

# WEAKLY SUPERVISED CAUSAL REPRESENTATION LEARNING

Johann Brehmer<sup>\*</sup>,<sup>1</sup> Pim de Haan<sup>\*</sup>,<sup>1,2</sup> Phillip Lippe,<sup>2</sup> and Taco Cohen<sup>1</sup>

<sup>1</sup>Qualcomm AI Research<sup>†</sup>      <sup>2</sup>QUVA Lab, University of Amsterdam  
 {jbrehmer, pim, tacos}@qti.qualcomm.com; p.lippe@uva.nl

## ABSTRACT

Learning high-level causal representations together with a causal model from unstructured low-level data such as pixels is impossible from observational data alone. We prove under mild assumptions that this representation is identifiable in a weakly supervised setting. This requires a dataset with paired samples before and after random, unknown interventions, but no further labels. Finally, we show that we can infer the representation and causal graph reliably in a simple synthetic domain using a variational autoencoder with a structural causal model as prior.

## 1 INTRODUCTION

The dynamics of many systems can be described in terms of some high-level variables and causal relations between them. Often, these causal variables are not known but only observed in some unstructured, low-level representation, such as the pixels of a camera feed. Learning the causal representations together with the causal structure between them is a challenging problem and important for instance for applications in robotics and autonomous driving (Schölkopf et al., 2021). Without prior assumptions on the data-generating process or supervision, it is impossible to identify the causal variables and their causal structure uniquely (Eberhardt, 2016; Locatello et al., 2019).

In this work, we show that a weak form of supervision is sufficient to identify both the causal representations and the structural causal model between them. We consider a setting in which we have access to data pairs, representing the system before and after a randomly chosen unknown intervention. Neither labels on the intervention targets nor active control of the interventions are necessary for our identifiability theorem, making this setting useful for offline learning. We prove that with this form of weak supervision, latent causal models (LCMs)—structural causal models (SCMs) together with a decoder from the causal factors to the data space—are identifiable up to a relabelling and elementwise reparameterizations of the causal variables.

We then introduce a practical implementation for LCM inference by using an SCM as a prior in a variational autoencoder (VAE). In a range of experiments, we show that LCMs can learn the true causal variables and the causal structure from unstructured data.

**Related work** Our work builds on the approach of Locatello et al. (2020) to *disentangled representation learning*. The authors introduce a similar weakly supervised setting where observations are collected before and after unknown interventions. In contrast to our work, however, they focus on disentangled representations, i. e. independent factors of variation with a trivial causal graph, which our work subsumes as a special case. Other relevant works on disentangled representation learning

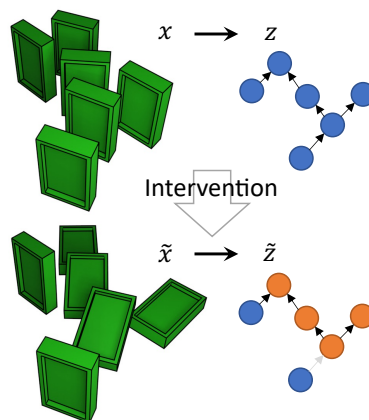


Figure 1: We learn to represent pixels  $x$  as causal variables  $z$ . The bottom shows the effect of intervening on the orange variable. We prove that variables and causal model can be identified from samples  $(x, \tilde{x})$ .

<sup>\*</sup>Equal contribution

<sup>†</sup>Qualcomm AI Research is an initiative of Qualcomm Technologies, Inc.

and (nonlinear) independent component analysis include Hyvärinen & Oja (2000); Shu et al. (2020); Khemakhem et al. (2020); Hälvä et al. (2021) and Lachapelle et al. (2022).

The problem of *causal representation learning* has been gaining attention lately, for a recent review see Schölkopf et al. (2021). Lu et al. (2021) learn causal representations by observing similar causal models in different environments. Von Kügelgen et al. (2021) use the weakly supervised setting to study self-supervised learning, using a known but non-trivial causal graph between content and style factors. Lippe et al. (2022b) learn causal representations from time-series data from labelled interventions, assuming that causal effects are not instantaneous but can be temporally resolved. Yang et al. (2021) propose to train a VAE with an SCM prior, but require the true causal variables as labels. To the best of our knowledge, our work is the first to provide identifiability guarantees for arbitrary, unknown causal graphs in this weakly supervised setting.

## 2 IDENTIFIABILITY OF LATENT CAUSAL MODELS FROM WEAK SUPERVISION

In this section, we show theoretically that causal variables and causal mechanisms are identifiable from weak supervision. In Sec. 3 we will then demonstrate how we can learn causal models in practice by training a causally structured VAE.

**Setup** We begin by defining latent causal models and the weakly supervised setting. Here, we only provide informal definitions and assume familiarity with common concepts from causality; see Appendix A.1 for a complete and precise treatment. In Appendix A.3, we discuss limitations of the setup and possible generalizations.

We describe the causal structure between latent variables as an Structural Causal Model (SCM). An SCM  $\mathcal{C}$  describes the relation between causal variables  $Z_1, \dots, Z_n$  with domains  $Z_i$  and noise variables  $\epsilon_1, \dots, \epsilon_n$  with domains  $E_i$  along a directed acyclic graph  $G$ . Causal mechanisms  $f_i : E_i \times \prod_{j \in \text{pa}_i} Z_j \rightarrow Z_i$  describe how the value of a causal variable is determined from the associated noise variables as well as the values of its parents in the graph. Finally, an SCM includes a probability measure for the noise variables.

An SCM entails a unique solution  $s : E \rightarrow Z$  defined by successively applying the causal mechanisms. We require the structural equations to be pointwise diffeomorphic,  $s$  is thus also diffeomorphic. It also entails an observational distribution  $p_{\mathcal{C}}(z)$  (Markov with respect to the graph of the SCM), which is the pushforward of  $p_{\mathcal{E}}$  through the solution.

A perfect intervention  $(I; (f_i)_{i \in I})$  modifies an SCM by replacing for a subset of the causal variables, called the intervention target set  $I = \{1, \dots, n\}$ , the causal mechanism  $f_i$  with a new mechanism  $f_i : E_i \rightarrow Z_i$ , which does not depend on the parents. The intervened SCM has a new solution  $s_I : E \rightarrow Z$ . We define interventions to be atomic if the number of targeted variables is one or zero.

We will reason about generative models in a data space  $X$ , in which the causal structure is latent. Also including a distribution of interventions, we define LCMs:

**Definition 1** (Latent causal model (LCM)). A latent causal model  $\mathcal{M} = \langle \mathcal{C}; X; g; I; p_I \rangle$  consists of

- an acyclic SCM  $\mathcal{C}$  that admits a faithful distribution,
- an observation space  $X$ ,
- a decoder  $g : Z \rightarrow X$  that is diffeomorphic onto its image,
- a set  $I$  of interventions on  $\mathcal{C}$ , and
- a probability measure  $p_I$  over  $I$ .

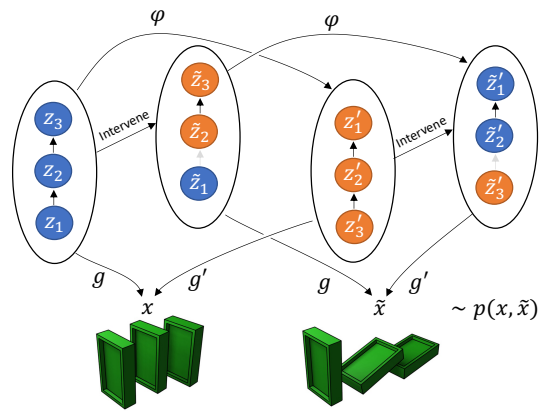


Figure 2: In LCM  $\mathcal{M}$ ,  $z_i$  denotes whether the  $i$ -th stone from the front is standing. Intervening on the second variable leads to  $\tilde{z}$ . The decoder  $g$  renders  $z, \tilde{z}$  as pixels  $x, \tilde{x}$ . LCM  $\mathcal{M}'$  has an equivalent representation in which  $z'_i$  denotes whether the  $i$ -th stone from the back has fallen. In Theorem 1, we prove that if and only if two causal models have the same pixel distribution  $p(x, \tilde{x})$ , there exists an LCM isomorphism  $\varphi$ : an element-wise reparametrization of the causal variables plus a permutation of the ordering that commutes with interventions and causal mechanisms.

We define two LCMs as equivalent if all of their components are equal up to a permutation of the causal variables and elementwise diffeomorphic reparameterizations of each variable, see Fig. 2.

**Definition 2** (LCM isomorphism (informal)). *Let  $\mathcal{M} = \langle h; C; X; g; l; p_{\mathcal{I}} \rangle$  and  $\mathcal{M}' = \langle h'; C'; X'; g'; l'; p'_{\mathcal{I}'} \rangle$  be two LCMs with identical observation space. An LCM isomorphism between them is a graph isomorphism  $\gamma : G(C) \rightarrow G(C')$  together with elementwise diffeomorphisms for noise and causal variables that tell us how to reparameterize them, such that the structure functions, noise distributions, decoder, intervention set, and intervention distribution of  $\mathcal{M}'$  are compatible with the corresponding elements of  $\mathcal{M}$  reparameterized through the graph isomorphism and elementwise diffeomorphisms.  $\mathcal{M}$  and  $\mathcal{M}'$  are equivalent,  $\mathcal{M} \sim \mathcal{M}'$ , if and only if there is an LCM isomorphism between them.*

Following Locatello et al. (2020), we define a generative process of pre and post interventional data:<sup>1</sup>

**Definition 3** (Weakly supervised generative process). *Consider an LCM  $\mathcal{M}$  where the underlying SCM has continuous noise spaces  $E_i$ , independent probabilities  $p_{E_i}$ , and admits a solution  $s$ . We define the weakly supervised generative process of data pairs  $(x; \mathbf{x}) \sim p_{\mathcal{M}}^{\mathbf{x}}(x; \mathbf{x})$  as follows:*

$$\begin{aligned} p_{E_i} &: & z &= s(\cdot); & x &= g(z); \\ l \sim p_{\mathcal{I}} &: & \mathbf{z} &= s_l(\cdot); & \mathbf{x} &= g(\mathbf{z}); \end{aligned} \quad (1)$$

Here we parameterize stochastic interventions on  $l$  by  $\sim_{-i} \sim p_{E_i}$  for  $i \in l$  and  $\sim_{-i} = \sim_i$  for  $i \notin l$ .

**Identifiability result** The main theoretical result of this paper is that an LCM  $\mathcal{M}$  can be identified from  $p(x; \mathbf{x})$  up to a relabeling and elementwise transformations of the causal variables:

**Theorem 1** (Identifiability of  $\mathbb{R}$ -valued LCMs from weak supervision). *Let  $\mathcal{M} = \langle h; C; X; g; l; p_{\mathcal{I}} \rangle$  and  $\mathcal{M}' = \langle h'; C'; X'; g'; l'; p'_{\mathcal{I}'} \rangle$  be LCMs with the following properties:*

- The LCMs have an identical observation space  $X$ .
- The SCMs  $C$  and  $C'$  both consist of  $n$  real-valued endogeneous causal variables and corresponding exogenous noise variables, i. e.  $E_i = Z_i = Z'_i = E'_i = \mathbb{R}$ .
- The intervention sets  $l$  and  $l'$  consist of all atomic, perfect interventions,  $l = \{f; \dots; f; z_0; g; \dots; f; z_n; g\}$  and similar for  $l'$ .
- The intervention distribution  $p_{\mathcal{I}}$  and  $p'_{\mathcal{I}'}$  have full support.

Then the following two statements are equivalent:

1. The LCMs entail equal weakly supervised distributions,  $p_{\mathcal{M}}^{\mathbf{x}}(x; \mathbf{x}) = p_{\mathcal{M}'}^{\mathbf{x}}(x; \mathbf{x})$ .
2. The LCMs are equivalent,  $\mathcal{M} \sim \mathcal{M}'$ .

Let us summarize the key steps of our proof, which we provide in its entirety in Appendix A.2. The direction 2  $\Rightarrow$  1 follows from the definition of equivalence. The direction 1  $\Rightarrow$  2 is proven constructively along the following steps:

1. We begin by defining a diffeomorphism  $\gamma = g'^{-1} \circ g : Z \rightarrow Z'$  and note that if  $z; \mathbf{z} \sim p_C^{\mathbf{z}}(z; \mathbf{z})$ , the weakly supervised distribution of causal variables of model  $C$ , then  $\gamma(z); \gamma(\mathbf{z}) \sim p_{C'}^{\mathbf{z}'}(\gamma(z); \gamma(\mathbf{z}))$ . The distribution over  $z; \mathbf{z}$  is a mixture, where each intervention  $l = f; i; g$  gives a mixture component; each component is supported on different a  $(n + 1)$ -dimensional submanifold. Therefore there exists a bijection between the components  $\gamma : [n] \rightarrow [n]$  that maps intervention targets  $l$  in  $\mathcal{M}$  to intervention targets  $l' = \gamma(l)$  in  $\mathcal{M}'$ . Furthermore, because the joint distribution  $z; \mathbf{z}$  is preserved by  $\gamma$ , first mapping with  $\gamma$ , then intervening,  $z; \mathbf{z} \sim_{l'} z'; \mathbf{z}' \sim_{l'} \mathcal{E}'$ , equals  $z; \mathbf{z} \sim_{\mathcal{E}} z'; \mathbf{z}' \sim_{\mathcal{E}'}$ .
2. Because  $l = f; i; g$  is a hard intervention, for the order  $z; \mathbf{z} \sim_{l'} z'; \mathbf{z}' \sim_{l'} \mathcal{E}'$ ,  $z'_{i_0}$  is independent of  $z'$ . Thus in both orders,  $z'_{i_0}$  is independent of  $z$ . This means that for the path through  $\mathcal{E}$ , the intervention sample  $z_i$  is transformed into  $z'_{i_0}$  independently of  $z$ . For  $\mathbb{R}$ -valued variables, this statistical independence implies that the transformation is constant in  $z$ , and thus  $\gamma(z)_{i_0}$  is constant in  $z_j$  for  $j \neq i$ .  $\gamma$  is therefore an elementwise reparameterization.
3. Using this, it is easy to show that  $\gamma$  is a causal graph isomorphism and that it is compatible with the causal mechanisms. This proves LCM equivalence  $\mathcal{M} \sim \mathcal{M}'$ .

<sup>1</sup>This construction is closely related to twinned SCMs (Bongers et al., 2021, Def. 2.17), typically used to compute counterfactual queries  $p(z_{\setminus i} | z_i, \mathbf{z}_i)$ . We instead focus on the joint distribution of pre-intervention and post-intervention data.

### 3 EXPERIMENTS

The identifiability result in Sec. 2 suggests that one can identify an LCM by finding the maximum-likelihood solution of  $p(x; \mathcal{X})$ . We now investigate whether we can indeed learn LCMs in practice. The datasets, models, and training setup are described in more detail in Appendix B.

**Relaxation of the weakly supervised setting** In realistic systems, the “before” and “after” state of a system are recorded some time apart. We do not expect that the causal variables unaffected by an intervention remain *perfectly* invariant during this time. Similarly, we want to allow for causal descendants of the intervention targets to slightly change from the exact value predicted by the weakly supervised process in Eq. (1). We model this by applying a Gaussian convolution with small variance  $\sigma^2$  to  $p(z; z|I)$  on all dimensions of  $z$  that are not intervened on, resulting in a continuous density. This relaxation means that there is a gap between the requirements of our identifiability theorem and the experimental setting. In particular, relaxing the exact manifold in the weakly supervised data space to a “fuzzy” one renders our argument for the identifiability of noise encodings and intervention targets invalid; see [Brehmer & Cranmer \(2020\)](#) for a related discussion. Nevertheless, we empirically find that we can still reliably identify LCMs in this more realistic setup.

**LCM implementation** In practice, we implement LCMs as VAEs ([Kingma & Welling, 2014](#)). In the simplest version, the causal graph  $G$  is fixed. The SCM and interventional distributions define the relaxed prior  $p(z; z|I)$ . Rather than with a diffeomorphism  $g$ , we map causal latents to observed variables through a stochastic encoder  $q(z|x)$  and decoder  $p(x|z)$ . We then maximize the following lower bound on the weakly supervised log likelihood:

$$\log p(x; \mathcal{X}) \approx \mathbb{E}_{\substack{z \sim q(z|x) \\ \tilde{z} \sim q(\tilde{z}|\tilde{x})}} \log \mathbb{E}_{I \sim p_I(I)} [p(z; z|I)] + \log p(x|z) + \log p(x|\tilde{z}) - \log q(z|x) - \log q(\tilde{z}|\tilde{x}) \quad (2)$$

For our experiments with small graphs, we instantiate one such LCM per graph structure, train them on the VAE loss in Eq. (2), and select the model with the lowest validation loss. To scale to larger graphs, it will be beneficial to learn a belief over graphs and to sample adjacency matrices and intervention targets from suitable priors, for instance using Gumbel-Softmax methods ([Jang et al., 2017](#); [Brouillard et al., 2020](#)) or other approaches ([Lippe et al., 2022a](#)); we will explore this in future work.

Following common practice in causal discovery ([Zheng et al., 2018](#); [Brouillard et al., 2020](#); [Lippe](#)

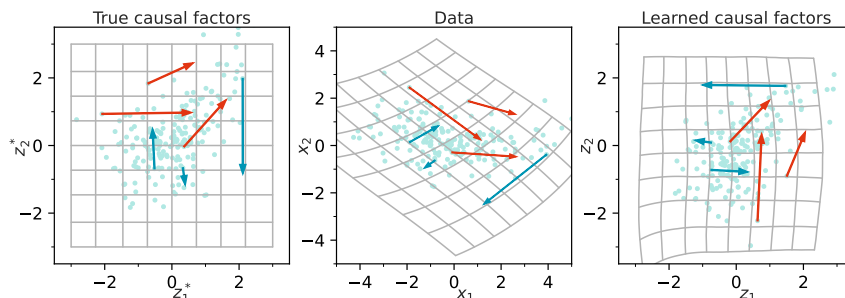


Figure 3: 2D toy data with graph  $z_2 \rightarrow z_1$ . The grey grids show the map between true causal factors, data, and latent causal factors learned by the LCM, with graph  $z_2 \rightarrow z_1$ . The mint dots indicate the observational data distribution, the arrows from  $z$  to  $\tilde{z}$  show interventions targeting  $z_1^*$  (red) or  $z_2^*$  (blue). The fact that axis-aligned lines in the true latent space are mapped to axis-aligned lines in the learned latent space implies that the disentanglement succeeded.

Table 1: Experiment results, comparing LCMs to unstructured  $\beta$ -VAEs and disentanglement VAEs (dVAE), best results in bold. For Causal3DIdent, we average over six datasets with different graphs, see Appendix B. We measure disentanglement with the DCI score  $D$ .  $D = 1$  implies disentanglement succeeds only with our method. LCMs learn the correct causal graph most of the time, visible as a structural Hamming distance SHD = 0. The quality of intervention inference is evaluated with the intervention negative log posterior ( $-\log p_I$ ).

Method	$D$	SHD	$-\log p_I$
<i>2D toy data</i>			
LCM	<b>0.99</b>	<b>0.0</b>	<b>0.28</b>
dVAE	0.35	n/a	0.33
$\beta$ -VAE	0.00	n/a	n/a
<i>Causal3DIdent</i>			
LCM	<b>0.98</b>	<b>0.17</b>	<b>0.20</b>
dVAE	0.57	n/a	1.98
$\beta$ -VAE	0.38	n/a	n/a

et al., 2022a), we incentivize learning the sparsest graph compatible with the data distribution by adding a regularization term proportional to the number of edges in the graph to the loss.

**Baselines** We compare LCMs to an unstructured  $\beta$ -VAE that treats  $X$  and  $X$  as i. i. d. and uses a standard Gaussian prior. We also compare to a disentanglement VAE, which models our weakly supervised process for a trivial causal graph (i. e. independent factors of variation), similar to the method proposed by Locatello et al. (2020).

**2D toy experiment** We first test LCMs in a toy experiment with  $X = Z = \mathbb{R}^2$ . Training data is generated from a nonlinear SCM with the graph  $Z_1 \dashv Z_2$  and mapped to the data space through a randomly initialized normalizing flow.

An LCM trained in the weakly supervised setting is able to reconstruct the causal factors and the causal graph accurately up to a permutation of the two variables, as shown in the Fig. 3. It fits the weakly supervised data distribution with a better log likelihood than the acausal baselines. In Tbl. 1 we quantify the quality of the learned representations with the DCI disentanglement score  $D$  (Eastwood & Williams, 2018). We find that our LCM is able to disentangle the causal factors almost perfectly, while the baselines, which assume independent factors of variation, fail as expected. Finally, we test whether the learned LCM correctly infers the interventions by computing the intervention posterior  $p(I|X; \mathbf{x})$  and evaluating it for the true intervention  $I^*$ , again finding better results for the causal LCM than for the baselines.

**Causal3DIdent** Next, we test LCMs on an adaptation of the Causal3DIdent dataset (von Kugelgen et al., 2021), which contains images of three-dimensional objects under variable positions and lighting conditions. We consider three latent variables representing object hue, the spotlight hue, and the position of the spotlight. We consider six versions of this dataset, each with a different causal graph, randomly initialized nonlinear structure functions, and heteroskedastic noise. These are mapped to images with a resolution of  $64 \times 64$ , see Fig. 4 for examples.

An LCM with convolutional encoder and decoder is again able to learn the causal variables and the causal structure between them, finding the correct causal graph in all but one settings<sup>2</sup>. The results in Tbl. 1 and Fig. 4 show that the learned representations are more disentangled than those learned by methods that do not account for causal structure and that the LCM can infer interventions more reliably than the baselines.

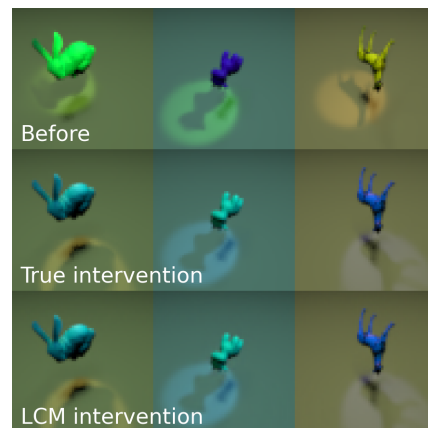


Figure 4: Causal3DIdent samples before (top row) and after (middle row) interventions, and post-intervention samples generated from the LCM under the intervention inferred from the data, indicating we correctly learned to intervene.

## 4 DISCUSSION

We have presented a method for causal representation learning in a weakly supervised setting, in which data consist of a system before and after a random intervention. We have proven that in this setting both the causal variables and the causal structure between them is identifiable up to permutations and elementwise reparameterizations of the causal variables. This extends the results by Locatello et al. (2020) from independent factors of variation (trivial causal graphs) to arbitrary causal structures. Our identifiability result relies on a few key assumptions, including that interventions are perfect, that all atomic interventions may be observed, and that the causal variables are real-valued. We discuss these requirements and their potential relaxation in Appendix A.3.

In practice, LCMs can be implemented in a variational autoencoder. We demonstrated in first experiments that this lets us reliably learn causal variables and nonlinear causal structure from unstructured pixel data.

<sup>2</sup>In the one configuration where our algorithm did not learn the correct graph, it instead learned a supergraph of the true graph, with one additional edge. We attribute this to insufficient regularization.

**Acknowledgments** We want to thank Dominik Neuenfeld and Frank Rösler for useful discussions and Gabriele Cesa, Yang Yang, and Yunfan Zhang for helping with our experiments.

## REFERENCES

- Blender Online Community. *Blender - a 3D modelling and rendering package*. Blender Foundation, Stichting Blender Foundation, Amsterdam, 2021. <http://www.blender.org>.
- Stephan Bongers, Patrick Forré, Jonas Peters, and Joris M. Mooij. Foundations of Structural Causal Models with Cycles and Latent Variables. *Annals of Statistics*, 49(5):2885–2915, 2021. doi: 10.1214/21-AOS2064.
- Johann Brehmer and Kyle Cranmer. Flows for simultaneous manifold learning and density estimation. In *Advances in Neural Information Processing Systems*, volume 33, pp. 442–453, 2020.
- Philippe Brouillard, Sébastien Lachapelle, Alexandre Lacoste, Simon Lacoste-Julien, and Alexandre Drouin. Differentiable Causal Discovery from Interventional Data. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020.
- Yuri Burda, Roger B. Grosse, and Ruslan Salakhutdinov. Importance weighted autoencoders. In Yoshua Bengio and Yann LeCun (eds.), *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*, 2016.
- Keenan Crane. Keenan’s 3D Model Repository. <https://www.cs.cmu.edu/~kmcraane/Projects/ModelRepository/>, 2021.
- Brian Curless and Marc Levoy. A Volumetric Method for Building Complex Models from Range Images. In *Proceedings of the 23rd Annual Conference on Computer Graphics and Interactive Techniques, SIGGRAPH ’96*, pp. 303–312, New York, NY, USA, 1996. Association for Computing Machinery. ISBN 0897917464. doi: 10.1145/237170.237269.
- Cian Eastwood and Christopher K I Williams. A framework for the quantitative evaluation of disentangled representations. *International Conference on Learning Representations*, February 2018.
- Frederick Eberhardt. Green and grue causal variables. *Synthese*, 193(4):1029–1046, 2016.
- Brendan Fong. Causal theories: A categorical perspective on bayesian networks. *arXiv preprint arXiv:1301.6201*, 2013.
- Tobias Fritz. A synthetic approach to markov kernels, conditional independence and theorems on sufficient statistics. *Advances in Mathematics*, 370:107239, 2020. ISSN 0001-8708. doi: <https://doi.org/10.1016/j.aim.2020.107239>. URL <https://www.sciencedirect.com/science/article/pii/S0001870820302656>.
- Hermanni Hälvä, Sylvain Le Corff, Luc Lehéricy, Jonathan So, Yongjie Zhu, Elisabeth Gassiat, and Aapo Hyvarinen. Disentangling Identifiable Features from Noisy Data with Structured Nonlinear ICA. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, 2021.
- Aapo Hyvärinen and Erkki Oja. Independent component analysis: algorithms and applications. *Neural Networks*, 13:411–430, 2000.
- Eric Jang, Shixiang Gu, and Ben Poole. Categorical Reparameterization with Gumbel-Softmax. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*, 2017.
- Ilyes Khemakhem, Diederik Kingma, Ricardo Monti, and Aapo Hyvarinen. Variational Autoencoders and Nonlinear ICA: A Unifying Framework. In Silvia Chiappa and Roberto Calandra (eds.), *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, volume 108 of *Proceedings of Machine Learning Research*, pp. 2207–2217. PMLR, 2020.

- Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. In Yoshua Bengio and Yann LeCun (eds.), *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.
- Diederik P. Kingma and Max Welling. Auto-Encoding Variational Bayes. In *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*, 2014.
- Venkat Krishnamurthy and Marc Levoy. Fitting Smooth Surfaces to Dense Polygon Meshes. In *Proceedings of the 23rd Annual Conference on Computer Graphics and Interactive Techniques, SIGGRAPH '96*, pp. 313–324, New York, NY, USA, 1996. Association for Computing Machinery. ISBN 0897917464. doi: 10.1145/237170.237270.
- Sebastien Lachapelle, Pau Rodriguez, Rémi Le, Yash Sharma, Katie E Everett, Alexandre Lacoste, and Simon Lacoste-Julien. Disentanglement via Mechanism Sparsity Regularization: A New Principle for Nonlinear ICA. In *First Conference on Causal Learning and Reasoning*, 2022.
- Phillip Lippe, Taco Cohen, and Efstratios Gavves. Efficient Neural Causal Discovery without Acyclicity Constraints. In *International Conference on Learning Representations*, 2022a.
- Phillip Lippe, Sara Magliacane, Sindy Löwe, Yuki M. Asano, Taco Cohen, and Efstratios Gavves. CIT-RIS: Causal Identifiability from Temporal Intervened Sequences. *arXiv preprint arXiv:2202.03169*, 2022b.
- 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 Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 4114–4124. PMLR, 2019.
- Francesco Locatello, Ben Poole, Gunnar Rätsch, Bernhard Schölkopf, Olivier Bachem, and Michael Tschannen. Weakly-supervised disentanglement without compromises. In *International Conference on Machine Learning*, pp. 6348–6359. PMLR, 2020.
- Chaochao Lu, Yuhuai Wu, José Miguel Hernández-Lobato, and Bernhard Schölkopf. Nonlinear invariant risk minimization: A causal approach. *arXiv preprint arXiv:2102.12353*, 2021.
- Martin Edward Newell. *The Utilization of Procedure Models in Digital Image Synthesis*. PhD thesis, The University of Utah, 1975. AAI7529894.
- Judea Pearl. *Causality : models, reasoning, and inference*. Cambridge University Press, Cambridge, U.K. New York, 2000. ISBN 978-0521895606.
- Emil Praun, Adam Finkelstein, and Hugues Hoppe. Lapped Textures. In *Proceedings of ACM SIGGRAPH 2000*, pp. 465–470, July 2000.
- Emily Riehl. *Category Theory in Context*. Aurora: Dover Modern Math Originals. Dover Publications, 2017. ISBN 9780486820804.
- Bernhard Schölkopf, Francesco Locatello, Stefan Bauer, Nan Rosemary Ke, Nal Kalchbrenner, Anirudh Goyal, and Yoshua Bengio. Toward causal representation learning. *Proceedings of the IEEE*, 109(5):612–634, 2021.
- Rui Shu, Yining Chen, Abhishek Kumar, Stefano Ermon, and Ben Poole. Weakly Supervised Disentanglement with Guarantees. In *International Conference on Learning Representations*, 2020.
- Stack Exchange user zhw. Lebesgue measure-preserving differentiable function. Mathematics Stack Exchange, 2016. URL <https://math.stackexchange.com/q/1755585>. (version: 2016-04-23).
- Greg Turk and Marc Levoy. Zippered Polygon Meshes from Range Images. In *Proceedings of the 21st Annual Conference on Computer Graphics and Interactive Techniques, SIGGRAPH '94*, pp. 311–318, New York, NY, USA, 1994. Association for Computing Machinery. ISBN 0897916670. doi: 10.1145/192161.192241.

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. In *Thirty-Fifth Conference on Neural Information Processing Systems*, 2021.

Mengyue Yang, Furui Liu, Zhitang Chen, Xinwei Shen, Jianye Hao, and Jun Wang. CausalVAE: Disentangled representation learning via neural structural causal models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9593–9602, 2021.

Xun Zheng, Bryon Aragam, Pradeep Ravikumar, and Eric P. Xing. DAGs with NO TEARS: Continuous Optimization for Structure Learning. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pp. 9492–9503, 2018.

## A IDENTIFIABILITY RESULT

### A.1 DEFINITIONS

Here we define objects and relations that were not formally defined in the main body of the paper, but are necessary to make Thm. 1 precise and to prove it.

We use the following notation:

- $[n] = \{1, \dots, n\}$
- $\text{pa}_i^C \subseteq [n]$  the set of parent nodes of node  $i$  in graph  $G(C)$ .
- $\text{desc}_i^C \subseteq [n]$  the set of descendant nodes of node  $i$  in graph  $G(C)$ , excluding  $i$  itself.
- $\text{anc}_i^C \subseteq [n]$  the set of ancestor nodes of node  $i$  in graph  $G(C)$ , excluding  $i$  itself.
- $\text{nonanc}_i^C = [n] \setminus (\text{anc}_i^C \cup \{i\})$  the set of non-ancestor nodes of node  $i$  in graph  $G(C)$ , excluding  $i$  itself.
- Given measure  $\rho$  on space  $A$  and measurable function  $f : A \rightarrow B$ ,  $f_*\rho$  is the push-forward measure on  $B$ .

We describe causal structure with SCMs.

**Definition 4** (Structural causal model (SCM)). *An SCM is a tuple  $C = (Z; E; F; p_E)$  consisting of the following:*

- domains  $Z = \{Z_1, \dots, Z_n\}$  of causal (endogenous) variables  $Z_1, \dots, Z_n$ ;
- domains  $E = \{E_1, \dots, E_n\}$  of noise (exogenous) variables  $E_1, \dots, E_n$ ;
- a directed acyclic graph  $G(C)$ , whose nodes are the causal variables and edges represent causal relations between the variables;
- causal mechanisms  $F = \{f_1, \dots, f_n\}$  with  $f_i : E_i \times \prod_{j \in \text{pa}_i} Z_j \rightarrow Z_i$ ; and
- a probability measure  $p_E(\cdot) = p_{E_1}(E_1) p_{E_2}(E_2) \dots p_{E_n}(E_n)$  with full support that admits a continuous density.

Additionally, we assume that  $\delta_i : \mathbb{R}^{\text{pa}_i} \times \mathbb{R}^{\text{pa}_i} \rightarrow E_i \times Z_i$  is a diffeomorphism.

We will need to reason about vectors being “equal up to permutation and elementwise reparameterizations”. We formalize this in the following definition:

**Definition 5** ( $\mathfrak{S}_n$ -diagonal). *Let  $\sigma : [n] \rightarrow [n]$  be a bijection (that is, a permutation). Let  $\tau : \prod_{i=1}^n X_i \rightarrow \prod_{i=1}^n Y_i$  be a function between product spaces. Then  $\tau$  is  $\sigma$ -diagonal if there exist functions, called components,  $\tau_i : X_i \rightarrow Y_{\sigma(i)}$  such that  $\delta_i : \mathbb{R}^{\text{pa}_i} \times \mathbb{R}^{\text{pa}_i} \rightarrow E_i \times Z_i$  is a diffeomorphism.*

This lets us define isomorphisms between SCMs:

**Definition 6** (Isomorphism of SCMs). *Let  $C = (Z; E; F; p_E)$  and  $C' = (Z'; E'; F'; p_{E'})$  be SCMs. An isomorphism  $\tau : C \rightarrow C'$  consists of*

1. a graph isomorphism  $\sigma : G(C) \rightarrow G(C')$  that tells us how to identify corresponding variables in the two models and which preserves parents:  $\text{pa}_{\sigma(i)}^{C'} = \sigma(\text{pa}_i^C)$  and



2. *-diagonal diffeomorphisms for noise and endogenous variables that tell us how to reparameterize them*  $\iota_{\mathcal{E}} : E \rightarrow E'$  and  $\iota_{\mathcal{Z}} : Z \rightarrow Z'$ , where  $\iota_{\mathcal{E}}$  must be measure preserving  $p_{\mathcal{E}^0} = \iota_{\mathcal{E}*} p_{\mathcal{E}}$ . For notational simplicity, we will drop the subscript in  $\iota_{\mathcal{Z}}$  and use the symbol  $\iota$  to refer both to the SCM isomorphism and the noise isomorphism.

The elementwise diffeomorphisms are required to make the following diagrams commute  $\forall i; i' = \iota(i)$ :

$$\begin{array}{ccc} Z_{\text{pa}_i} & E_i & \xrightarrow{(\iota_{\text{pa}_i}, \iota_{E_i})} Z'_{\text{pa}_{i^0}} & E_{i^0} \\ \downarrow f_i & & & \downarrow f_{i^0} \\ Z_i & & \xrightarrow{\iota_i} & Z'_{i^0} \end{array} \quad (3)$$

Intuitively, this says that if we apply a causal mechanism  $f_i$  and then reparameterize the causal variable  $i$  using  $\iota_i$ , we get the same thing as first reparameterizing the parents and noise variable of variable  $i$ , and then applying the causal mechanism  $f_{i^0}$ .

To reason about interventions, we equip SCMs with intervention distributions in the following definition.

**Definition 7** (Intervention structural causal model (ISCM)). *An intervention structural causal model (ISCM) is a tuple  $D = \langle \mathcal{C}; I; p_{\mathcal{I}} \rangle$  of*

1. *an acyclic SCM  $C = \langle \mathcal{Z}; E; F; p_{\mathcal{E}} \rangle$  that admits a faithful distribution, meaning that conditional independence of causal variables  $Z$  implies  $d$ -separation (Pearl, 2000).*
2. *a set  $I$  of interventions on  $C$ , where each intervention  $(I; (f_i)_{i \in I}) \in I$  consist of*
  - (a) *a subset  $I \subseteq \mathcal{V}; \dots; \text{ng}$  of the causal variables, called the intervention target set, and*
  - (b) *for each  $i \in I$ , a new causal mechanism  $f_i : E_i \rightarrow Z_i$  which replaces the original mechanism and which does not depend on the parents.*

*We define intervention set  $I$  to be atomic if the number of targeted variables is one or zero.*
3. *a probability measure  $p_{\mathcal{I}}$  over  $I$ .*

We can extend the notion of isomorphism from SCMs to ISCMs.

**Definition 8** (Isomorphism of ISCMs). *Let  $D = \langle \mathcal{C}; I; p_{\mathcal{I}} \rangle$  and  $D' = \langle \mathcal{C}'; I'; p'_{\mathcal{I}^0} \rangle$  be ISCMs. An ISCM isomorphism is an SCM isomorphism  $\iota : C \rightarrow C'$  with underlying graph isomorphism  $\iota : G(C) \rightarrow G(C')$  and a  $\iota_{\mathcal{E}}$ -diagonal diffeomorphism  $\iota_{\mathcal{E}} : E \rightarrow E'$  such that*

- *the graph isomorphism  $\iota$  induces a bijection of intervention sets*

$$\iota_{\mathcal{I}} : I \rightarrow I' : (I; (f_i)_{i \in I}) \mapsto (I'; (f'_{i^0})_{i^0 \in I})$$

- *for each intervention  $(I; (f_i)_{i \in I}) \in I$ , and each intervened on variable  $i \in I$ , the following diagram commutes:*

$$\begin{array}{ccc} E_i & \xrightarrow{f_{E,i}} & E'_{(i)} \\ \downarrow \tilde{f}_i & & \downarrow \tilde{f}'_{(i)} \\ Z_i & \xrightarrow{\iota_i} & Z'_{(i)} \end{array} \quad (4)$$

- *$\iota_{\mathcal{E}}$  is measure preserving, i. e.  $p_{\mathcal{E}^0} = (\iota_{\mathcal{E}})_* p_{\mathcal{E}}$ .*
- *the bijection  $\iota_{\mathcal{I}} : I \rightarrow I'$  preserves the distribution over interventions:  $p_{\mathcal{I}} = p'_{\mathcal{I}^0}$ .*

Latent Causal Models (LCMs), defined in Def. 1, add a map to the data space to an ILCM. We can lift ISCM isomorphisms to LCM isomorphisms by requiring that these decoders must respect the ISCM isomorphism.

**Definition 9** (Isomorphism of LCMs). *Let  $M = \langle \mathcal{C}; X; g; I; p_{\mathcal{I}} \rangle$  and  $M' = \langle \mathcal{C}'; X'; g'; I'; p'_{\mathcal{I}^0} \rangle$  be LCMs with identical observation space  $X = X'$ . An LCM isomorphism of LCM is an ISCM isomorphism  $\iota : D \rightarrow D'$  such that the decoders respect the SCM isomorphism, so this diagram must commute:*

$$\begin{array}{ccc} Z & \xrightarrow{\iota} & Z' \\ \searrow g & & \swarrow g' \\ & X & \end{array} \quad (5)$$

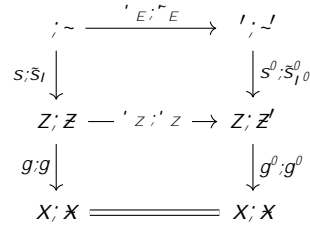


Figure 5: An illustration of the spaces and maps in our definitions and proof. When LCMs  $\mathcal{M}; \mathcal{M}'$  are isomorphic, all squares in the diagram should commute. Additionally, all maps should preserve the weakly supervised distributions on the variables and all horizontal maps should be  $\mathcal{I}$ -diagonal. Note that the latent variables  $(\mathcal{Z}; \mathcal{Z}')$  can differ up to a diffeomorphism, but the  $\mathcal{X}$  variables are actually observed, so must be identically equal. From that equality, the other horizontal maps are uniquely defined.

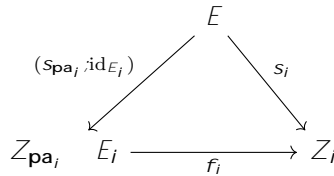
*Remark 1.* By defining objects and isomorphisms, we have defined a groupoid of SCMs, a groupoid of ISCMs and a groupoid of LCMs, as the isomorphisms are composed and inverted in an obvious way.

**Definition 10** (Equivalence). *We call two SCMs, ISCMs, or LCMs equivalent if an isomorphism exists between them.*

Informally, two SCMs, ISCMs, or LCMs are equivalent if there is a  $\mathcal{I}$ -diagonal map between their causal variables (i. e. the causal variables are equal up to permutation and elementwise diffeomorphisms), there is a  $\mathcal{I}$ -diagonal map between their noise encodings, and all other structure (decoders, intervention sets, intervention distributions) is compatible with these reparameterizations.

Next, we define the solution function of an SCM or ISCM, which maps from noise variables to causal variables by repeatedly applying the causal mechanisms.

**Definition 11** (Solution). *Given an ISCM  $D = \langle \mathcal{H}; \mathcal{I}; p_{\mathcal{I}} \rangle$ , the solution function  $s : E \rightarrow \mathcal{Z}$  is the unique function such that for all  $i \in [n]$ , the following diagram commutes (Bongers et al., 2021)*



Or equations, we have that  $s(\cdot)_i = f_i(\cdot; s(\text{pa}_i))$ . Similarly, intervention  $(\mathcal{I}; (f_i)_{i \in \mathcal{I}}) \in \mathcal{I}$  yields a solution function  $s_{\mathcal{I}} : E \rightarrow \mathcal{Z}$  with the modified causal mechanisms.

For example, with two variables with  $\mathcal{Z}_1 \perp \mathcal{Z}_2$ , the solution is given by:

$$s : E \rightarrow \mathcal{Z} : \begin{pmatrix} \mathcal{Z}_1 \\ \mathcal{Z}_2 \end{pmatrix} = \begin{pmatrix} f_1(\cdot) \\ f_2(\cdot; f_1(\cdot)) \end{pmatrix}$$

Since we require causal mechanisms to be pointwise diffeomorphic, the solution function is a diffeomorphism as well.

Pushing the noise distribution of an SCM through the solution function finally gives us the (observable) distribution entailed by an SCM or ISCM. In an ISCM or LCM we can define several other (observable or interventional) distributions.

**Definition 12** (Distributions). *Given an LCM  $\mathcal{M} = \langle \mathcal{H}; \mathcal{X}; g; \mathcal{I}; p_{\mathcal{I}} \rangle$ , we have the following generative process:*

$$\begin{array}{l}
 p_{\mathcal{E}}; \quad \quad \quad z = s(\cdot); \quad \quad \quad x = g(z); \quad \quad \quad e = s^{-1}(z) \\
 \mathcal{I} \quad p_{\mathcal{I}}; \quad \quad \sim \quad p_{\mathcal{E}}(\cdot; \mathcal{I}); \quad \quad z = s_{\mathcal{I}}(\cdot); \quad \quad x = g(z); \quad \quad e = s^{-1}(z); \quad (6)
 \end{array}$$

where  $p(\cdot; \mathcal{I}; i \in \mathcal{I}) = p_{\mathcal{E}}(\cdot)$  and  $p(\cdot; \mathcal{I}; i \notin \mathcal{I}) = \delta(\cdot; i)$  is the Dirac measure.

Then we define the following weakly supervised distributions:

- The weakly supervised noise distribution with interventions:  $p_C^{\mathcal{E}, \mathcal{I}}(\cdot; \cdot; l)$ .
- The weakly supervised causal distribution with interventions:  $p_C^{\mathcal{Z}, \mathcal{I}}(z; z; l)$ .
- The weakly supervised observational distribution with interventions:  $p_{\mathcal{M}}^{\mathcal{X}, \mathcal{I}}(x; x; l)$ .

These distributions are given by appropriate pushforwards of the noise distributions through the transformations in Eq. (6).

By marginalizing over  $l$ , we get  $p_C^{\mathcal{E}}; p_C^{\mathcal{Z}}; p_C^{\mathcal{O}}; p_{\mathcal{M}}^{\mathcal{X}}$  respectively.

The relationships between all the maps can be found in Fig. 5.

## A.2 IDENTIFIABILITY PROOF

First, we prove two auxiliary lemmata.

**Lemma 1.** *Let  $f : [0; 1] \rightarrow [0; 1]$  be differentiable and Lebesgue measure preserving. Then either  $f(x) = x$  or  $f(x) = 1 - x$ .*

*Proof.* We follow the proof from Stack Exchange user zhw (2016). Let  $\mu$  be the Lebesgue measure. Measure preservation means that for any measurable subset  $U \subseteq [0; 1]$ ,  $\mu(U) = \mu(f^{-1}(U))$ .

First, note that  $f$  is surjective, because otherwise the image of  $f$  is a proper subinterval  $[a; b] \subsetneq [0; 1]$  and  $\mu(f^{-1}([a; b])) = \mu([0; 1]) = 1 > \mu([a; b]) = b - a$ , which contradicts measure-preservation.

Define the open ball  $B(x; r) = \{y \in [0; 1] \mid |y - x| < r\}$ . Suppose that  $f'(0) = 0$  for some  $x \in [0; 1]$ . Then there exists an  $r > 0$  such that  $f(B(x; r)) \subseteq B(f(x); r/4)$ , and thus  $B(x; r) \subseteq f^{-1}(B(f(x); r/4))$ . Therefore,  $\mu(B(x; r)) \leq \mu(f^{-1}(B(f(x); r/4))) = \mu(B(f(x); r/4)) = r/4 = r/2$ , contradicting measure preservation. Hence  $f'(x) \neq 0$  on  $[0; 1]$ .

By the Darboux theorem,  $f'$  is either strictly positive or strictly negative on the interval and thus  $f$  is either strictly increasing or decreasing and thus a bijection. Assume that it is strictly increasing, then  $\mu(x \in [0; 1] \mid x = f^{-1}(f([0; x]))) = \mu(f([0; x])) = f(x) - f(0) = f(x)$ . Similarly, if it is strictly decreasing, we find  $f(x) = 1 - x$ .  $\square$

**Lemma 2.** *Let  $A = C = \mathbb{R}$  and  $B = \mathbb{R}^n$ . Let  $f : A \times B \rightarrow C$  be differentiable. Define differentiable measures  $p_A$  on  $A$  and  $p_C$  on  $C$ . Let  $g : B \rightarrow [0; 1]$ ,  $f(\cdot; b) : A \rightarrow C$  be measure-preserving. Then  $f$  is constant in  $B$ .*

*Proof.* Let  $P_A : A \rightarrow [0; 1]$ ;  $P_C : C \rightarrow [0; 1]$  be the diffeomorphic cumulative density functions. Then  $P_A^{-1}$  and  $P_C^{-1}$  are measure-preserving maps from the uniform distribution on  $[0; 1]$ . Now write  $g : [0; 1] \rightarrow B$ ;  $f(\cdot; b) : A \rightarrow C$  such that this diagram of measure-preserving differentiable maps commutes:

$$\begin{array}{ccc}
 A & \xrightarrow{f(\cdot; b)} & C \\
 P_A \searrow & & \nearrow P_C^{-1} \\
 [0; 1] & \xrightarrow{g(\cdot; b)} & [0; 1]
 \end{array}$$

Then  $g$  is differentiable and  $g : B \rightarrow [0; 1]$  measure-preserving  $[0; 1] \rightarrow [0; 1]$ . By the previous Lemma 1, the only differentiable measure-preserving functions  $[0; 1] \rightarrow [0; 1]$  are id and  $1 - \text{id}$ . As  $g$  is continuous in  $B$ , it can not vary between id and  $1 - \text{id}$  and thus  $g$ , and consequently  $f$  are constant in  $B$ .  $\square$

We can interpret this lemma in terms of statistical independence. Starting from a product measure on  $A \times B$ , the requirements of the lemma correspond to  $a \perp\!\!\!\perp b$  and  $c \perp\!\!\!\perp b$ . The lemma thus defines a sense in which for real-valued variables, statistical independence implies functional independence (the converse is always true).

Now in the remainder of this subsection, we prove the main theorem.

**Theorem 1** (Identifiability of  $\mathbb{R}$ -valued LCMs from weak supervision). *Let  $\mathcal{M} = \langle h; C; X; g; l; p_{\mathcal{I}} \rangle$  and  $\mathcal{M}' = \langle h'; C'; X'; g'; l'; p'_{\mathcal{I}'} \rangle$  be LCMs with the following properties:*

- The SCMs  $C$  and  $C'$  both consist of  $n$  real-valued endogeneous variables, i. e.  $E_i = Z_i = Z'_i = E'_i = \mathbb{R}$ .
- The intervention sets  $I$  and  $I'$  consist of the empty intervention and all atomic interventions,  $I = \{f; \bar{f}z_0g; \dots; \bar{f}z_rgg\}$  and similar for  $I'$ .
- The intervention distribution  $p_{\mathcal{I}}$  and  $p'_{\mathcal{I}}$  have full support.

Then the following two statements are equivalent:

1. The weakly supervised distributions entailed by the LCMs are equal,  $p_{\mathcal{M}}(x; \mathbf{x}) = p_{\mathcal{M}^0}(x; \mathbf{x})$ .
2. The LCMs are equivalent,  $\mathcal{M} \equiv \mathcal{M}'$ .

*Proof.* “(2)  $\Rightarrow$  (1)”: If the LCMs are equivalent, then the fact that  $\tau_{\mathcal{E}}$  and  $\tau'_{\mathcal{E}}$  are measure preserving and that diagrams (3) and (4) commute, implies that  $p_{\mathcal{C}^0}^{\mathcal{Z}^0} = (\tau_{\mathcal{E}}; \tau'_{\mathcal{E}})_* p_{\mathcal{C}}^{\mathcal{Z}}$ . Then because diagram (5) commutes, the weakly supervised distributions coincide,  $p_{\mathcal{M}^0}^{\mathcal{X}} = p_{\mathcal{M}}^{\mathcal{X}}$ .

“(1)  $\Rightarrow$  (2)”: Conversely, if the weakly supervised distributions coincide,  $p_{\mathcal{M}^0}^{\mathcal{X}} = p_{\mathcal{M}}^{\mathcal{X}}$ , the images of  $g: Z \rightarrow X; g': Z' \rightarrow X$  coincide,

$$g = g'^{-1} \circ g': Z \rightarrow Z' \quad (7)$$

is a diffeomorphism, and  $\tau_{\mathcal{E}}$  preserves the weakly supervised distribution over causal variables:  $p_{\mathcal{C}^0}^{\mathcal{Z}^0} = (\tau_{\mathcal{E}}; \tau'_{\mathcal{E}})_* p_{\mathcal{C}}^{\mathcal{Z}}$ .

LCM equivalence then follows from showing that  $\tau_{\mathcal{E}}: D \rightarrow D'$  is an ISCM isomorphism, where  $D = hC; I; p_{\mathcal{I}}; i$  and  $D' = hC'; I'; p'_{\mathcal{I}}; i$  be the ISCMs inherent to  $\mathcal{M}$  and  $\mathcal{M}'$ . We show this in the following steps:

1. For each intervention  $I$  in  $D$ , there is a corresponding intervention  $I'$  in  $D'$ , given by a permutation  $\sigma: [n] \rightarrow [n]$ , such that  $\tau_{\mathcal{E}}$  preserves the interventional distribution.
2. The diffeomorphism  $\tau_{\mathcal{E}}$  is  $\sigma$ -diagonal.
3. The permutation  $\sigma$  preserved the ancestry structure of graphs  $G(C)$  and  $G(C')$ .
4. The diffeomorphism  $\tau_{\mathcal{E}}: E \rightarrow E'$  of noise variables is  $\sigma$ -diagonal.
5. The causal mechanisms are compatible with  $\tau_{\mathcal{E}}$ .

**Step 1: Interventions preserved** Remember that the diffeomorphism  $\tau_{\mathcal{E}}: Z \rightarrow Z'$  is such that  $p_{\mathcal{C}^0}^{\mathcal{Z}^0} = (\tau_{\mathcal{E}}; \tau'_{\mathcal{E}})_* p_{\mathcal{C}}^{\mathcal{Z}}$ . For atomic interventions  $I \not\subseteq J \supseteq I$ , consider the intersection of the supports of the weakly supervised distribution for interventions on  $I$  and  $J$ :  $U = \text{supp } p_{\mathcal{C}}^{\mathcal{Z}; \mathcal{I}}(z; z \setminus I) \setminus \text{supp } p_{\mathcal{C}}^{\mathcal{Z}; \mathcal{I}}(z; z \setminus J) \subseteq Z \subseteq Z'$ . Note that  $U$  has zero measure in  $p_{\mathcal{C}}^{\mathcal{Z}; \mathcal{I}}(U \setminus I) = p_{\mathcal{C}}^{\mathcal{Z}; \mathcal{I}}(U \setminus J) = 0$ . The distribution is thus a discrete mixture on  $(z; \bar{z})$  of non-overlapping distributions.

The diffeomorphism  $(\tau_{\mathcal{E}}; \tau'_{\mathcal{E}})$  must map between these mixtures. Thus there exists a bijection  $\sigma: I \rightarrow I'$ , also inducing a permutation  $\sigma: [n] \rightarrow [n]$ , such that

$$p_{\mathcal{C}^0}^{\mathcal{Z}^0; \mathcal{I}^0} = (\tau_{\mathcal{E}}; \tau'_{\mathcal{E}})_* p_{\mathcal{C}}^{\mathcal{Z}; \mathcal{I}}:$$

**Step 2:  $\tau_{\mathcal{E}}$  is  $\sigma$ -diagonal** This measure preservation lets us define two equal distributions on  $Z \subseteq \mathcal{Z}' \subseteq I$ , namely  $(\text{id}_{\mathcal{Z}}; \tau'_{\mathcal{E}}; \text{id}_I)_* p_{\mathcal{C}}^{\mathcal{Z}; \mathcal{I}}$  and  $(\tau_{\mathcal{E}}^{-1}; \text{id}_{\mathcal{Z}^0}; \tau'_{\mathcal{E}})_* p_{\mathcal{C}^0}^{\mathcal{Z}^0; \mathcal{I}^0}$ . In particular, these must then have equal conditionals  $p(z' \setminus j \setminus z; I)$ . Thus, for any  $U \subseteq \mathcal{Z}'; z \supseteq Z; I \supseteq I$ ,

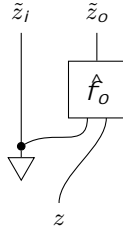
$$p_{\mathcal{C}^0}^{\mathcal{Z}^0; \mathcal{I}^0}(z' \setminus j \setminus z; I) = p_{\mathcal{C}}^{\mathcal{Z}; \mathcal{I}}(z \setminus \tau_{\mathcal{E}}^{-1}(U) \setminus j; I)$$

The conditional probability  $p_{\mathcal{C}}^{\mathcal{Z}; \mathcal{I}}(z \setminus j \setminus z; I)$  can be interpreted as a stochastic map  $Z \rightarrow Z'$ . The above relation can then be written as a commuting diagram of stochastic maps,  $\mathcal{I} \supseteq I; I' = \sigma(I)$ :

$$\begin{array}{ccc} Z & \xrightarrow{p_{\mathcal{C}}^{\mathcal{Z}; \mathcal{I}}(z \setminus j \setminus z; I)} & Z' \\ \downarrow \tau_{\mathcal{E}} & & \downarrow \tau'_{\mathcal{E}} \\ Z & \xrightarrow{p_{\mathcal{C}^0}^{\mathcal{Z}^0; \mathcal{I}^0}(z^0 \setminus j^0 \setminus z^0; I^0)} & Z' \end{array} \quad (8)$$

where we treat  $\tau_{\mathcal{E}}: Z \rightarrow Z'$  as a deterministic stochastic map.

For any variable  $i \in [n]$ , write the other nodes as  $o = [n] \setminus i$ . Let  $I = i$ . Then  $\rho_C^{\mathcal{Z}; \mathcal{I}}(z_j | z; I)$  can be written as a string diagram of stochastic maps:

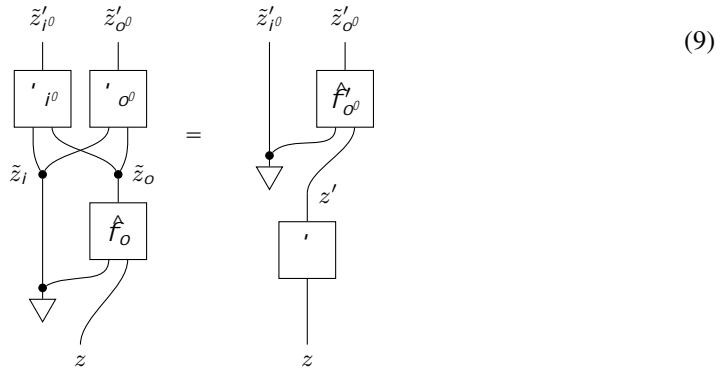


This string diagram represents a conditional probability distribution  $\rho(z_i; z_o | z)$  and is read from the bottom to the top. String diagrams map formally to a generative process (Fritz, 2020) and have been used previously in the context of causal models (Fong, 2013). In this case, the diagram maps to:

$$z_i \quad \rho(z_i); \quad z_o = \hat{f}_o(z_i; z)$$

where  $\rho(z_i)$  is the interventional distribution and the deterministic map  $\hat{f}_o : \mathcal{E}_i \times \mathcal{Z} \rightarrow \mathcal{E}_o$  can be constructed from the inverse solution  $s^{-1} : \mathcal{Z} \rightarrow \mathcal{E}$  and the causal mechanisms. Each box in a string diagram of stochastic maps denotes a stochastic map and each line to a measurable space. The triangle is the stochastic map  $\rho_i$  (the star denoting the one-point space; maps from which correspond to probability distributions over the codomain). The  $\hat{f}_o$  represents copying a variable.

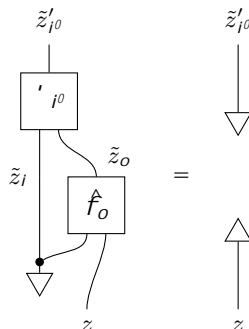
The above commuting diagram (8) can then be written as the equality of the following two string diagrams, where  $(I) = I' = i$ ;  $\mathcal{O}' = [n] \setminus i$ . We write  $\rho' : \mathcal{Z} \rightarrow \mathcal{Z}'$  as the pair  $\rho'_{i0} : \mathcal{Z} \rightarrow \mathcal{Z}'_{i0}$ ;  $\rho'_{o0} : \mathcal{Z} \rightarrow \mathcal{Z}'_{o0}$  obtained by projecting the output of  $\rho$  to the partition  $\mathcal{Z}' = \mathcal{Z}'_{i0} \times \mathcal{Z}'_{o0}$ :



This should be read as the equality of the two conditional probability distributions  $\rho(z'_{i0}, z'_{o0} | z)$  generated in the following way:

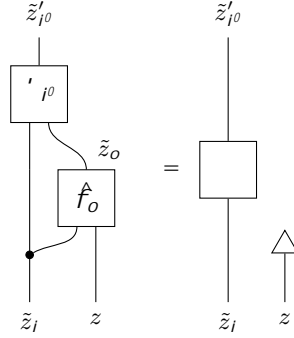
$$\begin{aligned} \text{Left: } & z_i \quad \rho(z_i); \quad z_o = \hat{f}_o(z_i; z); \quad z'_{i0} = \rho'_{i0}(z_i; z_o); \quad z'_{o0} = \rho'_{o0}(z_i; z_o); \\ \text{Right: } & z' = \rho'(z); \quad z'_{i0} \quad \rho'(z'_{i0}); \quad z'_{o0} = \rho'_{o0}(z'_{i0}; z'); \end{aligned}$$

The string diagram equality (9) implies equality when we disregard outputs  $\mathcal{Z}'_{o0}$ :

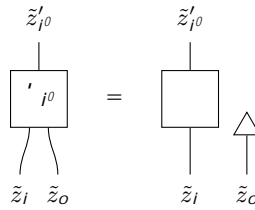


where the upwards pointing triangle represents discarding a variable.

Using Lemma 2, and the fact that  $\mathcal{E}_i = \mathcal{E}'_{j^0} = \mathbb{R}$ , the composed differentiable function  $\mathcal{E}_i \times \mathcal{Z} \rightarrow \mathcal{E}'_{j^0}$  is constant in  $\mathcal{Z}$ . Thus we have a deterministic function  $\mathcal{E}_i \rightarrow \mathcal{E}'_{j^0}$  such that:

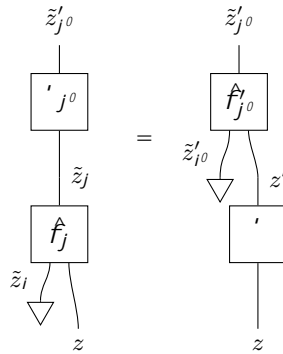


The deterministic function  $\mathcal{E}_i \times \mathcal{Z} \rightarrow \mathcal{E}'_{j^0}$  is surjective and both the left- and right-hand side can be seen as first applying this function (though the output is discarded on the right hand side), which implies there exists a function  $\mathcal{E}_i \rightarrow \mathcal{E}'_{j^0}$  such that



In words, the function  $'_{j^0} : \mathcal{Z}_i \times \mathcal{Z}_O \rightarrow \mathcal{Z}'_{j^0}$  is constant in  $\mathcal{Z}_O$ . This holds for all  $i$  and thus  $'$  is  $\mathcal{Z}$ -diagonal.

**Step 3: Ancestry preserved** Let  $i \neq j \in [n]$ ,  $i' = (i)$ ,  $j' = (j)$ , and  $l = fig$ . Writing  $'$  as  $\mathcal{Z}$ -diagonal, the commuting diagram (8) for the  $j'$  component of  $z'$ , can be written as the following string diagram:



The left hand side is a deterministic map  $\mathcal{Z} \rightarrow \mathcal{E}'_{j^0}$  if and only if  $\hat{f}_j$  is constant in  $\mathcal{E}_i$  which by faithfulness is the case if and only if  $i \notin \text{anc}_j$ . The same holds on the right hand side, so  $\delta i \neq j \in [n]$ ,  $i \in \text{anc}_j^{\mathcal{E}} \iff (i) \in \text{anc}_{(j)}$ .

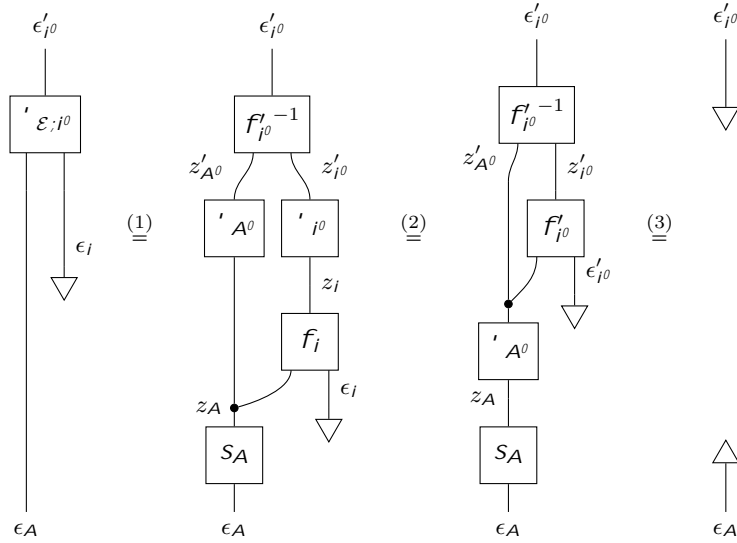
**Step 4: Noise map diagonal** Define  $'_{\mathcal{E}} = s'^{-1} \circ ' : S : E \rightarrow E'$ . Note that  $'_{\mathcal{E}}(\cdot)_{j^0}$  only depends on  $i$  and  $\text{anc}_i$ , because  $S(\cdot)_{\text{anc}_i, i}$  and  $s'^{-1}(z')_{j^0}$  only depend on ancestors,  $'$  is  $\mathcal{Z}$ -diagonal and preserves ancestry.

The map  $'$  is measure-preserving. Thus  $\delta i$  and writing  $A = \text{anc}_i$ , the conditional  $p(z_i | z_A) = p(z_i | z_{\text{pa}_i})$ , interpreted as a stochastic map, is preserved by  $'$ . We can express this as another

commuting diagram, in which the two paths from  $E_A$  to  $E'_{j^0}$  must be equal:

$$\begin{array}{ccccc}
 E_A & \xrightarrow{S_A} & Z_A & \xrightarrow{p(z_i|z_{\text{pa}_i})} & Z_{A;i} \\
 & & \downarrow \epsilon'_A & & \downarrow \epsilon'_{A;i} \\
 & & Z'_{A^0} & \xrightarrow{p(z_{i^0}|z_{\text{pa}_{i^0}})} & Z'_{A^0;i^0} \xrightarrow{f'_{i^0}} E'_{j^0}
 \end{array}$$

where  $f'_{j^0}{}^{-1}(z') = f(z'_{\text{pa}_{j^0}})^{-1}(z'_{j^0})$ . Then we have:



where the first equality follows from the definition of  $\epsilon'_{j^0}$ , the second equality from the commuting diagram above and the third equality from the fact that  $f'_{j^0}$  and  $f'_{j^0}{}^{-1}$  cancel. Then, again using Lemma 2, the map on the left hand side must be constant in  $\epsilon'_A$ . The noise encoding is thus also  $\epsilon'_A$ -diagonal.

**Step 5: Equivalence** Consider for a variable  $i$  and with  $i' = (i)$  the following commuting diagram of deterministic maps. Note that we write the causal mechanism  $f_i$  as a function of all ancestors, so it is constant in the non-parents. Because of faithfulness, it is non-constant in the parents. Since  $\epsilon'_A$  preserves ancestors,  $f'_{j^0}$  is well-typed.

$$\begin{array}{ccc}
 E & \xrightarrow{\epsilon'_E} & E' \\
 \downarrow (S_{\text{anc}_i}; \text{id}_{E_i}) & & \downarrow (S_{\text{anc}_{i^0}}; \text{id}_{E'_{i^0}}) \\
 Z_{\text{anc}_i} & \xrightarrow{\epsilon'_i} & Z'_{\text{anc}_{i^0}} \\
 \downarrow f_i & & \downarrow f'_{i^0} \\
 Z_i & \xrightarrow{\epsilon'_{Z;i}} & Z'_{j^0}
 \end{array}$$

The composition of the left vertical maps is equal to  $S_i$ , the composition of the right vertical maps to  $S'_{j^0}$ . Therefore and because of the definition of  $\epsilon'_E$ , the outer and the top square commute. Then, because  $(S_{\text{anc}_i}; \text{id}_{E_i})$  is surjective, the bottom square also commutes (Riehl, 2017, Lemma 1.6.21).

Then for  $Z_j \supseteq \text{pa}_i^C$ , we have that

$$Z_j \supseteq \text{pa}_i^C \iff f_i \text{ not constant in } Z_j \iff f'_{i^0} \text{ not constant in } Z'_{j^0} \iff Z'_{j^0} \supseteq \text{pa}_{i^0}^{C'}$$

And thus  $\epsilon'_A$  not only preserves ancestry, but also parenthood and is thus a graph isomorphism  $\epsilon'_A : G(C) \cong G(C')$ . Diagram (3) commutes, and we have established an SCM isomorphism  $\epsilon'_A : C \cong C'$ .

To have this also be an ISCM isomorphism, we need diagram (4) to commute and the distribution

over interventions to be preserved. For the first, use the fact that all maps in (4) are isomorphisms to simply define  $\tilde{\cdot}_\varepsilon$  so that the diagram commutes. The second follows directly from the assumptions. Hence  $\tilde{\cdot} : D \rightarrow D'$  is an ISCM isomorphism,  $D \cong D'$ , and—together with the arguments in the beginning of this proof—finally  $\mathcal{M} \cong \mathcal{M}'$ .  $\square$

### A.3 LIMITATIONS & GENERALIZATION

Our identifiability result relies on a few assumptions. Here we discuss some key requirements of Thm. 1 and whether they can be relaxed.

**Diffeomorphic causal mechanisms** In Def. 4, we require causal mechanisms to be pointwise diffeomorphisms from noise variables to causal variables. Under some mild smoothness assumptions, any SCM can be brought into this form by elementwise redefinitions of the variables, without affecting the observational or interventional distributions. (However, such a redefinition may change counterfactual distributions.)

**Perfect interventions** Our proof of Thm. 1 requires perfect interventions, i. e. intervened-upon mechanisms not depending on any causal variables. This is arguably the biggest mismatch between our assumptions and many real-world systems.

**Diffeomorphic decoder** Definition 1 and Thm. 1 assume that the map from causal variables to observed data is given by a deterministic, diffeomorphic decoder. However, our practical implementation in a VAE uses a stochastic decoder and allows for noisy data. Our experiments provide empirical evidence for identifiability in this setting. We believe that it may be possible to extend Thm. 1 to stochastic decoders, similarly to Khemakhem et al. (2020). We plan to study this extension in future work.

**Real-valued causal variables** Theorem 1 assumes real-valued causal and noise variables,  $Z_i = E_i = \mathbb{R}$ . We can easily extend this to intervals  $(a; b) \subseteq \mathbb{R}$ , as these are isomorphic to  $\mathbb{R}$ . However, the extension to arbitrary continuous spaces or  $\mathbb{R}^n$  is not straightforward. The main reason is that our proof relies on Lemma 2, which does not generalize.

Let us provide a counterexample for identifiability with circle  $S^1$ -valued causal variables.

**Example 1** ( $S^1$ -valued non-identifiable LCMs). Consider an LCM  $\mathcal{M} = \langle \mathcal{C}; X; g; I; p_{\mathcal{I}} \rangle$  with the following components:

- The SCM  $\mathcal{C}$  consists of two circle-valued variables  $z_1; z_2 \in S^1$  with noise variables  $\varepsilon_1; \varepsilon_2 \in S^1$ . We parameterize  $S^1$  as  $[0; 2\pi)$  with addition defined modulo  $2\pi$ .
- The causal graph is  $z_1 \rightarrow z_2$ .
- The causal mechanisms are  $f_1(\varepsilon_1) = \varepsilon_1$  and  $f_2(\varepsilon_2; z_1) = \varepsilon_2 + z_1$ .
- The solution function is  $s(\varepsilon_1; \varepsilon_2) = (\varepsilon_1; \varepsilon_2 + \varepsilon_1)$ .
- The noise variables are distributed as  $\varepsilon_1 \sim U$ , uniformly, and  $\varepsilon_2 \sim q$ , which we require to not be invariant under translations (so in particular not uniform). For example, one can take the von Mises distribution  $\log q(\varepsilon_2) = \cos(\varepsilon_2) + \text{const}$ .
- The observation space is  $X$  and the decoder  $g : S^1 \times S^1 \rightarrow X$  is diffeomorphic.
- The intervention set  $I$  consists of the empty intervention, atomic interventions on  $z_1$  with  $z_1 \sim U$ , and atomic interventions on  $z_2$  with  $z_2 \sim U$ . Each of these interventions has probability  $\frac{1}{3}$  in  $p_{\mathcal{I}}$ .

Note that the SCM is faithful, as  $z_1 \not\perp z_2$  in the observational distribution, because  $q$  is not translationally invariant. The LCM entails the weakly supervised causal distribution

$$p_{\mathcal{C}}^z(z; z) = U(z_1) q(z_2 - z_1) \frac{1}{3} (\delta_{z_1 - z_1} (z_2 - z_2) + \frac{1}{3} U(z_1) (\delta_{z_2 - z_2} (z_1 + z_1) + \frac{1}{3} (\delta_{z_1 - z_1} U(z_2)) \quad (10)$$

with Dirac delta  $\delta$ . The weakly supervised data distribution is then given by  $p_{\mathcal{M}}^x = (g_*; g_*)p_{\mathcal{C}}^z$ .

Now consider a second LCM  $\mathcal{M}' = \langle \mathcal{C}'; X; g'; I'; p_{\mathcal{I}'} \rangle$ :



- The SCM  $C'$  consists of two circle-valued variables  $Z'_1, Z'_2 \in S^1$  with noise variables  $U'_1, U'_2 \in S^1$ .
- The causal graph is trivial and the causal mechanisms are given by the identity,  $F'_i(Z'_i) = U'_i$ .
- The noise variables are distributed as  $U'_1 \sim U$  and  $U'_2 \sim q$ .
- The observation space is  $X$  and the decoder  $g' : S^1 \times S^1 \rightarrow X$  is given by the diffeomorphism  $g'(z') = g(s(z'))$ , where  $s$  is the solution function of  $C$ .
- The intervention set  $I'$  consists of empty interventions, atomic interventions on  $Z'_1$  with  $Z'_1 \sim U$ , and atomic interventions on  $Z'_2$  with  $Z'_2 \sim U$ . Each of these interventions has probability  $\frac{1}{3}$  in  $p_{I'}$ .

We find a weakly supervised causal distribution

$$p_{C'}^{Z'}(z'; z') = U(Z'_1) q(Z'_2) \frac{1}{3} (Z'_1 = z'_1) (Z'_2 = z'_2) + \frac{1}{3} U(Z'_1) (Z'_2 = z'_2) + \frac{1}{3} (Z'_1 = z'_1) U(Z'_2) \quad (11)$$

Clearly, two LCMs are not equivalent, because their graphs are non-isomorphic. Yet, if we define

$$\gamma : Z \rightarrow Z' : (z_1; z_2) \mapsto (z_1; z_2 - z_1)$$

then the weakly supervised distribution of the causal variables is preserved:

$$\begin{aligned} ((\gamma; \gamma) \circ p_{C'}^{Z'}) (z'; z') &= p_C^Z((Z_1; Z_2 + Z_1); (z'_1; z'_2 + z'_1)) \\ &= U(Z_1) q(Z_2 + Z_1 - z'_1) \frac{1}{3} (Z_1 = z'_1) (Z_2 + Z_1 - z'_1 = z'_2 + z'_1) \\ &\quad + \frac{1}{3} U(Z_1) (Z_2 + Z_1 - z'_1 = z'_2 + z'_1) (Z_1 = z'_1) + \frac{1}{3} (Z_1 = z'_1) U(Z_2 + Z_1) \\ &= U(Z_1) q(Z_2) \frac{1}{3} (Z_1 = z'_1) (Z_2 = z'_2) \\ &\quad + \frac{1}{3} U(Z_1) (Z_2 = z'_2) + \frac{1}{3} (Z_1 = z'_1) U(Z_2) \\ &= p_{C'}^{Z'}(z'; z') \end{aligned}$$

where we use that the density  $U$  is constant. Also, because  $\gamma = s^{-1}$  and  $g'(z') = g(s(z'))$ , we have that  $p_{\mathcal{M}}^X = p_{\mathcal{M}'}^X$ .

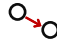
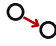

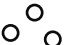
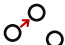
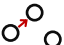
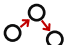
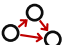
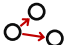
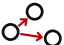
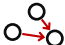
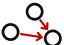
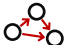
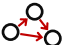
So these two models with their non-isomorphic graph structures have identical weakly-supervised distributions on the observables  $X; \mathcal{X}$ . They therefore provide a counter-example for a straightforward generalization of Thm. 1 to causal variables with arbitrary continuous domains.

Why does identifiability fail in this example? It is because the interventional distributions in  $\mathcal{M}$  have an accidental symmetry not expected by the graph structure, namely translational invariance. This makes it possible to fit the weakly supervised distribution with a simpler causal graph. This is related to faithfulness, but the standard definition of faithfulness only concerns observational distributions (and in this sense both  $\mathcal{M}$  and  $\mathcal{M}'$  are faithful). Because of this accidental symmetry, steps 1 and 3 of our proof do not hold any more.

We have circumvented such issues in Thm. 1 by requiring that all causal and noise variables are  $\mathbb{R}$ -valued. In this setting, functional dependence implies statistical dependence, as formalized in Lemma 2, and the counterexample does not work. We conjecture that it is possible to generalize Thm. 1 to arbitrary continuous domains under mild additional assumptions, but leave this for future work.

Finally, we believe that such accidental symmetries are unlikely in the sense that under an appropriate measure over LCMs, non-identifiable LCMs have zero measure. We find it likely that this issue will not occur frequently in practical systems (unless these are finetuned to exhibit exactly this behaviour). Overall, we conjecture that identifiability from weak supervision can be generalized beyond the real-valued case presented in Thm. 1.

Table 2: Experiment results. For each experiment, we show the true causal graph underlying the data-generating process. We then show the results from our LCMs and compare to unstructured -VAE and disentanglement VAE (dVAE) baselines. We show the learned causal graph, the structural Hamming distance SHD between the learned and the true graph, the DCI disentanglement score ( $D$ ), and the intervention negative log posterior ( $-\log p_I$ ). Best results in bold.

Dataset	True graph	Method	Learned graph	SHD	$D$	$\log p_I$	
2D toy data		LCM		<b>0</b>	<b>0.99</b>	<b>0.28</b>	
		dVAE	–	n/a	0.35	0.33	
		$\beta$ -VAE	–	n/a	0.00	n/a	
Causal3DIdent		LCM		<b>0</b>	<b>1.00</b>	<b>0.16</b>	
		dVAE	–	n/a	<b>1.00</b>	<b>0.16</b>	
		$\beta$ -VAE	–	n/a	0.55	n/a	
		LCM		<b>0</b>	<b>0.99</b>	0.23	
		dVAE	–	n/a	0.84	<b>0.22</b>	
		$\beta$ -VAE	–	n/a	0.68	n/a	
		LCM		1	<b>0.96</b>	<b>0.19</b>	
		dVAE	–	n/a	0.20	2.88	
		$\beta$ -VAE	–	n/a	0.06	n/a	
		LCM		<b>0</b>	<b>0.96</b>	<b>0.21</b>	
		dVAE	–	n/a	0.46	4.31	
		$\beta$ -VAE	–	n/a	0.44	n/a	
		LCM		<b>0</b>	<b>0.98</b>	<b>0.21</b>	
		dVAE	–	n/a	0.62	0.24	
		$\beta$ -VAE	–	n/a	0.35	n/a	
		LCM		<b>0</b>	<b>0.96</b>	<b>0.18</b>	
		dVAE	–	n/a	0.32	4.06	
		$\beta$ -VAE	–	n/a	0.24	n/a	
	Average		LCM		<b>0.17</b>	<b>0.98</b>	<b>0.20</b>
			dVAE		n/a	0.57	1.98
			$\beta$ -VAE		n/a	0.38	n/a

## B EXPERIMENT DETAILS

### B.1 2D TOY EXPERIMENT

In our first experiment, we generate latent data in  $Z = \mathbb{R}^2$  from a nonlinear SCM with graph  $z_1 \rightarrow z_2$ . In particular, we have that  $z_1 \sim N(z_1; 0; 1^2)$  and  $z_2 \sim N(z_1; 0.3z_1^2 + 0.6z_1; 0.8^2)$ . This latent data is mapped through the data space with a randomly initialized coupling flow with five affine coupling layers interspersed with random permutations of the dimensions. For the weakly supervised setting we use a uniform intervention prior over  $f; ; fz_1g; fz_2gg$ . We generate  $10^5$  training samples,  $10^5$  additional training samples for the models used to compute the DCI metrics,  $10^4$  validation samples, and  $10^4$  evaluation samples.

The learned LCM consists of an SCM prior, an encoder, and a decoder. In the SCM, the graph is fixed (we “learn” the graph by training multiple LCMs with different fixed graphs and then selecting the model with the best validation loss). Each causal mechanism is implemented as an MLP of the parents that outputs the parameters of an affine transformation from a standard normal noise variable to a causal variable. The encoder and decoder are diagonal Gaussians, with mean and standard deviations output by an MLP. For each MLP, we use two hidden layers with 100 units each and ReLU activations.

The disentanglement VAE baseline uses the same setup, except with a trivial graph. The  $\beta$ -VAE uses



Figure 6: Effect of varying the learned causal factors on the image in the “chain” subset of the Causal3DIdent dataset. We encode a single test images (middle column) into the three learned latents, vary each of these causal factors independently, and show the reconstructed images. The LCM (top) learns a representation that is quite disentangled:  $z_1$  corresponds to the spotlight position,  $z_2$  to the spotlight hue, and  $z_3$  to the object hue. In contrast, the acausal dVAE baseline strongly entangles these factors in its learned representation.



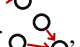

the same encoder and decoder, but uses a standard normal prior rather than an SCM and treats  $x$  and  $x$  as two i. i. d. samples from the same distribution.

All models are trained on the VAE loss in Eq. (2) plus a regularization term (0.1 times the number of edges in the graph). We train for  $10^5$  steps using the Adam optimizer Kingma & Ba (2015) with a batch size of 100. The learning rate is initially  $10^{-3}$  and is annealed with a cosine schedule. We estimate the model log likelihood using importance weighting with importance sampling (a la IWAE Burda et al. (2016)), using 10 samples at validation time and 100 samples at test time. We use a manifold “fuzziness” of  $\epsilon = 0.1$ . For each method and each graph, we train three models with different random seeds and select the model with the best validation log likelihood.

## B.2 CAUSAL3DIDENT

In the Causal3DIdent experiments we consider six different datasets, each generated from a different causal graph, SCM, and decoder. The six causal graphs we consider are:

- the trivial graph  $\circ \circ$
- single edge  $\circ \rightarrow \circ$

- the chain ,
- the fork ,
- the collider , and
- the full graph .

For each of these subsets, we randomly generate a nonlinear SCM with heteroskedastic noise: for each causal mechanism, we randomly initialize an MLP that outputs the scale and shift of an affine transformation as a function of the causal parents. We choose an MLP initialization scheme that emphasizes nontrivial, nonlinear causal effects. We then identify a random permutation of the three causal variables with three high-level concepts in the Causal3DIdent dataset: the object hue, the spotlight hue, and the spotlight position. We use the following causal graphs:

- single edge: object hue  $\perp$  spotlight position;
- chain: spotlight position  $\perp$  spotlight hue  $\perp$  object hue;
- fork: spotlight hue  $\perp$  spotlight position, object hue;
- collider: spotlight hue  $\perp$  object hue  $\perp$  spotlight position;
- full graph: spotlight hue  $\perp$  object hue  $\perp$  spotlight position, spotlight hue  $\perp$  spotlight position.

Since all of these properties are defined on a range  $[0;2)$ , we apply an elementwise  $\arctanh$  transform and rescaling to our variables such that they populate a subset of  $[0;2)$ . This also avoids topological issues. Next, we generate images in  $64 \times 64$  resolution following the procedure described in von Kügelgen et al. (2021). We use Blender (Blender Online Community, 2021) to generate 3D rendered images based on the previously defined causal variables. To increase diversity of the six datasets, we render each dataset with a different object: Teapot (Newell, 1975), Armadillo (Krishnamurthy & Levoy, 1996), Hare (Turk & Levoy, 1994), Cow (Crane, 2021), Dragon (Curless & Levoy, 1996), and Horse (Praun et al., 2000). We generate  $10^5$  training samples,  $10^4$  validation samples, and  $10^4$  evaluation samples.

The learned LCMs consist again of an SCM prior, which is the same as in the 2D toy experiment, an encoder, and a decoder. For the encoder and decoder we use a convolutional architecture with four residual blocks, using downsampling via average-pooling and bilinear upsampling, respectively. We do not use BatchNorm, as we found that that can lead to practical issues when images in a batch are very similar.

Our training setup is as in the 2D toy experiment, except that we use a batch size of 64, train for  $2.3 \cdot 10^5$  steps, and use an initial learning rate of  $3 \cdot 10^{-4}$ . We find it beneficial to begin training with a lower weight of the KL divergence in the VAE loss,  $\beta = 0.01$ , and increasing this until the final value of  $\beta = 0.1$  during the first half of training. We initialize the manifold “fuzziness” parameter to 0.2 and anneal it to 0.01 over the first half of training. For each method and each graph, we train three models with different random seeds and select the model with the best validation log likelihood.

We report our results in Tbl. 2 and visualize the disentanglement properties of the learned representations in Fig. 6.