

# REPRESENTATION LEARNING FOR DISTRIBUTIONAL PERTURBATION EXTRAPOLATION

Anonymous authors  
Paper under double-blind review

## ABSTRACT

We consider the problem of modelling the effects of perturbations such as gene knockdowns or drug combinations on low-level measurements like RNA sequencing data. Specifically, given data collected under some perturbations, we aim to predict the distribution of measurements for new perturbations. To address this challenging extrapolation task, we posit that perturbations act additively in a suitable, unknown embedding space. More precisely, we formulate the generative process underlying the observed data as a latent variable model, in which perturbations amount to mean shifts in latent space. We prove that the representation and perturbation effects are identifiable up to affine transformation and use this to characterize the class of unseen perturbations for which we obtain extrapolation guarantees. To estimate the model from data, we propose the perturbation distribution autoencoder (PDAE) which is trained by maximising the distributional similarity between true and predicted perturbation distributions. The trained model can then be used to predict previously unseen perturbation distributions. Preliminary empirical evidence suggests that PDAE compares favourably to CPA (Lotfollahi et al., 2023) and other baselines at predicting the effects of unseen perturbations.

## 1 INTRODUCTION

Due to technological progress, large-scale perturbation data is becoming more abundant across several scientific fields. This is particularly the case for single-cell biology, where advancements in gene editing, sequencing, and mass spectrometry have led to the collection of vast transcriptomic or proteomic databases for various drug and gene perturbations (Dixit et al., 2016; Jinek et al., 2012; Norman et al., 2019; Wang et al., 2009; Weinstein et al., 2013). However, the exponential number of possible combinations of perturbations renders exhaustive experimentation prohibitive. Observations are thus typically only available for a subset of perturbations of interest, e.g., some single and double gene knockdowns or certain dosages of drugs. This necessitates models capable of extrapolating to unseen combinations of perturbations, e.g., new multi-gene knockdowns or dosage combinations.

**Prior Work.** Several recent works leverage machine learning for biological perturbation modelling, e.g., to generalize to new cell types (Bunne et al., 2023; Lotfollahi et al., 2019), unseen combinations of perturbations (Lotfollahi et al., 2023), or entirely new perturbations by leveraging the molecular structure of the involved compounds (Hetzl et al., 2022; Qi et al., 2024; Yu & Welch, 2022) or prior knowledge about gene-gene interactions (Kamimoto et al., 2023; Roohani et al., 2024). A common theme is the use of representation learning techniques such as autoencoders (Hinton & Salakhutdinov, 2006; Kingma & Welling, 2014; Rumelhart et al., 1986) to embed observations in a latent space, in which the effects of perturbations are assumed to take on a simpler (e.g., additive) form. However, existing studies are purely empirical and lack theoretical underpinning. Despite promising results, the capabilities and fundamental limitations of existing methods thus remain poorly understood.

**Overview and Contributions.** In this work, we present a principled, theoretically-grounded approach for perturbation extrapolation. Given the unpaired nature of the available data (each cell is only measured under one experimental condition), we consider the task of predicting population-level effects of perturbations, which we formalize as a distributional regression problem (§ 2). We then postulate a generative model (§ 3) which, similar to prior works, assumes that perturbations act as mean shifts in a suitable latent space, see Fig. 1 for an overview. We analyse this model

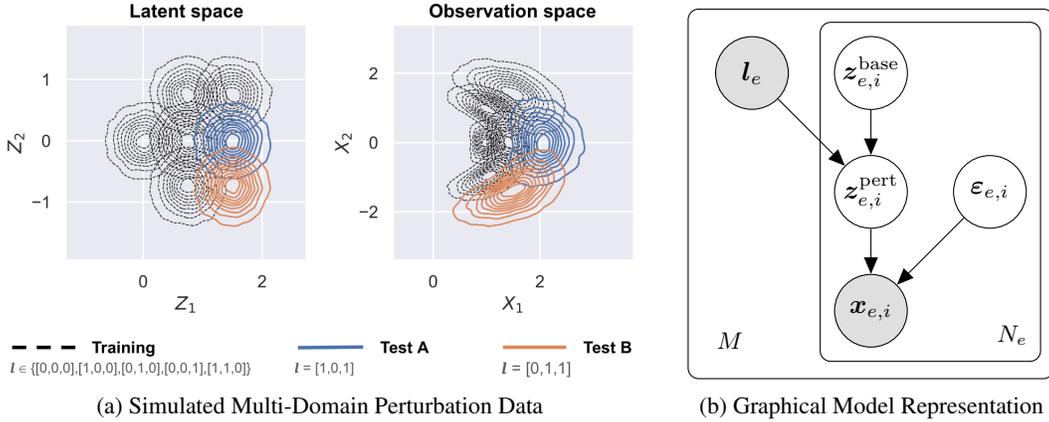


Figure 1: **Task Description and Assumed Data Generating Process.** (a) During training, we are given  $M = 5$  training data sets in observation space (right, grey), each of which is generated under a known combination of  $K = 3$  elementary perturbations. The corresponding perturbation labels  $\mathbf{l}_e$  are shown below the plots. During testing, we are given a new perturbation label and the task is to predict the corresponding distribution in observation space (right, blue and orange). We tackle this task by assuming that the effect of perturbations is linear additive in a suitable latent space (left). Both plots show kernel density estimates of the distributions. (b) Each dataset comprises a perturbation label  $\mathbf{l}_e$  and  $N_e$  observations  $\mathbf{x}_{e,i}$ . Perturbations act as mean-shifts on a latent basal state,  $\mathbf{z}_{e,i}^{\text{pert}} = \mathbf{z}_{e,i}^{\text{base}} + \mathbf{W}\mathbf{l}_e$ . A stochastic nonlinear decoder with noise  $\varepsilon_{e,i}$  then yields the observed  $\mathbf{x}_{e,i} = \mathbf{f}(\mathbf{z}_{e,i}^{\text{pert}}, \varepsilon_{e,i})$ . Shaded and white nodes indicate observed and unobserved/latent variables, respectively.

class theoretically (§ 4), proving that the latent representation and the relative training perturbation effects are identifiable up to affine transformation (Thm. 4.1). This result implies extrapolation guarantees for unseen perturbations that can be expressed as linear combinations of training perturbations (Thm. 4.6). Based on these insights, we devise an autoencoder-based estimation method (§ 5) which uses the energy score (Gneiting & Raftery, 2007) to assess the distributional similarity between predicted and ground-truth perturbation data. In preliminary simulations (§ 6), our approach compares favourably to the compositional perturbation autoencoder (CPA; Lotfollahi et al., 2023) in terms of mean prediction and distributional fit.

**Notation.** We write scalars as  $a$ , column-vectors as  $\mathbf{a}$ , and matrices as  $\mathbf{A}$ . We use uppercase for random variables and lowercase for their realizations. Equality in distribution is denoted by  $\stackrel{d}{=}$  and the pushforward of a distribution  $\mathbb{P}$  by a measurable function  $f$  is denoted by  $f_{\#}\mathbb{P}$ . The Euclidean (L2) norm is denoted by  $\|\cdot\|$ . Further, we use the shorthands  $[n] = \{1, \dots, n\}$  and  $[n]_0 = [n] \cup \{0\}$ .

## 2 PROBLEM SETTING: DISTRIBUTIONAL PERTURBATION EXTRAPOLATION

Let  $\mathbf{x} \in \mathbb{R}^{d_x}$  be an observation (e.g., omics data) that is obtained under one of several possible experimental conditions. We model these conditions as combinations of  $K$  elementary perturbations, each of which we assume can be encoded by a real number. Further, let  $\mathbf{l} \in \mathbb{R}^K$  be a perturbation label that indicates if, or how much of, each perturbation was applied before collecting the corresponding  $\mathbf{x}$ .

*Example 2.1 (Gene perturbations).* For data arising from gene knockouts, a perturbation can be represented by a binary  $\mathbf{l} \in \{0, 1\}^K$ , where  $K$  denotes the number of potential targets and  $l_k = 1$  if and only if target  $k$  was subject to a knockout experiment. For example,  $\mathbf{l} = (1, 0, 1, 0, 0)$  indicates a multi-gene knockout on targets one and three.

*Example 2.2 (Drug perturbations).* For data arising from applying varying amounts of  $K$  different drugs, perturbations can be represented by continuous, non-negative labels  $\mathbf{l} \in \mathbb{R}_+^K$ , where  $l_k$  indicates the (relative or absolute) amount of drug  $k$  that was administered.

We have access to  $M+1$  experimental datasets  $\mathcal{D}_0, \mathcal{D}_1, \dots, \mathcal{D}_M$ , each comprising a sample of  $N_e$  observations  $\mathbf{x}$  and a perturbation label  $\mathbf{l}$ , i.e., for all experiments or environments  $e \in [M]_0$ ,

$$\mathcal{D}_e = \left( \{\mathbf{x}_{e,i}\}_{i=1}^{N_e}, \mathbf{l}_e \right). \quad (2.1)$$

Given data in the form of (2.1), the task we consider is to predict the effects of new perturbations  $l_{\text{test}} \notin \{l_0, l_1, \dots, l_M\}$  without observing any data from this condition (“zero-shot”). In particular, we are interested in the distribution over observations resulting from  $l_{\text{test}}$ . That is, we aim to leverage the training domains (2.1) to learn a map

$$l \mapsto \mathbb{P}_{\mathbf{X}|l} \quad (2.2)$$

which extrapolates beyond the training support of  $l$ , i.e., the predictions should remain reliable for new inputs  $l_{\text{test}}$ . Since (2.2) targets the full conditional distribution—rather than, say, the conditional mean  $\mathbb{E}[\mathbf{X}|l]$ —it constitutes a (multi-variate) distributional regression task (Koenker, 2005; Koenker & Bassett Jr, 1978), also referred to as probabilistic forecasting (Gneiting & Raftery, 2007) or conditional generative modelling (Mirza, 2014; Sohn et al., 2015; Winkler et al., 2019). We therefore refer to our problem setting as *distributional perturbation extrapolation*.

Formally, extrapolation means that the value of the function (2.2) at  $l_{\text{test}}$  is determined by its values on the training support  $\{l_0, l_1, \dots, l_M\}$ . Intuitively, for this to be feasible,  $l_{\text{test}}$  must be somehow related to the training perturbations  $l_e$ . For example, given data resulting from individual perturbations, predict the effects of combinations thereof. This type of extrapolation to new combinations of inputs is also called compositional generalization (Goyal & Bengio, 2022; Lake et al., 2017). It is known to be challenging (Montero et al., 2022; 2021; Schott et al., 2022) and requires assumptions that sufficiently constrain the model class (Brady et al., 2023; 2025; Dong & Ma, 2022; Lachapelle et al., 2023; Lippl & Stachenfeld, 2024; Wiedemer et al., 2024a;b).

### 3 MODEL: PERTURBATIONS AS MEAN SHIFTS IN LATENT SPACE

We now specify a generative process for the observed data in (2.1). In so doing, we aim to strike a balance between imposing sufficient structure on (2.2) to facilitate extrapolation, while remaining flexible enough to model the complicated, nonlinear effects which perturbations may have on the distribution of observations. Similar to Lotfollahi et al. (2023), we therefore model the effect of perturbations in a latent space with perturbation-relevant latent variables  $z \in \mathbb{R}^{d_z}$ , which are related to the observations  $\mathbf{x}$  via a nonlinear (stochastic) mixing function or generator  $\mathbf{f}$ . The full generative process amounts to a hierarchical latent variable model, which additionally contains noise variables  $\varepsilon$  that capture other variation underlying the observations  $\mathbf{x}$ , and which is represented as a graphical model in Fig. 1b. Specifically, we posit for all  $e \in [M]_0 \cup \{\text{test}\}$  and all  $i \in [N_e]$ :

$$z_{e,i}^{\text{base}} \sim \mathbb{P}_{\mathbf{Z}}, \quad (3.1)$$

$$z_{e,i}^{\text{pert}} := z_{e,i}^{\text{base}} + \mathbf{W}l_e, \quad (3.2)$$

$$\varepsilon_{e,i} \sim \mathbb{Q}_{\varepsilon}, \quad (3.3)$$

$$\mathbf{x}_{e,i} := \mathbf{f}\left(z_{e,i}^{\text{pert}}, \varepsilon_{e,i}\right), \quad (3.4)$$

where  $(z_{e,i}^{\text{base}})_{e \in [M]_0, i \in [N_e]}$  are independent and identically distributed (i.i.d.) according to  $\mathbb{P}_{\mathbf{Z}}$ , and  $(\varepsilon_{e,i}^{\text{base}})_{e \in [M]_0, i \in [N_e]}$  are i.i.d. according to  $\mathbb{Q}_{\varepsilon}$  and jointly independent of  $(z_{e,i}^{\text{base}})_{e \in [M]_0, i \in [N_e]}$ .

The basal state  $z^{\text{base}}$  in (3.1) describes the unperturbed state of latent variables, which can, in principle, be affected by perturbations, and is distributed according to a base distribution  $\mathbb{P}_{\mathbf{Z}}$ . The perturbation matrix  $\mathbf{W} \in \mathbb{R}^{d_z \times K}$  in (3.2) captures the effect of the  $K$  elementary perturbations encoded in  $l$  on the latents and turns  $z^{\text{base}}$  into perturbed latents  $z^{\text{pert}}$ . Since the same perturbation  $l_e$  is applied for all  $i \in [N_e]$ , all intra-dataset variability in  $z^{\text{pert}}$  is due to  $\mathbb{P}_{\mathbf{Z}}$ . The noise variables  $\varepsilon \in \mathbb{R}^{d_\varepsilon}$  in (3.3) capture all other variation in the observed data that is unaffected by perturbations. It is distributed according to a fixed, uninformative distribution  $\mathbb{Q}_{\varepsilon}$  such as a standard isotropic Gaussian. The noise serves as an additional input to the (stochastic) mixing function or generator  $\mathbf{f} : \mathbb{R}^{d_z} \times \mathbb{R}^{d_\varepsilon} \rightarrow \mathbb{R}^{d_x}$  in (3.4), which produces observations for the perturbed latent. This generative model allows us to model any conditional distribution  $\mathbb{P}_{\mathbf{X}|z}$  and is more flexible than, e.g., a Gaussian decoder with mean and covariance parametrised by  $\mathbf{f}$  as used by Lotfollahi et al. (2023).

For a given  $e$  and  $l_e$ , the generative process in (3.1)–(3.4) induces a distribution  $\mathbb{P}_e$  over observations  $\mathbf{x}$ , which we also denote by  $\mathbb{P}_{\mathbf{X}|l_e}$ , defined as the push-forward of  $\mathbb{P}_{\mathbf{Z}}$  and  $\mathbb{Q}_{\varepsilon}$  through (3.2) and (3.4), such that:

$$\forall e \in [M]_0 : \quad (\mathbf{x}_{e,i})_{i \in [N_e]} \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}_e. \quad (3.5)$$

## 4 THEORY: IDENTIFIABILITY AND EXTRAPOLATION GUARANTEES

In this section, we present our theoretical analysis for the model class introduced in § 3.

**Identifiability.** We first study identifiability, that is, the question under what assumptions and up to what ambiguities certain parts of the postulated generative process can be provably recovered assuming access to the full distributions. As established by the following results, our model class is identifiable up to affine transformation, provided that the training perturbations are sufficiently diverse, the dimension of the latent space is known and some additional technical assumptions hold.

**Theorem 4.1** (Affine identifiability for Gaussian latents). *For  $M \in \mathbb{Z}_{\geq 0}$ , let  $l_0, \dots, l_M \in \mathbb{R}^K$  be  $M+1$  perturbation labels. Let  $f, \tilde{f} : \mathbb{R}^{d_Z} \rightarrow \mathbb{R}^{d_X}$ ,  $\mathbf{W}, \tilde{\mathbf{W}} \in \mathbb{R}^{d_Z \times K}$ , and  $\mathbb{P}, \tilde{\mathbb{P}}$  be distributions on  $\mathbb{R}^{d_Z}$  such that the models  $(f, \mathbf{W}, \mathbb{P})$  and  $(\tilde{f}, \tilde{\mathbf{W}}, \tilde{\mathbb{P}})$  induce the same observed distributions, i.e.,*

$$\forall e \in [M]_0 : \quad f(\mathbf{Z} + \mathbf{W}l_e) \stackrel{d}{=} \tilde{f}(\tilde{\mathbf{Z}} + \tilde{\mathbf{W}}l_e), \quad \text{where } \mathbf{Z} \sim \mathbb{P} \text{ and } \tilde{\mathbf{Z}} \sim \tilde{\mathbb{P}}. \quad (4.1)$$

Assume further that:

- (i) **[invertibility]**  $f$  and  $\tilde{f}$  are  $C^2$ -diffeomorphisms onto their respective images;
- (ii) **[Gaussianity]**  $\mathbf{Z}$  and  $\tilde{\mathbf{Z}}$  are non-degenerate multi-variate Gaussians, i.e.,  $\mathbb{P} = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  and  $\tilde{\mathbb{P}} = \mathcal{N}(\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}})$  for some  $\boldsymbol{\mu}, \tilde{\boldsymbol{\mu}} \in \mathbb{R}^{d_Z}$  and positive-definite  $\boldsymbol{\Sigma}, \tilde{\boldsymbol{\Sigma}} \in \mathbb{R}^{d_Z \times d_Z}$ ;
- (iii) **[sufficient diversity]** the matrix  $\tilde{\mathbf{W}}\mathbf{L} \in \mathbb{R}^{d_Z \times M}$ , where  $\mathbf{L} \in \mathbb{R}^{K \times M}$  is the matrix with columns  $(l_e - l_0)$  for  $e \in [M]$ , has full row rank, i.e.,  $\text{rank}(\tilde{\mathbf{W}}\mathbf{L}) = d_Z$ .

Then the latent representation and the effects of the observed perturbation combinations relative to  $l_0$  (as captured by  $\mathbf{W}\mathbf{L}$ ) are identifiable up to affine transformation in the following sense:

$$\forall z : \quad \tilde{f}^{-1} \circ f(z) = \mathbf{A}z + \mathbf{b}, \quad (4.2)$$

$$\tilde{\mathbf{W}}\mathbf{L} = \mathbf{A}\mathbf{W}\mathbf{L}, \quad (4.3)$$

where  $\mathbf{A} := \tilde{\boldsymbol{\Sigma}}^{\frac{1}{2}}\boldsymbol{\Sigma}^{-\frac{1}{2}}$  and  $\mathbf{b} := \tilde{\boldsymbol{\mu}} - \mathbf{A}\boldsymbol{\mu} + (\tilde{\mathbf{W}} - \mathbf{A}\mathbf{W})l_0$ .

**Corollary 4.2** (Affine recovery of the perturbation matrix). *If, in addition to the assumptions of Thm. 4.1,  $\mathbf{L} \in \mathbb{R}^{K \times M}$  has full row rank (i.e.,  $\text{rank}(\mathbf{L}) = K \leq M$ ), then the perturbation matrix  $\mathbf{W}$  is identifiable up to affine transformation in the sense that*

$$\tilde{\mathbf{W}} = \mathbf{A}\mathbf{W}, \quad (4.4)$$

for  $\mathbf{A} := \tilde{\boldsymbol{\Sigma}}^{\frac{1}{2}}\boldsymbol{\Sigma}^{-\frac{1}{2}}$ . In this case, the expression for  $\mathbf{b}$  in (4.2) simplifies to  $\mathbf{b} = \tilde{\boldsymbol{\mu}} - \mathbf{A}\boldsymbol{\mu}$ .

The proofs of Thm. 4.1 and Cor. 4.2 are provided in Appx. B.1 and B.2.

**Discussion.** Thm. 4.1 can be interpreted as follows. Fix a set of perturbation labels  $(l_0, \dots, l_M)$  and a data generating process parametrised by  $(f, \mathbf{W}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$ . Then, for all  $\tilde{\boldsymbol{\mu}}$  and  $\tilde{\boldsymbol{\Sigma}}$ , and for all  $\tilde{\mathbf{W}}$  such that (4.3) holds for  $\mathbf{A} = \tilde{\boldsymbol{\Sigma}}^{\frac{1}{2}}\boldsymbol{\Sigma}^{-\frac{1}{2}}$ , there exists a unique  $\tilde{f}$ , characterized by (4.2), such that  $(\tilde{f}, \tilde{\mathbf{W}}, \tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}})$  gives rise to the same observed distributions of  $\mathbf{X}_e$ . At the same time, any  $(\tilde{f}, \tilde{\mathbf{W}}, \tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}})$  for which this holds for all  $e \in [M]_0$  is of this form. In other words, the mean  $\boldsymbol{\mu}$  and covariance  $\boldsymbol{\Sigma}$  of the basal state are completely unidentifiable, but the mixing function  $f$  and the relative shift matrix  $\mathbf{W}\mathbf{L}$  are identifiable up to affine transformation—provided that the observed training perturbation conditions are sufficiently diverse, as formalised by assumption (iii).

**Remark 4.3** (Sufficient diversity). The matrix product  $\mathbf{W}\mathbf{L} \in \mathbb{R}^{d_Z \times M}$  captures the relative effects of the observed perturbations. Specifically,  $(\mathbf{W}\mathbf{L})_{je}$  corresponds to the shift in the  $j^{\text{th}}$  latent  $Z_j$  resulting from  $l_e$ , relative to a reference condition  $l_0$ . Moreover, assumption (iii) of Thm. 4.1 implies

$$\min \left\{ \text{rank}(\tilde{\mathbf{W}}), \text{rank}(\mathbf{L}) \right\} \geq \text{rank}(\tilde{\mathbf{W}}\mathbf{L}) = d_Z. \quad (4.5)$$

Hence, sufficient diversity requires at least  $d_Z$  elementary perturbations whose associated shift vectors  $w_k \in \mathbb{R}^{d_Z}$  are linearly independent, and we must observe at least  $d_Z$  perturbation conditions  $l_e$  other than  $l_0$  such that the relative perturbation vectors  $(l_e - l_0) \in \mathbb{R}^K$  are linearly independent.

*Remark 4.4* (Choice of reference). Since the environments are unordered, the choice of reference environment is arbitrary. Here, we choose  $e = 0$  as reference without loss of generality. Intuitively, if a perturbation is always present (e.g.,  $l_{e,1} = 1$  for all  $e$ ), then its effects cannot be discerned from the basal state. Therefore, only the effects of the relative perturbations ( $l_e - l_0$ ) can be recovered. In practice, we often have access to an unperturbed, purely observational control condition with  $l_0 = \mathbf{0}$ .

*Remark 4.5* (Deterministic vs noisy mixing.). The mixing function  $f$  in Thm. 4.1 is deterministic, i.e., does not take a separate noise variable  $\varepsilon$  as input, cf. (3.4). In principle, noise can be appended to  $z^{\text{pert}}$  as additional dimensions that are not influenced by perturbations. However, this increases  $d_Z$  and thus makes it harder to satisfy sufficient diversity (see Remark 4.3). Alternatively, the setting of additive noise,  $f(Z) + \varepsilon$ , can be reduced to the noiseless case (Khemakhem et al., 2020).

**From Identifiability to Extrapolation.** Since we aim to make distributional predictions for new perturbations  $l_{\text{test}}$ , identifiability is only of intermediary interest. The following result highlights the usefulness of the affine identifiability established in Thm. 4.1 for extrapolation. In particular, it shows that this allows us to uniquely predict the observable effects of certain unseen perturbations—specifically, those which can be expressed as linear combinations of the observed perturbations.

**Theorem 4.6** (Extrapolation to span of relative perturbations). *Under the same setting and assumptions as in Thm. 4.1, let  $l_{\text{test}} \in \mathbb{R}^K$  be an unseen perturbation label such that*

$$(l_{\text{test}} - l_0) \in \text{span} \left( \{l_e - l_0\}_{e \in [M]} \right). \quad (4.6)$$

*Then the effect of  $l_{\text{test}}$  is uniquely identifiable in the sense that*

$$\mathbf{X}_{\text{test}} = f(Z + \mathbf{W}l_{\text{test}}) \stackrel{d}{=} \tilde{f} \left( \tilde{Z} + \tilde{\mathbf{W}}l_{\text{test}} \right) = \tilde{\mathbf{X}}_{\text{test}}. \quad (4.7)$$

The proof of Thm. 4.6 is provided in Appx. B.3.

*Remark 4.7* (Additive vs linear perturbations). For our affine identifiability result (Thm. 4.1), it is not necessary that the mean shifts  $\mathbf{W}l_e$  are linear in  $l_e$ . If we replace  $\mathbf{W}l_e$  and  $\tilde{\mathbf{W}}l_e$  in (4.1) with arbitrary shift vectors  $c_e, \tilde{c}_e \in \mathbb{R}^{d_Z}$ , the same result can be shown to hold with  $\mathbf{W}L$  and  $\tilde{\mathbf{W}}L$  replaced by  $C$  and  $\tilde{C}$ , defined as the matrices with columns  $(c_e - c_0)$  and  $(\tilde{c}_e - \tilde{c}_0)$ , respectively. That is, the relative shift vectors are identifiable up to affine transformation, regardless of whether they are linear in  $l$ . This has implications, e.g., for the CPA model of Lotfollahi et al. (2023) which includes element-wise nonlinear dose-response functions applied to  $l$ , see Appx. C.4 for details. However, linearity is leveraged in the proof of our extrapolation result (Thm. 4.6) where it is used in (4.6) to establish a link between  $l_{\text{test}}$  and the training perturbations. Since only  $l_{\text{test}}$  is observed at test time, the above argument thus cannot easily be extended to the extrapolation setting, as this would require establishing a link between  $c_{\text{test}}$  and the training shifts, all of which are unobserved.

## 5 ESTIMATION METHOD: PERTURBATION DISTRIBUTION AUTOENCODER

To leverage the extrapolation guarantees of Thm. 4.6, we seek to estimate the parts of the generative process that are identifiable according to Thm. 4.1 from the available data in (2.1). To this end, we build on an autoencoder framework and adapt it for multi-domain perturbation modelling and distributional regression. Our method, the *perturbation distribution autoencoder* (PDAE), comprises an encoder, a perturbation matrix, and a (stochastic) decoder, trained to maximise the similarity between pairs of true and simulated perturbation distributions, see Fig. 2 for an overview.

**Encoder.** The encoder  $\hat{g} : \mathbb{R}^{d_x} \rightarrow \mathbb{R}^{d_z}$  maps observations  $x$  to the space of perturbation-relevant latents  $z$ . Ideally, it should invert the stochastic mixing function in (3.4) in the sense of recovering the perturbed latent state  $z_{e,i}^{\text{pert}}$  in (3.2) from observation  $x_{e,i}$ . We therefore denote the encoder output by

$$\hat{z}_{e,i}^{\text{pert}} := \hat{g}(x_{e,i}), \quad (5.1)$$

and refer to it as estimated perturbed latent. This is a key difference to the CPA method of Lotfollahi et al. (2023), which seeks an encoder that maps to the latent basal state, regardless of the domain  $e$ .

**Perturbation Model.** If the encoder recovers the perturbed latents up to affine transformation, the additivity of perturbation effects assumed in (3.2) allows us to map between the latent distributions

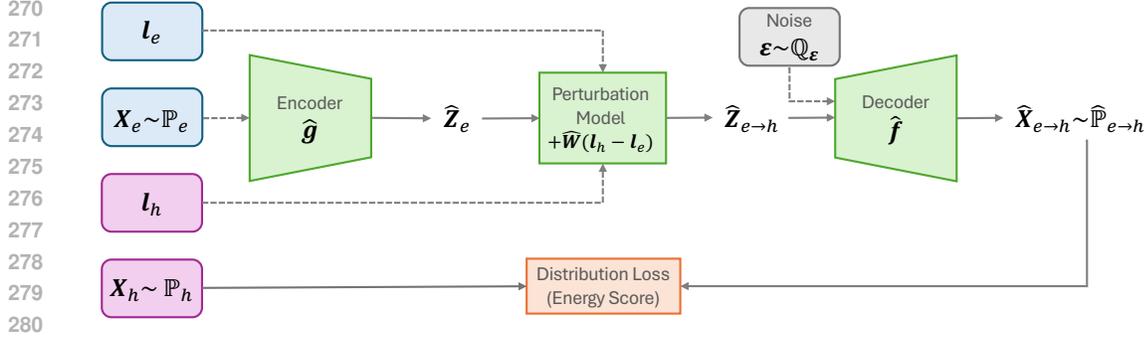


Figure 2: **Overview of the Perturbation Distribution Autoencoder (PDAE).** The distribution of a target perturbation condition  $h$  (purple) is simulated by encoding, perturbing, and decoding data from a source condition  $e$  (blue). Dashed arrows indicate model inputs and green boxes model components with learnable parameters, trained to maximise the similarity (orange) between the empirical true and simulated target distributions for all pairs of training domains  $(e, h)$ . At test time, the target perturbation label  $l_h$  is replaced with an unseen  $l_{\text{test}}$ .

underlying different perturbation conditions. Specifically, we use a perturbation model parametrised by a perturbation matrix  $\widehat{\mathbf{W}} \in \mathbb{R}^{d_z \times K}$  to create synthetic perturbed latents from domain  $h$  as:

$$\widehat{\mathbf{z}}_{e \rightarrow h, i}^{\text{pert}} := \widehat{\mathbf{z}}_{e, i}^{\text{pert}} + \widehat{\mathbf{W}}(\mathbf{l}_h - \mathbf{l}_e). \quad (5.2)$$

which can be interpreted as undoing the effects of perturbation  $l_e$  (i.e., mapping back to the latent basal state) and then simulating perturbation  $l_h$ . (For  $h = e$ , this has no effect and  $\widehat{\mathbf{z}}_{e \rightarrow h, i}^{\text{pert}} = \widehat{\mathbf{z}}_{e, i}^{\text{pert}}$ .)

**Decoder.** The decoder  $\widehat{\mathbf{f}} : \mathbb{R}^{d_z} \times \mathbb{R}^{d_e} \rightarrow \mathbb{R}^{d_x}$  maps estimated latents  $\widehat{\mathbf{z}}$  and noise  $\varepsilon \sim \mathbb{Q}_\varepsilon$  back to observations. When viewed as a function of  $\widehat{\mathbf{z}}$  only, it is stochastic and induces the distribution  $\widehat{\mathbf{f}}(\widehat{\mathbf{z}}, \cdot)_{\#} \mathbb{Q}_\varepsilon$  from which we can sample synthetic observations,

$$\widehat{\mathbf{X}}_{e \rightarrow h, i} = \widehat{\mathbf{f}}\left(\widehat{\mathbf{z}}_{e \rightarrow h, i}^{\text{pert}}, \varepsilon\right) \quad \text{where} \quad \varepsilon \sim \mathbb{Q}_\varepsilon. \quad (5.3)$$

**Simulating Perturbation Distributions.** Given a distribution  $\mathbb{P}_e$  and the corresponding perturbation label  $l_e$ , our model facilitates sampling synthetic observations for another perturbation condition with label  $l_h$ . We denote the resulting distribution by  $\widehat{\mathbb{P}}_{e \rightarrow h}$ , formally defined as the distribution of

$$\widehat{\mathbf{f}}\left(\widehat{\mathbf{g}}(\mathbf{X}_e) + \widehat{\mathbf{W}}(\mathbf{l}_h - \mathbf{l}_e), \varepsilon\right) \quad \text{where} \quad \mathbf{X}_e \sim \mathbb{P}_e \quad \text{and} \quad \varepsilon \sim \mathbb{Q}_\varepsilon. \quad (5.4)$$

**Learning Objective.** To learn  $(\widehat{\mathbf{g}}, \widehat{\mathbf{f}}, \widehat{\mathbf{W}})$ , we propose to minimise the pairwise distribution loss

$$\mathcal{L}\left(\widehat{\mathbf{g}}, \widehat{\mathbf{f}}, \widehat{\mathbf{W}}; \{(\mathbb{P}_e, \mathbf{l}_e)\}_{e \in [M]_0}\right) = \sum_{e, h \in [M]_0} d\left(\widehat{\mathbb{P}}_{e \rightarrow h}, \mathbb{P}_h\right) \quad (5.5)$$

where  $\widehat{\mathbb{P}}_{e \rightarrow h}$  depends on  $(\mathbb{P}_e, \mathbf{l}_e, \mathbf{l}_h)$  and the model parameters via (5.4); and  $d$  is a measure of dissimilarity between distributions. Here, we use the negative expected energy score for  $d$ , i.e.,

$$d\left(\widehat{\mathbb{P}}_{e \rightarrow h}, \mathbb{P}_h\right) = -\mathbb{E}_{\mathbf{X}_h \sim \mathbb{P}_h} \left[ \text{ES}_\beta\left(\widehat{\mathbb{P}}_{e \rightarrow h}, \mathbf{X}_h\right) \right], \quad (5.6)$$

where  $\text{ES}_\beta$  denotes the energy-score (Gneiting & Raftery, 2007), defined for  $\beta \in (0, 2)$  as

$$\text{ES}_\beta(\mathbb{P}, \mathbf{x}) = \frac{1}{2} \mathbb{E}_{\mathbf{X}, \mathbf{X}' \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}} \|\mathbf{X} - \mathbf{X}'\|^\beta - \mathbb{E}_{\mathbf{X} \sim \mathbb{P}} \|\mathbf{X} - \mathbf{x}\|^\beta. \quad (5.7)$$

It is a strictly proper scoring rule, meaning that the expected energy score  $\mathbb{E}_{\mathbf{X}}[\text{ES}_\beta(\mathbb{P}, \mathbf{X})]$  is maximised if and only if  $\mathbf{X} \sim \mathbb{P}$ , see Appx. C.1 for details on probabilistic forecasting and scoring rules. Combined with its computational simplicity, this property makes the negative expected energy score a popular loss function for distributional regression (Shen & Meinshausen, 2024a;b).

**Corollary 5.1.** *The objective in (5.5) is minimised if and only if  $\mathbb{P}_h = \widehat{\mathbb{P}}_{e \rightarrow h}$  for all  $e, h \in [M]_0$ .*

**Training.** Since we only have access to empirical distributions, we approximate the expectations in (5.6) and (5.7) with Monte Carlo samples based on the available data (2.1). Moreover, for a fixed encoder, the optimal perturbation matrix is available in closed form and given by the ordinary least squares solution to regressing the domain-specific means of  $\widehat{z}_{e,i}^{\text{pert}}$  on the corresponding perturbation labels  $l_e$ . We, therefore, write  $\widehat{W}$  as a function of  $\widehat{g}$  and optimise (5.5) with respect to the parameters of  $\widehat{g}$  and  $\widehat{f}$  using stochastic gradient descent (Kingma & Ba, 2015).

**Prediction.** To simulate the distribution for a new perturbation label  $l_{\text{test}}$ , we use our model to compute the synthetic perturbed test latents  $\widehat{z}_{e \rightarrow \text{test}, i}^{\text{pert}}$  for all  $e \in [M]_0$  and all  $i \in [N_e]$  via (5.2), and then sample the corresponding synthetic test observations  $\widehat{X}_{e \rightarrow \text{test}, i}$  according to (5.3). In other words, our estimate of  $\mathbb{P}_{\text{test}}$  is the pooled version of the domain-specific empirical synthetic distributions,

$$\widehat{\mathbb{P}}_{\text{test}} = \frac{1}{M+1} \sum_{e \in [M]_0} \widehat{\mathbb{P}}_{e \rightarrow \text{test}} \quad (5.8)$$

## 6 EXPERIMENTS

We present preliminary empirical evidence that our approach can outperform existing methods at distributional perturbation extrapolation. As we focus on a simple, controlled setting with synthetic data, our results should be viewed as proof of concept, rather than as comprehensive empirical study.

**Data.** For ease of visualisation, we consider  $d_Z = d_X = 2$ -dimensional latents and observations. The base distribution  $\mathbb{P}_Z$  is a zero-mean, isotropic Gaussian with standard deviation  $\sigma = 0.25$ . We consider  $K = 3$  elementary perturbations with associated shift vectors  $w_1 = (1, 0)^\top$ ,  $w_2 = (0, 1)^\top$ , and  $w_3 = (1, 1)^\top$ . We create  $M+1 = 4$  training domains with labels  $l_0 = (0, 0, 0)^\top$ ,  $l_1 = (1, 0, 0)^\top$ ,  $l_2 = (0, 1, 0)^\top$ , and  $l_3 = (0, 0, 1)^\top$ , and test on  $l_{\text{test}}^{\text{ID}} = (1, 1, 0)^\top$  and  $l_{\text{test}}^{\text{OOD}} = (1, 0, 1)^\top$ . By construction,  $l_{\text{test}}^{\text{ID}}$  results in the same mean shift of  $(1, 1)^\top$  as  $l_3$ , whereas  $l_{\text{test}}^{\text{OOD}}$  results in a different shift not seen during training. We, therefore, refer to the respective test cases as in-distribution (ID) and out-of-distribution (OOD) relative to the decoder inputs seen during training. To generate observations, we use the complex exponential  $x = f(z) = e^{z_1}(\cos z_2, \sin z_2)$  as a deterministic nonlinear mixing function, which was also used to generate Fig. 1a (where both test cases are partially OOD). The resulting datasets are in shown in Fig. 3 (left) in Appx. A.

**Methods.** We compare our approach with CPA (Lotfollahi et al., 2023) and the following baselines:

*Pooled Mean:* pool all training observations, then output the mean;

*Pseudobulked Mean:* pool only data arising from individual perturbations involved in the combination to be predicted (e.g.,  $l_1$  and  $l_2$  for  $l_{\text{test}}^{\text{ID}}$ ), then output the mean;

*Linear Regression:* linearly regress the environment-specific means of observations  $\mu_x^e$  on  $l_e$  and use the resulting model to predict the test means  $\mu_x^{\text{test}}$  from  $l_{\text{test}}$ .

The former two were used by Lotfollahi et al. (2023); we propose the latter as an additional baseline.

**Metrics.** To assess distributional fit, we use the energy distance (ED; Székely & Rizzo, 2013), i.e., twice the normalized negative expected energy score, and the maximum mean discrepancy (MMD; Gretton et al., 2012) with Gaussian kernel and bandwidth chosen by the median heuristic, see Appx. C.2 for more on measures of distributional similarity. Since some methods only predict the mean, we also report the L2 norm of the difference between predicted and true mean,  $\|\mu_x - \widehat{\mu}_x\|$ .

**Experimental Details.** We generate  $N_e = 2^{14}$  observations for each domain. Both CPA and PDAE use 4-hidden layer MLPs with 64 hidden units as encoders and decoders and are trained for 2000 epochs using a batch size of  $2^{12}$ . For CPA, all other hyperparameters are set to their default values. For PDAE, we set  $\beta = 1$  in (5.7) and use a learning rate of 0.005.

**Results.** The quantitative results are summarized in Tab. 1, see also Fig. 3 in Appx. A for qualitative results. For the ID test setting, PDAE performs best, achieving near-perfect distributional fit. CPA outperforms linear regression but does substantially worse than PDAE at both mean and distributional prediction. For the OOD test setting, all methods perform much worse, with PDAE yielding the least bad performance. This failure for the OOD case is expected. Despite learning a good representation and perturbation model (as evident from the strong ID test performance), the decoder did not encounter inputs similar to the perturbed test latents during training and has thus not learnt which part of the observation space to map them to.

Table 1: **Results on Simulated Data.** For all metrics, lower is better. Best results highlighted in bold. Brackets indicate distributional similarities calculated using only the mean, i.e., a sample size of one.

Method	In-Distribution Test			Out-of-Distribution Test		
	ED	MMD	$\ \mu_x - \hat{\mu}_x\ _2$	ED	MMD	$\ \mu_x - \hat{\mu}_x\ _2$
Pooled Mean	(1.23)	(4.85)	0.82	(2.80)	(5.98)	1.74
Pseudobulked Mean	(1.49)	(5.43)	0.96	(2.45)	(5.56)	1.57
Linear Regression	(0.81)	(3.42)	0.60	(1.22)	(3.35)	0.88
CPA (Lotfollahi et al., 2023)	0.17	0.57	0.36	3.09	5.02	2.15
PDAE (Ours)	<b>0.001</b>	<b>0.005</b>	<b>0.03</b>	<b>0.45</b>	<b>1.33</b>	<b>0.61</b>

## 7 DISCUSSION

**Single-Cell vs Population-Level Effects.** Since cells are typically destroyed during measurement, each unit is only observed for one perturbation condition. As a result, we do not have access to paired data, which would be required to establish ground-truth single-cell level effects. Instead, we can only observe the effects of perturbations at the population level. We, therefore, formulate the perturbation extrapolation task, our theoretical guarantees, and our learning objective in distributional terms. This contrasts with some prior works, which, despite a lack of ground truth training data, pursue the seemingly infeasible task of making counterfactual predictions at the single-cell level. While our model can, in principle, also make such predictions, we stress that they are not falsifiable from empirical data and that our guarantees do not extend to the single-cell resolution. For this reason, we consider distributional perturbation extrapolation a more meaningful and feasible task formulation.

**Relation to Causal Models.** In the field of causal inference, experimental data resulting from perturbations is modelled via interventions in an underlying causal model. For example, in the structural causal model (SCM; Pearl, 2009) framework, interventions modify a set of assignments, which determine each variable from its direct causes and unexplained noise. In general, our model for the effect of perturbations in latent space (§ 3) differs from how interventions are treated in SCMs. However, as detailed in Appx. D, if we restrict our attention to linear SCMs, then the class of shift interventions can be viewed as a special case of our model with  $K = d_Z$  elementary perturbations and a particular choice of perturbation matrix  $\mathbf{W}$ . In this sense, our approach may be interpreted as causal representation learning (Schölkopf et al., 2021) with a linear latent causal model (Buchholz et al., 2024; Squires et al., 2023) and generalized shift interventions (Zhang et al., 2024).

**Is Gaussianity Necessary?** In Thm. 4.1, Gaussianity of  $\mathbb{P}_Z$  is assumed to prove identifiability. While this is a sufficient condition, it may not be necessary and can possibly be relaxed. Indeed, our method does not explicitly enforce any particular latent distribution. Moreover, preliminary empirical evidence suggests that training on simulated data generated from a Laplacian or Uniform basal state distribution also yields identifiable latent spaces and perturbation effects.

**Open Problems.** We consider the requirement for the decoder to generalize to new inputs the biggest open problem for distributional perturbation extrapolation. While this issue is absent from our theory, where  $\mathbb{P}_Z$  is Gaussian and thus has full support, it can pose major challenges in practice when learning from finite data. As discussed at the end of § 6, the issue of decoder extrapolation is orthogonal to learning the correct representation and perturbation model and thus appears fundamental. Besides addressing decoder extrapolation, future work should evaluate the proposed approach on more complex, noisy real-world data, extend our theoretical results to partially identifiable settings, and pursue extensions that incorporate covariates and nonlinear dose-response functions.

### MEANINGFULNESS STATEMENT

We consider a “meaningful representation of life” an embedding of biological data that facilitates drawing non-trivial inferences, such as predicting the effects of new (combinations of) interventions or generalizing to new cell types or species. In the present work, we focus on a class of representations which we show to be provably meaningful in this sense and propose a principled estimation method to learn such representations from multi-domain perturbation data.

## REFERENCES

- 432  
433  
434 Ludwig Baringhaus and Carsten Franz. On a new multivariate two-sample test. *Journal of Multi-*  
435 *variate Analysis*, 88(1):190–206, 2004. [Cited on p. 16.]  
436  
437 Mikhail Belkin and Partha Niyogi. Laplacian eigenmaps and spectral techniques for embedding and  
438 clustering. In *Advances in Neural Information Processing Systems*, volume 14, 2001. [Cited on  
439 p. 17.]  
440 Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new  
441 perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(8):1798–1828,  
442 2013. [Cited on p. 17.]  
443 Jack Brady, Roland S Zimmermann, Yash Sharma, Bernhard Schölkopf, Julius von Kügelgen, and  
444 Wieland Brendel. Provably learning object-centric representations. In *International Conference*  
445 *on Machine Learning*, pp. 3038–3062. PMLR, 2023. [Cited on p. 3.]  
446 Jack Brady, Julius von Kügelgen, Sebastien Lachapelle, Simon Buchholz, Thomas Kipf, and  
447 Wieland Brendel. Interaction asymmetry: A general principle for learning composable abstrac-  
448 tions. In *International Conference on Learning Representations*, 2025. [Cited on p. 3.]  
449 Simon Buchholz, Goutham Rajendran, Elan Rosenfeld, Bryon Aragam, Bernhard Schölkopf, and  
450 Pradeep Ravikumar. Learning linear causal representations from interventions under general non-  
451 linear mixing. In *Advances in Neural Information Processing Systems*, volume 36, 2024. [Cited  
452 on p. 8.]  
453 Charlotte Bunne, Stefan G Stark, Gabriele Gut, Jacobo Sarabia Del Castillo, Mitch Levesque,  
454 Kjong-Van Lehmann, Lucas Pelkmans, Andreas Krause, and Gunnar Rätsch. Learning single-  
455 cell perturbation responses using neural optimal transport. *Nature methods*, 20(11), 2023. [Cited  
456 on p. 1.]  
457 Lawrence Cayton. Algorithms for manifold learning, 2005. Technical Report, University of Cali-  
458 fornia at San Diego. [Cited on p. 17.]  
459 Atray Dixit, Oren Parnas, Biyu Li, Jenny Chen, Charles P Fulco, Livnat Jerby-Arnon, Nemanja D  
460 Marjanovic, et al. Perturb-seq: dissecting molecular circuits with scalable single-cell RNA pro-  
461 filing of pooled genetic screens. *cell*, 167(7):1853–1866, 2016. [Cited on p. 1.]  
462 Kefan Dong and Tengyu Ma. First steps toward understanding the extrapolation of nonlinear models  
463 to unseen domains. In *International Conference on Learning Representations*, 2022. [Cited on  
464 p. 3.]  
465 Tilmann Gneiting and Adrian E Raftery. Strictly proper scoring rules, prediction, and estimation.  
466 *Journal of the American Statistical Association*, 102(477):359–378, 2007. [Cited on p. 2, 3, 6,  
467 and 17.]  
468 Anirudh Goyal and Yoshua Bengio. Inductive biases for deep learning of higher-level cognition.  
469 *Proceedings of the Royal Society A*, 478(2266):20210068, 2022. [Cited on p. 3.]  
470 Arthur Gretton, Karsten M Borgwardt, Malte J Rasch, Bernhard Schölkopf, and Alexander Smola.  
471 A kernel two-sample test. *The Journal of Machine Learning Research*, 13(1):723–773, 2012.  
472 [Cited on p. 7 and 17.]  
473 Leon Hetzel, Simon Boehm, Niki Kilbertus, Stephan Günemann, Fabian Theis, et al. Predicting  
474 cellular responses to novel drug perturbations at a single-cell resolution. In *Advances in Neural*  
475 *Information Processing Systems*, volume 35, pp. 26711–26722, 2022. [Cited on p. 1.]  
476 Geoffrey E Hinton and Ruslan R Salakhutdinov. Reducing the dimensionality of data with neural  
477 networks. *science*, 313(5786):504–507, 2006. [Cited on p. 1 and 17.]  
478 Martin Jinek, Krzysztof Chylinski, Ines Fonfara, Michael Hauer, Jennifer A Doudna, and Em-  
479 manuelle Charpentier. A programmable dual-RNA-guided DNA endonuclease in adaptive bacte-  
480 rial immunity. *science*, 337(6096):816–821, 2012. [Cited on p. 1.]  
481 Kenji Kamimoto, Blerta Stringa, Christy M Hoffmann, Kunal Jindal, Lilianna Solnica-Krezel, and  
482 Samantha A Morris. Dissecting cell identity via network inference and in silico gene perturbation.  
483 *Nature*, 614(7949):742–751, 2023. [Cited on p. 1.]  
484  
485

- 486 Ilyes Khemakhem, Diederik Kingma, Ricardo Monti, and Aapo Hyvarinen. Variational autoen-  
487 coders and nonlinear ica: A unifying framework. In *International Conference on Artificial Intel-*  
488 *ligence and Statistics*, pp. 2207–2217, 2020. [Cited on p. 5.]
- 489 Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International*  
490 *Conference on Learning Representations*, 2015. [Cited on p. 7.]
- 491 Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In *International Conference*  
492 *on Learning Representations*, 2014. [Cited on p. 1.]
- 493 Roger Koenker. *Quantile regression*. Cambridge University Press, 2005. [Cited on p. 3.]
- 494 Roger Koenker and Gilbert Bassett Jr. Regression quantiles. *Econometrica: Journal of the Econo-*  
495 *metric Society*, pp. 33–50, 1978. [Cited on p. 3.]
- 496 Sébastien Lachapelle, Divyat Mahajan, Ioannis Mitliagkas, and Simon Lacoste-Julien. Additive  
497 decoders for latent variables identification and cartesian-product extrapolation. In *Advances in*  
498 *Neural Information Processing Systems*, volume 36, 2023. [Cited on p. 3.]
- 499 Brenden M Lake, Tomer D Ullman, Joshua B Tenenbaum, and Samuel J Gershman. Building  
500 machines that learn and think like people. *Behavioral and brain sciences*, 40:e253, 2017. [Cited  
501 on p. 3.]
- 502 Guillaume Lample, Neil Zeghidour, Nicolas Usunier, Antoine Bordes, Ludovic Denoyer, and  
503 Marc’Aurelio Ranzato. Fader networks: Manipulating images by sliding attributes. In *Advances*  
504 *in Neural Information Processing Systems*, volume 30, 2017. [Cited on p. 19.]
- 505 Samuel Lippl and Kim Stachenfeld. When does compositional structure yield compositional gener-  
506 alization? a kernel theory. *arXiv preprint arXiv:2405.16391*, 2024. [Cited on p. 3.]
- 507 Mohammad Lotfollahi, F Alexander Wolf, and Fabian J Theis. scgen predicts single-cell perturba-  
508 tion responses. *Nature methods*, 16(8):715–721, 2019. [Cited on p. 1.]
- 509 Mohammad Lotfollahi, Anna Klimovskaia Susmelj, Carlo De Donno, Leon Hetzel, Yuge Ji, Ignacio  
510 L Ibarra, Sanjay R Srivatsan, Mohsen Naghipourfar, Riza M Daza, Beth Martin, Jay Shendure,  
511 Jose L McFaline-Figueroa, Pierre Boyeau, F Alexander Wolf, Nafissa Yakubova, Stephan Gün-  
512 nemann, Cole Trapnell, David Lopez-Paz, and Fabian J Theis. Predicting cellular responses to  
513 complex perturbations in high-throughput screens. *Molecular Systems Biology*, 19(6):e11517,  
514 2023. [Cited on p. 1, 2, 3, 5, 7, 8, 18, and 19.]
- 515 James E Matheson and Robert L Winkler. Scoring rules for continuous probability distributions.  
516 *Management Science*, 22(10):1087–1096, 1976. [Cited on p. 16.]
- 517 Mehdi Mirza. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014. [Cited  
518 on p. 3.]
- 519 Milton L. Montero, Jeffrey Bowers, Rui Ponte Costa, Casimir JH Ludwig, and Gaurav Malhotra.  
520 Lost in latent space: Examining failures of disentangled models at combinatorial generalisation.  
521 In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in*  
522 *Neural Information Processing Systems*, 2022. [Cited on p. 3.]
- 523 Milton Llera Montero, Casimir JH Ludwig, Rui Ponte Costa, Gaurav Malhotra, and Jeffrey Bow-  
524 ers. The role of disentanglement in generalisation. In *International Conference on Learning*  
525 *Representations*, 2021. [Cited on p. 3.]
- 526 Thomas M Norman, Max A Horlbeck, Joseph M Replogle, Alex Y Ge, Albert Xu, Marco Jost,  
527 Luke A Gilbert, and Jonathan S Weissman. Exploring genetic interaction manifolds constructed  
528 from rich single-cell phenotypes. *science*, 365(6455):786–793, 2019. [Cited on p. 1.]
- 529 Judea Pearl. *Causality: Models, Reasoning, and Inference*. Cambridge University Press, 2nd edition,  
530 2009. [Cited on p. 8 and 19.]
- 531 Xiaoning Qi, Lianhe Zhao, Chenyu Tian, Yueyue Li, Zhen-Lin Chen, Peipei Huo, Runsheng Chen,  
532 Xiaodong Liu, Baoping Wan, Shengyong Yang, et al. Predicting transcriptional responses to novel  
533 chemical perturbations using deep generative model for drug discovery. *Nature Communications*,  
534 15(1):1–19, 2024. [Cited on p. 1.]
- 535 Yusuf Roohani, Kexin Huang, and Jure Leskovec. Predicting transcriptional outcomes of novel  
536 multigene perturbations with gears. *Nature Biotechnology*, 42(6):927–935, 2024. [Cited on p. 1.]

- 540 David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by back-  
541 propagating errors. *nature*, 323(6088):533–536, 1986. [Cited on p. 1 and 17.]
- 542 Lawrence K Saul and Sam T Roweis. Think globally, fit locally: Unsupervised learning of low  
543 dimensional manifolds. *Journal of Machine Learning Research*, 4:119–155, 2003. [Cited on  
544 p. 17.]
- 545 Bernhard Schölkopf and Alexander J Smola. *Learning with kernels: support vector machines,  
546 regularization, optimization, and beyond*. MIT press, 2002. [Cited on p. 17.]
- 547 Bernhard Schölkopf, Francesco Locatello, Stefan Bauer, Nan Rosemary Ke, Nal Kalchbrenner,  
548 Anirudh Goyal, and Yoshua Bengio. Toward causal representation learning. *Proceedings of  
549 the IEEE*, 109(5):612–634, 2021. [Cited on p. 8.]
- 550 Lukas Schott, Julius von Kügelgen, Frederik Träuble, Peter Vincent Gehler, Chris Russell, Matthias  
551 Bethge, Bernhard Schölkopf, Francesco Locatello, and Wieland Brendel. Visual representation  
552 learning does not generalize strongly within the same domain. In *International Conference on  
553 Learning Representations*, 2022. [Cited on p. 3.]
- 554 Dino Sejdinovic, Bharath Sriperumbudur, Arthur Gretton, and Kenji Fukumizu. Equivalence of  
555 distance-based and rkhs-based statistics in hypothesis testing. *The Annals of Statistics*, pp. 2263–  
556 2291, 2013. [Cited on p. 17.]
- 557 Xinwei Shen and Nicolai Meinshausen. Distributional principal autoencoders. *arXiv preprint  
558 arXiv:2404.13649*, 2024a. [Cited on p. 6 and 18.]
- 559 Xinwei Shen and Nicolai Meinshausen. Engression: extrapolation through the lens of distributional  
560 regression. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 2024b.  
561 [Cited on p. 6.]
- 562 Kihyuk Sohn, Honglak Lee, and Xinchen Yan. Learning structured output representation using  
563 deep conditional generative models. In *Advances in Neural Information Processing Systems*,  
564 volume 28, 2015. [Cited on p. 3.]
- 565 Chandler Squires, Anna Seigal, Salil Bhate, and Caroline Uhler. Linear causal disentanglement via  
566 interventions. In *40th International Conference on Machine Learning*, 2023. [Cited on p. 8.]
- 567 Gábor J Székely and Maria L Rizzo. Energy statistics: A class of statistics based on distances.  
568 *Journal of Statistical Planning and Inference*, 143(8):1249–1272, 2013. [Cited on p. 7 and 17.]
- 569 Joshua B Tenenbaum, Vin de Silva, and John C Langford. A global geometric framework for  
570 nonlinear dimensionality reduction. *science*, 290(5500):2319–2323, 2000. [Cited on p. 17.]
- 571 Zhong Wang, Mark Gerstein, and Michael Snyder. RNA-seq: a revolutionary tool for transcrip-  
572 tomics. *Nature reviews genetics*, 10(1):57–63, 2009. [Cited on p. 1.]
- 573 John N Weinstein, Eric A Collisson, Gordon B Mills, Kenna R Shaw, Brad A Ozenberger, Kyle  
574 Ellrott, Ilya Shmulevich, Chris Sander, and Joshua M Stuart. The cancer genome atlas pan-cancer  
575 analysis project. *Nature genetics*, 45(10):1113–1120, 2013. [Cited on p. 1.]
- 576 Thaddäus Wiedemer, Jack Brady, Alexander Panfilov, Attila Juhos, Matthias Bethge, and Wieland  
577 Brendel. Provable compositional generalization for object-centric learning. In *International Con-  
578 ference on Learning Representations*, 2024a. [Cited on p. 3.]
- 579 Thaddäus Wiedemer, Prasanna Mayilvahanan, Matthias Bethge, and Wieland Brendel. Composi-  
580 tional generalization from first principles. In *Advances in Neural Information Processing Systems*,  
581 volume 36, 2024b. [Cited on p. 3.]
- 582 Christina Winkler, Daniel Worrall, Emiel Hoogeboom, and Max Welling. Learning likelihoods with  
583 conditional normalizing flows. *arXiv preprint arXiv:1912.00042*, 2019. [Cited on p. 3.]
- 584 Hengshi Yu and Joshua D Welch. Perturbnet predicts single-cell responses to unseen chemical and  
585 genetic perturbations. *BioRxiv*, pp. 2022–07, 2022. [Cited on p. 1.]
- 586 Jiaqi Zhang, Kristjan Greenewald, Chandler Squires, Akash Srivastava, Karthikeyan Shanmugam,  
587 and Caroline Uhler. Identifiability guarantees for causal disentanglement from soft interventions.  
588 *Advances in Neural Information Processing Systems*, 36, 2024. [Cited on p. 8.]
- 589  
590  
591  
592  
593

# Appendix

## A ADDITIONAL RESULTS

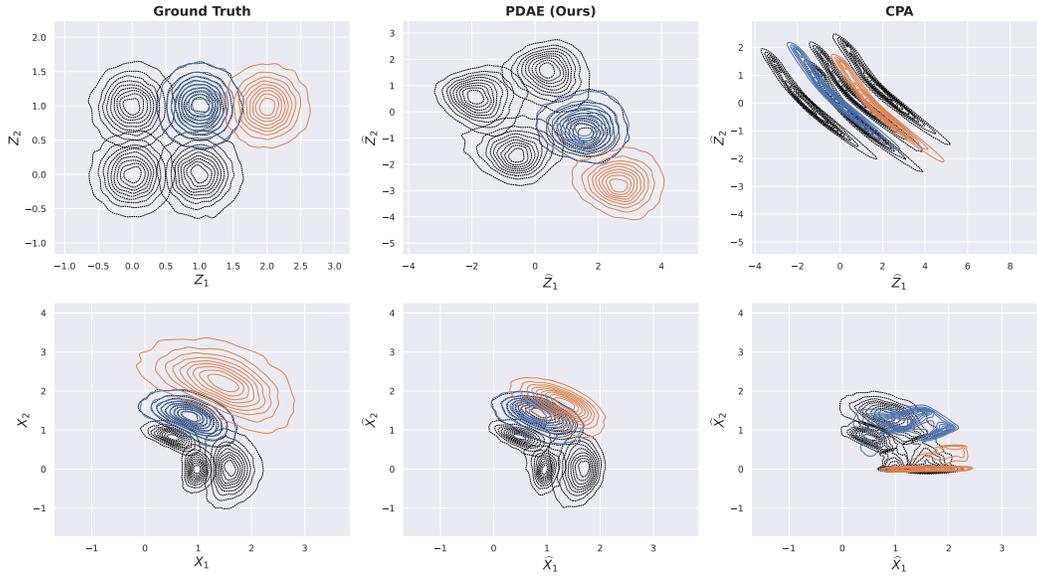


Figure 3: **Comparison of PDAE and CPA on Synthetic Data.** Shown are the results of our experiment described in § 6. Rows correspond to latent space (top) and observation space (bottom). Columns show the ground truth data (left), PDAE predictions (center), and CPA predictions (right). Training domains are shown in grey, the in-distribution (ID) test domain (which overlaps with one of the training domains) in blue, and the out-of-distribution (OOD) test domain in orange. All plots show kernel density estimates of the distributions. As can be seen, PDAE recovers an affine transformation of the true latents (top, center) leading to accurate distributional predictions for the training and ID test domain (bottom, center). However, the OOD test domain is mapped to a part of the latent space not seen during training (top, center). As a result, the corresponding decoder output does not accurately match the true OOD distribution (bottom left vs center). CPA appears to learn a latent space, in which all perturbed latent distributions are co-linear (top right), and the predicted distributions do not match the ground truth particularly well, particularly for the test conditions (bottom left vs right). However, recall that—unlike PDAE—CPA is not trained for distributional reconstruction, see Appx. C.4 for details.

## B PROOFS

### B.1 PROOF OF THM. 4.1

**Theorem 4.1** (Affine identifiability for Gaussian latents). *For  $M \in \mathbb{Z}_{\geq 0}$ , let  $\mathbf{l}_0, \dots, \mathbf{l}_M \in \mathbb{R}^K$  be  $M+1$  perturbation labels. Let  $\mathbf{f}, \tilde{\mathbf{f}} : \mathbb{R}^{d_z} \rightarrow \mathbb{R}^{d_x}$ ,  $\mathbf{W}, \tilde{\mathbf{W}} \in \mathbb{R}^{d_z \times K}$ , and  $\mathbb{P}, \tilde{\mathbb{P}}$  be distributions on  $\mathbb{R}^{d_z}$  such that the models  $(\mathbf{f}, \mathbf{W}, \mathbb{P})$  and  $(\tilde{\mathbf{f}}, \tilde{\mathbf{W}}, \tilde{\mathbb{P}})$  induce the same observed distributions, i.e.,*

$$\forall e \in [M]_0 : \quad \mathbf{f}(\mathbf{Z} + \mathbf{W}\mathbf{l}_e) \stackrel{d}{=} \tilde{\mathbf{f}}(\tilde{\mathbf{Z}} + \tilde{\mathbf{W}}\mathbf{l}_e), \quad \text{where } \mathbf{Z} \sim \mathbb{P} \text{ and } \tilde{\mathbf{Z}} \sim \tilde{\mathbb{P}}. \quad (4.1)$$

Assume further that:

- (i) **[invertibility]**  $\mathbf{f}$  and  $\tilde{\mathbf{f}}$  are  $C^2$ -diffeomorphisms onto their respective images;
- (ii) **[Gaussianity]**  $\mathbf{Z}$  and  $\tilde{\mathbf{Z}}$  are non-degenerate multi-variate Gaussians, i.e.,  $\mathbb{P} = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  and  $\tilde{\mathbb{P}} = \mathcal{N}(\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}})$  for some  $\boldsymbol{\mu}, \tilde{\boldsymbol{\mu}} \in \mathbb{R}^{d_z}$  and positive-definite  $\boldsymbol{\Sigma}, \tilde{\boldsymbol{\Sigma}} \in \mathbb{R}^{d_z \times d_z}$ ;
- (iii) **[sufficient diversity]** the matrix  $\tilde{\mathbf{W}}\mathbf{L} \in \mathbb{R}^{d_z \times M}$ , where  $\mathbf{L} \in \mathbb{R}^{K \times M}$  is the matrix with columns  $(\mathbf{l}_e - \mathbf{l}_0)$  for  $e \in [M]$ , has full row rank, i.e.,  $\text{rank}(\tilde{\mathbf{W}}\mathbf{L}) = d_z$ .

Then the latent representation and the effects of the observed perturbation combinations relative to  $\mathbf{l}_0$  (as captured by  $\mathbf{W}\mathbf{L}$ ) are identifiable up to affine transformation in the following sense:

$$\forall \mathbf{z} : \quad \tilde{\mathbf{f}}^{-1} \circ \mathbf{f}(\mathbf{z}) = \mathbf{A}\mathbf{z} + \mathbf{b}, \quad (4.2)$$

$$\tilde{\mathbf{W}}\mathbf{L} = \mathbf{A}\mathbf{W}\mathbf{L}, \quad (4.3)$$

where  $\mathbf{A} := \tilde{\boldsymbol{\Sigma}}^{\frac{1}{2}}\boldsymbol{\Sigma}^{-\frac{1}{2}}$  and  $\mathbf{b} := \tilde{\boldsymbol{\mu}} - \mathbf{A}\boldsymbol{\mu} + (\tilde{\mathbf{W}} - \mathbf{A}\mathbf{W})\mathbf{l}_0$ .

*Proof.* Let  $p_e$  and  $\tilde{p}_e$  denote the densities of

$$\mathbf{Z}_e := \mathbf{Z} + \mathbf{W}\mathbf{l}_e \quad (B.1)$$

and

$$\tilde{\mathbf{Z}}_e := \tilde{\mathbf{Z}} + \tilde{\mathbf{W}}\mathbf{l}_e, \quad (B.2)$$

respectively. Due to (4.1),  $\mathbf{f}$  and  $\tilde{\mathbf{f}}$  have the same image. Thus, by the invertibility assumption (i), the function  $\mathbf{h} := \tilde{\mathbf{f}}^{-1} \circ \mathbf{f} : \mathbb{R}^{d_z} \rightarrow \mathbb{R}^{d_z}$  is a well-defined  $C^2$ -diffeomorphism. The change of variable formula applied to

$$\tilde{\mathbf{Z}}_e \stackrel{d}{=} \mathbf{h}(\mathbf{Z}_e) \quad (B.3)$$

then yields for all  $e$  and all  $\mathbf{z}$ :

$$p_e(\mathbf{z}) = \tilde{p}_e(\mathbf{h}(\mathbf{z})) |\det \mathbf{J}_{\mathbf{h}}(\mathbf{z})|, \quad (B.4)$$

where  $\mathbf{J}_{\mathbf{h}}(\mathbf{z})$  denotes the Jacobian of  $\mathbf{h}$ . By taking logarithms of (B.4) and contrasting domain  $e$  with a reference domain with  $e = 0$ , the determinant terms cancel and we obtain for all  $e$  and all  $\mathbf{z}$ :

$$\log p_e(\mathbf{z}) - \log p_0(\mathbf{z}) = \log \tilde{p}_e(\mathbf{h}(\mathbf{z})) - \log \tilde{p}_0(\mathbf{h}(\mathbf{z})). \quad (B.5)$$

Next, denote the densities of  $\mathbf{Z}$  and  $\tilde{\mathbf{Z}}$  by  $p$  and  $\tilde{p}$ , respectively. From (B.1) and (B.2) it then follows that, for all  $e$  and all  $\mathbf{z}$ ,  $p_e$  and  $\tilde{p}_e$  can respectively be expressed in terms of  $p$  and  $\tilde{p}$  as follows:

$$p_e(\mathbf{z}) = p(\mathbf{z} - \mathbf{W}\mathbf{l}_e), \quad (B.6)$$

$$\tilde{p}_e(\mathbf{z}) = \tilde{p}(\mathbf{z} - \tilde{\mathbf{W}}\mathbf{l}_e). \quad (B.7)$$

By substituting these expressions in (B.5), we obtain for all  $e$  and all  $\mathbf{z}$ :

$$\log p(\mathbf{z} - \mathbf{W}\mathbf{l}_e) - \log p(\mathbf{z} - \mathbf{W}\mathbf{l}_0) = \log \tilde{p}(\mathbf{h}(\mathbf{z}) - \tilde{\mathbf{W}}\mathbf{l}_e) - \log \tilde{p}(\mathbf{h}(\mathbf{z}) - \tilde{\mathbf{W}}\mathbf{l}_0). \quad (B.8)$$

By the Gaussianity assumption (ii), the contrast of log-densities in (B.8) takes the following form for all  $e$  and all  $\mathbf{z}$ :

$$(\mathbf{l}_e - \mathbf{l}_0)^\top \mathbf{W}^\top \boldsymbol{\Sigma}^{-1} (\mathbf{z} - \boldsymbol{\mu}) - \frac{1}{2} (\mathbf{l}_e - \mathbf{l}_0)^\top \mathbf{W}^\top \boldsymbol{\Sigma}^{-1} \mathbf{W} (\mathbf{l}_e + \mathbf{l}_0) \quad (\text{B.9})$$

$$= (\mathbf{l}_e - \mathbf{l}_0)^\top \widetilde{\mathbf{W}}^\top \widetilde{\boldsymbol{\Sigma}}^{-1} (\mathbf{h}(\mathbf{z}) - \widetilde{\boldsymbol{\mu}}) - \frac{1}{2} (\mathbf{l}_e - \mathbf{l}_0)^\top \widetilde{\mathbf{W}}^\top \widetilde{\boldsymbol{\Sigma}}^{-1} \widetilde{\mathbf{W}} (\mathbf{l}_e + \mathbf{l}_0) \quad (\text{B.10})$$

We aim to show that the two representations are related by an affine transformation, i.e., that all second-order derivatives of  $\mathbf{h}$  are zero everywhere. Taking gradients w.r.t.  $\mathbf{z}$  yields for all  $e$  and all  $\mathbf{z}$ :

$$(\mathbf{l}_e - \mathbf{l}_0)^\top \mathbf{W}^\top \boldsymbol{\Sigma}^{-1} = (\mathbf{l}_e - \mathbf{l}_0)^\top \widetilde{\mathbf{W}}^\top \widetilde{\boldsymbol{\Sigma}}^{-1} \mathbf{J}_h(\mathbf{z}). \quad (\text{B.11})$$

Let  $\mathbf{L} \in \mathbb{R}^{K \times M}$  be the matrix with columns  $\mathbf{l}_e - \mathbf{l}_0$  for  $e \in [M]$ . Then, for all  $\mathbf{z}$ :

$$\mathbf{L}^\top \mathbf{W}^\top \boldsymbol{\Sigma}^{-1} = \mathbf{L}^\top \widetilde{\mathbf{W}}^\top \widetilde{\boldsymbol{\Sigma}}^{-1} \mathbf{J}_h(\mathbf{z}). \quad (\text{B.12})$$

Differentiating once more w.r.t.  $\mathbf{z}$  yields for all  $\mathbf{z}$ :

$$\mathbf{0} = \mathbf{L}^\top \widetilde{\mathbf{W}}^\top \widetilde{\boldsymbol{\Sigma}}^{-1} \mathbf{H}_h(\mathbf{z}), \quad (\text{B.13})$$

where the 3-tensor  $\mathbf{H}_h(\mathbf{z}) \in \mathbb{R}^{d_z \times d_z \times d_z}$  denotes the Hessian of  $\mathbf{h}$ , i.e., for all  $i, j \in [d_z]$  and all  $\mathbf{z}$

$$\mathbf{0} = \mathbf{L}^\top \widetilde{\mathbf{W}}^\top \widetilde{\boldsymbol{\Sigma}}^{-1} \frac{\partial^2}{\partial z_i \partial z_j} \mathbf{h}(\mathbf{z}). \quad (\text{B.14})$$

By assumption (iii), the matrix  $\mathbf{L}^\top \widetilde{\mathbf{W}}^\top$  has full column rank and thus a left inverse, i.e., there exists  $\mathbf{V} \in \mathbb{R}^{d_z \times M}$  such that  $\mathbf{V} \mathbf{L}^\top \widetilde{\mathbf{W}}^\top = \mathbf{I}_{d_z}$ . Multiplying (B.14) on the left by  $\widetilde{\boldsymbol{\Sigma}} \mathbf{V}$  then yields for all  $i, j \in [d_z]$  and all  $\mathbf{z}$ :

$$\frac{\partial^2}{\partial z_i \partial z_j} \mathbf{h}(\mathbf{z}) = \mathbf{0}. \quad (\text{B.15})$$

This implies that  $\mathbf{h}$  must be affine, i.e., there exist  $\mathbf{A} \in \mathbb{R}^{d_z \times d_z}$  and  $\mathbf{b} \in \mathbb{R}^{d_z}$  such that for all  $\mathbf{z}$ :

$$\mathbf{h}(\mathbf{z}) = \mathbf{A} \mathbf{z} + \mathbf{b}. \quad (\text{B.16})$$

Further, since  $\mathbf{h}$  is invertible,  $\mathbf{A}$  must be invertible.

Recall that  $\widetilde{\mathbf{Z}}_e \stackrel{d}{=} \mathbf{h}(\mathbf{Z}_e) = \mathbf{A} \mathbf{Z}_e + \mathbf{b}$ . It follows from (B.1), (B.2) and assumption (ii) that for all  $e$ :

$$\mathcal{N}(\widetilde{\boldsymbol{\mu}} + \widetilde{\mathbf{W}} \mathbf{l}_e, \widetilde{\boldsymbol{\Sigma}}) = \mathcal{N}(\mathbf{A}(\boldsymbol{\mu} + \mathbf{W} \mathbf{l}_e) + \mathbf{b}, \mathbf{A} \boldsymbol{\Sigma} \mathbf{A}^\top). \quad (\text{B.17})$$

By equating the covariances, we find

$$\mathbf{A} \boldsymbol{\Sigma} \mathbf{A}^\top = \widetilde{\boldsymbol{\Sigma}} \quad \implies \quad \mathbf{A} = \widetilde{\boldsymbol{\Sigma}}^{\frac{1}{2}} \boldsymbol{\Sigma}^{-\frac{1}{2}}, \quad (\text{B.18})$$

where the matrix square roots exist and are unique since  $\boldsymbol{\Sigma}$  and  $\widetilde{\boldsymbol{\Sigma}}$  are symmetric and positive definite. By equating the means, we obtain for all  $e$ :

$$\widetilde{\boldsymbol{\mu}} + \widetilde{\mathbf{W}} \mathbf{l}_e = \mathbf{A}(\boldsymbol{\mu} + \mathbf{W} \mathbf{l}_e) + \mathbf{b} \quad \iff \quad (\widetilde{\mathbf{W}} - \mathbf{A} \mathbf{W}) \mathbf{l}_e = \mathbf{A} \boldsymbol{\mu} + \mathbf{b} - \widetilde{\boldsymbol{\mu}}. \quad (\text{B.19})$$

By contrasting (B.19) for all  $e \in [M]$  with  $e = 0$  as before, we obtain

$$(\widetilde{\mathbf{W}} - \mathbf{A} \mathbf{W}) \mathbf{L} = \mathbf{0} \quad \iff \quad \widetilde{\mathbf{W}} \mathbf{L} = \mathbf{A} \mathbf{W} \mathbf{L}. \quad (\text{B.20})$$

Finally, by choosing  $e = 0$  in (B.19) we obtain the desired expression for  $\mathbf{b}$ ,

$$\mathbf{b} = \widetilde{\boldsymbol{\mu}} - \mathbf{A} \boldsymbol{\mu} + (\widetilde{\mathbf{W}} - \mathbf{A} \mathbf{W}) \mathbf{l}_0. \quad (\text{B.21})$$

This completes the proof.  $\square$

## B.2 PROOF OF COR. 4.2

**Corollary 4.2** (Affine recovery of the perturbation matrix). *If, in addition to the assumptions of Thm. 4.1,  $\mathbf{L} \in \mathbb{R}^{K \times M}$  has full row rank (i.e.,  $\text{rank}(\mathbf{L}) = K \leq M$ ), then the perturbation matrix  $\mathbf{W}$  is identifiable up to affine transformation in the sense that*

$$\widetilde{\mathbf{W}} = \mathbf{A}\mathbf{W}, \quad (4.4)$$

for  $\mathbf{A} := \widetilde{\Sigma}^{\frac{1}{2}}\Sigma^{-\frac{1}{2}}$ . In this case, the expression for  $\mathbf{b}$  in (4.2) simplifies to  $\mathbf{b} = \widetilde{\boldsymbol{\mu}} - \mathbf{A}\boldsymbol{\mu}$ .

*Proof.* If  $\text{rank}(\mathbf{L}) = K$ , then  $\mathbf{L}$  has a right inverse, i.e., there exists  $\mathbf{K} \in \mathbb{R}^{M \times K}$  such that  $\mathbf{L}\mathbf{K} = \mathbf{I}_K$ . Right multiplication of (B.20) by  $\mathbf{K}$  then yields

$$\widetilde{\mathbf{W}} = \mathbf{A}\mathbf{W} \quad (B.22)$$

Finally, substitution of (B.22) into (B.21) yields

$$\mathbf{b} = \widetilde{\boldsymbol{\mu}} - \mathbf{A}\boldsymbol{\mu}. \quad (B.23)$$

□

## B.3 PROOF OF THM. 4.6

**Theorem 4.6** (Extrapolation to span of relative perturbations). *Under the same setting and assumptions as in Thm. 4.1, let  $\mathbf{l}_{\text{test}} \in \mathbb{R}^K$  be an unseen perturbation label such that*

$$(\mathbf{l}_{\text{test}} - \mathbf{l}_0) \in \text{span}\left(\{\mathbf{l}_e - \mathbf{l}_0\}_{e \in [M]}\right). \quad (4.6)$$

Then the effect of  $\mathbf{l}_{\text{test}}$  is uniquely identifiable in the sense that

$$\mathbf{X}_{\text{test}} = \mathbf{f}(\mathbf{Z} + \mathbf{W}\mathbf{l}_{\text{test}}) \stackrel{d}{=} \widetilde{\mathbf{f}}\left(\widetilde{\mathbf{Z}} + \widetilde{\mathbf{W}}\mathbf{l}_{\text{test}}\right) = \widetilde{\mathbf{X}}_{\text{test}}. \quad (4.7)$$

*Proof.* With  $\mathbf{h} = \widetilde{\mathbf{f}}^{-1} \circ \mathbf{f}$ , the condition in (4.7) is equivalent to

$$\mathbf{h}(\mathbf{Z} + \mathbf{W}\mathbf{l}_{\text{test}}) \stackrel{d}{=} \widetilde{\mathbf{Z}} + \widetilde{\mathbf{W}}\mathbf{l}_{\text{test}}. \quad (B.24)$$

By Thm. 4.1, we have  $\mathbf{h}(\mathbf{z}) = \mathbf{A}\mathbf{z} + \mathbf{b}$  for

$$\mathbf{A} = \widetilde{\Sigma}^{\frac{1}{2}}\Sigma^{-\frac{1}{2}}, \quad (B.25)$$

$$\mathbf{b} = \widetilde{\boldsymbol{\mu}} - \mathbf{A}\boldsymbol{\mu} + (\widetilde{\mathbf{W}} - \mathbf{A}\mathbf{W})\mathbf{l}_0. \quad (B.26)$$

Together with  $\mathbf{Z} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$ , this lets us compute the distribution of the LHS of (B.24) as:

$$\mathbf{h}(\mathbf{Z} + \mathbf{W}\mathbf{l}_{\text{test}}) = \mathbf{A}(\mathbf{Z} + \mathbf{W}\mathbf{l}_{\text{test}}) + \mathbf{b} \sim \mathcal{N}\left(\mathbf{A}\boldsymbol{\mu} + \mathbf{A}\mathbf{W}\mathbf{l}_{\text{test}} + \mathbf{b}, \mathbf{A}\Sigma\mathbf{A}^\top\right). \quad (B.27)$$

Similarly, since  $\widetilde{\mathbf{Z}} \sim \mathcal{N}\left(\widetilde{\boldsymbol{\mu}}, \widetilde{\Sigma}\right)$ , the distribution of the RHS of (B.24) is given by

$$\widetilde{\mathbf{Z}} + \widetilde{\mathbf{W}}\mathbf{l}_{\text{test}} \sim \mathcal{N}\left(\widetilde{\boldsymbol{\mu}} + \widetilde{\mathbf{W}}\mathbf{l}_{\text{test}}, \widetilde{\Sigma}\right). \quad (B.28)$$

From (B.25), it follows directly that

$$\mathbf{A}\Sigma\mathbf{A}^\top = \widetilde{\Sigma}. \quad (B.29)$$

To complete the proof, it thus remains to show that the means of (B.27) and (B.28) are equal, i.e.,

$$\mathbf{A}\boldsymbol{\mu} + \mathbf{A}\mathbf{W}\mathbf{l}_{\text{test}} + \mathbf{b} = \widetilde{\boldsymbol{\mu}} + \widetilde{\mathbf{W}}\mathbf{l}_{\text{test}}. \quad (B.30)$$

Starting from the LHS of (B.30), by substituting the expression for  $\mathbf{b}$  from (B.26) we obtain:

$$\mathbf{A}\boldsymbol{\mu} + \mathbf{A}\mathbf{W}\mathbf{l}_{\text{test}} + \mathbf{b} = \mathbf{A}\boldsymbol{\mu} + \mathbf{A}\mathbf{W}\mathbf{l}_{\text{test}} + \widetilde{\boldsymbol{\mu}} - \mathbf{A}\boldsymbol{\mu} + (\widetilde{\mathbf{W}} - \mathbf{A}\mathbf{W})\mathbf{l}_0 \quad (B.31)$$

$$= \widetilde{\boldsymbol{\mu}} + \mathbf{A}\mathbf{W}(\mathbf{l}_{\text{test}} - \mathbf{l}_0) + \widetilde{\mathbf{W}}\mathbf{l}_0 \quad (B.32)$$

Next, by (4.6), i.e., the assumption that  $(l_{\text{test}} - l_0)$  lies in the span of  $\{l_e - l_0\}_{e \in [M]}$ , there exists  $\alpha \in \mathbb{R}^M$  such that

$$l_{\text{test}} - l_0 = \sum_{e \in [M]} \alpha_e (l_e - l_0) = L\alpha \quad (\text{B.33})$$

where, as before,  $L \in \mathbb{R}^{K \times M}$  is the matrix with columns  $(l_e - l_0)$ . Substituting (B.33) into (B.32) yields

$$A\mu + AWl_{\text{test}} + b = \tilde{\mu} + AWL\alpha + \tilde{W}l_0. \quad (\text{B.34})$$

Finally, it follows from Thm. 4.1 that

$$AWL = \tilde{W}L. \quad (\text{B.35})$$

which upon substitution into (B.34) yields the desired equality from (B.30)

$$A\mu + AWl_{\text{test}} + b = \tilde{\mu} + \tilde{W}L\alpha + \tilde{W}l_0 \quad (\text{B.36})$$

$$= \tilde{\mu} + \tilde{W}(L\alpha + l_0) \quad (\text{B.37})$$

$$= \tilde{\mu} + \tilde{W}l_{\text{test}}. \quad (\text{B.38})$$

This completes the proof.  $\square$

## C ADDITIONAL BACKGROUND MATERIAL AND RELATED WORK

Since the problem of interest (§ 2) involves making distributional predictions for new perturbation conditions, we review some basics of probabilistic forecasting (Appx. C.1) and measuring the similarity between two distributions (Appx. C.2), in our case typically between an empirical distribution and its predicted counterpart. We then turn to representation learning with encoder-decoder architectures (Appx. C.3), in particular a recent approach that also targets distributional reconstruction. Finally, we cover some prior efforts on perturbation modelling and extrapolation (Appx. C.4), which we draw inspiration from.

### C.1 PROBABILISTIC FORECASTING AND SCORING RULES

Let  $\Omega$  be a sample space,  $\mathcal{A}$  a  $\sigma$ -algebra of subsets of  $\Omega$ , and  $\mathcal{P}$  a convex class of probability measures on  $(\Omega, \mathcal{A})$ . A *probabilistic prediction* or *probabilistic forecast* is a mapping into  $\mathcal{P}$ , which outputs predictive distributions  $\mathbb{P}$  over outcomes  $x \in \Omega$ . Probabilistic forecasting can thus be viewed as a distributional generalization point prediction (i.e., deterministic forecasting), which maps directly into  $\Omega$ .

To evaluate, compare, or rank different forecasts, it is useful to assign them a numerical score reflecting their quality. A *scoring rule* is a function  $S : \mathcal{P} \times \Omega \rightarrow \mathbb{R}$  that assigns a score  $S(\mathbb{P}, x)$  to forecast  $\mathbb{P}$  if event  $x$  materializes, with higher scores corresponding to better forecasts—akin to (negative) loss or cost functions for point predictions. If  $x$  is distributed according to  $\mathbb{Q}$ , we denote the *expected score* by  $S(\mathbb{P}, \mathbb{Q}) = \mathbb{E}_{x \sim \mathbb{Q}}[S(\mathbb{P}, x)]$ . The scoring rule is called *proper* if  $S(\mathbb{Q}, \mathbb{Q}) \geq S(\mathbb{P}, \mathbb{Q})$  for all  $\mathbb{P}$ , and *strictly proper* if equality holds if and only if  $\mathbb{P} = \mathbb{Q}$ .

**CRPS.** For continuous scalar random variables (i.e.,  $\Omega = \mathbb{R}$ ), a popular scoring rule is the *continuous ranked probability score* (CRPS; Matheson & Winkler, 1976). When  $\mathcal{P}$  is the space of Borel probability measures on  $\mathbb{R}$  with finite first moment, it is strictly proper and given by:<sup>1</sup>

$$\text{CRPS}(\mathbb{P}, x) = \frac{1}{2} \mathbb{E}_{X, X' \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}} |X - X'| - \mathbb{E}_{X \sim \mathbb{P}} |X - x|. \quad (\text{C.1})$$

If  $\mathbb{P}$  is a point mass, the negative CRPS reduces to the absolute error loss function; it can thus be viewed as a generalization thereof to probabilistic forecasts.

<sup>1</sup>The original definition by Matheson & Winkler (1976) is  $\text{CRPS}(F, x) = -\int_{-\infty}^{\infty} (F(y) - 1\{y \geq x\})^2 dy$  where  $F$  is the CDF of  $\mathbb{P}$ , but the simpler form in (C.1) has been shown to be equivalent (Baringhaus & Franz, 2004, Lemma 2.2).

**Energy Score.** Gneiting & Raftery (2007) propose the *energy score* as a multi-variate generalization of the CRPS for vector-valued  $\mathbf{x} \in \Omega = \mathbb{R}^{d_x}$ . For  $\beta \in (0, 2)$ , it is defined by

$$\text{ES}_\beta(\mathbb{P}, \mathbf{x}) = \frac{1}{2} \mathbb{E}_{\mathbf{X}, \mathbf{X}' \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}} \|\mathbf{X} - \mathbf{X}'\|_2^\beta - \mathbb{E}_{\mathbf{X} \sim \mathbb{P}} \|\mathbf{X} - \mathbf{x}\|_2^\beta, \quad (\text{C.2})$$

where  $\|\cdot\|_2$  denotes the Euclidean (L2) norm.  $\text{ES}_\beta$  is strictly proper w.r.t.  $\mathcal{P}_\beta$ , the set of Borel probability measures  $\mathbb{P}$  for which  $\mathbb{E}_{\mathbf{X} \sim \mathbb{P}} \|\mathbf{X}\|_2^\beta$  is finite (Gneiting & Raftery, 2007).

## C.2 ASSESSING DISTRIBUTIONAL SIMILARITY

**Energy Distance.** The expected energy score  $\text{ES}_\beta(\mathbb{P}, \mathbb{Q}) = \mathbb{E}_{\mathbf{Y} \sim \mathbb{Q}}[\text{ES}_\beta(\mathbb{P}, \mathbf{Y})]$  constitutes a measure of similarity between  $\mathbb{P}$  and  $\mathbb{Q}$  and is closely linked to the *energy distance* (Székely & Rizzo, 2013):

$$\begin{aligned} \text{ED}_\beta(\mathbb{P}, \mathbb{Q}) &= 2\mathbb{E}_{\mathbf{X} \sim \mathbb{P}, \mathbf{Y} \sim \mathbb{Q}} \|\mathbf{X} - \mathbf{Y}\|_2^\beta - \mathbb{E}_{\mathbf{X}, \mathbf{X}' \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}} \|\mathbf{X} - \mathbf{X}'\|_2^\beta - \mathbb{E}_{\mathbf{Y}, \mathbf{Y}' \stackrel{\text{i.i.d.}}{\sim} \mathbb{Q}} \|\mathbf{Y} - \mathbf{Y}'\|_2^\beta \\ &= 2(\text{ES}_\beta(\mathbb{Q}, \mathbb{Q}) - \text{ES}_\beta(\mathbb{P}, \mathbb{Q})) \geq 0 \end{aligned} \quad (\text{C.3})$$

with equality if and only if  $\mathbb{P} = \mathbb{Q}$ , since  $\text{ES}_\beta$  is a strictly proper scoring rule.

**Maximum Mean Discrepancy (MMD).** Another well-known distance between probability measures that is rooted in kernel methods (Schölkopf & Smola, 2002) is the *maximum mean discrepancy* (MMD; Gretton et al., 2012), which for a positive definite kernel  $k : \Omega \times \Omega \rightarrow \mathbb{R}$  is given by

$$\text{MMD}_k^2(\mathbb{P}, \mathbb{Q}) = \mathbb{E}_{\mathbf{X}, \mathbf{X}' \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}} [k(\mathbf{X}, \mathbf{X}')] - 2\mathbb{E}_{\mathbf{X} \sim \mathbb{P}, \mathbf{Y} \sim \mathbb{Q}} [k(\mathbf{X}, \mathbf{Y})] + \mathbb{E}_{\mathbf{Y}, \mathbf{Y}' \stackrel{\text{i.i.d.}}{\sim} \mathbb{Q}} [k(\mathbf{Y}, \mathbf{Y}')] . \quad (\text{C.5})$$

**Energy Distance as a Special Case of MMD.** As shown by Sejdinovic et al. (2013), if  $k$  in (C.5) is chosen to be the positive-definite *distance kernel*<sup>2</sup>

$$k_\beta(\mathbf{X}, \mathbf{Y}) = \frac{1}{2} \left( \|\mathbf{X}\|_2^\beta + \|\mathbf{Y}\|_2^\beta - \|\mathbf{X} - \mathbf{Y}\|_2^\beta \right) \quad (\text{C.6})$$

then the energy distance is recovered as a special case of MMD,

$$\text{ED}_\beta(\mathbb{P}, \mathbb{Q}) = 2 \text{MMD}_{k_\beta}^2(\mathbb{P}, \mathbb{Q}) . \quad (\text{C.7})$$

## C.3 REPRESENTATION LEARNING

Many modern data sources of interest contain high-dimensional and unstructured observations  $\mathbf{x}$ , such as audio, video, images, or text. Representation learning aims to transform such data into a more compact, lower-dimensional set of features  $\mathbf{z}$  (the representation) which preserves most of the relevant information while making it more easily accessible, e.g., for use in downstream tasks (Bengio et al., 2013).<sup>3</sup> For example, a representation of images of multi-object scenes could be a list of the contained objects, along with their size, position, colour, etc. A key assumption underlying this endeavour is the so-called *manifold hypothesis* which posits that high-dimensional natural data tends to lie near a low-dimensional manifold embedded in the high-dimensional ambient space; this idea is also at the heart of several (nonlinear) dimension reduction techniques (Belkin & Niyogi, 2001; Cayton, 2005; Saul & Roweis, 2003; Tenenbaum et al., 2000).

**Autoencoder (AE).** An autoencoder (AE; Hinton & Salakhutdinov, 2006; Rumelhart et al., 1986) is a pair of functions  $(g, f)$ , consisting of an *encoder*  $g : \mathbb{R}^{d_x} \rightarrow \mathbb{R}^{d_z}$  mapping observations  $\mathbf{X} \sim \mathbb{P}$  to their representation  $\mathbf{Z} := g(\mathbf{X})$ , and a *decoder*  $f : \mathbb{R}^{d_z} \rightarrow \mathbb{R}^{d_x}$  mapping representations  $\mathbf{Z}$  to

<sup>2</sup>induced by the negative-definite semi-metric  $\rho(\mathbf{X}, \mathbf{Y}) = \|\mathbf{X} - \mathbf{Y}\|_2^\beta$  on  $\Omega = \mathbb{R}^{d_x}$  (centered at the origin),

<sup>3</sup>Representation learning is thus closely related to the classical task of dimension reduction. The former usually refers to nonlinear settings, involves some form of machine learning, and tends to be more focused on usefulness in terms of downstream tasks, rather than, say, explained variance.

918 their reconstructions in observation space,  $\widehat{\mathbf{X}} := \mathbf{f}(\mathbf{Z}) = \mathbf{f}(\mathbf{g}(\mathbf{X}))$ . Typically,  $d_Z < d_X$ , such  
 919 that there is a bottleneck and perfect reconstruction is not feasible. Autoencoders are usually trained  
 920 with a (point-wise, mean) reconstruction objective, i.e., the aim is to minimise the mean squared  
 921 error

$$922 \quad \mathcal{L}_{\text{AE}}(\mathbf{g}, \mathbf{f}; \mathbb{P}) := \mathbb{E}_{\mathbf{X} \sim \mathbb{P}} \left\| \mathbf{X} - \widehat{\mathbf{X}} \right\|_2^2 = \mathbb{E}_{\mathbf{X} \sim \mathbb{P}} \left\| \mathbf{X} - \mathbf{f}(\mathbf{g}(\mathbf{X})) \right\|_2^2, \quad (\text{C.8})$$

923 w.r.t. both  $\mathbf{g}$  and  $\mathbf{f}$ . As a result, the optimal decoder  $\mathbf{f}_{\text{AE}}^*$  for any fixed encoder choice  $\mathbf{g}$  is given by  
 924 the conditional mean

$$925 \quad \mathbf{f}_{\text{AE}}^*(\mathbf{z}; \mathbf{g}, \mathbb{P}) = \mathbb{E}_{\mathbf{X} \sim \mathbb{P}} [\mathbf{X} | \mathbf{g}(\mathbf{X}) = \mathbf{z}]. \quad (\text{C.9})$$

926  
 927 **Distributional Principal Autoencoder (DPA).** Since the objective of a standard AE is mean re-  
 928 construction, the distribution of reconstructions  $\widehat{\mathbf{X}}$  is typically not the same as the distribution of  $\mathbf{X}$   
 929 (unless the encoder is invertible, i.e. the compression is lossless and perfect reconstruction is feasi-  
 930 ble, which is usually not the case in practice). To address this, Shen & Meinshausen (2024a) pro-  
 931 posed the distributional principal autoencoder (DPA) which targets distributional (rather than mean)  
 932 reconstruction. A DPA also consists of an encoder-decoder pair  $(\mathbf{g}, \mathbf{f})$ . However, unlike in a stan-  
 933 dard AE, the DPA decoder  $\mathbf{f} : \mathbb{R}^{d_Z} \times \mathbb{R}^{d_\epsilon} \rightarrow \mathbb{R}^{d_X}$  is stochastic and receives an additional noise  
 934 term  $\epsilon$  as input, which is sampled from a fixed distribution  $\mathbb{Q}_\epsilon$  such as a standard isotropic Gaussian.  
 935 The DPA loss function is constructed such that for a fixed encoder  $\mathbf{g}$ , the optimal DPA decoder  $\mathbf{f}_{\text{DPA}}^*$   
 936 maps a given latent embedding  $\mathbf{z}$ , to the distribution of  $\mathbf{X}$ , given  $\mathbf{g}(\mathbf{X}) = \mathbf{z}$ , i.e.,

$$937 \quad \mathbf{f}_{\text{DPA}}^*(\mathbf{z}, \epsilon; \mathbf{g}) \stackrel{d}{=} (\mathbf{X} | \mathbf{g}(\mathbf{X}) = \mathbf{z}), \quad (\text{C.10})$$

938 where  $\stackrel{d}{=}$  denotes equality in distribution. This means that the decoder evaluated at  $\mathbf{z}$  should match  
 939 the distribution of all realizations of  $\mathbf{X}$  that are mapped by the encoder to  $\mathbf{z}$ . At the same time,  
 940 the DPA encoder aims to minimise the variability in the distributions in (C.10) by encoding the first  
 941  $d_Z$  ‘‘principal’’ components. As shown by Shen & Meinshausen (2024a), both of these goals are  
 942 achieved by the following DPA objective,

$$943 \quad \begin{aligned} 944 \quad \mathcal{L}_{\text{DPA}}(\mathbf{g}, \mathbf{f}; \mathbb{P}) &= \mathbb{E}_{\mathbf{X} \sim \mathbb{P}, \epsilon \sim \mathbb{Q}_\epsilon} \left\| \mathbf{X} - \mathbf{f}(\mathbf{g}(\mathbf{X}), \epsilon) \right\|_2^\beta - \frac{1}{2} \mathbb{E}_{\mathbf{X} \sim \mathbb{P}, \epsilon, \epsilon' \stackrel{\text{iid}}{\sim} \mathbb{Q}_\epsilon} \left\| \mathbf{f}(\mathbf{g}(\mathbf{X}), \epsilon) - \mathbf{f}(\mathbf{g}(\mathbf{X}), \epsilon') \right\|_2^\beta \\ 945 \quad &= -\mathbb{E}_{\mathbf{X} \sim \mathbb{P}} \left[ \text{ES}_\beta \left( \mathbf{f}(\mathbf{g}(\mathbf{X}), \cdot) \# \mathbb{Q}_\epsilon, \mathbf{X} \right) \right], \end{aligned} \quad (\text{C.11})$$

946 where  $\mathbf{f}(\mathbf{z}, \cdot) \# \mathbb{Q}_\epsilon$  denotes the pushforward distribution of  $\mathbb{Q}_\epsilon$  through the function  $\mathbf{f}(\mathbf{z}, \cdot)$ , i.e.,  
 947 the distribution of  $\mathbf{f}(\mathbf{z}, \epsilon)$  for a fixed  $\mathbf{z}$  when  $\epsilon \sim \mathbb{Q}_\epsilon$ . In other words, a DPA minimizes the negative  
 948 expected energy score between  $\mathbf{X}$  and the corresponding (stochastic) decoder output, conditional  
 949 on the encoding of  $\mathbf{X}$ . Due to this conditioning, the DPA objective differs from an energy distance by  
 950 a ‘‘normalization constant’’ which depends on the encoder and encourages capturing principal (i.e.,  
 951 variation-minimizing) components, rather than random latent dimensions.

#### 952 C.4 PERTURBATION MODELLING

953 **Compositional Perturbation Autoencoder (CPA).** Lotfollahi et al. (2023) propose the compo-  
 954 sitional perturbation autoencoder (CPA) as a model for compositional extrapolation of perturbation  
 955 data. Specifically, they assume the following model:

$$956 \quad \mathbf{z}^{\text{pert}} = \mathbf{z}^{\text{base}} + \mathbf{W}^{\text{pert}} \begin{pmatrix} h_1(l_1) \\ \dots \\ h_K(l_K) \end{pmatrix} + \sum_{j=1}^J \mathbf{W}_j^{\text{cov}} \mathbf{c}_j, \quad (\text{C.12})$$

957 where  $\mathbf{z}^{\text{base}} \sim \mathbb{P}_{\mathbf{Z}}$  denotes an unperturbed basal state; the matrix  $\mathbf{W}^{\text{pert}} \in \mathbb{R}^{d_Z \times K}$  encodes the  
 958 additive effect of each elementary perturbation;  $\{h_k : \mathbb{R} \rightarrow \mathbb{R}\}_{k \in [K]}$  are unknown, possibly non-  
 959 linear dose-response curves;  $\{\mathbf{c}_j \in \mathbb{R}^{K_j}\}_{j \in [J]}$  are observed one-hot vectors capturing  $J$  additional  
 960 discrete covariates, such as cell-types or species; and the matrices  $\{\mathbf{W}_j^{\text{cov}} \in \mathbb{R}^{d_Z \times K_j}\}_{j \in [J]}$  encode  
 961 additive covariate-specific effects. Further, the basal state  $\mathbf{z}^{\text{base}}$  is assumed to be independent of  
 962 the perturbation labels  $\mathbf{l}$  and covariates  $\mathbf{C} = (\mathbf{c}_1, \dots, \mathbf{c}_J)$ . Observations  $\mathbf{x}$  are then drawn from a  
 963 Gaussian whose mean and variance are determined by the perturbed latent state  $\mathbf{z}^{\text{pert}}$ ,

$$964 \quad \mathbf{x} \sim \mathcal{N} \left( \boldsymbol{\mu}(\mathbf{z}^{\text{pert}}), \sigma^2(\mathbf{z}^{\text{pert}}) \mathbf{I} \right). \quad (\text{C.13})$$

To fit this model, Lotfollahi et al. (2023) employ an adversarial autoencoder (Lample et al., 2017). First, an encoder  $g$  estimates the basal state

$$\widehat{\mathbf{z}}^{\text{base}} = g(\mathbf{x}). \quad (\text{C.14})$$

The estimated perturbed latent state  $\widehat{\mathbf{z}}^{\text{pert}}$  is then computed according to (C.12) using (C.14) and learnt estimates  $\widehat{\mathbf{W}}^{\text{pert}}$ ,  $\{\widehat{h}_k\}$ , and  $\{\widehat{\mathbf{W}}_l^{\text{cov}}\}$ . Finally, a (deterministic) decoder  $f$  uses  $\widehat{\mathbf{z}}^{\text{pert}}$  to compute estimates of the mean and variance in (C.13), i.e.,  $(\widehat{\boldsymbol{\mu}}, \widehat{\sigma}^2) = f(\widehat{\mathbf{z}}^{\text{pert}})$ . All learnable components of the model are trained by minimizing the (Gaussian) negative log-likelihood of the observed data  $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{l}_i, \mathbf{C}_i)\}_{i \in [N]}$ . To encourage the postulated independence of  $\widehat{\mathbf{z}}^{\text{base}}$  and  $(\mathbf{l}, \mathbf{C})$ , an additional adversarial loss is used, which minimizes the predictability of the latter from the former.

## D RELATION TO CAUSAL MODELS

In the field of causal inference, experimental data resulting from perturbations is typically modelled via *interventions* in an underlying causal model. In the *structural causal model* (SCM) framework of Pearl (2009), interventions modify a set of assignments, which determine each variable as a function of its direct causes and unexplained noise.

**Definition D.1** (Acyclic SCM). An acyclic structural causal model (SCM)  $\mathcal{M} = (\mathcal{S}, \mathbb{P}_{\mathcal{U}})$  with *endogenous* variables  $\mathbf{V} = \{V_1, \dots, V_n\}$  and *exogenous* variables  $\mathbf{U} = \{U_1, \dots, U_m\}$  consists of:

- (i) a set of *structural equations*

$$\mathcal{S} = \left\{ V_i := f_i \left( \mathbf{V}_{\text{pa}(i)}, U_i \right) \right\}_{i=1}^n, \quad (\text{D.1})$$

where  $f_i$  are deterministic functions;  $U_i \subseteq \mathbf{U}$ ; and  $\mathbf{V}_{\text{pa}(i)} \subseteq \mathbf{V} \setminus \{V_i\}$  is the set of *causal parents* of  $V_i$ , such that the *induced causal graph* with vertices  $[n]$  and edges  $j \rightarrow i$  iff.  $j \in \text{pa}(i)$  is acyclic;

- (ii) a joint distribution  $\mathbb{P}_{\mathcal{U}}$  over the exogenous variables.

The *induced distribution*  $\mathbb{P}_{\mathbf{V}}$  of  $\mathcal{M}$  is given by the push-forward of  $\mathbb{P}_{\mathcal{U}}$  via  $\mathcal{S}$ . An intervention replaces  $\mathcal{S}$  by a new set of assignments  $\mathcal{S}'$  such the graph induced by  $\mathcal{S}'$  is acyclic. The *interventional distribution* is the induced distribution  $\mathcal{M}' = (\mathcal{S}', \mathbb{P}_{\mathcal{U}})$ .<sup>4</sup>

Our assumed generative process for the effect of perturbations (§ 3) is, in general, different from interventions in an SCM. However, as we show next, certain types of SCMs and interventions are recovered as a special case of our model.

### D.1 SHIFT INTERVENTIONS IN LINEAR SCMS AS A SPECIAL CASE

Consider a linear SCM with  $d_Z$  endogenous variables  $\mathbf{Z}$  and  $d_U$  exogenous variables  $\mathbf{U}$  of the form

$$\mathbf{Z} := \mathbf{A}\mathbf{Z} + \mathbf{U}, \quad (\text{D.2})$$

where  $\mathbf{A}$  is a (lower-triangular) weighted adjacency matrix. The observational (i.e., non-intervened) distribution of  $\mathbf{Z}$  induced by (D.2) is most easily understood via the *reduced form* expression

$$\mathbf{Z} = (\mathbf{I} - \mathbf{A})^{-1}\mathbf{U} \quad (\text{D.3})$$

and thus given by  $\mathbb{P}_{\mathbf{Z}} = (\mathbf{I} - \mathbf{A})_{\#}^{-1}\mathbb{P}_{\mathcal{U}}$ . Now consider the class of *shift interventions* parametrised by constant shift vectors  $\mathbf{c}_e \in \mathbb{R}^{d_Z}$ , which modify the original SCM in (D.2) to

$$\mathbf{Z}_e := \mathbf{A}\mathbf{Z}_e + \mathbf{U} + \mathbf{c}_e. \quad (\text{D.4})$$

Analogous to (D.3), the reduced form of (D.4) is given by

$$\begin{aligned} \mathbf{Z}_e &= (\mathbf{I} - \mathbf{A})^{-1}\mathbf{U} + (\mathbf{I} - \mathbf{A})^{-1}\mathbf{c}_e \\ &= \mathbf{Z} + (\mathbf{I} - \mathbf{A})^{-1}\mathbf{c}_e, \end{aligned} \quad (\text{D.5})$$

<sup>4</sup>Interventions can also introduce new sources of randomness  $\mathbf{U}'$ , which, for sake of simplicity, we do not consider here.

1026 where the second equality follows from (D.3).

1027 Thus, if we take the shift vectors as perturbation labels (i.e.,  $K = d_Z$  and  $l_e = c_e$ ) and use a linear  
1028 perturbation model of the form

$$1030 \mathbf{Z}_e^{\text{pert}} = \phi(\mathbf{Z}_e^{\text{base}}, l_e) = \mathbf{Z}_e^{\text{base}} + \mathbf{W}l_e \quad (\text{D.6})$$

1031 with  $\mathbf{W} = (\mathbf{I} - \mathbf{A})^{-1}$ , then our perturbation model captures shift interventions in a linear SCM  
1032 with adjacency matrix  $\mathbf{A}$ .

1033 *Remark D.2.* The above argument does not require causal sufficiency: it still holds if  $\mathbb{P}_{\mathcal{U}}$  is not  
1034 factorized (e.g., due to hidden confounding).

1035 *Remark D.3.* When trained on data generated according to (D.4), the adjacency matrix associated  
1036 with the learnt perturbation matrix  $\widehat{\mathbf{W}}$  is given by  $\widehat{\mathbf{A}} = \mathbf{I} - \widehat{\mathbf{W}}^{-1}$ .

1037 *Remark D.4.* The correspondence between mean shift perturbations and shift interventions appears  
1038 to only hold for *linear* SCMs. Consider, for example, a nonlinear additive noise model,

$$1040 \mathbf{Z} := \mathbf{f}(\mathbf{Z}) + \mathbf{U}, \quad (\text{D.7})$$

1041 with reduced form given by

$$1042 \mathbf{Z} = \mathbf{g}(\mathbf{U}), \quad (\text{D.8})$$

1043 where  $\mathbf{g}$  is the inverse of the mapping  $z \mapsto z - \mathbf{f}(z)$ . For shift interventions in (D.7) to match our  
1044 perturbation model, we then must have

$$1045 \mathbf{Z}_e = \mathbf{g}(\mathbf{U} + c_e) = \mathbf{g}\left(\mathbf{g}^{-1}(\mathbf{Z}) + c_e\right) = \phi(\mathbf{Z}, l_e) \quad (\text{D.9})$$

1046 for suitable  $\phi$  and  $l_e$ . Thus, if  $\mathbf{f}$  is nonlinear, so is  $\mathbf{g}$  and therefore  $\phi$ . In other words, shift interven-  
1047 tions in a nonlinear SCM do not, in general, amount to mean shift perturbations (i.e., linear  $\phi$ ).

1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079