\mathbf{g}(\cdot)$  into generate model  $\mathbf{f}(\cdot)$  since mask layer will not influence the quality of disentangled representation  $\mathbf{z}$ .

$$\begin{aligned}
p_{\theta}(\mathbf{x}|\mathbf{u}) &= p_{\tilde{\theta}}(\mathbf{x}|\mathbf{u}), \\
\Rightarrow \int \int_{\mathbf{z}, \epsilon} p_{\theta}(\mathbf{x}|\mathbf{z}, \epsilon) p_{\theta}(\mathbf{z}, \epsilon|\mathbf{u}) d\mathbf{z} d\epsilon \\
&= \int \int_{\mathbf{z}, \epsilon} p_{\tilde{\theta}}(\mathbf{x}|\mathbf{z}, \epsilon) p_{\tilde{\theta}}(\mathbf{z}, \epsilon|\mathbf{u}) d\mathbf{z} d\epsilon, \\
\Rightarrow \int_{\mathbf{z}} p_{\theta}(\mathbf{x}|\mathbf{z}) p_{\theta}(\mathbf{z}|\mathbf{u}) d\mathbf{z} &= \int \int_{\mathbf{z}} p_{\tilde{\theta}}(\mathbf{x}|\mathbf{z}) p_{\tilde{\theta}}(\mathbf{z}|\mathbf{u}) d\mathbf{z}, \\
\Rightarrow \int_{\mathbf{x}'} p_{\theta}(\mathbf{x}|\mathbf{f}^{-1}(\mathbf{x}')) p_{\theta}(\mathbf{f}^{-1}(\mathbf{x}')|\mathbf{u}) |\det(J_{\mathbf{f}^{-1}}(\mathbf{x}'))| d\mathbf{x}' \\
&= \int_{\mathbf{x}'} p_{\theta}(\mathbf{x}|\tilde{\mathbf{f}}^{-1}(\mathbf{x}')) p_{\tilde{\theta}}(\tilde{\mathbf{f}}^{-1}(\mathbf{x}')|\mathbf{u}) |\det(J_{\tilde{\mathbf{f}}^{-1}}(\mathbf{x}'))| d\mathbf{x}'. \tag{17}
\end{aligned}$$

In determining function  $\mathbf{f}$ , there exist a Gaussian distribution  $p_{\xi}(\xi)$  which has infinitesimal variance. Then, the  $p_{\theta}(\mathbf{x}|\mathbf{f}^{-1}(\mathbf{x}'))$  can be written as  $p_{\xi}(\mathbf{x} - \mathbf{x}')$ . As the assumption (1) holds, this term is vanished. Then in our method, there exists the following equation:

$$\begin{aligned}
p_{\theta}(\mathbf{f}^{-1}(\mathbf{x}')|\mathbf{u}) |\det(J_{\mathbf{f}^{-1}}(\mathbf{x}'))| &= p_{\tilde{\theta}}(\tilde{\mathbf{f}}^{-1}(\mathbf{x}')|\mathbf{u}) |\det(J_{\tilde{\mathbf{f}}^{-1}}(\mathbf{x}'))|, \\
\Rightarrow \tilde{p}_{\theta}(\mathbf{x}) &= \tilde{p}_{\tilde{\theta}}(\mathbf{x}). \tag{18}
\end{aligned}$$

Adopting the definition of multivariate Gaussian distribution, we define

$$\lambda_s(\mathbf{u}) = \begin{bmatrix} \lambda_1^s(u_1) & & \\ & \ddots & \\ & & \lambda_n^s(u_n) \end{bmatrix}. \tag{19}$$

There exists the following equations:

$$\begin{aligned}
&\log |\det(J_{\mathbf{f}^{-1}}(\mathbf{x}))| - \log \mathbf{Q}(\mathbf{f}^{-1}(\mathbf{x})) + \log \mathbf{Z}(\mathbf{u}) \tag{20} \\
&+ \sum_{s=1}^2 \mathbf{T}_s(\mathbf{f}^{-1}(\mathbf{x})) \lambda_s(\mathbf{u}), \\
&= \log |\det(J_{\tilde{\mathbf{f}}^{-1}}(\mathbf{x}))| - \log \tilde{\mathbf{Q}}(\tilde{\mathbf{f}}^{-1}(\mathbf{x})) + \log \tilde{\mathbf{Z}}(\mathbf{u}) \\
&+ \sum_{s=1}^2 \tilde{\mathbf{T}}_s(\tilde{\mathbf{f}}^{-1}(\mathbf{x})) \tilde{\lambda}_s(\mathbf{u}), \tag{21}
\end{aligned}$$

where  $\mathbf{Q}$  denotes the base measure. In Gaussian distribution, it is  $\sigma(\mathbf{z})$ .

In learning process,  $\tilde{\mathbf{A}}$  is restricted as DAG. Thus, the  $\tilde{\mathbf{C}}$  exists which is full rank matrix. The item which is not related to  $u$  in Eq. 21 are cancelled out [31].

$$\sum_{s=1}^2 \mathbf{T}_s(\mathbf{f}^{-1}(\mathbf{x})) \lambda_s(\mathbf{u}) = \sum_{s=1}^2 \tilde{\mathbf{T}}_s(\tilde{\mathbf{f}}^{-1}(\mathbf{x})) \tilde{\lambda}_s(\mathbf{u}) + \mathbf{b}, \tag{22}$$

where  $\mathbf{b}$  is a vector related to  $\mathbf{u}$ .

In our model, there exist a deterministic relationship  $\mathbf{C}$  between  $\epsilon$  and  $\mathbf{z}$  where  $\mathbf{C} = (\mathbf{I} - \mathbf{A}^T)^{-1}$ . Thus we could get equivalent of Eq. 22 as follows,

$$\sum_{s=1}^2 \mathbf{T}_s(\mathbf{C}\mathbf{h}(\mathbf{x})) \lambda_s(\mathbf{u}) = \sum_{s=1}^2 \tilde{\mathbf{T}}_s(\tilde{\mathbf{C}}\tilde{\mathbf{h}}(\mathbf{x})) \tilde{\lambda}_s(\mathbf{u}) + \mathbf{b}', \tag{23}$$

where  $s$  denote the index of sufficient statistics of Gaussian distributions, indexing the mean (1) and the variance (2).

By assuming that the additional observation  $u_i$  is different, it is guaranteed that coefficients of the observations for different concepts are distinct. Thus, there exists an invertible matrix corresponding to additional information  $\mathbf{u}$ :

$$\mathbf{L} = \begin{bmatrix} \lambda_1(\mathbf{u}) & \\ & \lambda_2(\mathbf{u}) \end{bmatrix}. \tag{24}$$

Since the assumption that  $u_i \neq 0$  holds,  $\mathbf{L}$  is  $2n \times 2n$  invertible and full rank diagonal matrix. Then, function of  $\lambda$  in Eq. 22 and Eq. 23 are replaced by Eq. 24, we could get:

$$\mathbf{L}\mathbf{T}(\mathbf{f}^{-1}(\mathbf{x})) = \tilde{\mathbf{L}}\tilde{\mathbf{T}}(\tilde{\mathbf{f}}^{-1}(\mathbf{x})) + \mathbf{b}, \tag{25}$$

$$\mathbf{T}(\mathbf{f}^{-1}(\mathbf{x})) = \mathbf{B}_2 \tilde{\mathbf{T}}(\tilde{\mathbf{f}}^{-1}(\mathbf{x})) + \mathbf{b}_2, \tag{26}$$

$$\mathbf{B}_2 = \begin{bmatrix} \lambda_{1,1}(u_1)^{-1} \tilde{\lambda}_{1,1}(u_1) & & \\ & \ddots & \\ & & \lambda_{n,2}(u_n) \tilde{\lambda}_{n,2}(u_n) \end{bmatrix}. \tag{27}$$

We replace  $\mathbf{f}^{-1}$  with  $\mathbf{C}\mathbf{h}$  and we could get the equations as below:

$$\mathbf{L}\mathbf{T}(\mathbf{C}\mathbf{h}(\mathbf{x})) = \tilde{\mathbf{L}}\tilde{\mathbf{T}}(\tilde{\mathbf{C}}\tilde{\mathbf{h}}(\mathbf{x})) \Rightarrow \mathbf{T}(\mathbf{h}(\mathbf{x})) = \mathbf{B}_1 \tilde{\mathbf{T}}(\tilde{\mathbf{h}}(\mathbf{x})) + \mathbf{b}_1, \tag{28}$$

where  $\mathbf{B}_3 = \mathbf{C}\tilde{\mathbf{C}}^{-1}$  is invertible matrix which corresponds to  $\mathbf{C}$  and  $\mathbf{B}_1 = \mathbf{L}^{-1}\mathbf{B}_3^{-1}\tilde{\mathbf{L}}$ . The definition of  $\tilde{\mathbf{L}}$  on learning model migrates the definition of  $\mathbf{L}$  on ground truth.

Then we adopt the definitions following [11]. According to the Lemma 3 in [11], we are able to pick out a pair  $(\epsilon_i, \epsilon_i^2)$  such that,  $(\mathbf{T}'_i(z_i), \mathbf{T}'_i(z_i^2))$  are linearly independent. Then concat the two points into a vector, and denote the Jacobian matrix  $\mathbf{Q} = [J_{\mathbf{T}}(\epsilon), J_{\mathbf{T}}(\epsilon^2)]$ , and define  $\tilde{\mathbf{Q}}$  on  $\tilde{\mathbf{T}}(\tilde{\mathbf{h}} \circ \mathbf{C}\mathbf{f}(\epsilon))$  in the same manner. By differentiating Eq. 28, we get

$$\mathbf{Q} = \mathbf{B}_1 \tilde{\mathbf{Q}}. \tag{29}$$

Since the assumption (2) that Jacobian of  $\mathbf{f}^{-1}$  is full rank holds, it can prove that both  $\mathbf{Q}$  and  $\tilde{\mathbf{Q}}$  are invertible matrix. Thus from Eq. 29,  $\mathbf{B}_1$  is invertible matrix. Using the same way as shown in Eq. 29, it can prove that  $\mathbf{B}_2$  is invertible matrix.

Eq. 26 and Eq. 28 both hold. Combining the two results supports the identifiability result in CausalVAE.

## B.2. Extension of Definition 1

In most of scenarios, latent variable is a low dimensional representation of the observation, since we are not interested in all the information in observations.

Therefore, we usually have  $d > n$ . We called it the reduction of dimension. We add auxiliary term as  $\lambda(\mathbf{x}) = \{\lambda(\mathbf{u}), \lambda'\}$ . In our model, Only  $n$  components of the latent variable are modulated, and its density has the form:

$$p_{\theta}(\mathbf{z}|\mathbf{u}) = \frac{\mathbf{Q}(\mathbf{z})}{\mathbf{Z}(\mathbf{u})} \exp \sum_i^n \mathbf{T}_i(z_i) \lambda_i(u_i) \quad (30)$$

and the term  $e^{\sum_{n+1}^d \mathbf{T}(z_i) \lambda_i}$  is simply absorbed into  $\mathbf{Q}(\mathbf{z})$ . When we evaluate Eq. 21 by new definition (Eq. 30), the dimension of  $p(\mathbf{z}|\mathbf{u})$  is  $n$ , because the remaining part is cancelled out.

Assume that  $p_{\theta}(\mathbf{x}|\mathbf{u})$  equal to  $p_{\tilde{\theta}}(\mathbf{x}|\mathbf{u})$ . For all the observational pairs  $(\mathbf{x}, \mathbf{u})$ , let  $J_h$  denote the Jacobian matrix of the encoder function. Following the definition in Theorem 2 in i VAE [11],  $\mathbf{B}$  will be indexed by 4 indicates  $(i, l, a, b)$ , where  $1 < i < d$  and  $1 < l < s$  refer to the rows and  $1 < a < d$  and  $1 < b < s$  refer to the columns. We define a following equation:

$$\mathbf{v} = \tilde{\mathbf{C}} \circ \tilde{\mathbf{h}} \circ \mathbf{f}(\mathbf{z}). \quad (31)$$

The goal is to show that  $v_i(\mathbf{z})$  is a function of only one  $z_j$ . We denote by  $v_i^r := \frac{\partial v_i}{\partial z_r}$  and  $v_i^{rt} := \frac{\partial^2 v_i}{\partial z_r \partial z_t}$ . By differentiating Eq. 26 with respect to  $z_s$ , we could get:

$$T'_{i,l}(z_i) = \sum_{a=1}^d \sum_{b=1}^s B_{2,(i,l,a,b)} \tilde{T}'_{a,b}(v_a(\mathbf{z})) v_a^r(\mathbf{z}). \quad (32)$$

**Lemma 1** (from Lemma 9 in Khemakhem et al. [12]): Consider a distribution that follows a strongly exponential family. Its sufficient statistic  $\tilde{\mathbf{T}}$  is differentiable almost surely. Then  $\tilde{T}'_i \neq 0$  almost everywhere on  $\mathbb{R}$  for all  $1 \leq i \leq s$ .

For  $r > n$ ,  $T'_{i,l}(z_i) = 0$ , according to Lemma 1,  $\tilde{T}'_{a,b}(v_a(\mathbf{z})) \neq 0$ , since  $\mathbf{B}_2$  is an invertible matrix, we can conclude that  $v_a^r(\mathbf{z}) = 0$  for all  $a < n$  and  $r > n$ . Therefore, we can conclude that each of the first  $n$  components of  $\mathbf{v}$  is only a function of one different  $z_j$ . Thus, when  $d > n$ , we could get the same conclusion as Theorem 1.

## B.3. Identifiability of Causal Graph

Consider the identifiability analysis in Appendix B.1. For the framework of CausalVAE, in Causal Layer, the latent code  $\mathbf{z}$  is identified since  $\mathbf{B}_2$  is a diagonal matrix which corresponds to learnt  $\tilde{\mathbf{z}}$  and  $\mathbf{z}$ . Since the true  $\epsilon$  and learnt  $\tilde{\epsilon}$  are linearly related,  $\mathbf{B}_1$ ,  $\mathbf{C}$  and  $\tilde{\mathbf{C}}$  are in a linear equivalent

class. In other words,  $\mathbf{C}$  or  $\mathbf{A}$  is identifiable in Causal Layer up to a linear equivalent class.

In our work, strict identifiability is guaranteed by the non-linear mask layer. Details of the Mask Layer are shown in Section 3.2 in main text. The Mask Layer uses non-linear functions and additional supervision signal  $\mathbf{u}$  (non-Gaussian) to help the model to identify the true causal graph in a linear equivalent class.

## C. Implementation Details

We use one NVIDIA Tesla P40 GPU as our training and inference device.

For the implementation of CausalVAE and other baselines, we extend  $\mathbf{z}$  to matrix  $\mathbf{z} \in \mathbb{R}^{n \times k}$  where  $n$  is the number of concepts and  $k$  is the latent dimension of each  $\mathbf{z}_i$ . The corresponding prior or conditional prior distributions of CausalVAE and other baselines are also adjusted (this means that we extend the multivariate Gaussian to the matrix Gaussian).

The subdimensions  $k$  for each synthetic (pendulum, water) experiments are set to be 4, and 32 for CelebA experiments. The implementation of continuous DAG constraint  $H(\mathbf{A})$  follows the code of [34]<sup>3</sup>.

### C.1. Data Preprocessing

#### C.1.1 Synthetic Simulator

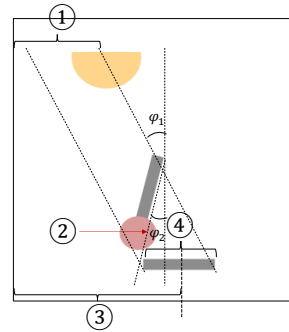


Figure 7. Generate Policy of Pendulum Simulator

Fig. 7 shows our policy of generating synthetic Pendulum data. The picture includes a pendulum. The angles of pendulum and the light are changing overtime, and projection laws are used to generate the shadows. Given the light POSITION and pendulum ANGLE, we get the angles  $\varphi_1$  and  $\varphi_2$ . Then the system can calculate the shadow POSITION and LENGTH using triangular functions. The causal graph of concepts is shown in Fig. 10 (a). In Pendulum generator, the image size is set to be  $96 \times 96$  with 4 channels. We generate about

<sup>3</sup><https://github.com/fishmoon1234/DAG-GNN>

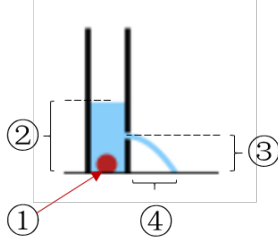


Figure 8. Generate Policy of Flow Simulator

7k images (6k for training and 1k for inference),  $\varphi_1$  and  $\varphi_2$  are ranged in around  $[-\frac{\pi}{4}, \frac{\pi}{4}]$ , and they are generated independently. For each image, we provide 4 labels, which include light position, pendulum angle, shadow position and shadow length. For light position, we use the value of center of semicircle (Fig.7 ①) as supervision signal. For the pendulum angle, we use the value of  $\varphi_2$  as supervision signal (Fig. 7 ②). For shadow position and shadow length, we use the length of Fig.7 ③ and Fig.7 ④ as supervision signal respectively.

Fig. 8 presents our policy of generating synthetic Flow data. Each image is of the  $96 \times 96 \times 4$  resolution, and consists of a cup of water and a ball. The original water level, the ball size (Fig.8①) and the location of hole (Fig.8③) vary over time. Given the ball size Fig. 8 and the original water level, we determine the WATER HEIGHT (Fig.8②). Then we generate WATER FLOW according to the Parabola law, where we additionally introduce a noise from  $\mathcal{N}(0, 0.01)$  to the gravitational acceleration. The causal graph of concepts is given in Fig. 10 (b). We consider four semantically meaningful concepts, BALLS SIZE, WATER HEIGHT, HOLE POSITION and WATER FLOW, whose supervised signals are given by the ball’s diameter (Fig.7 ①), the length of Fig. 7 ②, the length of Fig.7 ③ and Fig.7 ④ respectively. The sample size is 8k with 6k for training and 2k for testing.

### C.1.2 Data Preprocess of CelebA

CelebA dataset contains 20K human face images. We preprocess the original dataset by following two steps:

(1) We divided the whole dataset into training dataset 85% and test dataset 15%.

(2) We only focus on facial features and resize the picture to be squared ( $128 \times 128$  with 3 channels).

## C.2. Intervention Experiments

### C.2.1 Synthetic

In synthetic experiments, we train the model on synthetic data for 80 epochs, and use this model to generate latent

code of representations. The hyperparameters of baselines are defined as default.

For CausalVAE, we set the  $\alpha = 0.3$  and  $(\beta, \gamma) = (1, 1)$ . We use  $\mathcal{N}(\mathbf{u}, |\mathbf{u}|)$  as the condition prior  $p_{\theta}(\mathbf{z}|\mathbf{u})$ . In the implementation of CausalVAE,  $|\mathbf{z}_{\text{mean}}|$  is used as the variance of condition prior.

The details of the neural networks are shown in Table 2. We all follows the neural network design strategy of Khemakhem *et al.* [12] to satisfy Theorem 1 assumption (ii).

### C.2.2 CelebA

We also present the DO-experiments of CausalVAE and CausalGAN. In the training of the models, we use face labels (AGE, GENDER and BEARD).

For CausalVAE, we set the  $\alpha = 0.3$  and  $(\beta, \gamma) = (1, 1)$ . We use  $\mathcal{N}(\mathbf{u}, \mathbf{I})$  as the condition prior  $p_{\theta}(\mathbf{z}|\mathbf{u})$ . For all the baseline, default hyperparameters and one common encoder and decoder structure are employed. For CausalGAN, we use the publicly available code<sup>4</sup>.

For all the VAE-based methods, mean and variance of the distribution of the latent variable are learned during training, and the latent code  $z$  are sampled from Conditional Gaussian Distribution  $p_{\theta}(\mathbf{z}|\mathbf{u})$ . In all experiments, we rescale the variance of learned representation  $\mathbf{z}$  by multiplying a factor 0.1 to the original one.

Training epoches for the model is set to be 80, and our proposed CausalVAE has a pretrain step to learn causal graph  $\mathbf{A}$ , which takes 10 epoches.

The details of the neural networks are shown in Table 3.

### C.3. The Pretrain Step for Causal Graph Learning

In our model, we need to learn the latent representation  $\mathbf{z}$  and causal graph  $\mathbf{A}$  simultaneously, whose optimal solution is not easy to find. Thus we adopt a pretrain stage to learn the causal graph  $\mathbf{A}$  in the Mask Layer. We adopt the augmented Lagrangian to learn  $\mathbf{A}$  in CausalVAE from the labels  $\mathbf{u}$  in Mask Layer first. During the pretrain process, we truncate the gradient of other part of model and solve the optimization problem in Eq. 34 to learn  $\mathbf{A}$ .

The augmentation approach is widely used in causal discovery method, like NOTEARS [29], DAG-GNN [34]. The pretrain is a stage that learns the graph by optimizing the following objective functions:

$$\begin{aligned} \text{minimize } l_u &= \mathbb{E}_{q_D} \|\mathbf{u} - \mathbf{A}^T \mathbf{u}\|_2^2 \\ \text{subject to } H(\mathbf{A}) &= 0 \end{aligned} \quad (33)$$

Then, we define an augmented Lagrangian:

$$l_{pre} = l_u + \lambda H(\mathbf{A}) + \frac{c}{2} H^2(\mathbf{A}) \quad (34)$$

<sup>4</sup><https://github.com/mkocaoglu/CausalGAN>

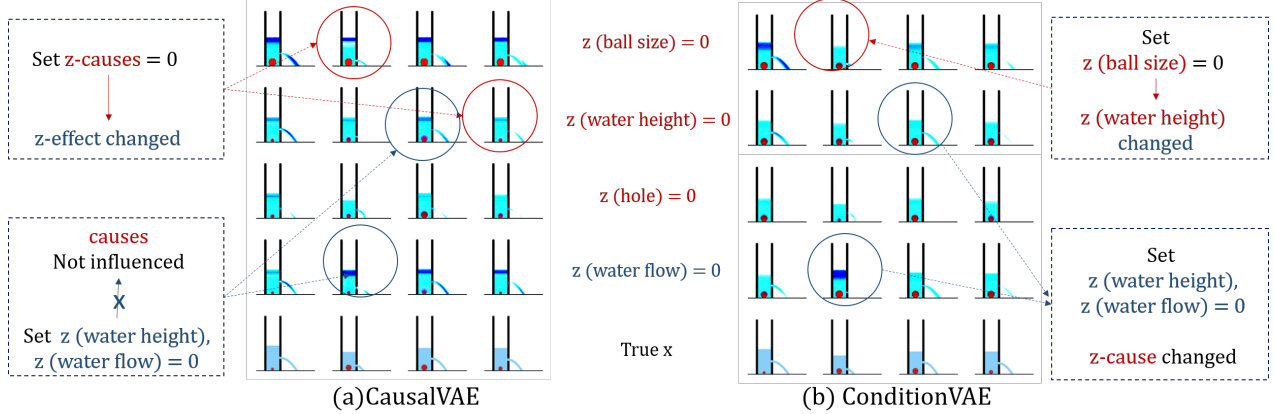


Figure 9. The results of Intervention experiments on the Flow dataset. Each row shows the result of controlling the BALL SIZE, WATER HEIGHT, HOLE, and WATER FLOW respectively. The bottom row is the original input image.

where  $\lambda$  is the Lagrangian multiplier and  $c$  is the penalty.

The following policy is used to update the  $\lambda$  and  $c$ :

$$\lambda_{s+1} = \lambda_s + c_s H(\mathbf{A}_s) \quad (35)$$

$$c_{s+1} = \begin{cases} c_s = \eta c_s, & \text{if } |H(\mathbf{A}_s)| > \gamma |H(\mathbf{A}_{s-1})| \\ c_s = c_s, & \text{otherwise} \end{cases}$$

where  $s$  is the iteration. In our experiments, we set  $\eta = 10$  and  $\gamma = \frac{1}{4}$ .

## D. Additional Experimental Results

In this section, we show more experimental results. Fig. 10 shows the causal graph among concepts in different dataset respectively. We here show results including experiments analyzing the properties of learned representation, intervening results and the learning process of the causal graph.

### D.1. The Property of Learned Representation

We test our method and baselines on both synthetic data and benchmark human face data. In the previous section, we already show the relationships between the learned representation  $\tilde{\mathbf{z}}$  and the target representation  $\mathbf{z}$  (related by a linear transformation formed as a diagonal matrix). In this section, we visualize it by scatter plot.

One of the important aspect of the generative model is that whether the learned representation aligns to the conditional prior we set. Our conditional prior is generated by the true label of each concept. The results show that the learned representations align to the expected representations. In figures, points are sampled from the joint distribution, and each color corresponds to one dimension.

The additional observations (labels) of Pendulum dataset and those of CelebA dataset are different. In Pendulum, the

labels are values within a fixed range. The labels in CelebA dataset are discrete (in  $\{-1, 1\}$ ). Thus the scatter plots are different.

The results show that the performance of our proposed method is better than all the baselines, including the supervised method and unsupervised method.

### D.2. The Learned Graph

We demonstrate the learning process of causal graph in this section. Fig. 14 shows the graph learned process of CelebA (BEARD). In this process, we initialize all the entries in  $\mathbf{A}$  as 0.5. After 5 epochs, the graph converges. We observe an almost correct graph in this group of concepts.

### D.3. Intervention Results

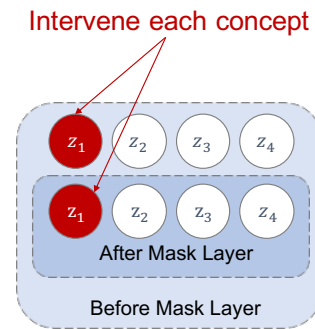


Figure 16. Intervention method

The intervention operations are as:

- For the learned model, we first put an random observed image  $\mathbf{x}$  into the encoder. In this process we could get  $\epsilon$  and  $\mathbf{z}$ .

encoder	decoder
4*96*96*900 fc. 1ELU	concepts*( 4* 300 fc. 1ELU )
900*300 fc. 1ELU	concepts* (300*300 fc. 1ELU)
300*2*concepts*k fc.	concepts*(300* 1024 fc. 1ELU)
-	concepts*(1024* 4*96*96 fc.)

Table 2. Network design of models trained on synthetic data.

encoder	decoder
-	(1*1 conv. 128 1LReLU(0.2), stride 1)
4*4 conv. 32 1LReLU (0.2), stride 2	(4*4 convtranspose. 64 1LReLU (0.2), stride 1)
4*4 conv. 64 1LReLU (0.2), stride 2	(4*4 convtranspose. 64 1LReLU (0.2), stride 2)
4*4 conv. 64 1LReLU(0.2), stride 2	(4*4 convtranspose. 32 1LReLU (0.2), stride 2)
4*4 conv. 64 1LReLU (0.2), stride 2	(4*4 convtranspose. 32 1LReLU (0.2), stride 2)
4*4 conv. 256 1LReLU (0.2), stride 2	(4*4 convtranspose. 32 1LReLU (0.2), stride 2)
1*1 conv. 3, stride 1	(4*4 convtranspose. 3 , stride 2)

Table 3. Network design of models trained on CelebA.

- Then for  $i$ -th concept, we fix the value of  $z_i$  and  $g_i(\mathbf{A}_i \circ \mathbf{z})$  as constants.
- Finally, we put the new  $\mathbf{z}$  into the decoder and get  $\mathbf{x}'$ .

Fig. 9 (a) demonstrates the intervention results of CausalVAE on Flow dataset. We see that when we intervene on the cause concept BALL SIZE, its child concepts WATER HEIGHT and WATER FLOW change correspondingly. Similarly, when the cause concept HOLE is intervened, its child concept WATER FLOW also changes. In contrast, intervening on effect concept WATER HEIGHT does not influence the causal concept BALL SIZE. Fig. 9(b) shows the results of ConditionVAE on Flow. We observe that when we intervene on BALL SIZE, WATER HEIGHT and WATER FLOW are affected as expected. However when we intervene on the effect concepts WATER HEIGHT and WATER FLOW, concept BALL SIZE is also influenced, which makes no sense. In general, the “do-intervention” of ConditionVAE performs worse than CausalVAE. The results support that by simply using a supervised model, one can not guarantee a causal disentangled representation.

The Fig. 17 demonstrates the result of CausalVAE on real world benchmark dataset CelebA (BEARD), with subfigures (a) (b) (c) (d) showing the intervention experiments on concepts of AGE, GENDER, BALD and BEARD respectively. The interventions perform well that when we intervened the cause concept GENDER, the BEARD changes correspondingly. Similarly, when the cause concept AGE is intervened, its child concept BALD also changes. In contrast, intervening effect concept BEARD does not influence the causal concepts GENDER and other unrelated concepts in Fig. 17 (d). Fig. 18 demonstrates the results of CausalGAN, with subfigures (a) (b) (c) (d) showing the intervention experiments on concepts CelebA (BEARD). We observe that when we intervene

GENDER, the BEARD are changed. But when we intervene BEARD, concept GENDER is also changed in third line as shown by Fig. 18 (d). In general, the ‘do-intervention’ of CausalGAN performs worse than CausalVAE.

The Fig. 19 demonstrates the result of CausalVAE on real world benchmark dataset CelebA (SMILE), with subfigures (a) (b) (c) (d) showing the intervention experiments on concepts of GENDER, SMILE, MOUTH OPEN and EYES OPEN respectively. The interventions perform well that when we intervened the cause concept GENDER, not only the appearance of GENDER but the eyes changed. When we intervened the cause concept SMILE, not only the appearance of SMILE but the MOUTH OPEN. In contrast, intervening effect concept MOUTH OPEN does not influence the causal concepts SMILE in Fig. 19 (d). Fig. 20 demonstrates the results of CausalGAN, with subfigures (a) (b) (c) (d) showing the intervention experiments on concepts CelebA (SMILE). We find that when we control SMILE, the mouth is changed, as shown in the second line of Fig. 20 (b). But we find sometimes the control of SMILE influence other unrelated concepts like GENDER (shown in first line of Fig. 20 (b)). In this concepts group, CausalGAN also shows relatively unstable intervention experiments compared to that of ours.

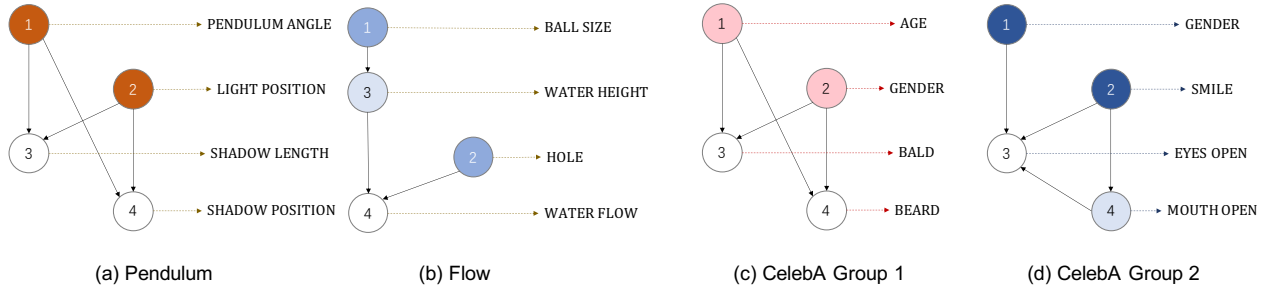


Figure 10. Causal graphs of three dataset. (a) shows the causal graph in pendulum dataset. The concepts are PENDULUM ANGLE, light POSITION, SHADOW POSITION and SHADOW LENGTH. (b) shows the causal graph in CelebA, on concepts AGE, GENDER and BEARD and BALD. (c) shows the causal graph in CelebA, on concepts GENDER, SMILE, EYES OPEN and MOUTH OPEN.

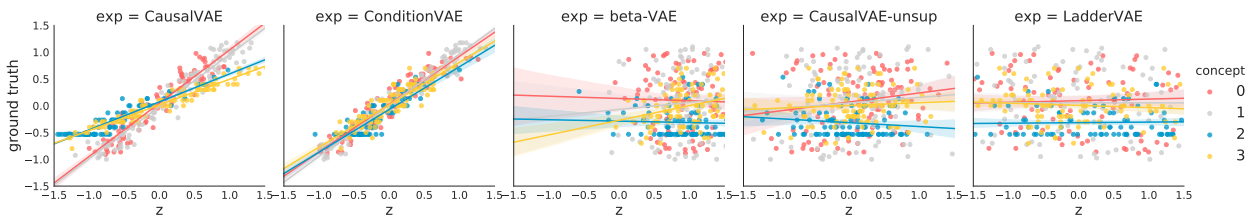


Figure 11. The figure shows the alignment of ground truth  $p(\mathbf{z}|\mathbf{u})$  and the learned latent factors  $q(\mathbf{z}|\mathbf{x}, \mathbf{u})$  on pendulum experiments. Although ConditionVAE is also the supervised method, our proposed CausalVAE shows a better performance.

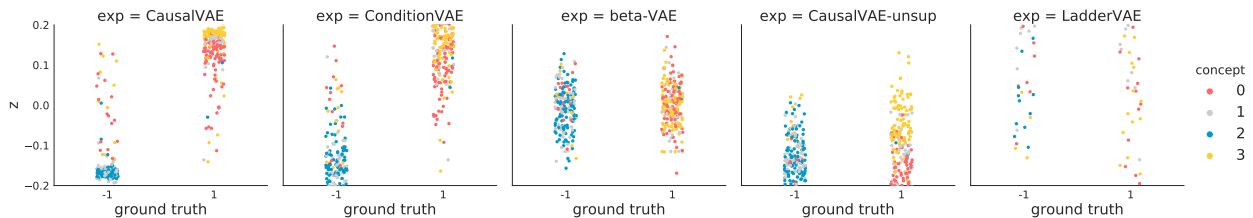


Figure 12. The figure shows the alignment of ground truth  $p(\mathbf{z}|\mathbf{u})$  and the learned latent factors  $q(\mathbf{z}|\mathbf{x}, \mathbf{u})$  on CelebA for the concepts (BEARD). The ground truth is a discrete distribution over  $\{-1, 1\}$ , and the color of the points indicates different dimensions. The factors learned by CausalVAE show the best alignment among all.

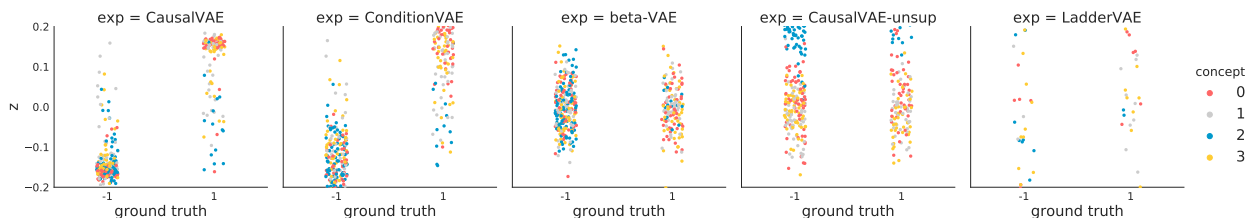


Figure 13. The figure shows the alignment between ground truth  $p(\mathbf{z}|\mathbf{u})$  and the learned latent factors  $q(\mathbf{z}|\mathbf{x}, \mathbf{u})$  on CelebA for 5 methods (CausalVAE, ConditionVAE,  $\beta$ -VAE, CausalVAE-unsup, LadderVAE from left to right). The ground truth is a distribution with mean taken from  $\{-1, 1\}$ , and the color of the points indicates different dimensions. The factors learned by CausalVAE show the best alignment among all. The concepts include: 1 GENDER; 2 SMILE; 3 EYES OPEN; 4 MOUTH OPEN.



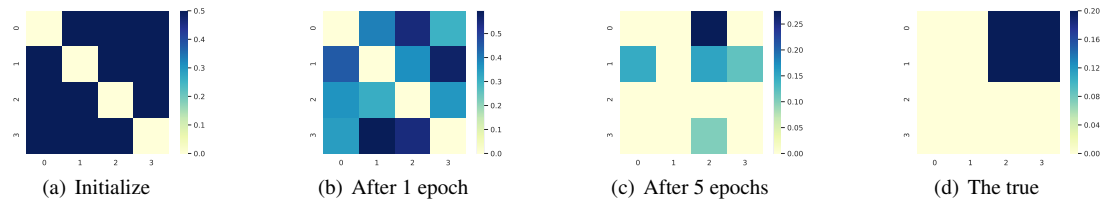


Figure 14. Learning process of causal graph  $A$  in CelebA (BEARD). The concepts include: 1 AGE; 2 GENDER; 3 BALD; 4 BEARD.

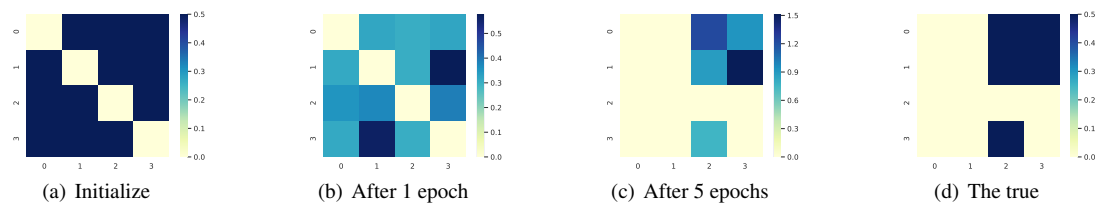


Figure 15. Learning process of causal graph  $A$  in CelebA (SMILE). The concepts include: 1 GENDER; 2 SMILE; 3 EYES OPEN; 4 MOUTH OPEN.

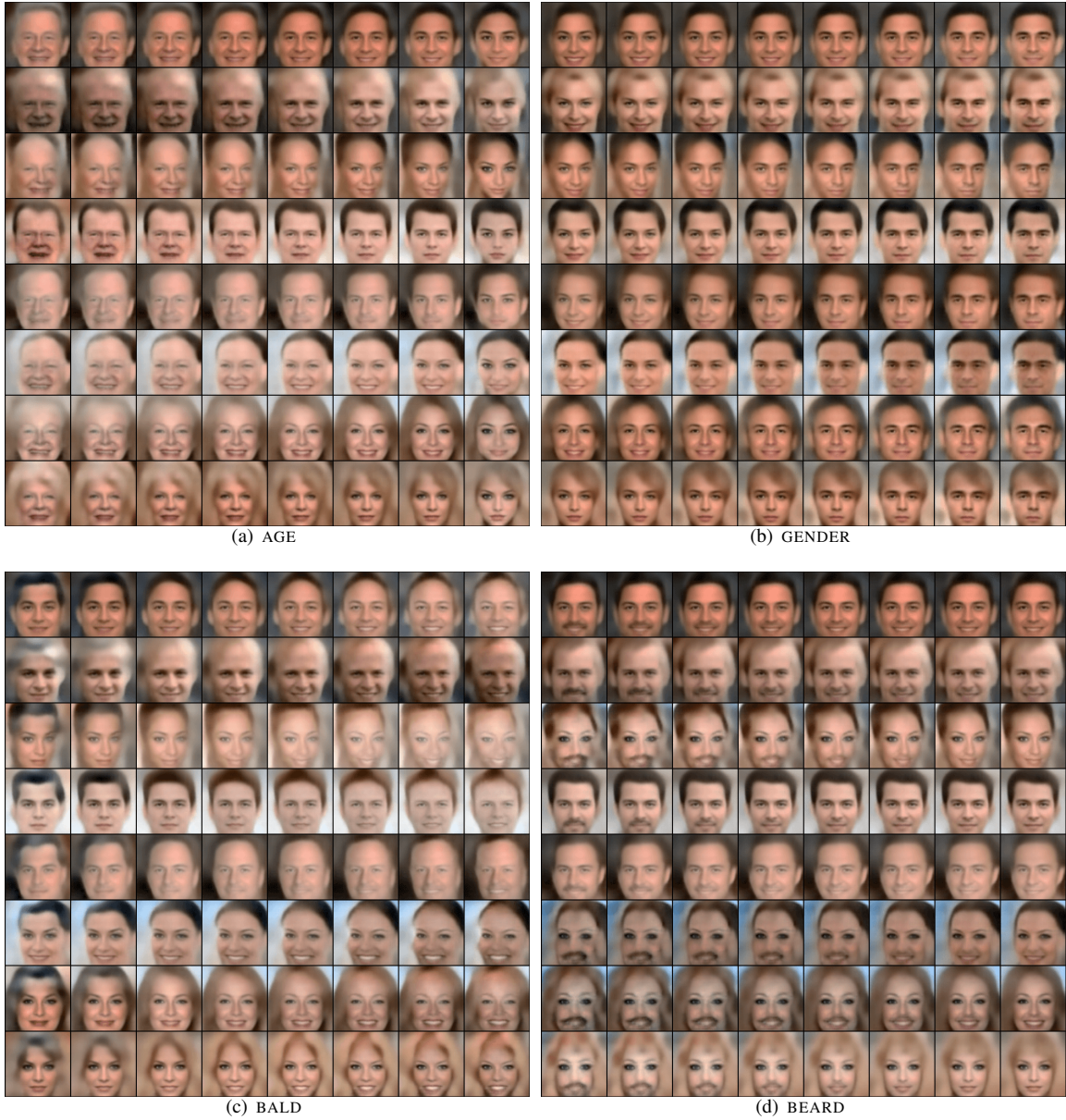


Figure 17. Results of CausalVAE model on CelebA (BEARD). The captions of the subfigures describe the controlled factors. From left to right, the pictures are results obtained by varying the value of the controlled factors.

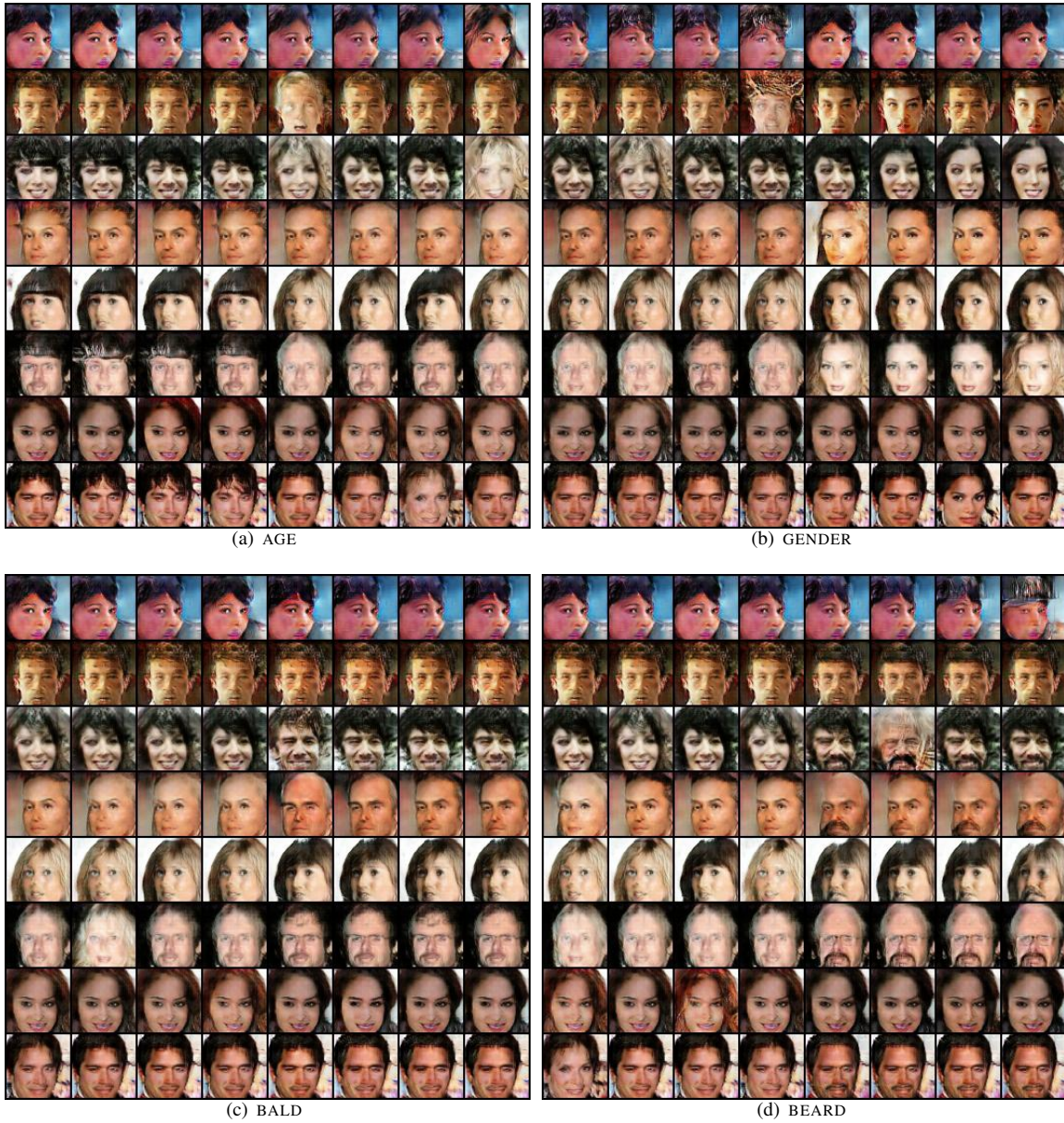


Figure 18. Results of CausalGAN [15] model on CelebA (BEARD). The captions of the subfigures describe the controlled factors. From left to right, the pictures are results obtained by varying the value of the controlled factors.

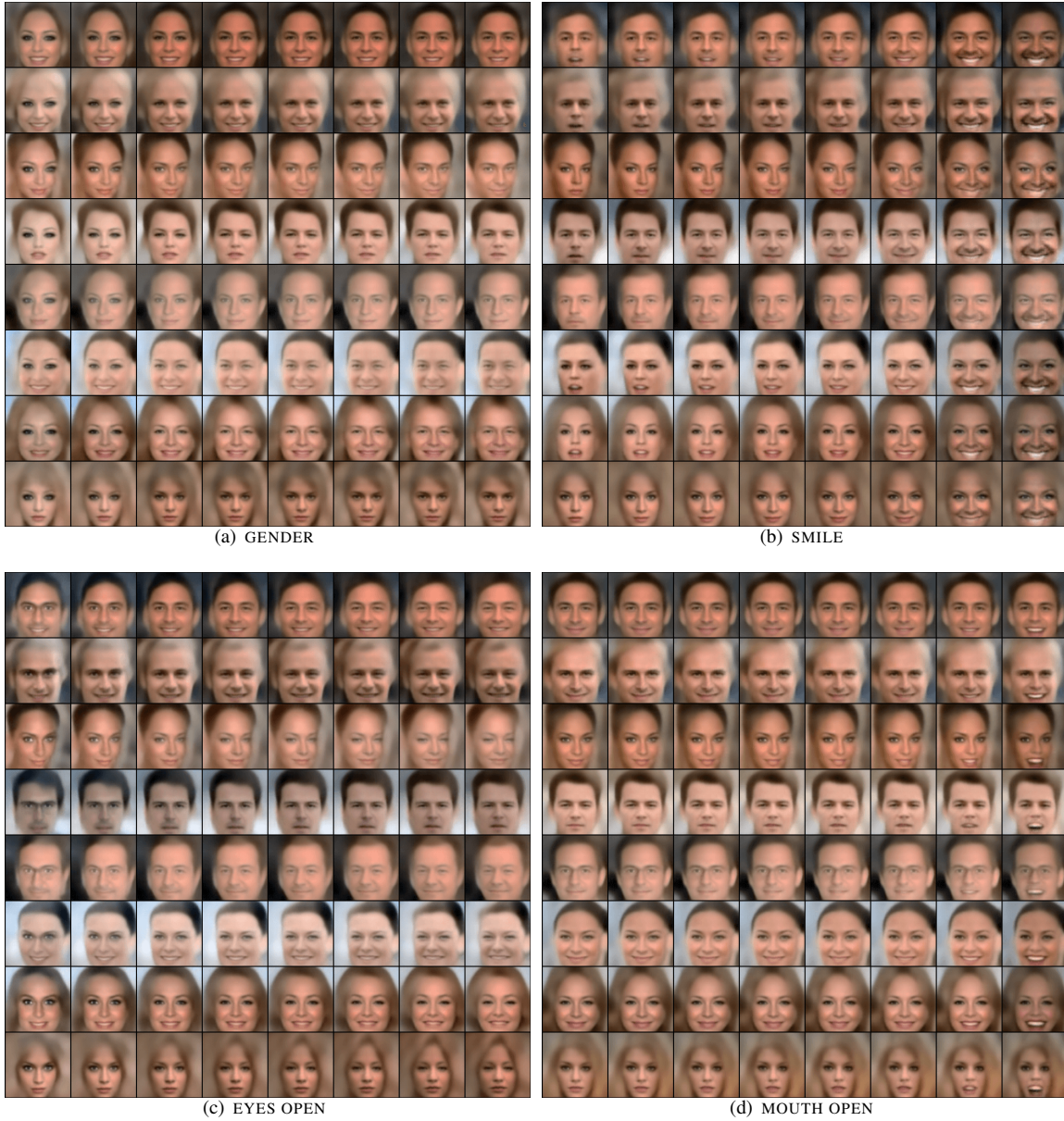


Figure 19. Results of CausalVAE model on CelebA (SMILE). The captions of the subfigures describe the controlled factors. From left to right, the pictures are results obtained by varying the value of the controlled factors.

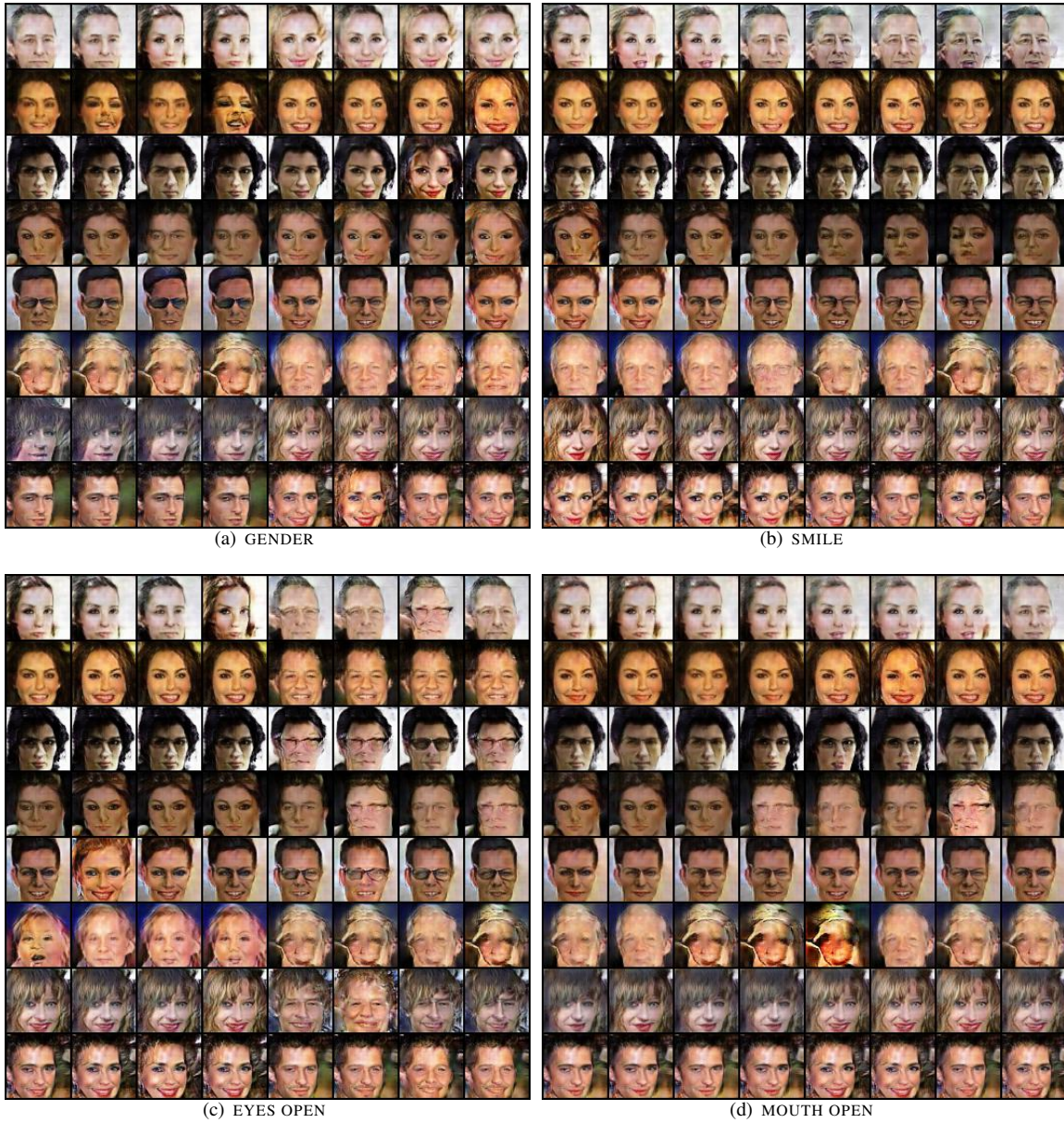


Figure 20. Results of CausalGAN model on CelebA (SMILE). The captions of the subfigures describe the controlled factors. From left to right, the pictures are results obtained by varying the value of the controlled factors.