

Provably Learning Disentangled & Compositional Models

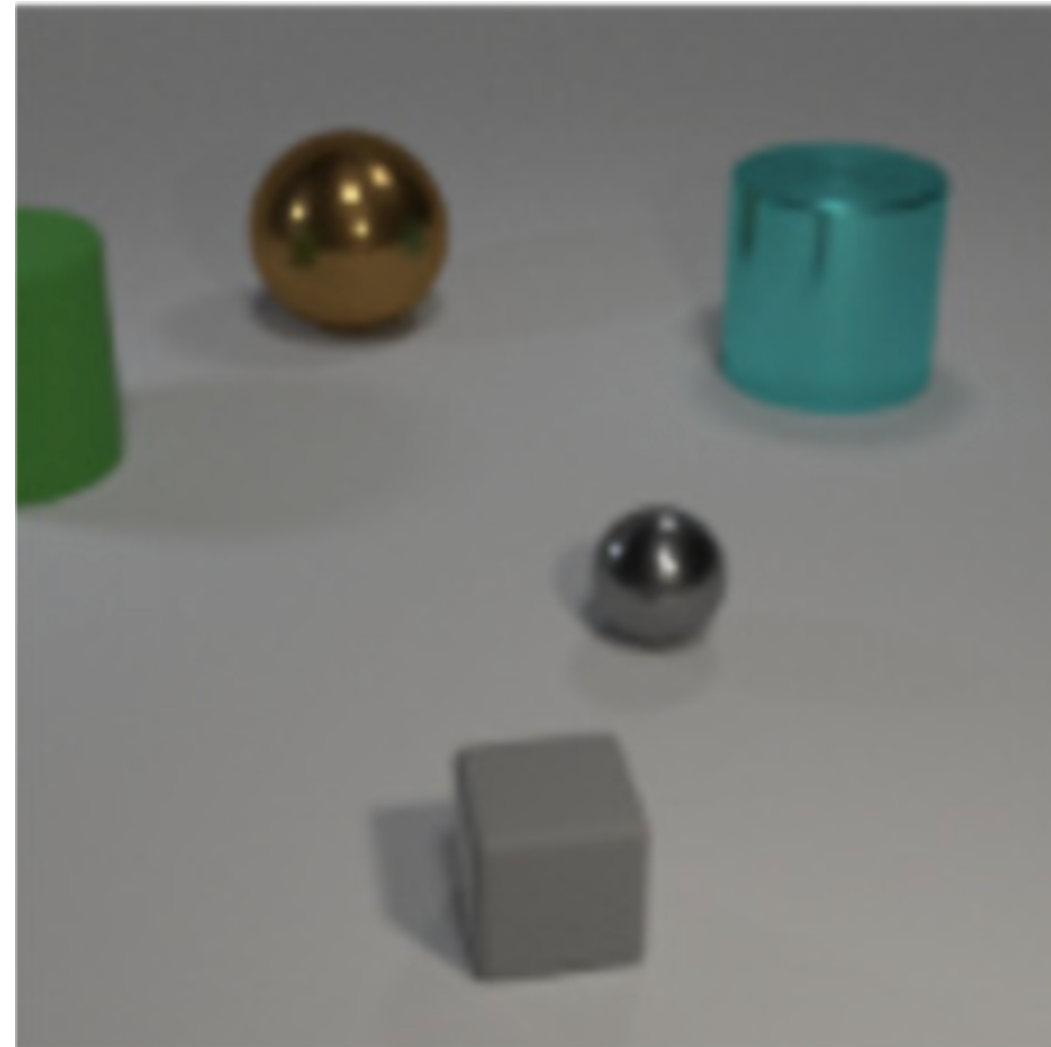
Prédoc 3 Examination

Student: Divyat Mahajan

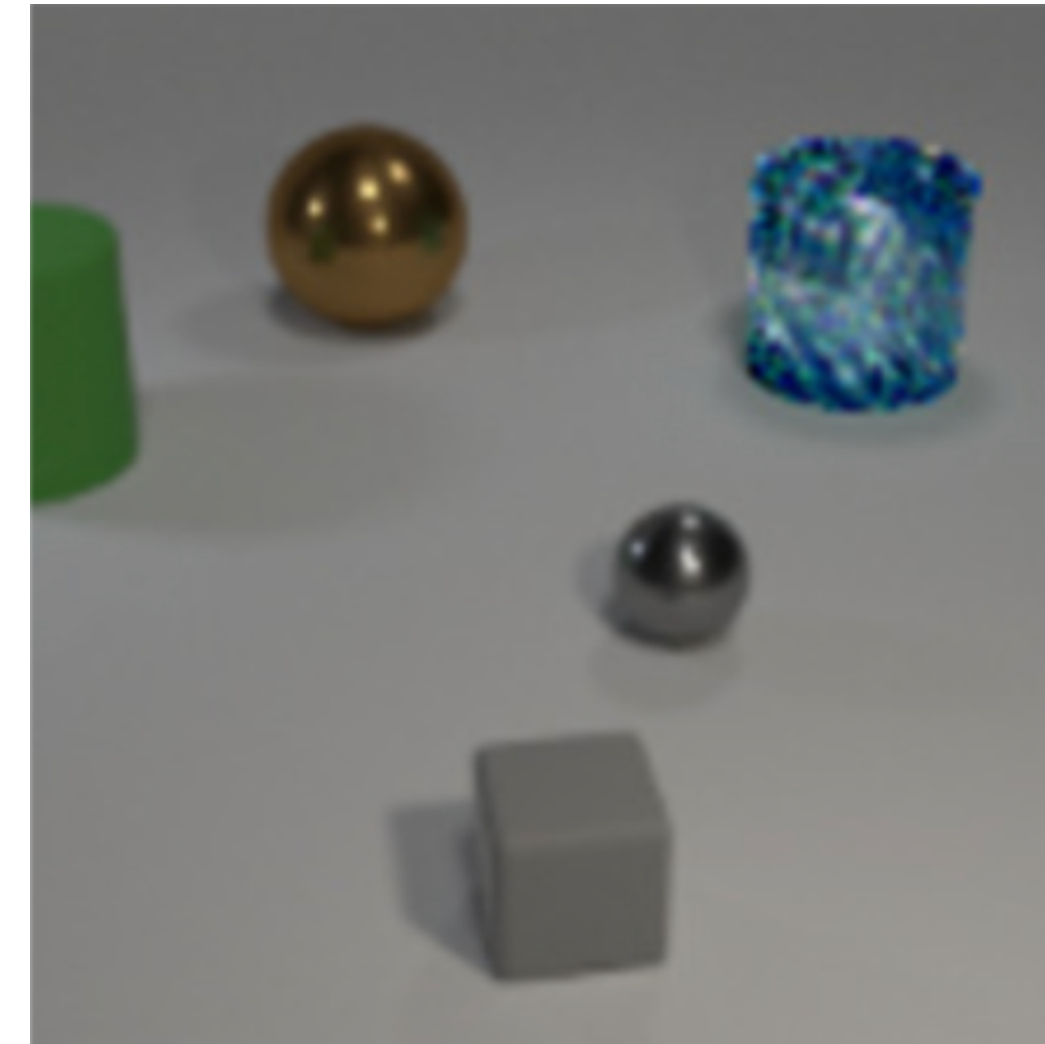
Committee: Ioannis Mitliagkas, Simon Lacoste-Julien, Yoshua Bengio

Covariate Shift

Training Samples



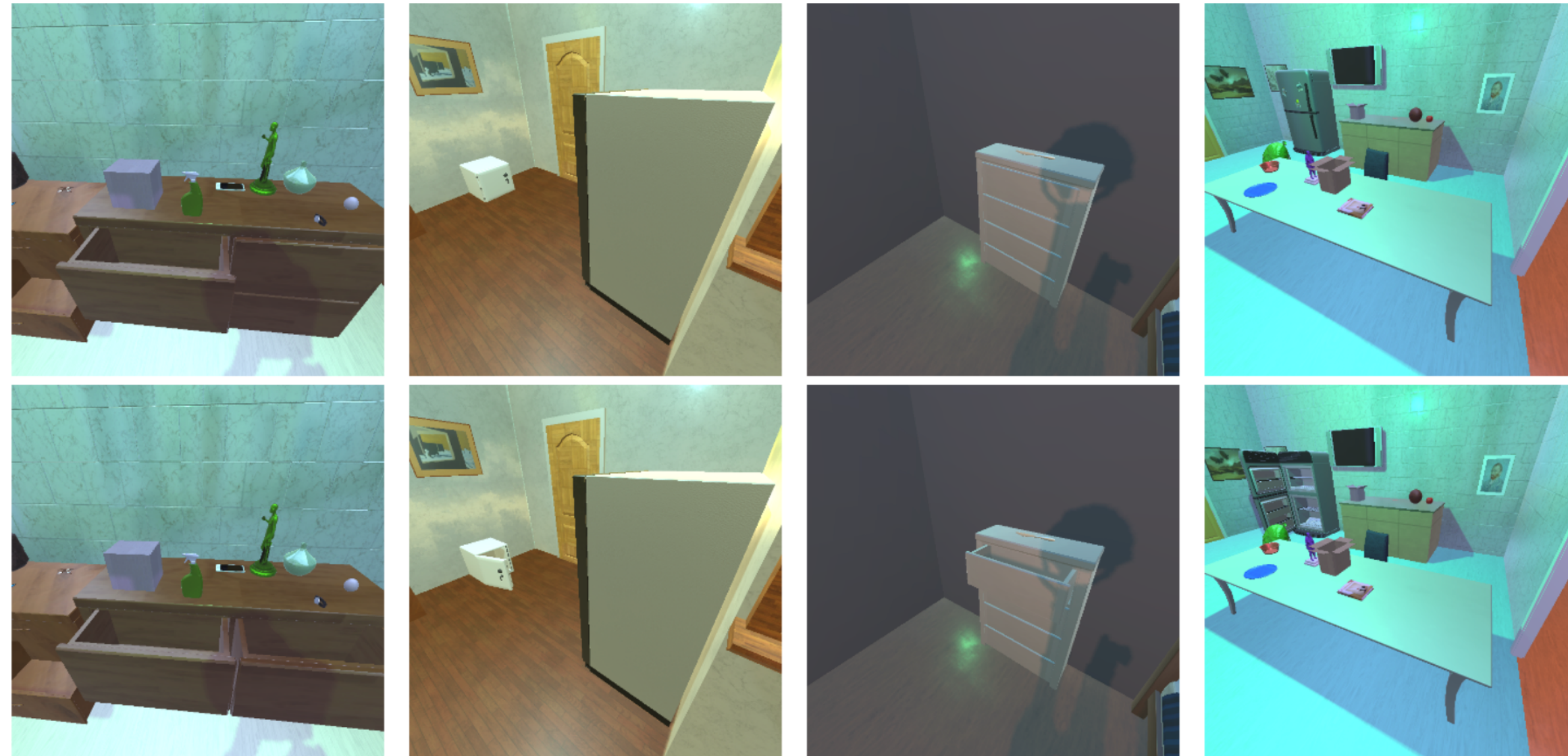
Test Samples



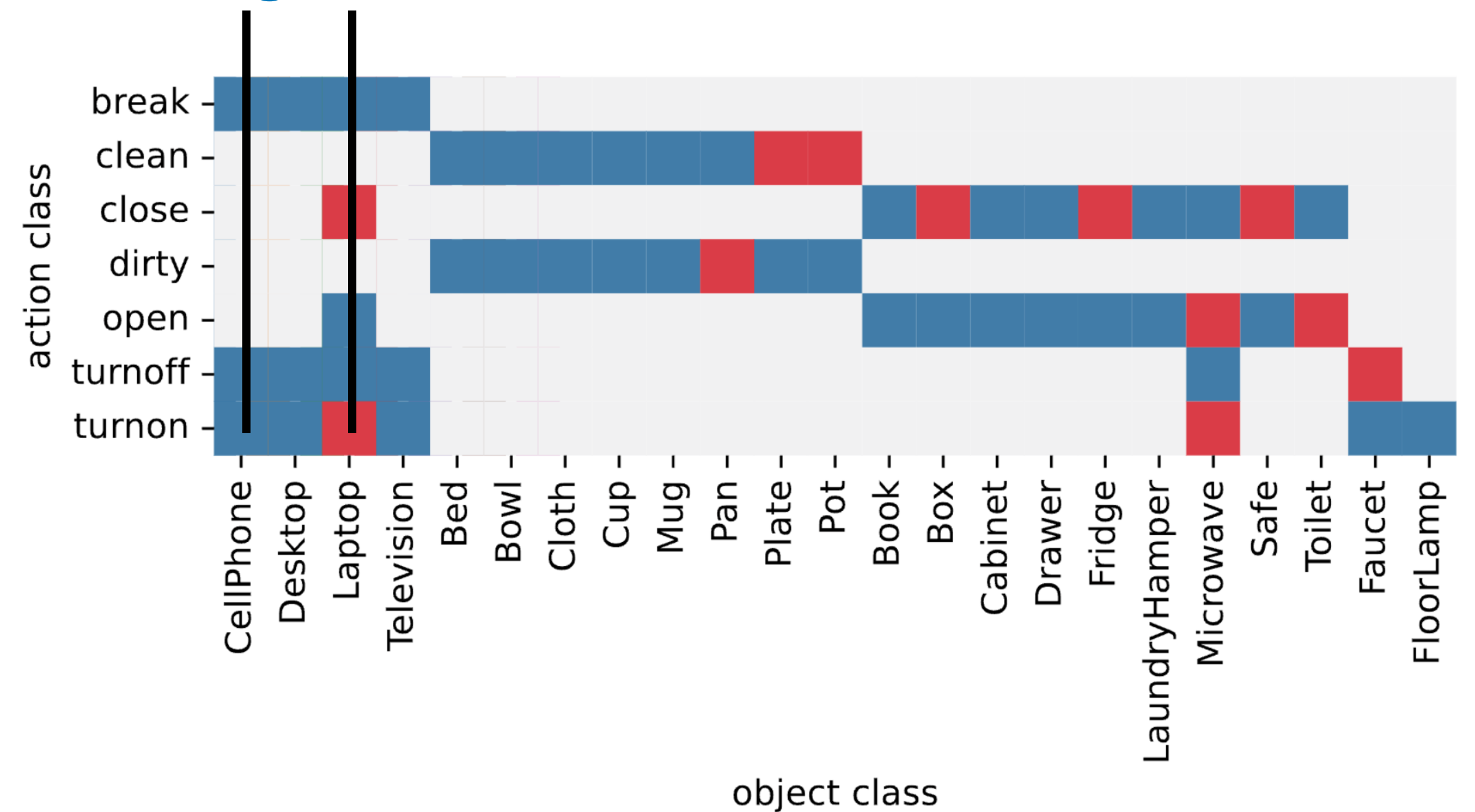
Source: Generalization and Robustness Implications in Object-Centric Learning by Dittadi et al. (2022)

Hypothesis: Learning disentangled representations can allow us to efficiently adapt to covariate shifts as it changes mechanisms in a sparse manner

Compositional Shift



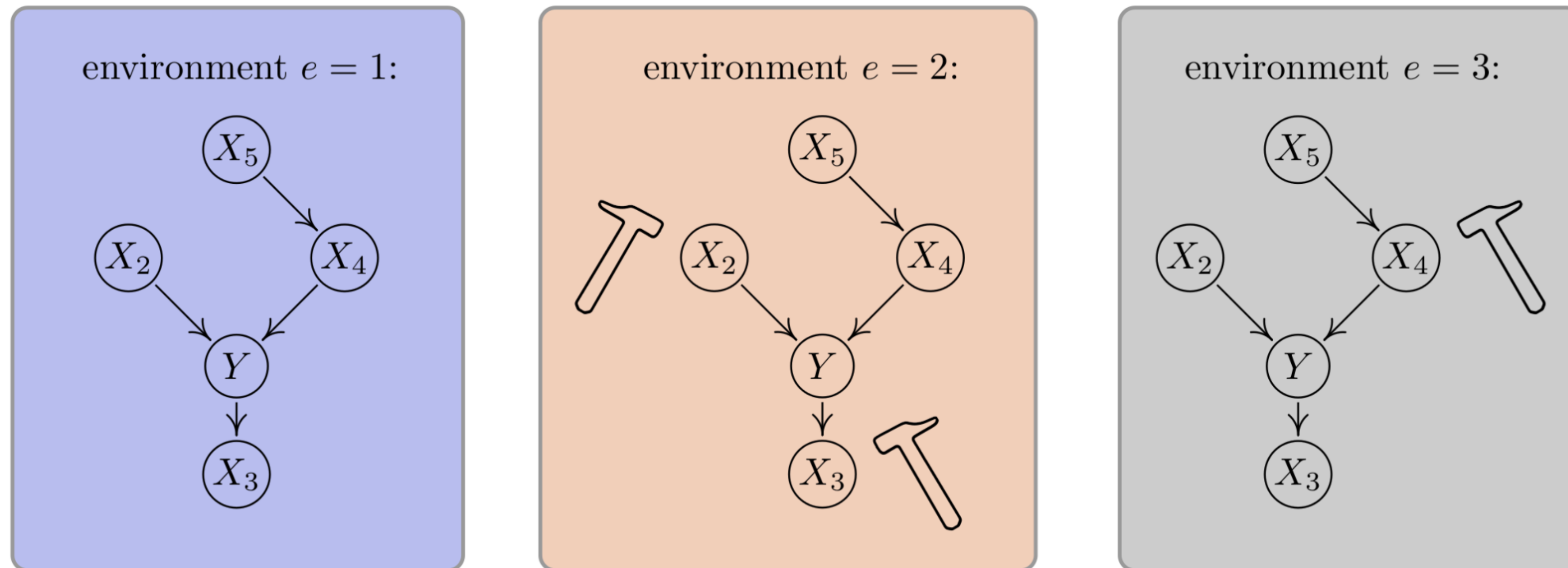
Training Test



Source: Causal Triplet by Liu et al. (2023)

Hypothesis: Learning disentangled representations can allow us to efficiently extrapolate to novel compositions

Distribution Shifts in SCMs

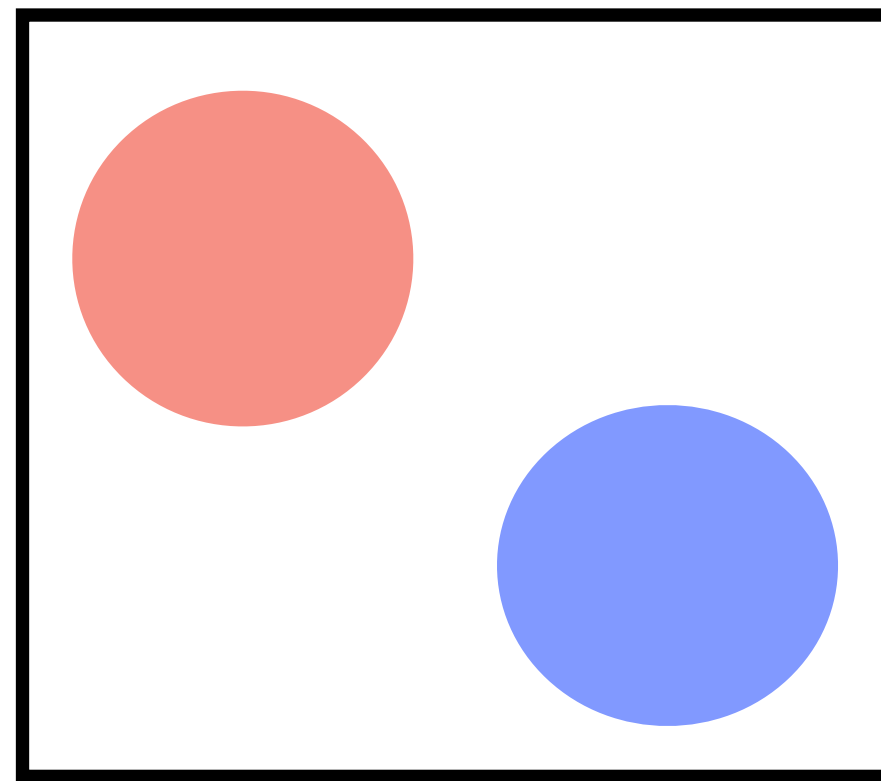


Source: Causal Inference using Invariant Prediction by Peters et al. (2015)

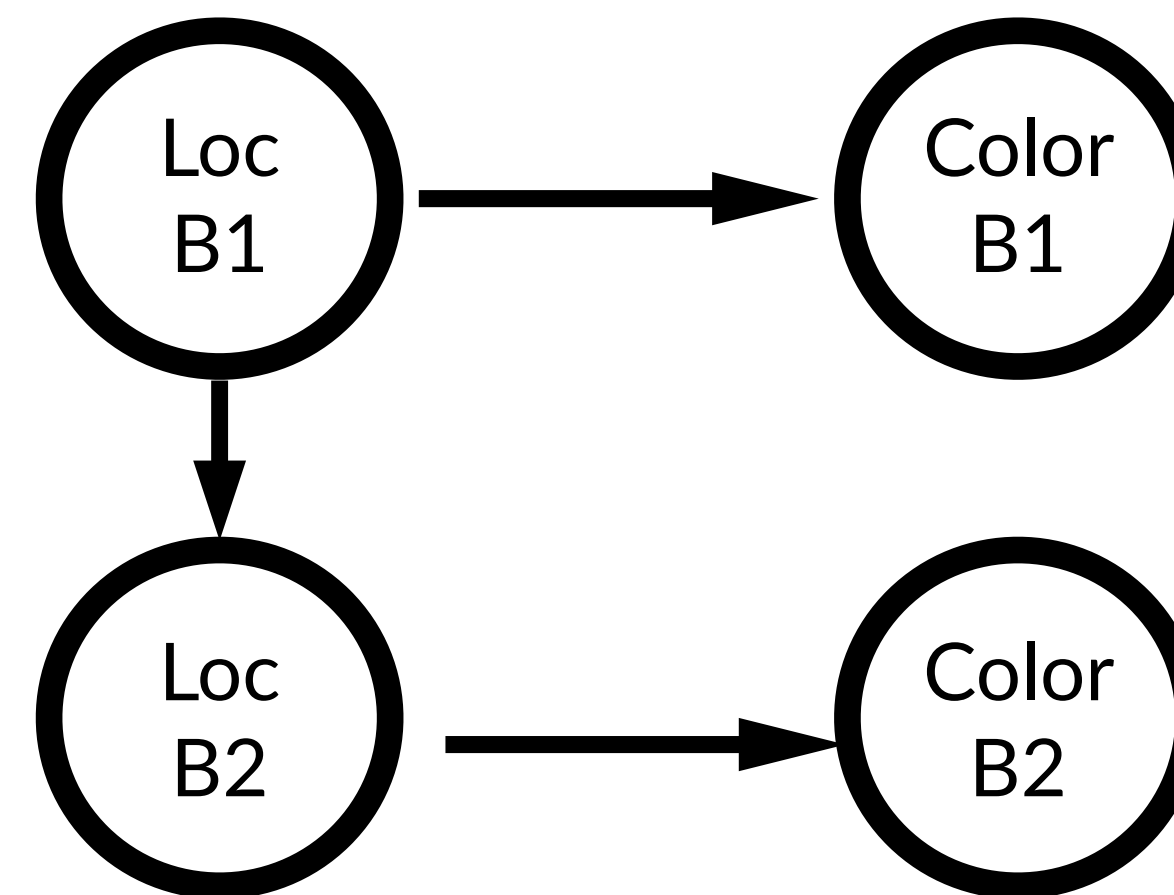
- **Independent Causal Mechanisms (ICM):** Changing one causal mechanism leads to no change in the other causal mechanisms
- **Sparse Mechanism Shift:** Effect of interventions is modular in structural causal models

Disentangled Representation Learning

Input



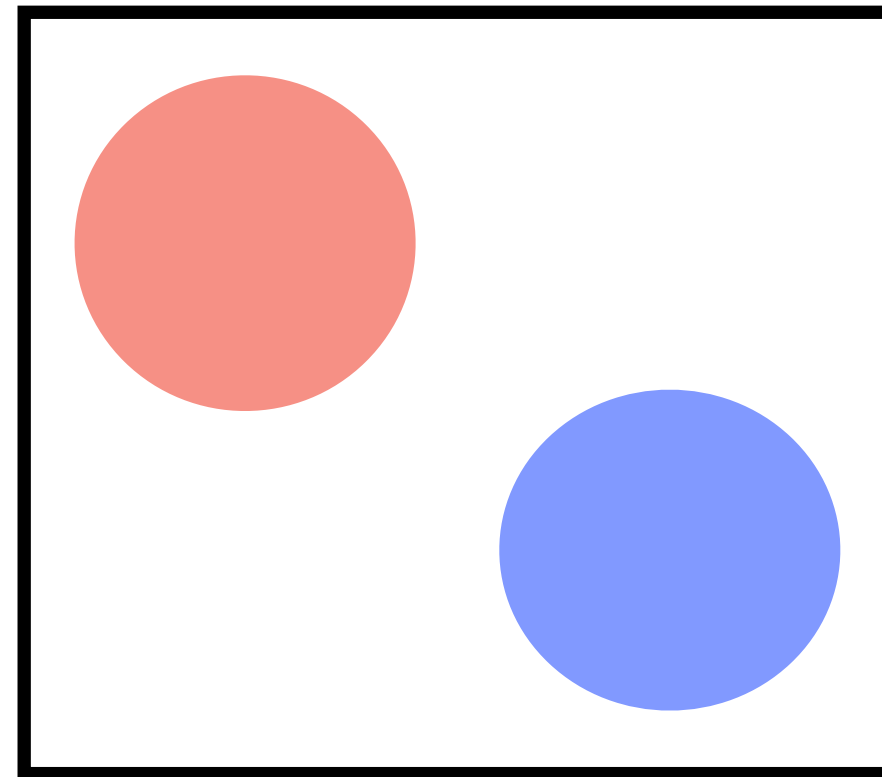
Latent Factors of Variations



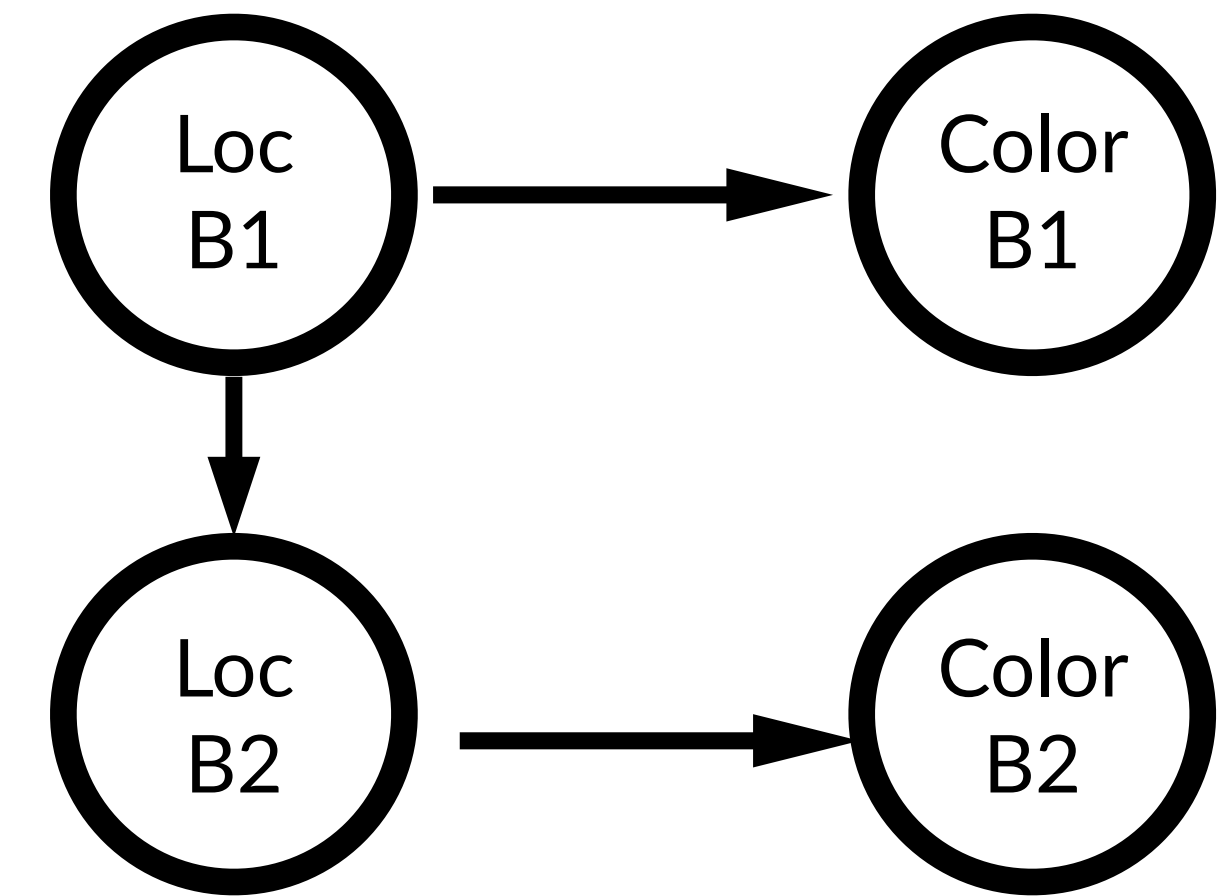
- **Setup:** $x = g(z)$ where $z \in \mathbb{R}^{d_z}$ are the latent (causal) factors of the data generation process (DGP), that are transformed to observations $x \in \mathbb{R}^{d_x}$
- **Goal:** Invert the DGP to get latent factors (z) from observations (x)

Indeterminacy in Latent Recovery

Input



Latent Factors of Variations

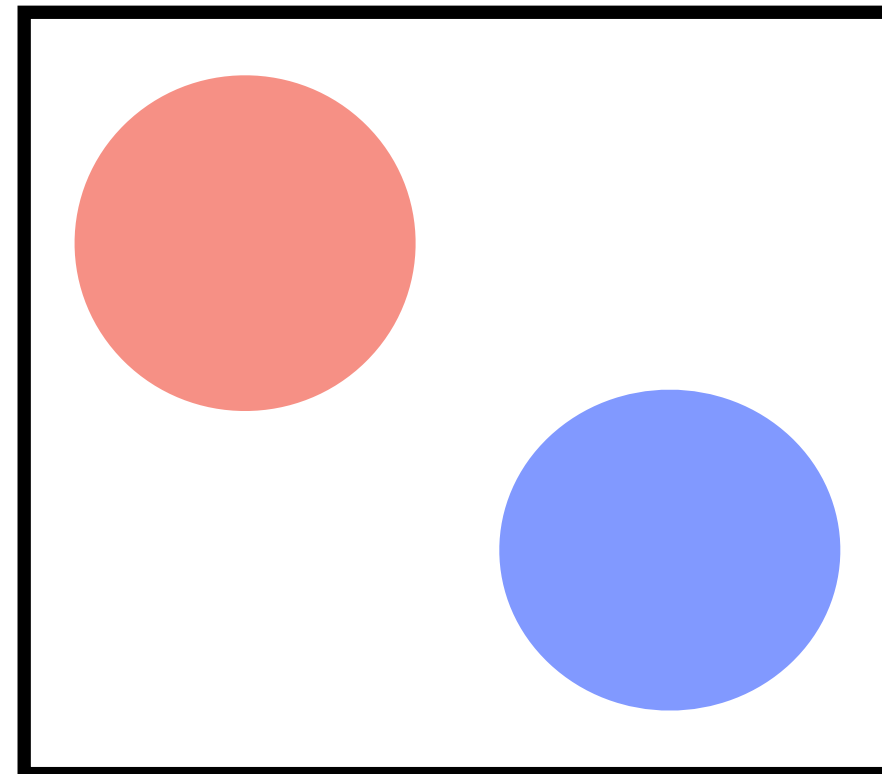


- **Reconstruction Objective:** Optimal encoder ($\hat{f} : \mathbb{R}^{d_x} \rightarrow \mathbb{R}^{d_z}$) and the optimal decoder ($\hat{g} : \mathbb{R}^{d_z} \rightarrow \mathbb{R}^{d_x}$) satisfy the following.

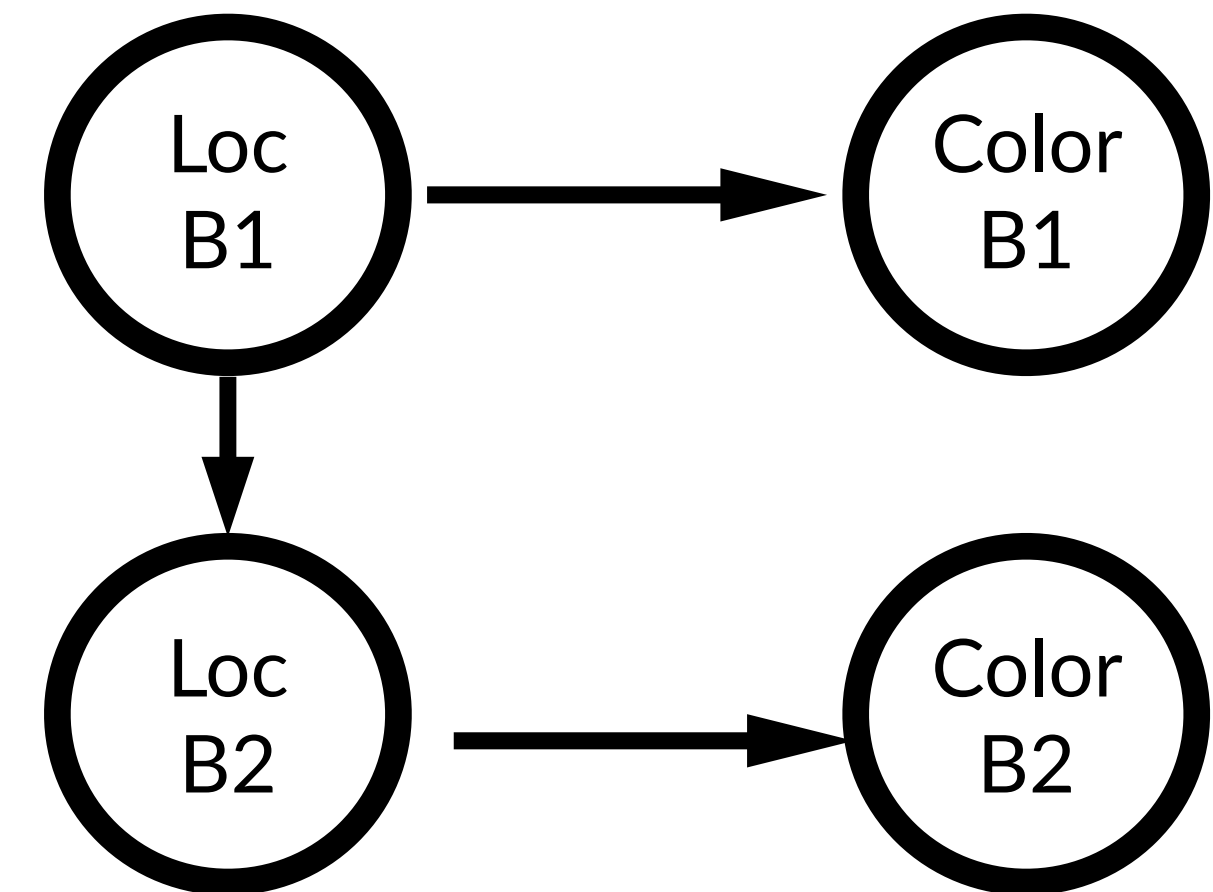
$$\mathbb{E}_{x \sim X} ||x - \hat{g}(\hat{f}(x))||^2 = 0 \quad \implies \quad \hat{z} = v(z) ; v(z) = \hat{g}^{-1} \circ g(z)$$

Indeterminacy in Latent Recovery

Input

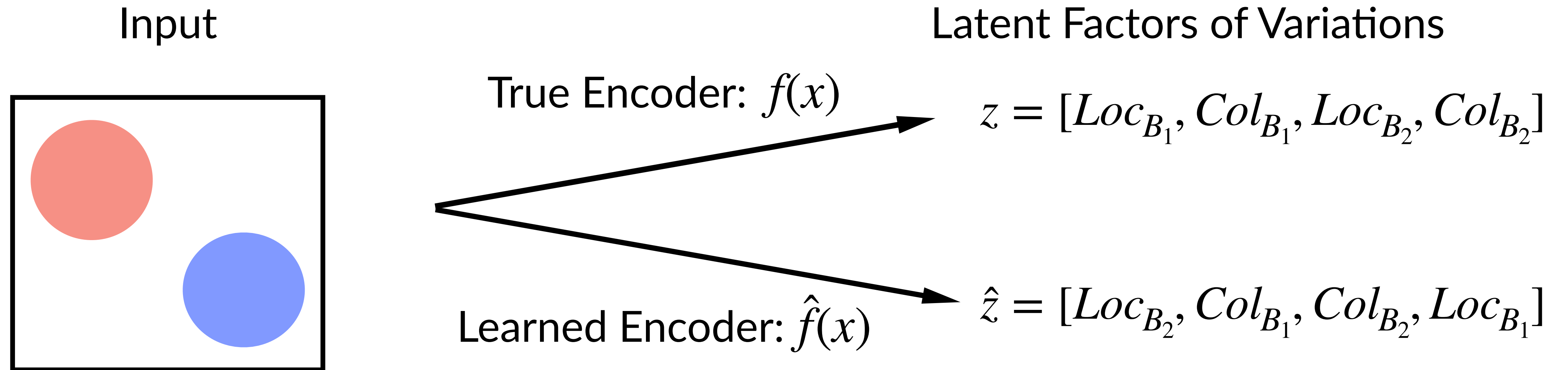


Latent Factors of Variations



- We need to constrain the indeterminacy in latent recovery; $\hat{z} = v(z) = \hat{g}^{-1} \circ g(z)$
- Both (\hat{f}, \hat{g}) and (f, g) explain data equally well, however, learned latents \hat{z} might be a complex transformation of true latents z

Latent Identification



- **Permutation & Scaling Identification:** $\hat{z} = \Pi \circ \Lambda z + b$ where Π is permutation matrix and Λ is invertible diagonal matrix
- **Local Disentanglement:** $\hat{z} = v(z)$ where Jacobian of v is permuted diagonal matrix

How to achieve identification guarantees?

$$\hat{z} = v(z) = \hat{g}^{-1} \circ g(z)$$

- Constraints on the mixing function (g) and learned decoder (\hat{g})
- Constrains on the latent distribution ($\mathbb{P}(Z)$) and enforcing learned latents ($\hat{z} = v(z)$) to satisfy them as well

How to achieve identification guarantees?

$$\hat{z} = v(z) = \hat{g}^{-1} \circ g(z)$$

- **Linear ICA:**
 - Constrain g, \hat{g} to be linear functions
 - Leads to linear identification as $v(z) = \hat{g}^{-1} \circ g(z)$ is a linear function
 - Constrain z, \hat{z} to have mutually independent components and *all components of z are non-gaussian*
 - Further restricts the linear $v(z)$ to permutation & scaling matrix.

Solving Non-linear ICA

Challenging Common Assumptions in the Unsupervised Learning of Disentangled Representations

Francesco Locatello^{1,2} Stefan Bauer² Mario Lucic³ Gunnar Rätsch¹ Sylvain Gelly³ Bernhard Schölkopf²
Olivier Bachem³

Unlike linear ICA, restricting Z, \hat{Z} to have mutually independent components is **not sufficient** to guarantee disentanglement for non-linear ICA!

Today's Talk

Disentanglement with Auxiliary Information

Towards Efficient Representation Identification in Supervised Learning

Kartik Ahuja*, **Divyat Mahajan***, Vasilis Syrgkanis, Ioannis Mitliagkas

Conference on Causal Learning and Reasoning [CleaR 2022]

Unsupervised Disentanglement & Cartesian-Product Extrapolation

Additive Decoders for Latent Variables Identification and Extrapolation

Sébastien Lachapelle*, **Divyat Mahajan***, Ioannis Mitliagkas, Simon Lacoste-Julien

Advances in Neural Information Processing Systems [NeurIPS 2023 (Oral)]

Extrapolation with Discrete Factors

Compositional Generalization with Additive Energy Models

Ongoing work in collaboration with Kartik Ahuja, Ioannis Mitliagkas, Mohammad Pezeshki, Pascal Vincent

Other Contributions

Causal Inference with Observational Data

Empirical Analysis of Model Selection for Heterogeneous Causal Effect Estimation

Divyat Mahajan, Ioannis Mitliagkas, Brady Neal, Vasilis Syrgkanis

International Conference on Learning Representations [ICLR 2024 (Spotlight)]

Disentanglement with Interventional Data

Interventional Causal Representation Learning

Kartik Ahuja, **Divyat Mahajan**, Yixin Wang, Yoshua Bengio

International Conference on Machine Learning [ICML 2023 (Oral)]

Benefits of Disentanglement for Downstream Tasks

Synergies between Disentanglement and Sparsity in Multi-Task Learning

Sébastien Lachapelle*, Tristan Deleu*, **Divyat Mahajan**, Ioannis Mitliagkas, Yoshua Bengio,

Simon Lacoste-Julien, Quentin Bertrand

International Conference on Machine Learning [ICML 2023]

Towards Efficient Representation Identification in Supervised Learning

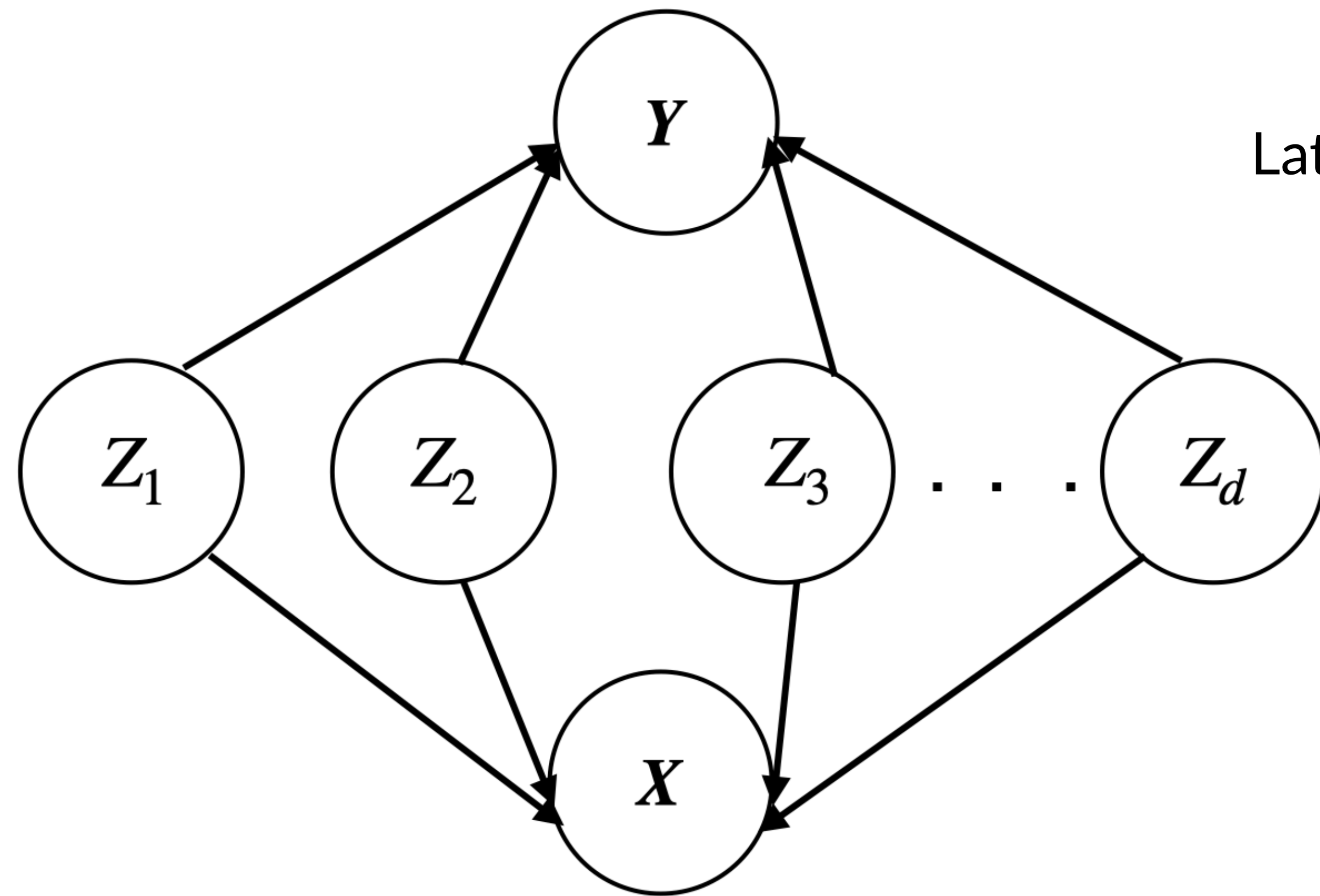
Kartik Ahuja*, Divyat Mahajan*, Vasilis Syrgkanis & Ioannis Mitliagkas

Conference on Causal Learning and Reasoning (CLeaR) 2022

*Equal contribution



Non Linear ICA with Auxiliary Information



$$Z = (Z_1, \dots, Z_d)$$

Latent Variable: Mutually independent & Non-Gaussian

$$Y \leftarrow \Gamma Z + N$$

Auxiliary Information: $Y \in \mathbf{R}^k$ & $Z \perp N$

$$X \leftarrow g(Z)$$

Observed non-linear mixing of latents, g is bijection

Independence Constrained ERM

Model: $W \circ \Phi$

$W \in \mathbb{R}^{d \times k}$: Linear Classifier

$\Phi \in \mathcal{H}_\Phi$: Non-Linear Representation

Empirical Risk Minimization (ERM):
$$\min_{W \in \mathbb{R}^{d \times k}, \Phi \in \mathcal{H}_\Phi} \sum_{i=1}^N \ell(W \circ \Phi(x_i), y_i)$$

IC-ERM:
$$\min_{W \in \mathbb{R}^{d \times k}, \Phi} \sum_{i=1}^N \ell(W \circ \Phi(x_i), y_i) \text{ s.t. Components of } \Phi(x) \text{ are i.i.d.}$$

Identification with IC-ERM

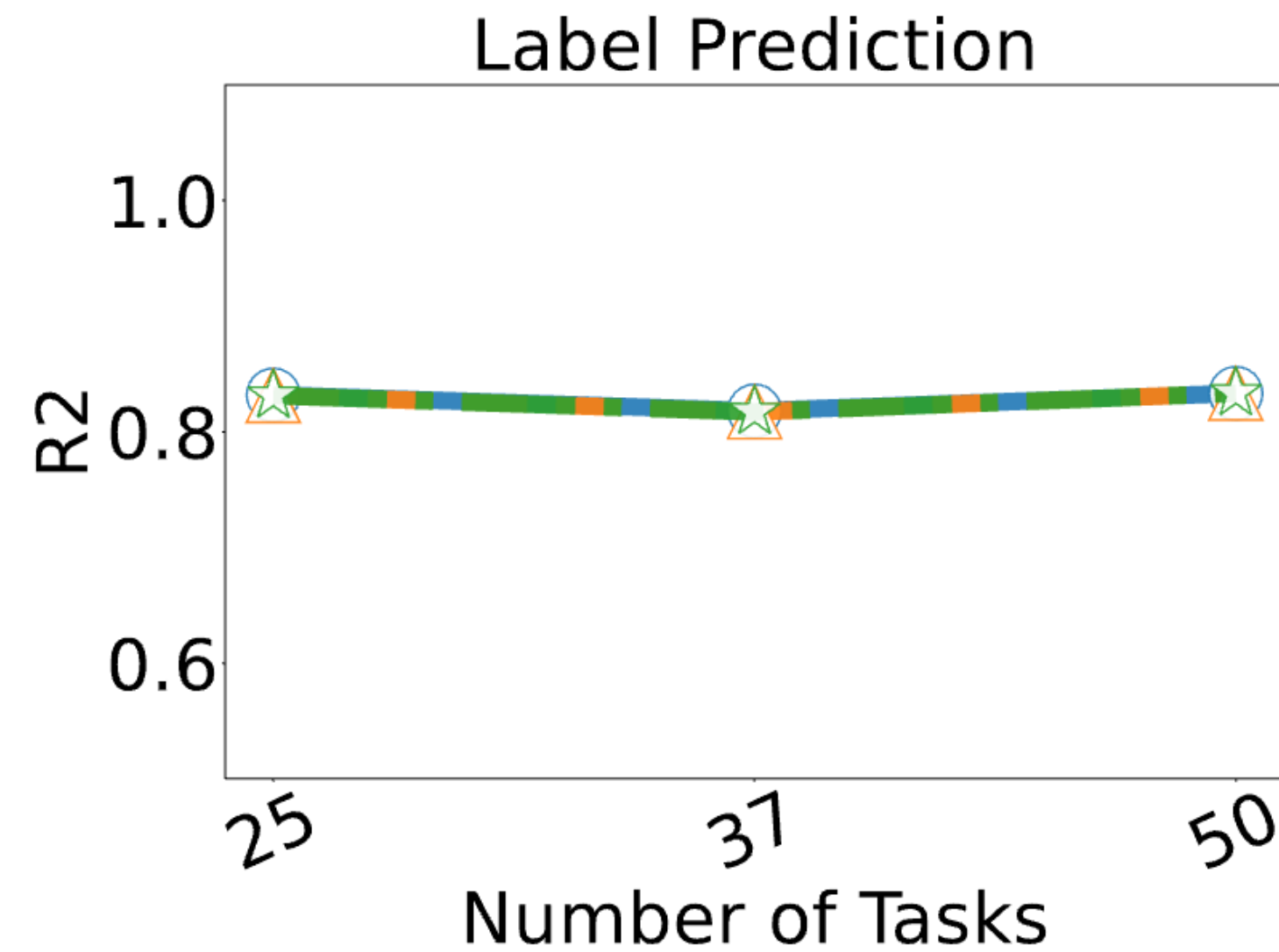
Assumption: Dimension of the label (k) is equal to the dimension of the latent (d)

Theorem (Informal): Under the above **assumption** as well as those on the data generation process (**mutual independence** of Z), we have the following:

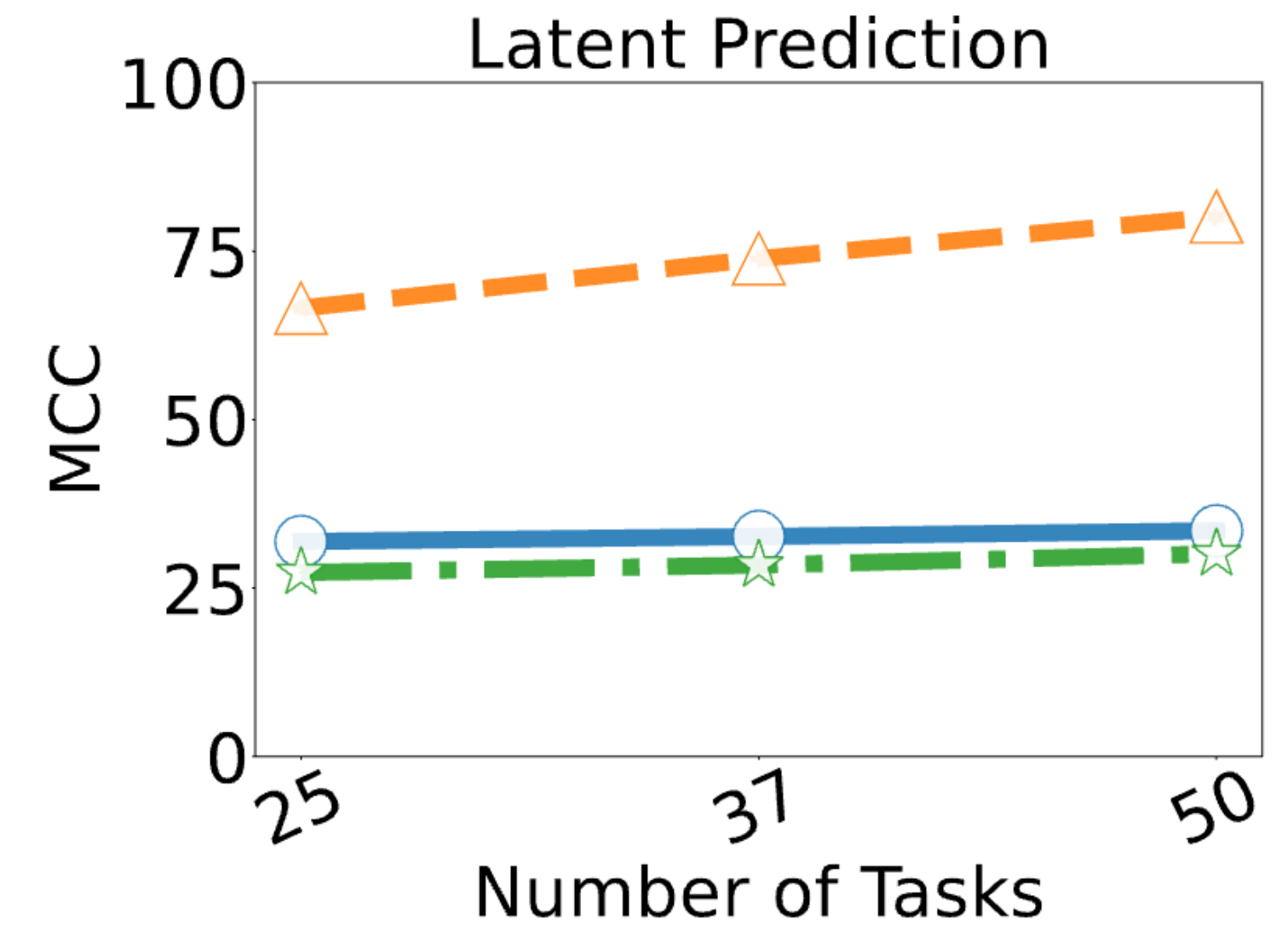
- **ERM:** Optimal solutions identify true latents up to **linear transformation**
- **IC-ERM:** Optimal solutions identify true latents up to **permutation & scaling**

Note: We also present identification results in the paper when $k < d$

Experiments



ERM ERM-ICA ERM-PCA



ERM ERM-ICA ERM-PCA

Results for regression task with latent dimension $d = 50$.

Disentanglement performance (MCC) improves as we observe more tasks.

Extending theory beyond mutual independence

Assumption: Latent Variables Z are mutually independent

Solution: Assume Γ to be sparse where $Y \leftarrow \Gamma Z + N$

Synergies between Disentanglement & Sparsity: Generalization & Identifiability in Multi-Task Learning

Sébastien Lachapelle*, Tristan Deleu*, Divyat Mahajan, Ioannis Mitliagkas, Yoshua Bengio, Simon Lacoste-Julien & Quentin Bertrand

*Equal contribution

International Conference on Machine Learning (ICML) 2023

Can we identify latents without auxiliary information?

Additive Decoders for Latent Variables Identification and Cartesian-Product Extrapolation

Sébastien Lachapelle*, Divyat Mahajan*, Ioannis Mitliagkas & Simon Lacoste-Julien

Neural Information Processing Systems (*NeurIPS*) 2023 (*Oral*)

*Equal contribution



Additive Decoders

$$x = g(z) = \sum_{B \in \mathcal{B}} g^{(B)}(z_B)$$

Observation
e.g. an image

Latent
Factors

Partition of $\{1, \dots, d_z\}$
e.g. $\mathcal{B} = \{\{1,2\}, \{3,4\}\}$

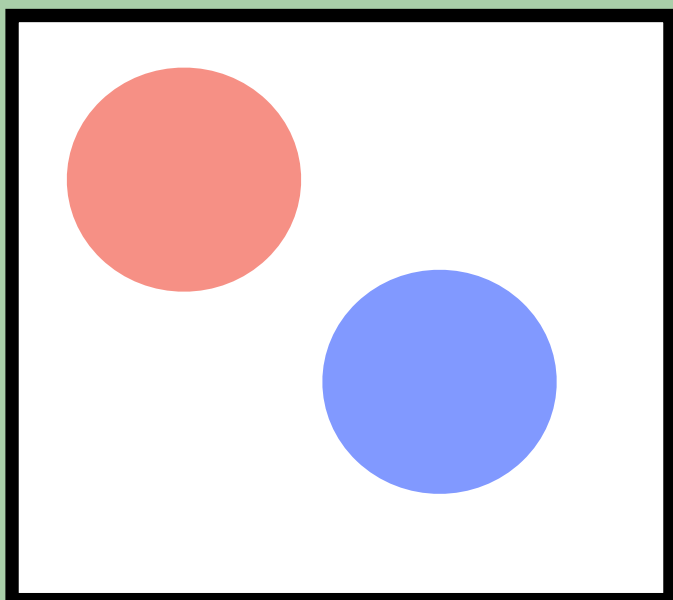
Sub-blocks of z

Example: Images of moving balls

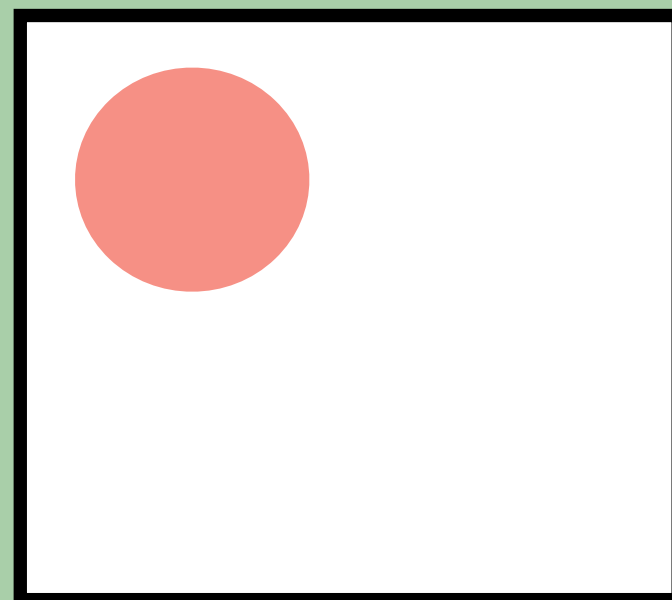
$$x = g(z)$$

$$g^{(B_1)}(z_1)$$

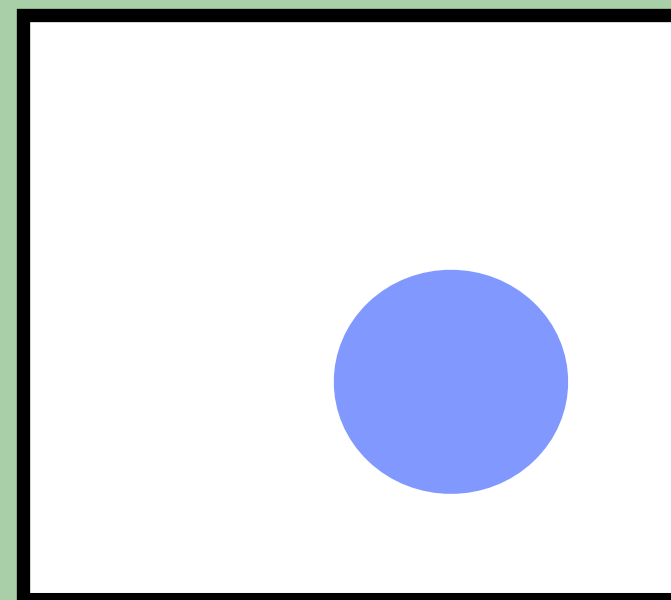
$$g^{(B_2)}(z_2)$$



=



+



$$\mathcal{B} = \{\{1,2\}, \{3,4\}\}$$

$$z_{B_1} = (z_1, z_2) \text{ Coordinates of } \text{red ball}$$

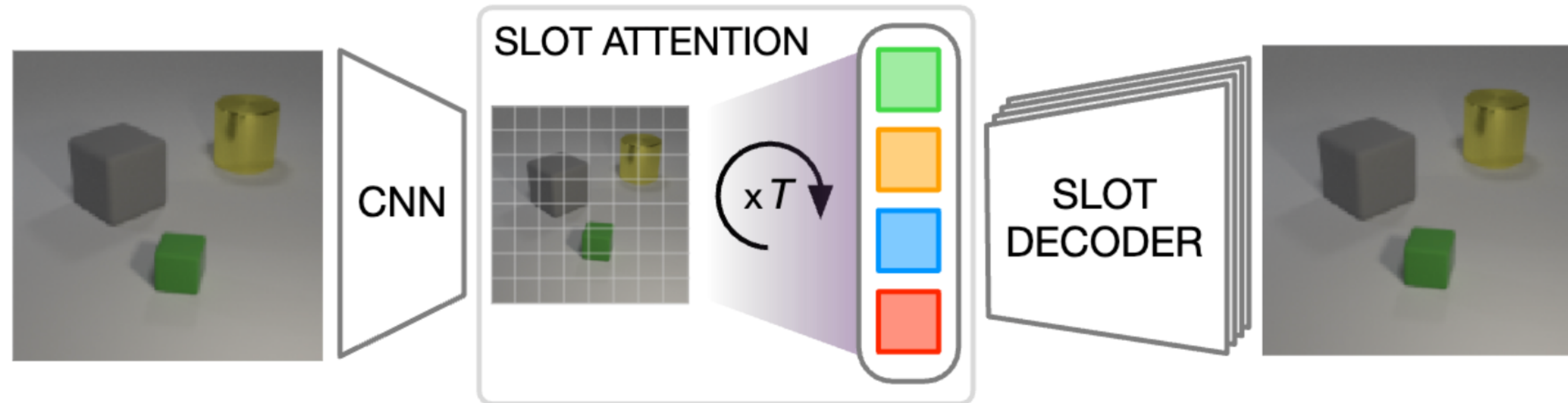
$$z_{B_2} = (z_3, z_4) \text{ Coordinates of } \text{blue ball}$$

$g^{(B)}$ Block-specific Decoder

Contribution

We introduce **additive decoders**: a simple architecture similar to object-centric decoders for which we can prove both **disentanglement** and **extrapolation** guarantees.

Decoder Architecture in Object-Centric Learning



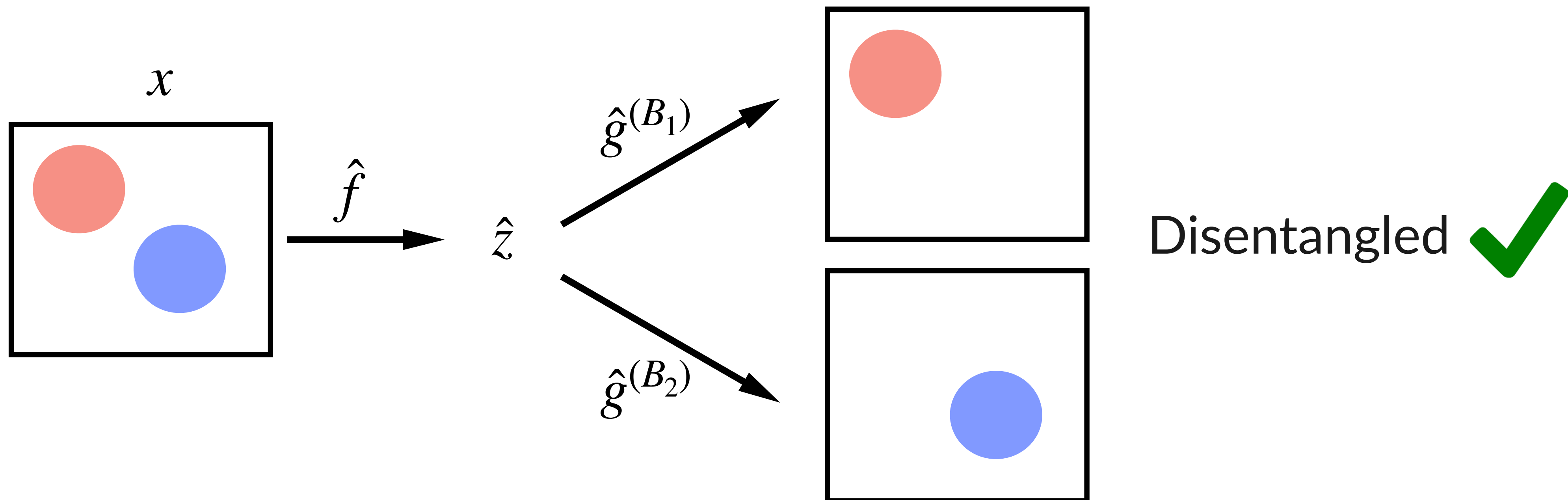
$$x = g(z) = \sum_{B \in \mathcal{B}} m^{(B)}(z) \odot g^{(B)}(z_B) \text{ where } m^{(B)} \text{ denote masking mechanism}$$

Object-centric learning approaches have shown impressive performance at disentanglement without using any weak supervision!

Block Disentanglement

- Learned Encoder: $\hat{f}(x)$
- Learned Additive Decoder: $\hat{g}(z) = \sum_{B \in \mathcal{B}} \hat{g}^{(B)}(z_B)$

If we optimise reconstruction loss perfectly, i.e., $\mathbb{E}[||x - \hat{g}(\hat{f}(x))||] = 0$, can we guarantee disentanglement of latent blocks?



Definition of Block Disentanglement

Learned decoder \hat{g} is disentangled w.r.t ground-truth decoder g if the learned block-specific decoders “imitate” the ground-truth ones

Definition of Block Disentanglement

Learned decoder \hat{g} is disentangled w.r.t ground-truth decoder g if the learned block-specific decoders “imitate” the ground-truth ones

Precisely, for all $B \in \mathcal{B}$ we have $v_{\pi(B)}(z) = \bar{v}_{\pi_B}(z_B)$

$$\hat{g}^{(B)}(z_B) = g^{(\pi(B))}(v_{\pi(B)}(z)) + c^{(B)}$$

Permutation that sends
blocks to blocks, i.e., $\pi(B) \in \mathcal{B}$

Invertible
Transformation

$$\sum_{B \in \mathcal{B}} c^{(B)} = 0$$

Local & Global Disentanglement

Local Disentanglement: $\pi(B)$ depends on z
Global Disentanglement: $\pi(B)$ independent of z

Local Disentanglement: $D_{i,j}^{\mathcal{V}_{\pi(B)}}(z) = 0 \quad \forall i \in \pi(B), j \notin B$

The permutation map can vary with samples in the case of local disentanglement
No unique mapping between the learned and true block decoders!

Latent Identification Guarantee

Theorem (Informal): Under the following assumptions

- Data Generation Process is additive, i.e, $x = \sum_{B \in \mathcal{B}} g^{(B)}(z_B)$
- Learned decoder is additive as well with total latent partitions as $|\mathcal{B}|$
- Ground-truth decoder is sufficiently non-linear (see paper)
- Block-specific decoders $g^{(B)}, \hat{g}^{(B)}$ are injective (for global disentanglement)

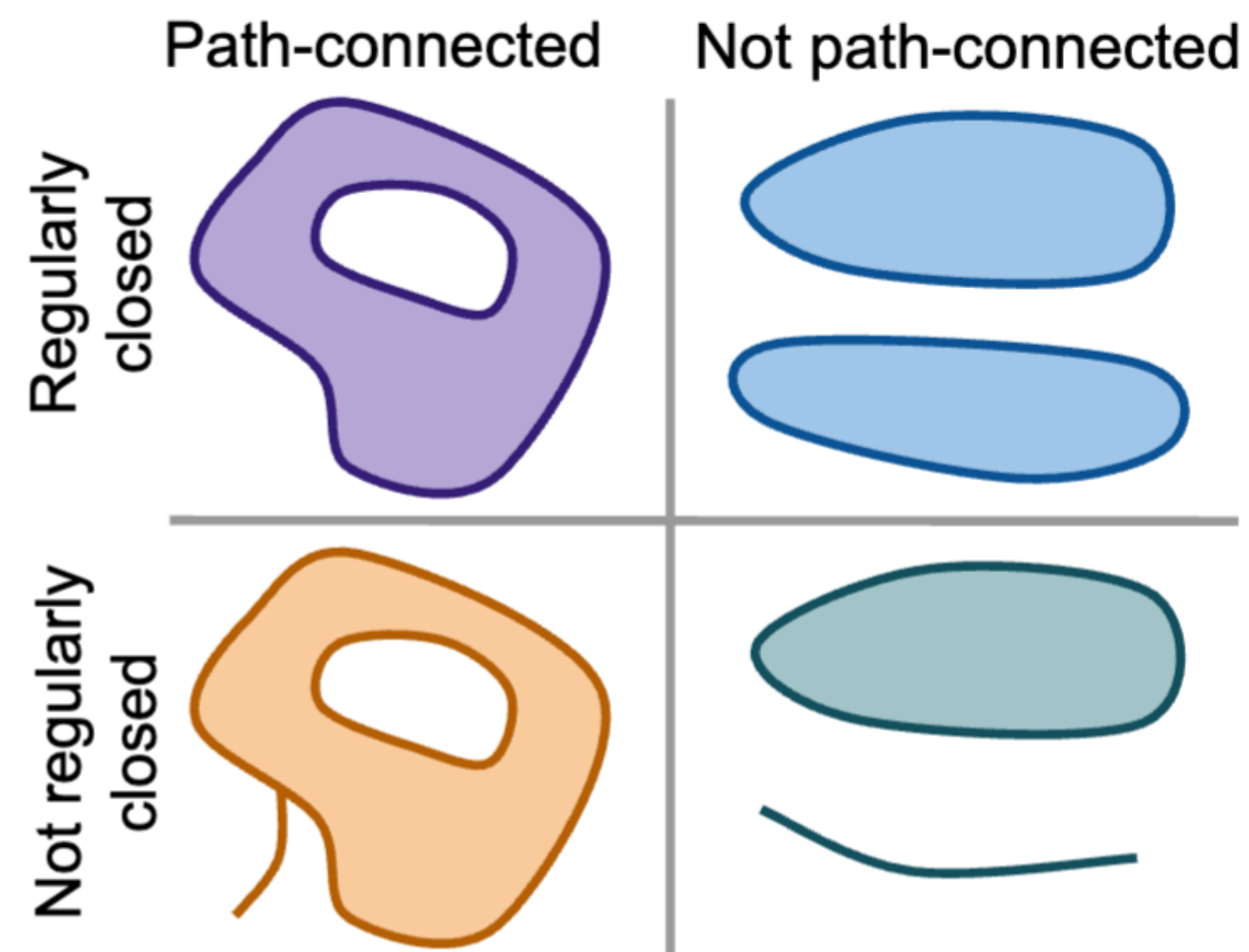
Then optimal reconstruction loss ($\mathbb{E}[||x - \hat{g}(\hat{f}(x))||]$) implies block disentanglement

We make no distributional assumptions on latent factors!

General Support for Latents

The assumptions made on the support of the distribution of latent factors

- Regularly Closed (For both local and global disentanglement)
 - Need this to define derivative uniquely over the support of training data
- Path-Connected (Only for global disentanglement)



Extrapolation with Additive Decoders

Cartesian Product Extrapolation

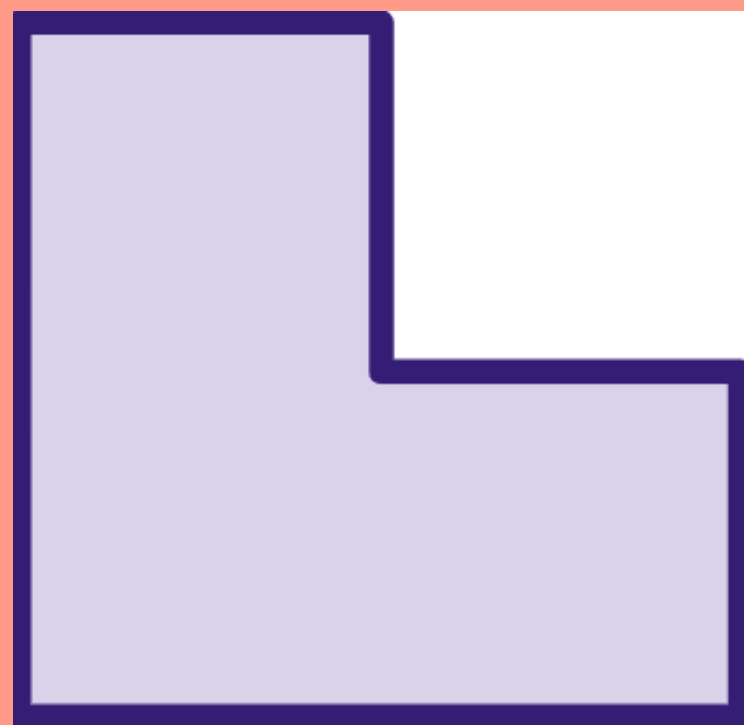
$\hat{\mathcal{Z}}^{train}$ = Support of learned latent factors observed during training

Cartesian Product Extrapolation:

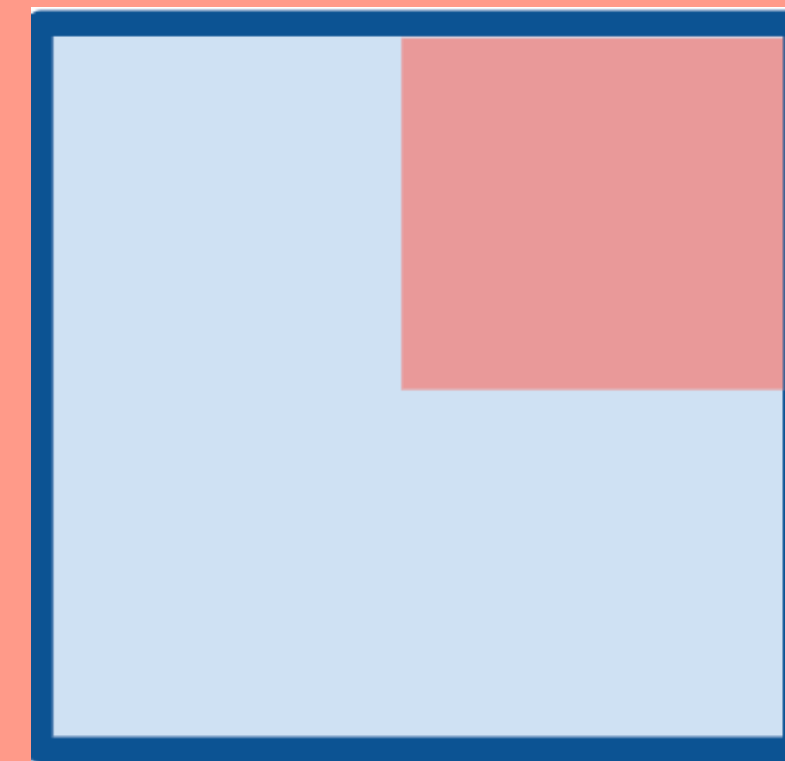
$$CPE_{\mathcal{B}}(\hat{\mathcal{Z}}^{train}) = \prod_{B \in \mathcal{B}} \hat{\mathcal{Z}}_B^{train} \text{ where } \hat{\mathcal{Z}}_B^{train} = \{\hat{z}_B \mid \hat{z} \in \hat{\mathcal{Z}}^{train}\}$$

Example:

$\hat{\mathcal{Z}}^{train}$



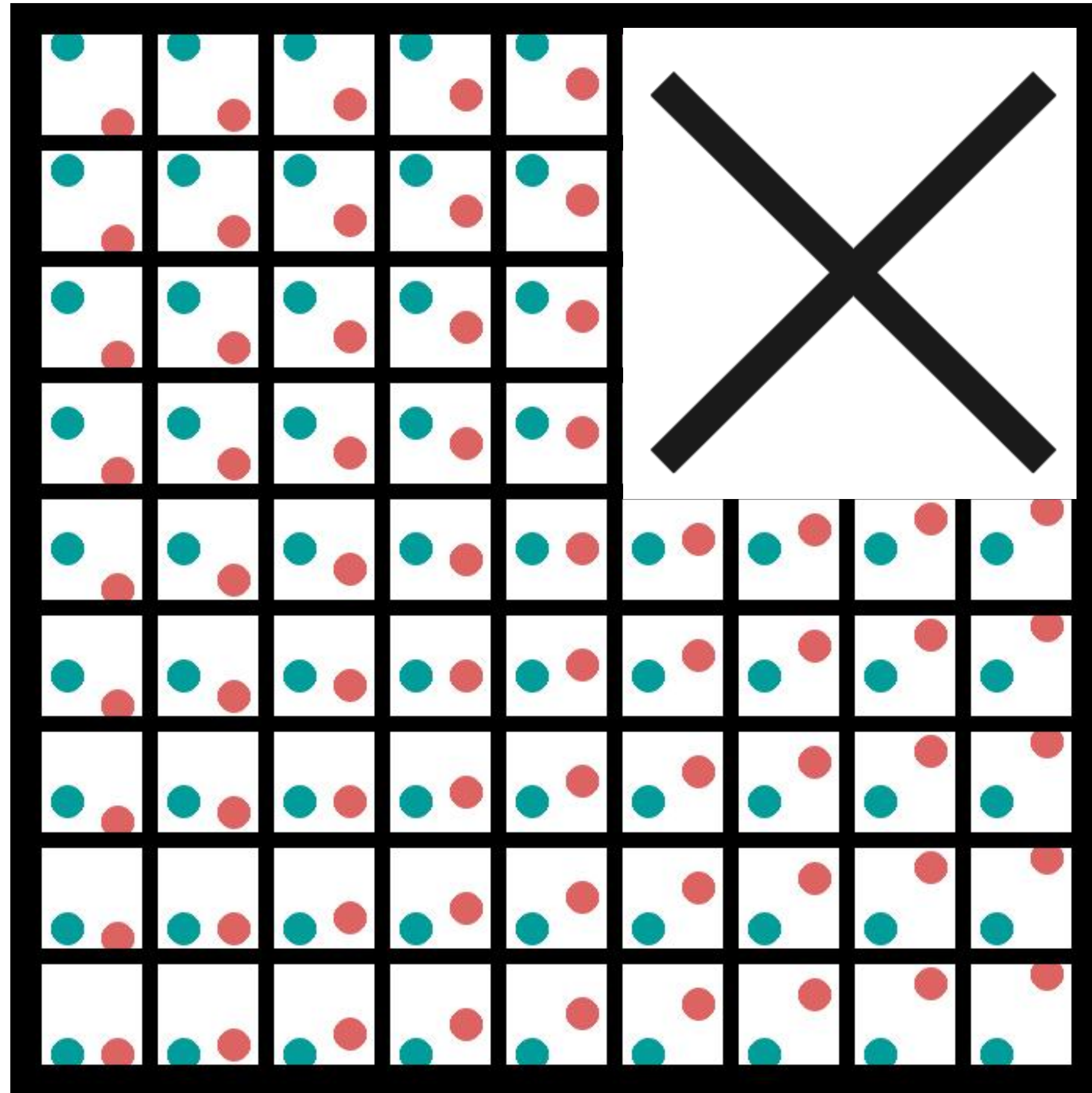
$CPE_{\mathcal{B}}(\hat{\mathcal{Z}}^{train})$



Corollary (Informal): Under same assumptions as previous theorem, the learned decoders imitate ground-truth decoders not only over $\hat{\mathcal{Z}}^{train}$ but also over $CPE_{\mathcal{B}}(\hat{\mathcal{Z}}^{train})$

Experiments

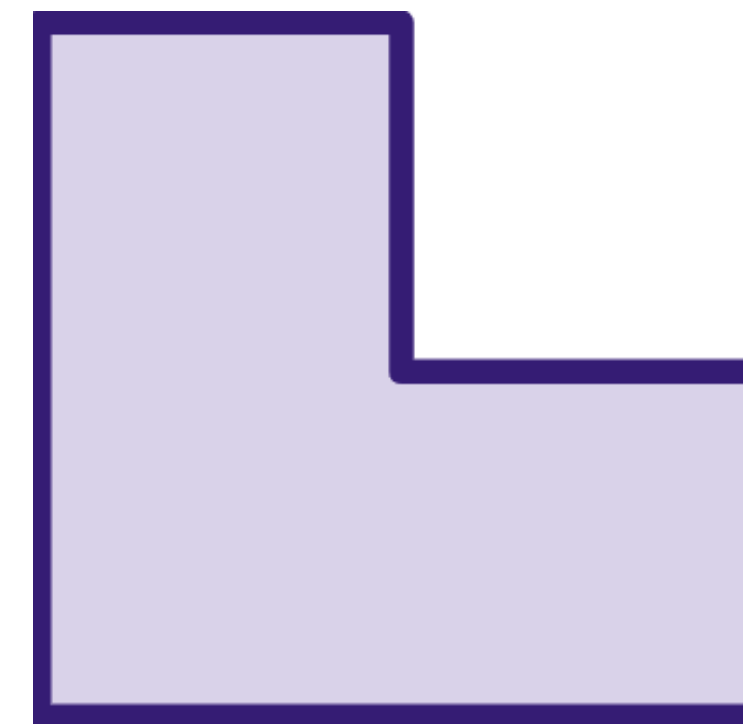
Extrapolation



Scalar Latent Dataset:

- Balls move only along y-axis
- Remove images where both balls have high y-coordinate to get L-shaped training support

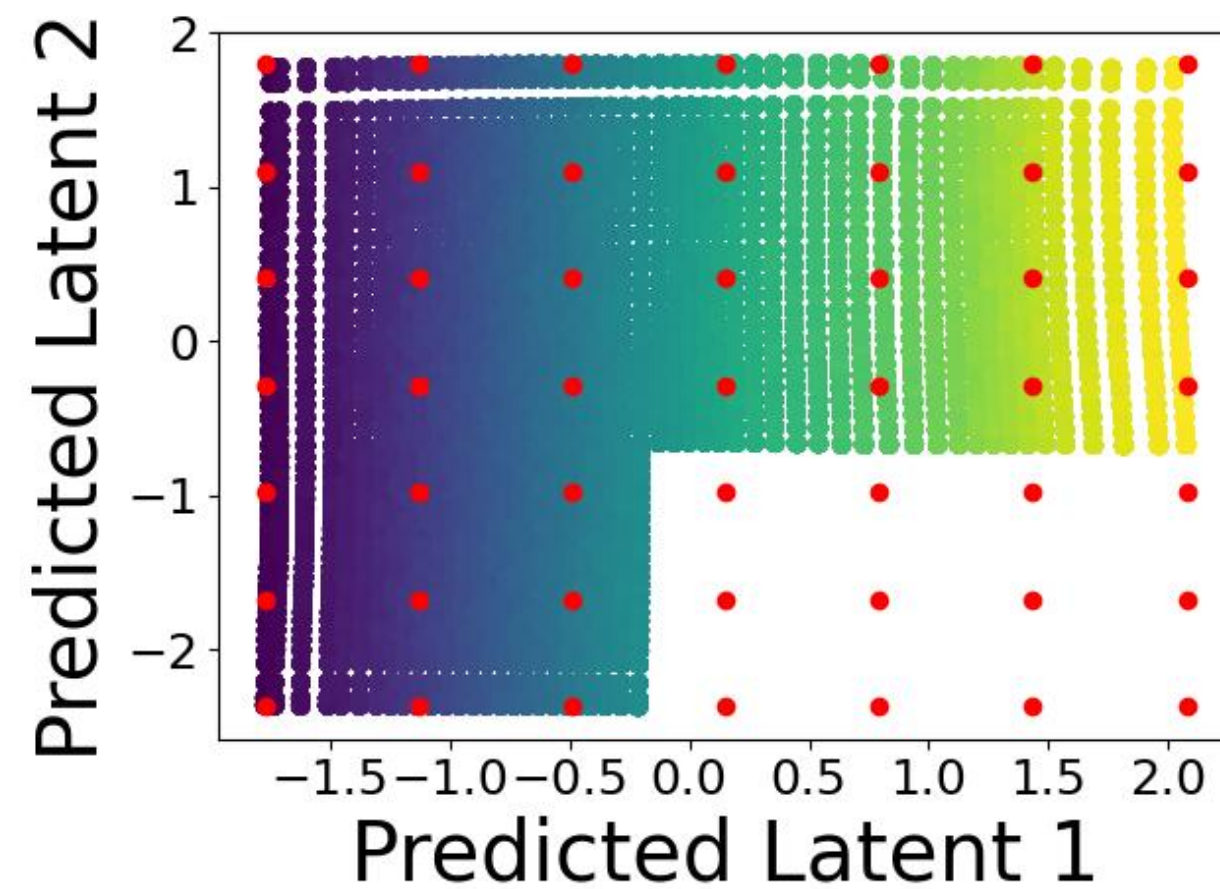
$\hat{\mathbb{Z}}^{train}$



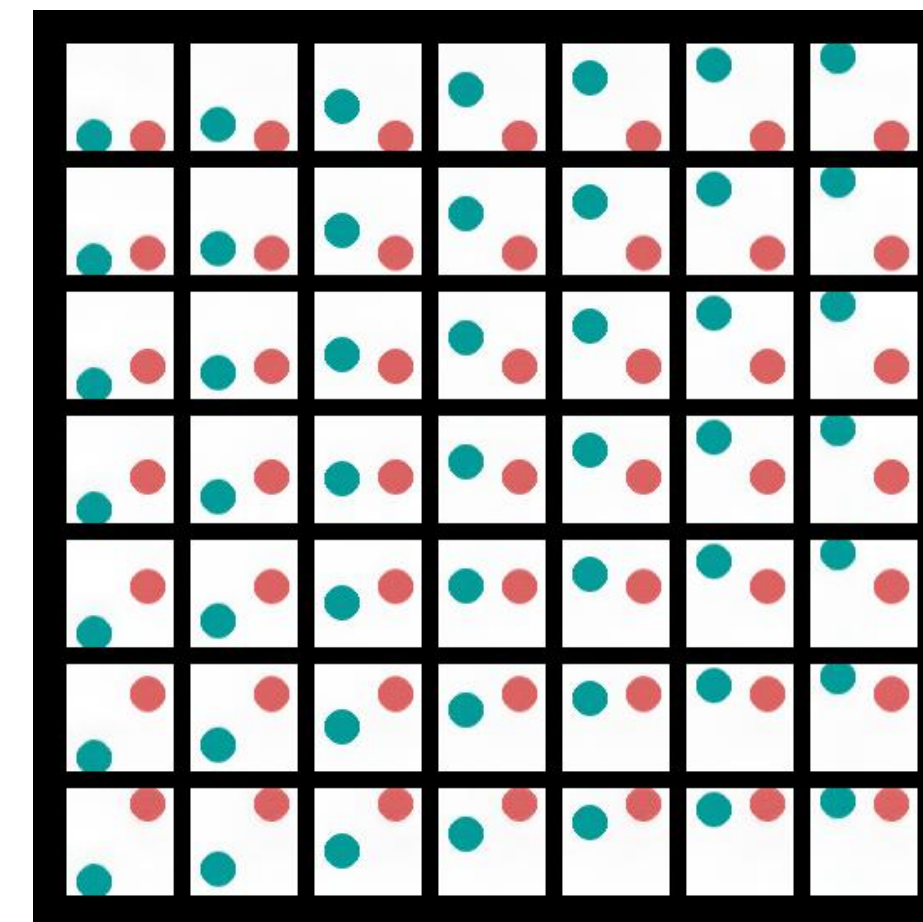
Extrapolation

Additive
Decoder

Learned Latent Space



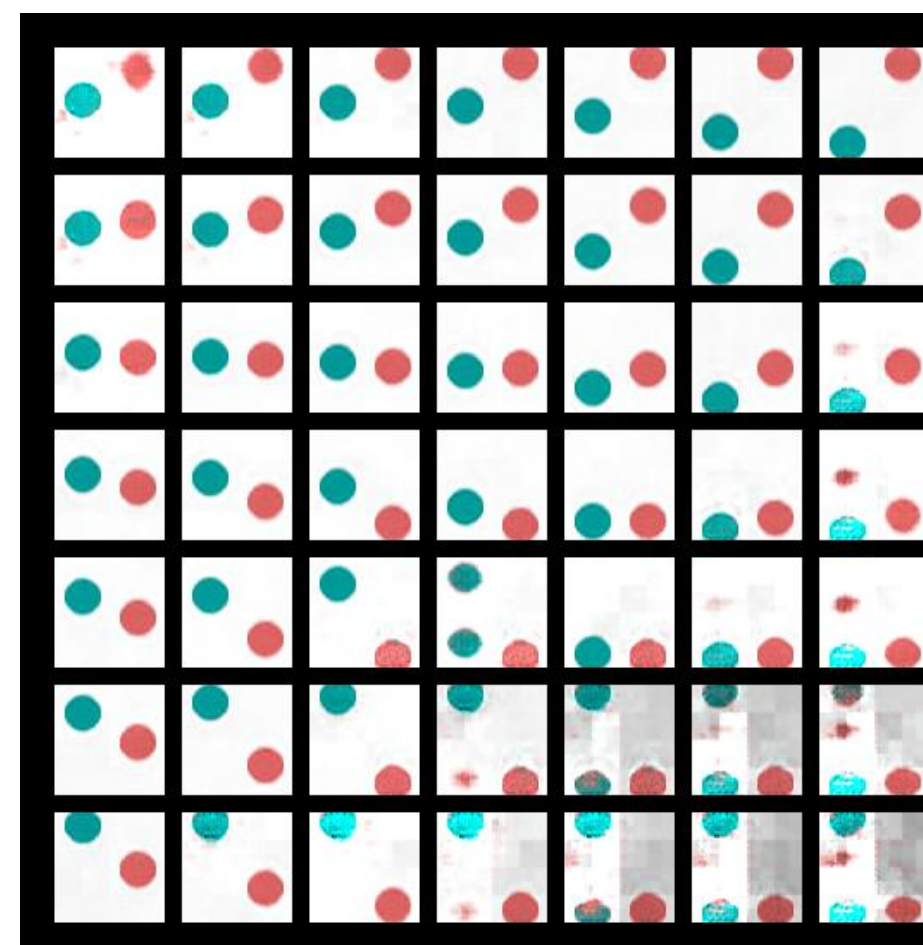
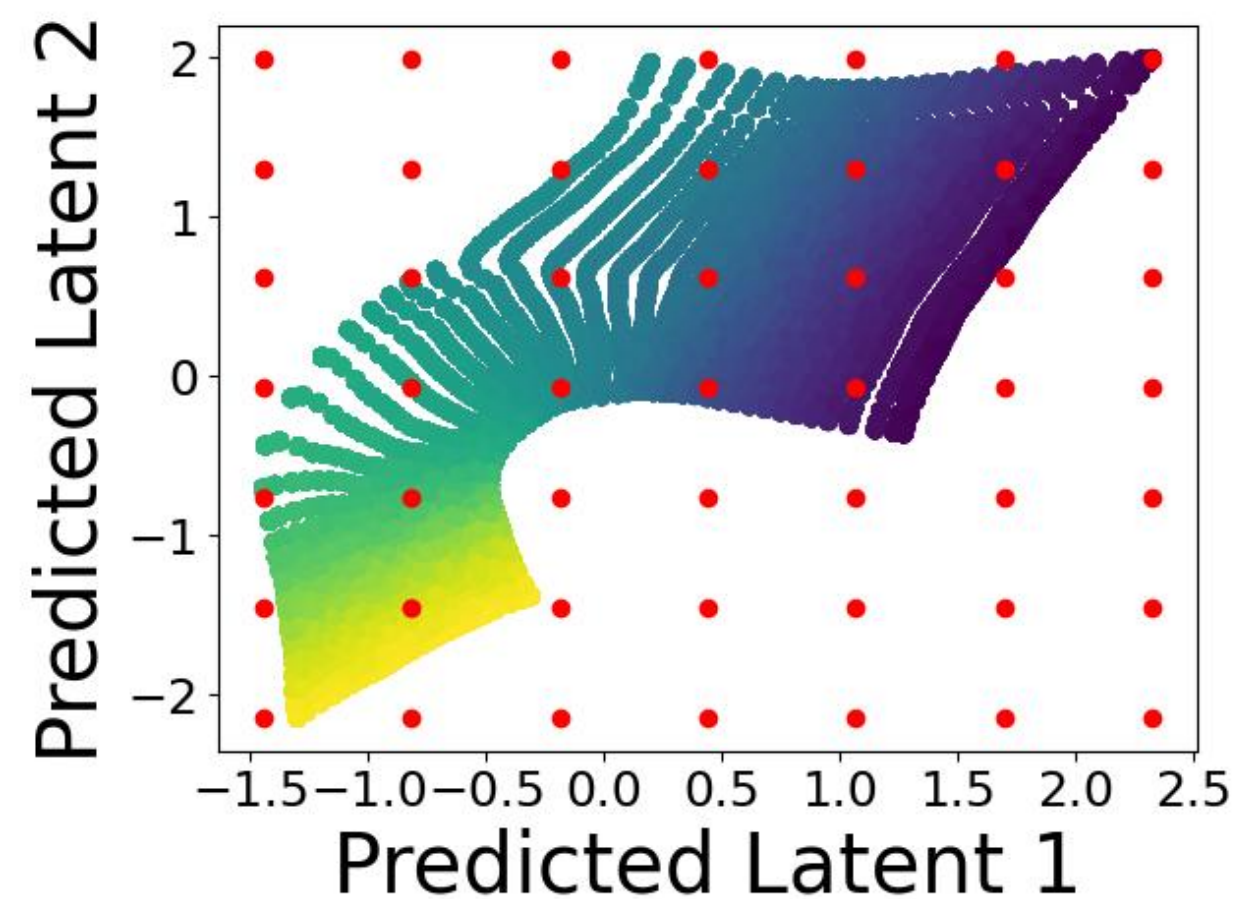
Generated Images



Disentangled

← Changing Latent 1 only
changes the blue ball

Non-Additive
Decoder



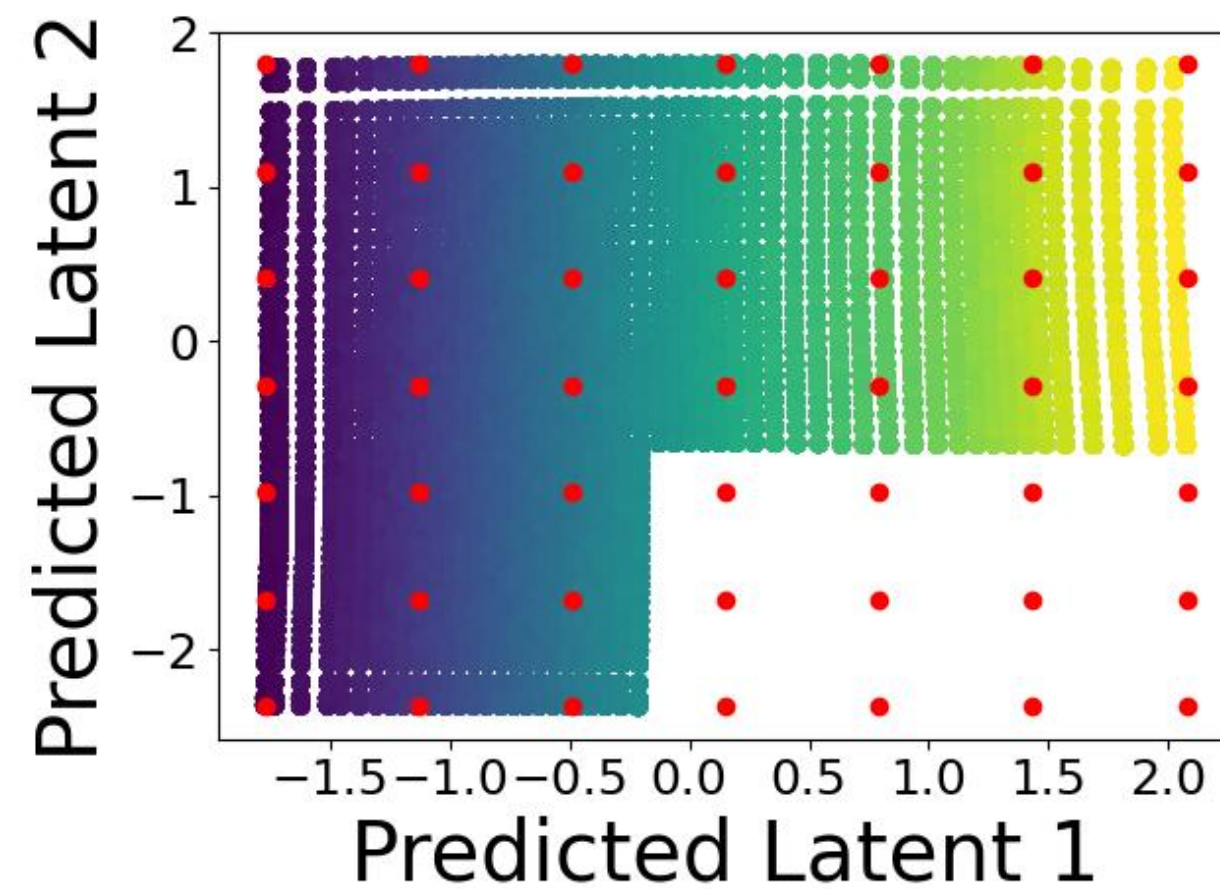
Entangled

← Changing Latent 1
changes both the balls

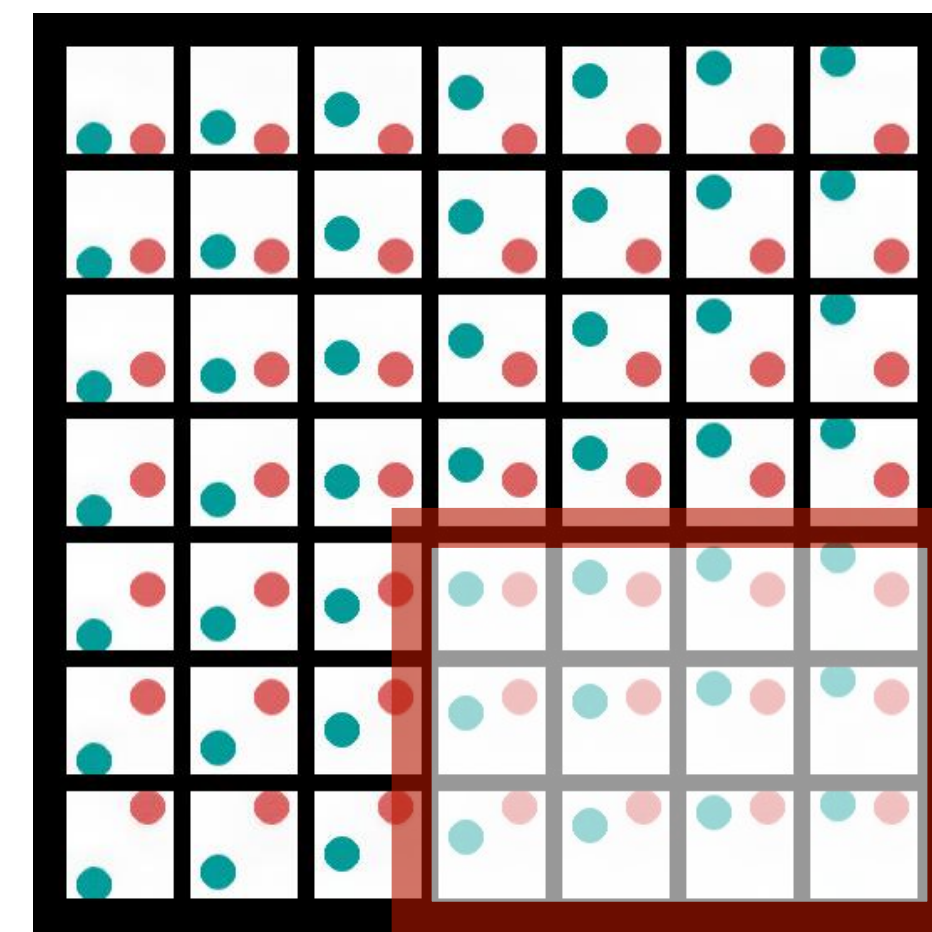
Extrapolation

Additive Decoder

Learned Latent Space



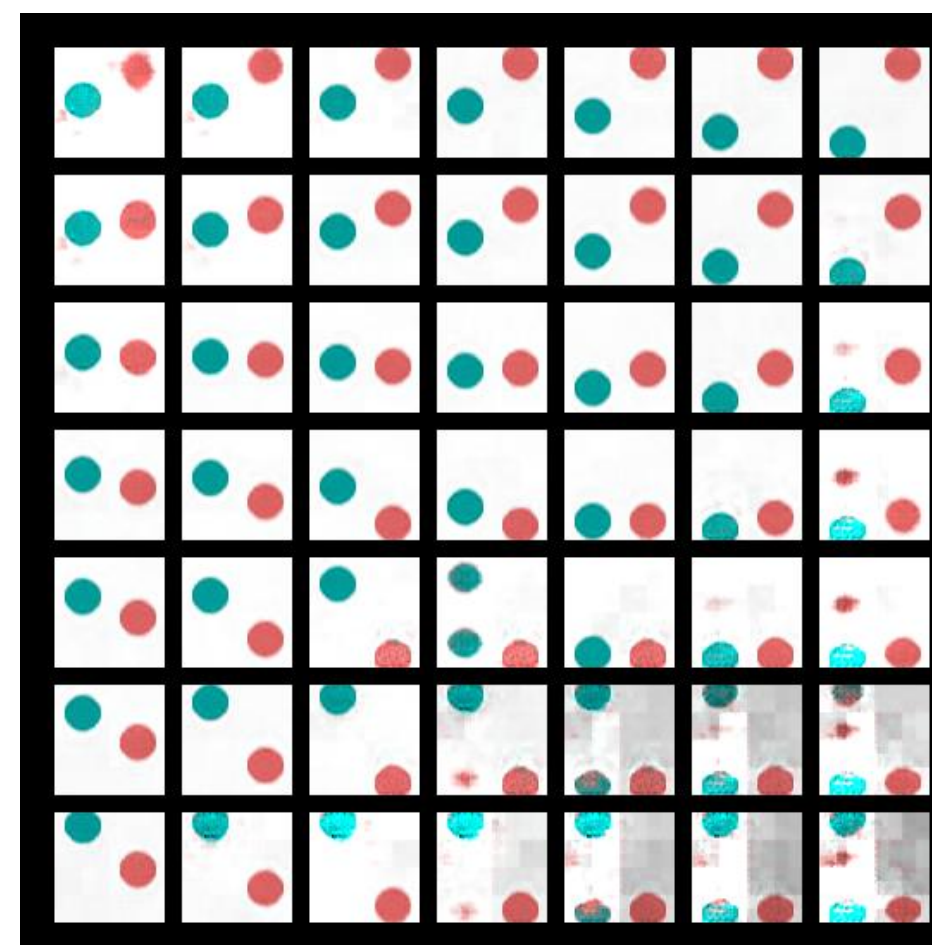
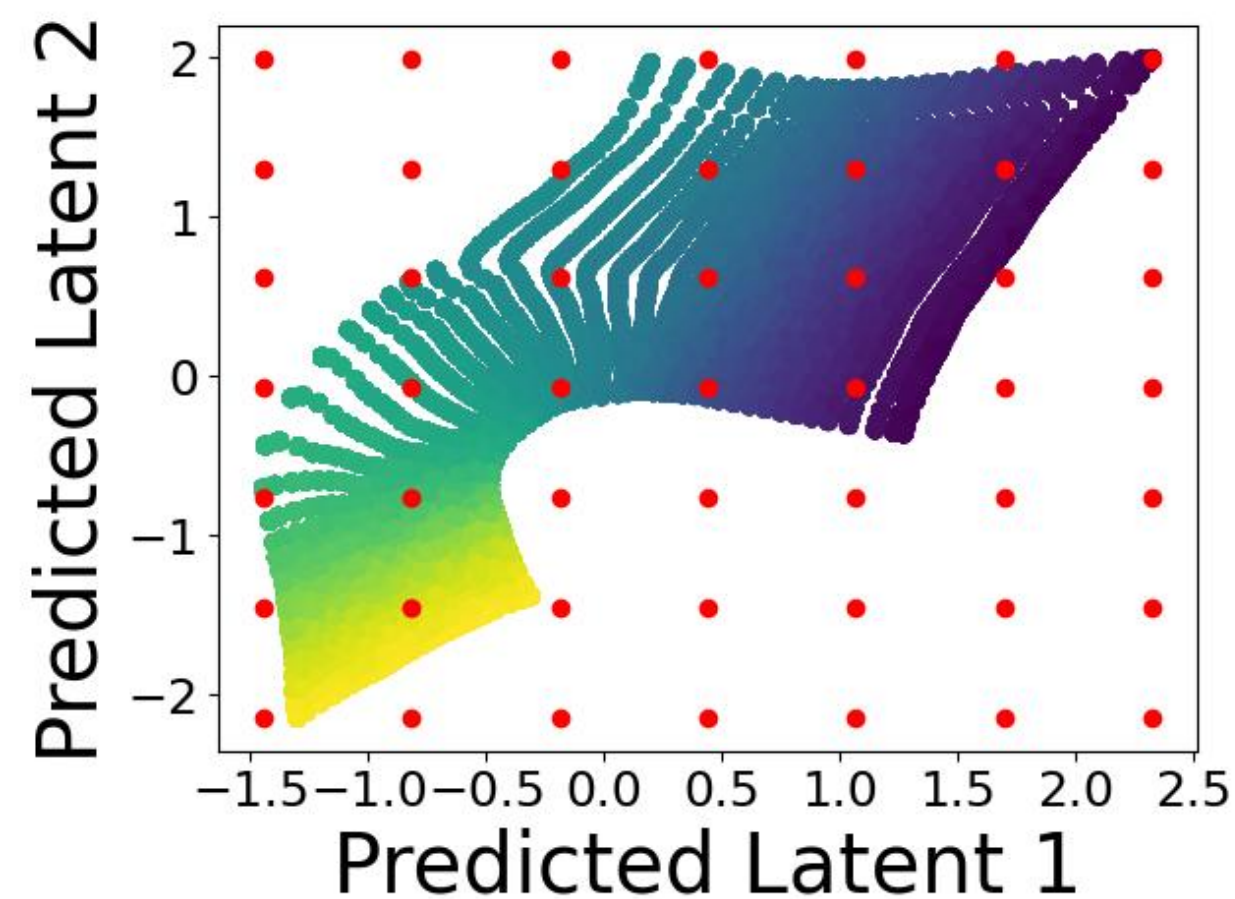
Generated Images



Disentangled

← These samples were never seen during training

Non-Additive Decoder



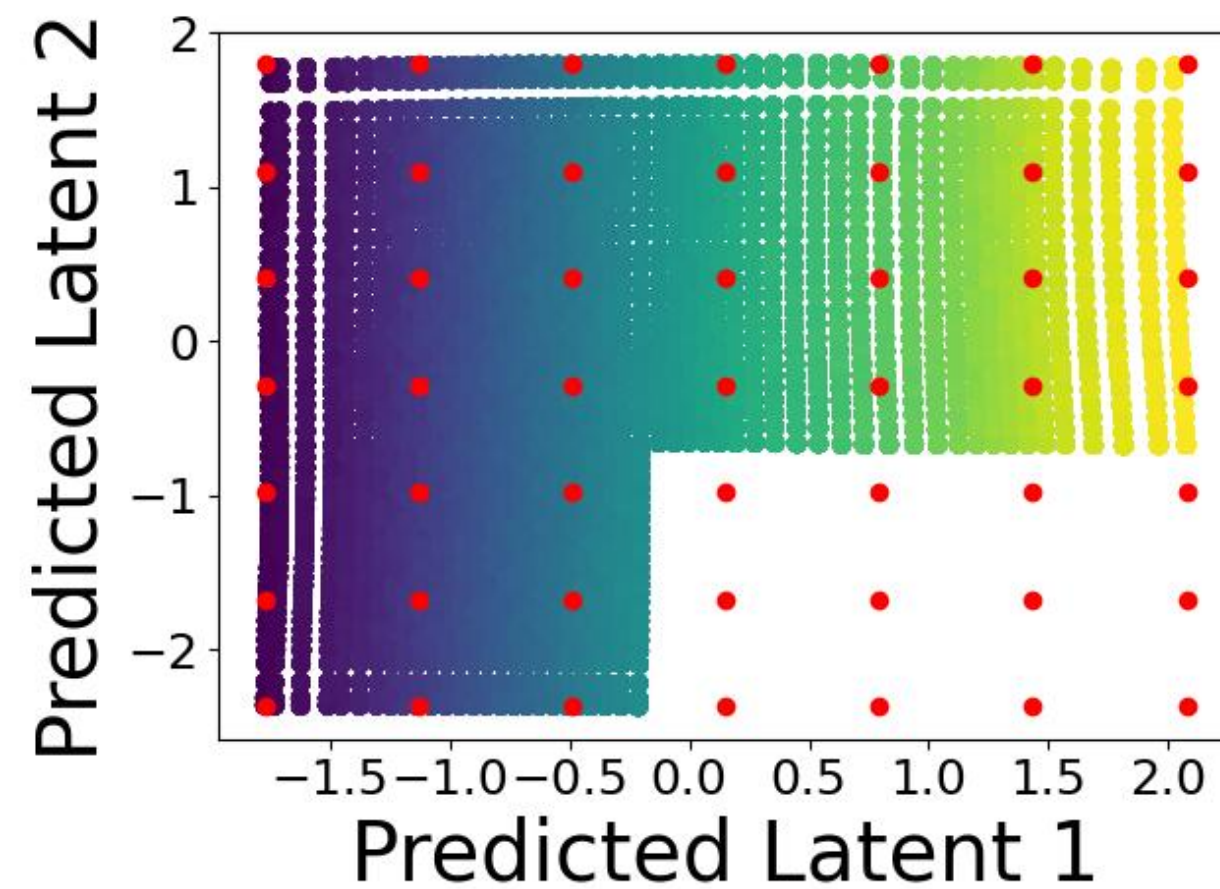
Entangled

Cannot generate unseen samples

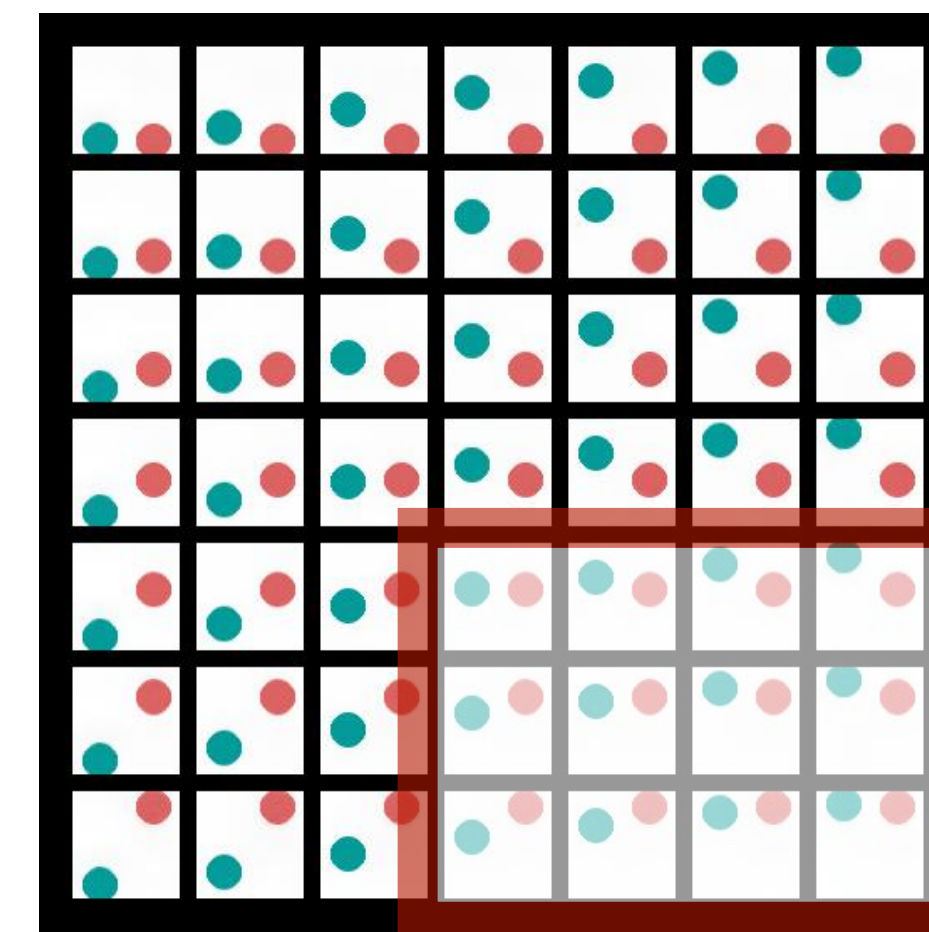
Extrapolation

Additive Decoder

Learned Latent Space



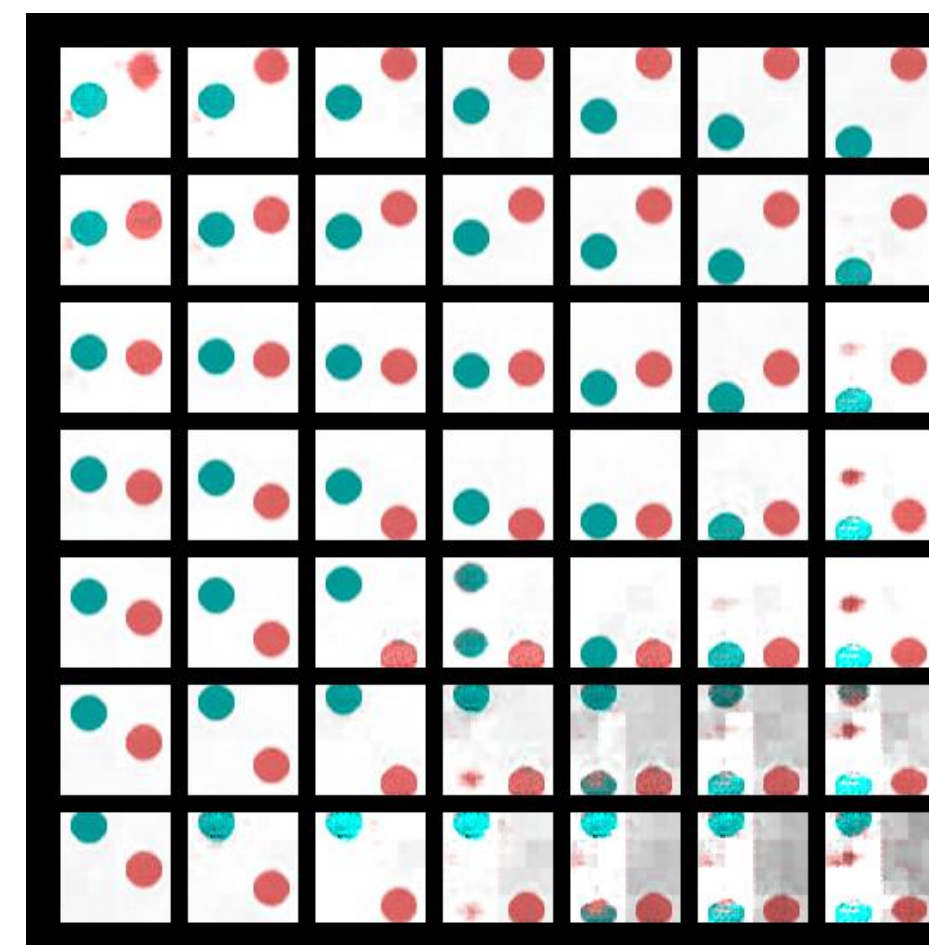
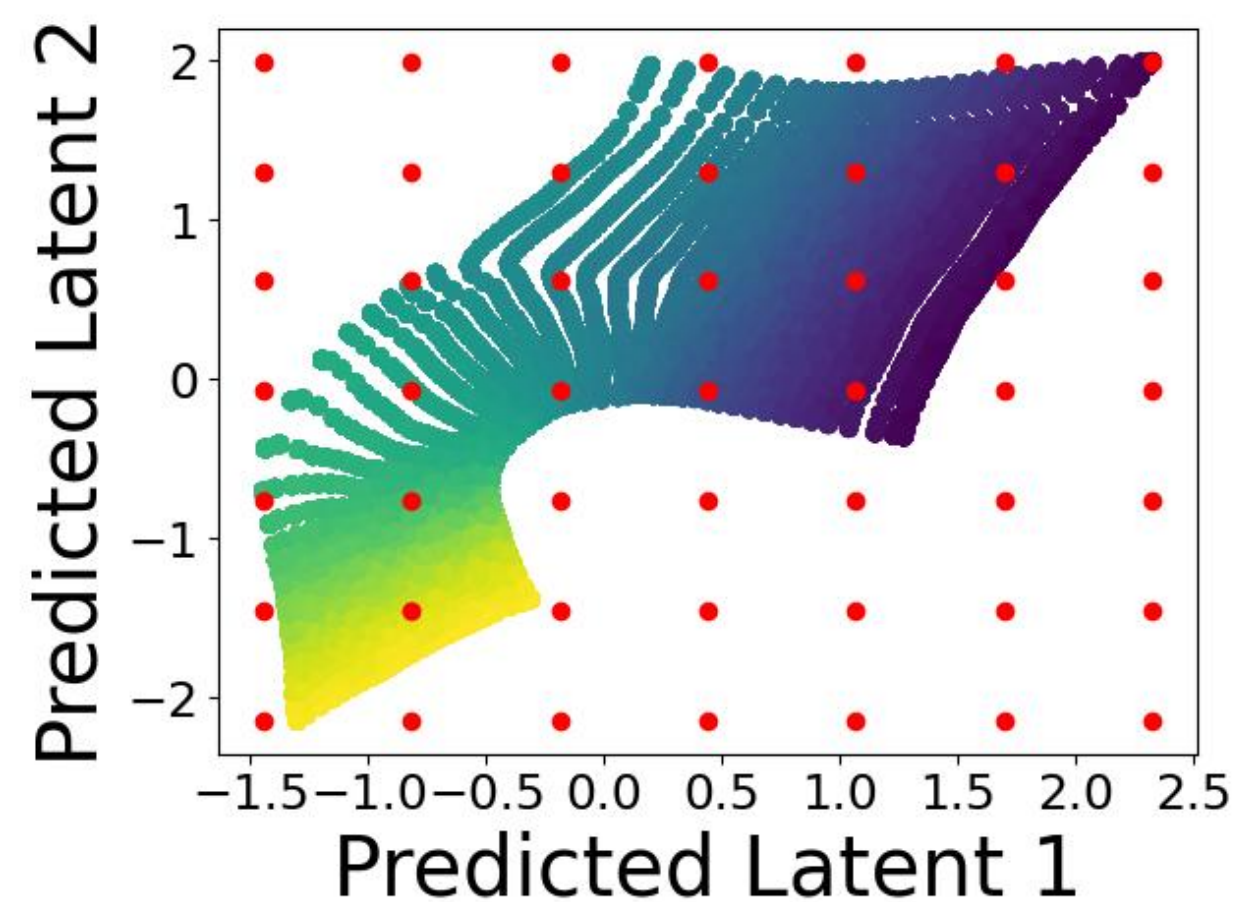
Generated Images



Disentangled

← These samples were never seen during training

Non-Additive Decoder

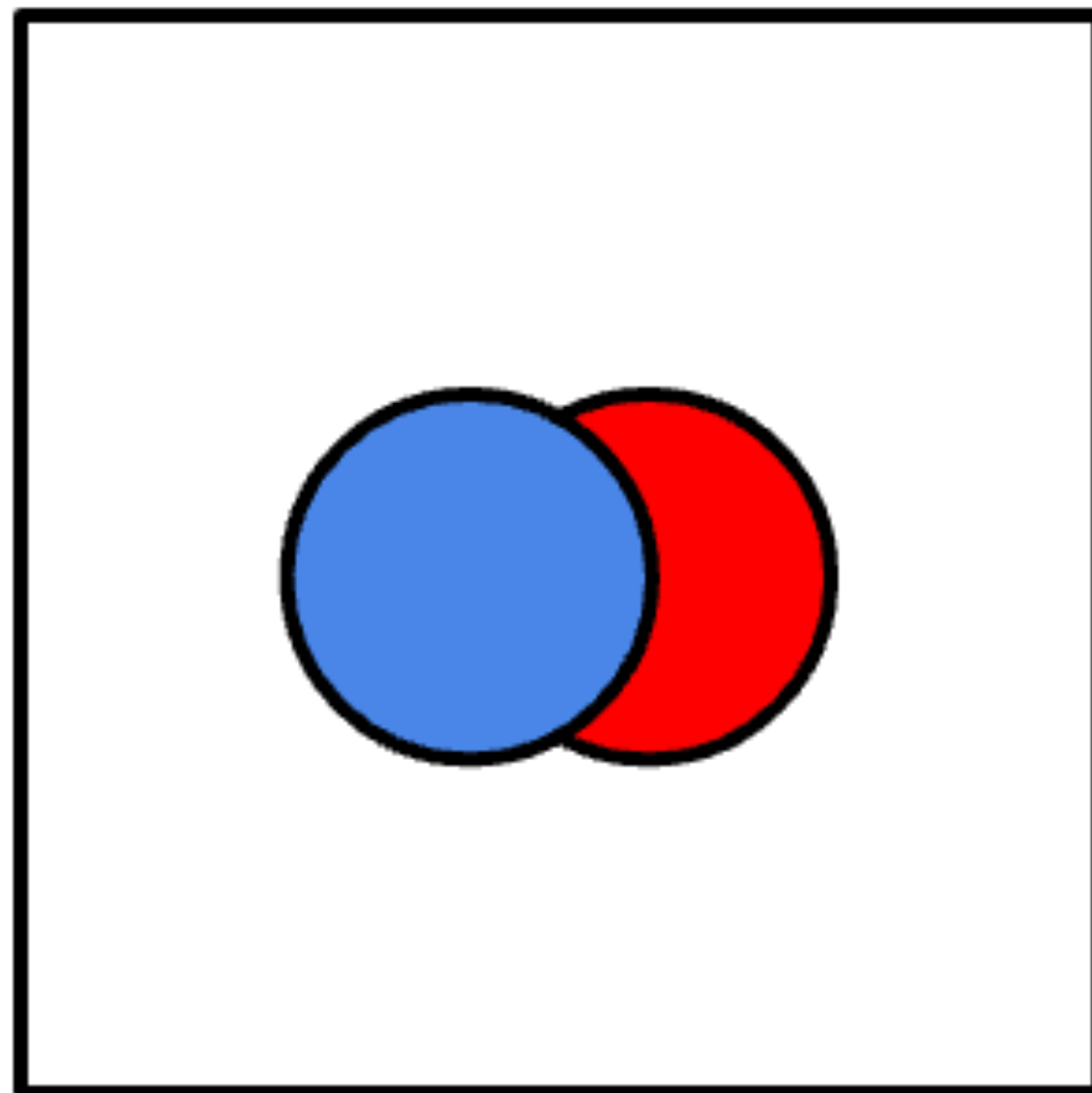


Entangled

Cannot generate unseen samples

Limitations of Additive Decoders

No interaction between latent factors



$$x = g(z_1) + g(z_2) \quad \times$$

$$x = m^1(z) \odot g(z_1) + m^2(z) \odot g(z_2) \quad \checkmark$$

Additive Decoders cannot model images with occlusions!

**Can we consider more expressive function classes
for provable extrapolation?**

Compositional Generalization with Additive Energy Models

Ongoing work with *Kartik Ahuja, Ioannis Mitliagkas, Mohammad Pezeshki, Pascal Vincent*



Additive Energy Models

$$p(x|z) = \frac{1}{B(z)} \exp\left(-\langle 1, \mathbf{E}(x, z) \rangle\right) \text{ where } \langle 1, \mathbf{E}(x, z) \rangle = \sum_{i=1}^{d_z} E_i(x, z_i)$$

Conditional distribution
of data given factors

Partition Function

Energy Function

Energy Function
for each component

- **Assumption:** The energy function can be decomposed as addition of energies with different components of z
- More expressive than additive decoders; can model interaction between components of z via the partition function $B(z) = \int \exp\left(-\langle 1, \mathbf{E}(x, z) \rangle\right) dx$

Contribution

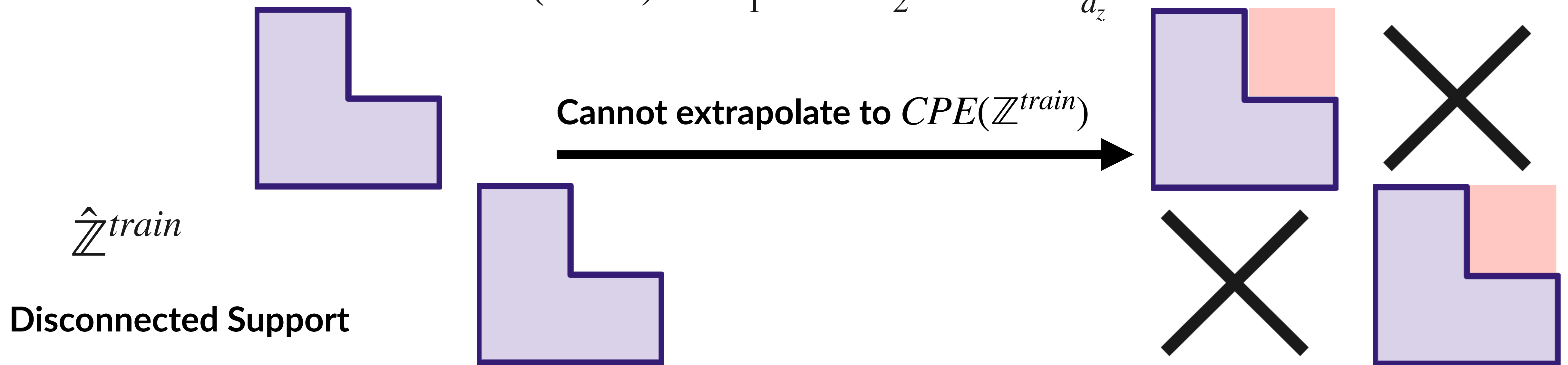
We prove extrapolation guarantees for **discrete factors** with **additive energy models**

Note: We assume the factors of variations z are observed to focus on extrapolation aspect of additive energy models.

Challenges with Disconnected Support

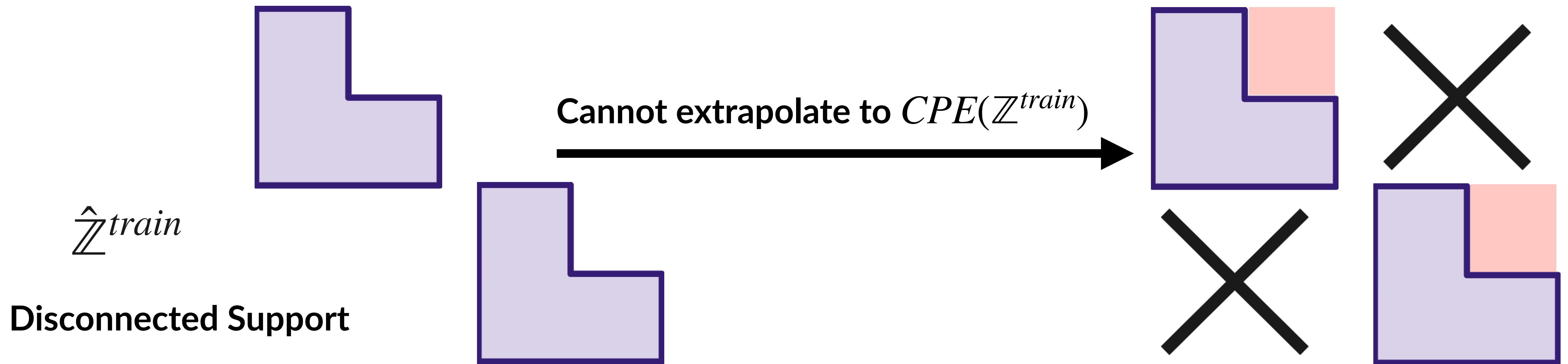
Lets revisit Cartesian-Product Extrapolation with additive functions $g(z) = \sum_{i=1}^{d_z} g_i(z_i)$ where each function $g_i : \mathbb{R} \rightarrow \mathbb{R}^{d_x}$ takes component z_i as input.

$$CPE(\mathbb{Z}^{train}) = \mathbb{Z}_1^{train} \times \mathbb{Z}_2^{train} \times \dots \times \mathbb{Z}_{d_z}^{train}$$

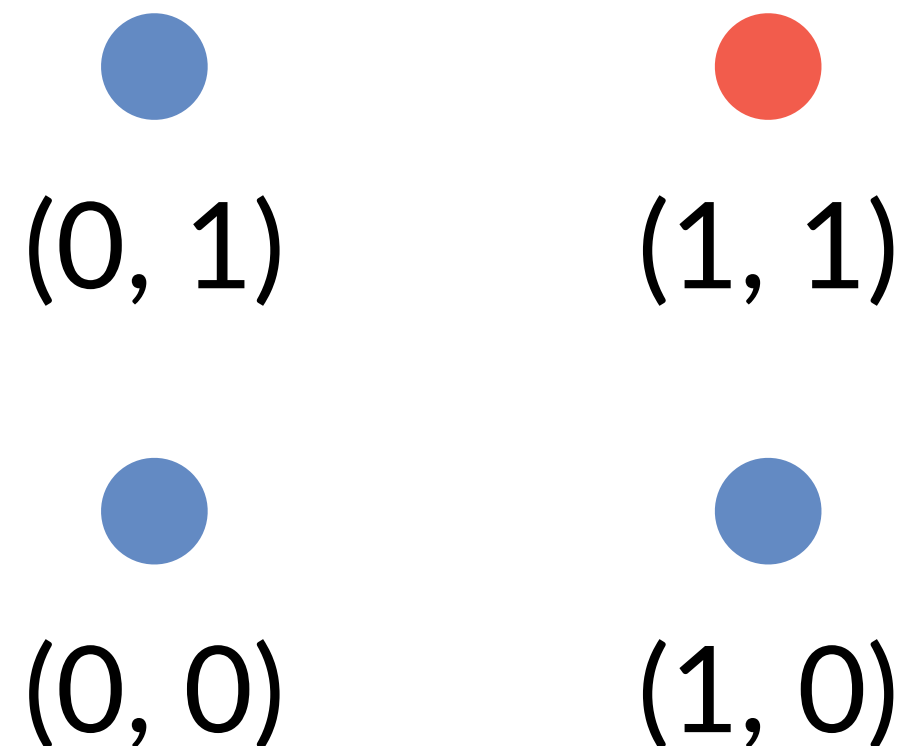


Challenges with Disconnected Support

- Disconnected support makes it hard to extrapolate to $CPE(\mathbb{Z}^{train})$
- This is a fundamental challenge when the factors z are discrete!



Affine Hull Extension



$$g(1,1) = g(1,0) + g(0,1) - g(0,0)$$

Extrapolation to novel discrete factors possible in this case!

● Training Factors ● Novel Factor

Affine Hull Extension of Training Support: $Aff(\mathbb{Z}^{train})$

Let each discrete component z_i takes one of m possible values.

Denote by $\tau(z_i) \in \mathbb{R}^m$ as the one-hot transformation of z_i ; $\tau(z) = [\tau(z_1), \dots, \tau(z_{d_z})] \in \mathbb{R}^{d_z \times m}$

Then $\forall z \in Aff(\mathbb{Z}^{train})$ we have $\tau(z) = \sum_{z \in \mathbb{Z}^{train}} \alpha_z \tau(z)$ where $\sum_{z \in \mathbb{Z}^{train}} \alpha_z = 1$

Can Affine Hull equal Cartesian Product?

- As we add more factors to \mathbb{Z}^{train} , then $Aff(\mathbb{Z}^{train})$ would increase as well
- Can we show that after enough samples $Aff(\mathbb{Z}^{train})$ spans the full grid $\times_{i=1}^{d_z} [m]$?

Theorem: Assume $d_z = 2$, i.e, $z = (z_1, z_2)$ where each z_i has m possible values.

If $|\mathbb{Z}^{train}| > 8c * m \log m$, then $Aff(\mathbb{Z}^{train}) = [m] \times [m]$ with probability $\geq 1 - \frac{1}{c}$

Affine Hull Extrapolation with Additive Functions

For all $z' \in \text{Aff}(\mathbb{Z}^{\text{train}})$, we have $g(z') = \sum_{z \in \mathbb{Z}^{\text{train}}} \alpha_z g(z)$ as $\tau(z') = \sum_{z \in \mathbb{Z}^{\text{train}}} \alpha_z \tau(z)$

True Function: $g(z) = \sum_{i=1}^{d_z} g_i(z_i)$ Learned Function: $\hat{g}(z) = \sum_{i=1}^{d_z} \hat{g}_i(z_i)$

Corollary: If $g(z) = \hat{g}(z) \forall z \in \mathbb{Z}^{\text{train}}$ then $g(z) = \hat{g}(z) \forall z \in \text{Aff}(\mathbb{Z}^{\text{train}})$

Affine Hull Extrapolation with Additive Energy Models

True Model:
$$p(x | z) = \frac{1}{B(z)} \exp\left(- \langle 1, \mathbf{E}(x, z) \rangle \right)$$

Learned Model:
$$p(x | z) = \frac{1}{\hat{B}(z)} \exp\left(- \langle 1, \hat{\mathbf{E}}(x, z) \rangle \right)$$

Theorem: If $p(x | z) = \hat{p}(x | z) \forall z \in \mathbb{Z}^{train}$ then $p(x | z) = \hat{p}(x | z) \forall z \in \text{Aff}(\mathbb{Z}^{train})$
under the assumption of invariant support of $p(x | z)$

Affine Hull Extrapolation for Discriminative Case

True Model:

$$p(z|x) = \text{Softmax}(\log p(x|z) + \log p(z)) \quad \text{where} \quad p(x|z) = \frac{1}{B(z)} \exp(-\langle 1, \mathbf{E}(x, z) \rangle)$$

Learned Model:

$$\hat{p}(z|x) = \text{Softmax}(\log \hat{p}(x|z) + \log p(z)) \quad \text{where} \quad \hat{p}(x|z) = \frac{1}{\hat{B}(z)} \exp(-\langle 1, \hat{\mathbf{E}}(x, z) \rangle)$$

Corollary: If $p(z|x) = \hat{p}(z|x) \forall z \in \mathbb{Z}^{train}$ then $p(z|x) = \hat{p}(z|x) \forall z \in \text{Aff}(\mathbb{Z}^{train})$
under the assumption of invariant support of $p(x|z)$

Inferring partition function $\hat{B}(z) = \int \exp(-\langle 1, \hat{\mathbf{E}}(x, z) \rangle) dx$ is challenging!

Affine Hull Extrapolation for Discriminative Case

Learned Model: $\hat{p}(z | x) = \text{Softmax} \left(- \langle 1, \hat{\mathbf{E}}(x, z) \rangle - \log \hat{M}(z) + \log p(z) \right)$

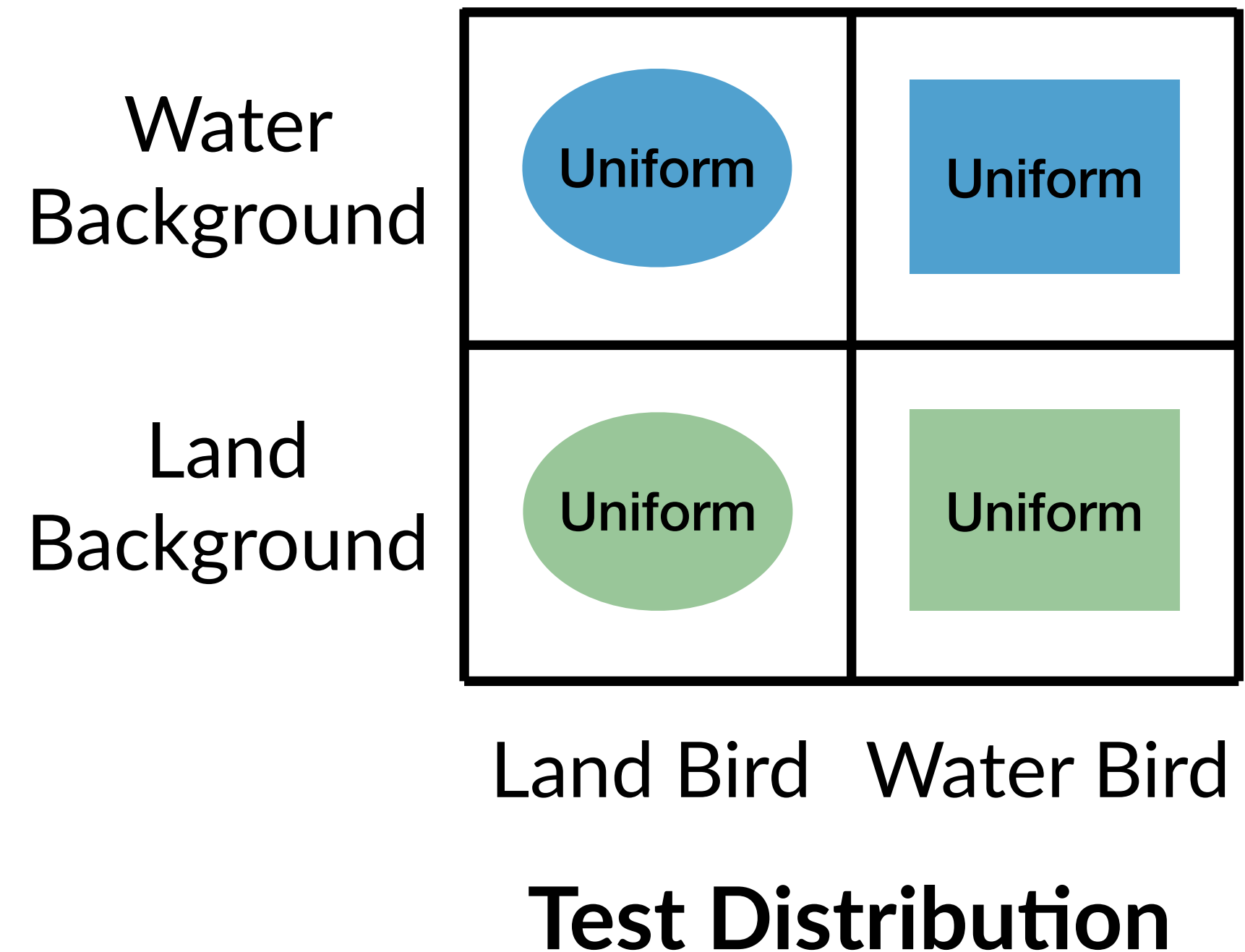
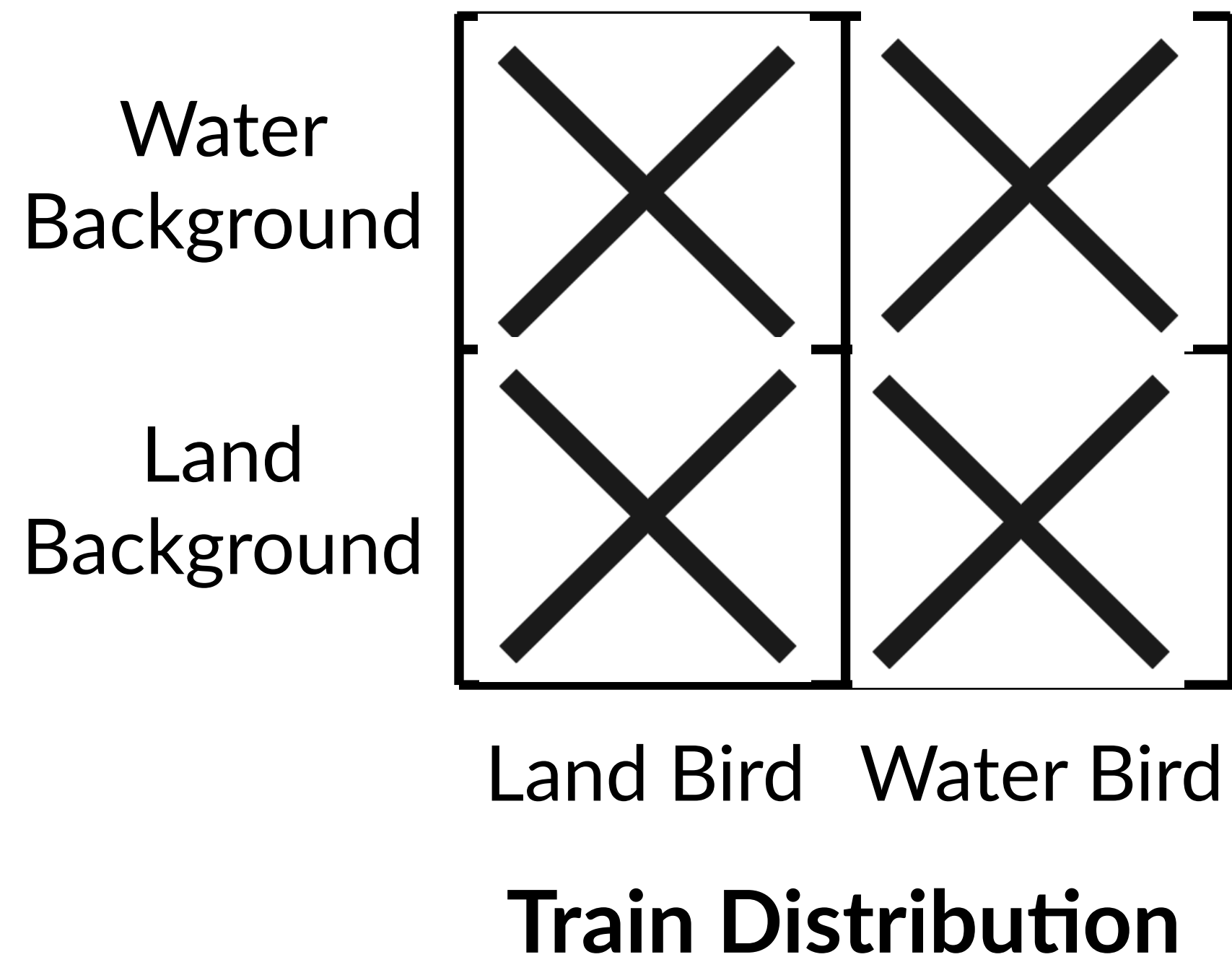
where $\hat{M}(z)$ is not constrained to be the partition function.

Theorem: If $p(z | x) = \hat{p}(z | x) \forall z \in \mathbb{Z}^{train}$ then $p(z | x) = \tilde{p}(z | x) \forall z \in \text{Aff}(\mathbb{Z}^{train})$ under the assumption of invariant support of $p(x | z)$, where $\tilde{p}(z | x)$ is defined as

$$\tilde{p}(z | x) = \text{Softmax} \left(- \langle 1, \hat{\mathbf{E}}(x, z) \rangle - \log \hat{Q}(z) + \log p(z) \right)$$

$$\hat{Q}(z) = \mathbb{E}_{x \sim p^{train}(x)} \left[\frac{\exp \left(- \langle 1, \hat{\mathbf{E}}(x, z) \rangle \right)}{\sum_{\tilde{z} \in \mathbb{Z}^{train}} \exp \left(- \langle 1, \hat{\mathbf{E}}(x, \tilde{z}) \rangle - \log \hat{M}(\tilde{z}) + \log p(\tilde{z}) \right)} \right]$$

Experiments: Compositional Distribution Shift



- Factors $z = (y, a)$ where y denotes the class label and a denotes the spurious attribute
- Compositional Shift: $\mathbb{Z}^{train} \neq \mathbb{Z}^{test}$ but $\mathbb{Z}^{test} = Aff(\mathbb{Z}^{train})$

Implementation of Proposed Approach

$$\hat{p}(z | x) = \textit{Softmax}\left(- \langle 1, \hat{\mathbf{E}}(x, z) \rangle - \log \hat{M}(z) + \log p(z) \right)$$



$$\hat{p}(z | x) = \textit{Softmax}\left(- \langle W_y, \phi(x) \rangle - \langle W_a, \phi(x) \rangle - \log \hat{M}(z) + \log p(z) \right)$$

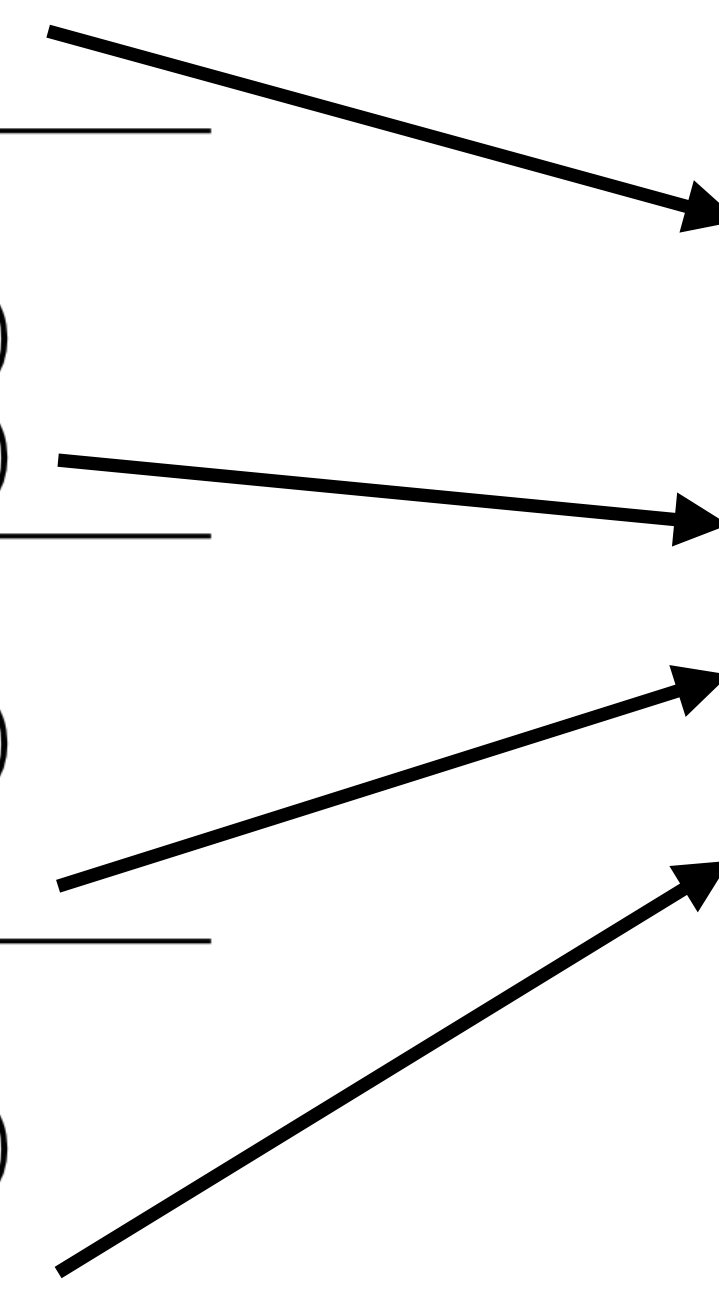
$\phi(x)$: Representations via pretrained ResNet-50 architecture

Learnable Parameters: W_y, W_a, \hat{M}

Learning Objective: $\min_{W_y, W_a, \hat{M}} \mathbb{E}_{p(y, a, x)} - \log \hat{p}(y, a | x)$

Results

Removed (y, a)	Method	Average Acc	Worst Group Acc
(0, 0)	ERM	0.76 (0.0)	0.69 (0.0)
(0, 0)	GroupDRO	0.86 (0.0)	0.78 (0.0)
(0, 0)	AddEnergy	0.88 (0.0)	0.86 (0.0)
(0, 1)	ERM	0.71 (0.0)	0.38 (0.0)
(0, 1)	GroupDRO	0.79 (0.01)	0.41 (0.05)
(0, 1)	AddEnergy	0.87 (0.0)	0.79 (0.01)
(1, 0)	ERM	0.81 (0.01)	0.1 (0.02)
(1, 0)	GroupDRO	0.92 (0.0)	0.74 (0.04)
(1, 0)	AddEnergy	0.88 (0.0)	0.85 (0.0)
(1, 1)	ERM	0.89 (0.0)	0.53 (0.0)
(1, 1)	GroupDRO	0.91 (0.0)	0.77 (0.04)
(1, 1)	AddEnergy	0.89 (0.0)	0.86 (0.0)



Better worst group accuracy than baselines

Results for the **Waterbirds** benchmark. The performance for both the metrics is denoted as mean \pm standard error over 3 random seeds on the test dataset

Future Work

Planned experiments for the current method

$$\hat{p}(z|x) = \textit{Softmax}\left(-\langle 1, \hat{\mathbf{E}}(x, z) \rangle - \log \hat{M}(z) + \log p(z)\right)$$

- Experiment on more complex subpopulation shift benchmarks
 - CivilComments, MultiNLI, NICO++, Causal Triplet
- Empirically verify the $O(m \log m)$ bound for the case of 2-dimensional factors
 - Design synthetic image datasets where we can control factors $z = (y, a)$

Future Work

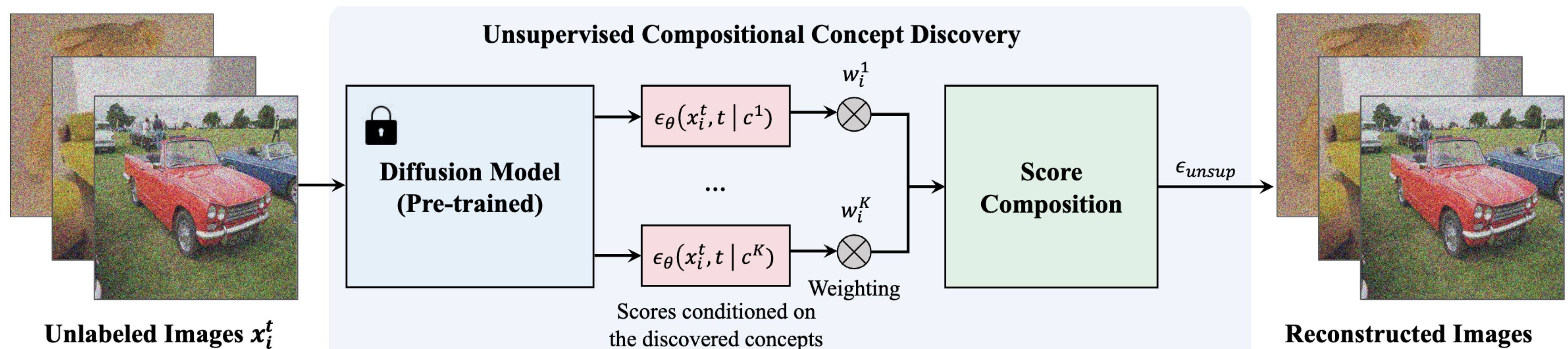
Extrapolation without labelled factors z ?

- A reasonable assumption is that we observe the class label (y) but not the spurious attributes $((a_1, \dots, a_{d_z-1}))$ where $z = (y, a_1, \dots, a_{d_z-1})$
 - Similar setup explored in recent works; XRM (Pezeshki et al.), ULA(Tsirigotis et al.)
 - However, their setup did not consider the extreme case of compositional shift
- Similar to disentanglement with weak supervision but our goal is to disentangle only the quantization of spurious features!

Future Work

Disentanglement with Additive Energy Models?

- Perhaps assumptions similar to additive decoders can help
 - $\nabla_x \log \hat{p}(x | \hat{z}) = \nabla_x \log p(x | z) \implies - \langle 1, \hat{\mathbf{E}}(x, \hat{z}) \rangle = - \langle 1, \mathbf{E}(x, z) \rangle$
- Empirical evidence for disentanglement with methods similar to additive energy models in recent work by Liu et al.



Thank You!

Backup Slides

Linear Identifiability of ERM

$$\begin{aligned}R(f) &= \mathbb{E}[\|Y - f(X)\|^2] \\&= \mathbb{E}[\|\Gamma Z + N - f \circ g(Z)\|^2] \\&= \mathbb{E}[\|\Gamma Z - f \circ g(Z)\|^2] + \mathbb{E}[\|N\|^2] - 2 * \mathbb{E}[(\Gamma Z - f \circ g(Z))^\top N] \\&= \mathbb{E}[\|\Gamma Z - f \circ g(Z)\|^2] + \mathbb{E}[\|N\|^2] \quad (\text{since } Z \perp N \text{ and } \mathbb{E}[N] = 0)\end{aligned}$$

Linear Identifiability of ERM

$$\begin{aligned}R(f) &= \mathbb{E}[\|Y - f(X)\|^2] \\&= \mathbb{E}[\|\Gamma Z + N - f \circ g(Z)\|^2] \\&= \mathbb{E}[\|\Gamma Z - f \circ g(Z)\|^2] + \mathbb{E}[\|N\|^2] - 2 * \mathbb{E}[(\Gamma Z - f \circ g(Z))^T N] \\&= \mathbb{E}[\|\Gamma Z - f \circ g(Z)\|^2] + \mathbb{E}[\|N\|^2] \quad (\text{since } Z \perp N \text{ and } \mathbb{E}[N] = 0)\end{aligned}$$

$$\begin{aligned}\text{Minimum risk } (R(f) = \mathbb{E}[\|N\|^2]) &\implies f = \Gamma Z \\ &\implies \hat{\Gamma}\hat{Z} = \Gamma Z\end{aligned}$$

Linear Identifiability of ERM

$$\begin{aligned}R(f) &= \mathbb{E}[\|Y - f(X)\|^2] \\&= \mathbb{E}[\|\Gamma Z + N - f \circ g(Z)\|^2] \\&= \mathbb{E}[\|\Gamma Z - f \circ g(Z)\|^2] + \mathbb{E}[\|N\|^2] - 2 * \mathbb{E}[(\Gamma Z - f \circ g(Z))^T N] \\&= \mathbb{E}[\|\Gamma Z - f \circ g(Z)\|^2] + \mathbb{E}[\|N\|^2] \quad (\text{since } Z \perp N \text{ and } \mathbb{E}[N] = 0)\end{aligned}$$

$$\begin{aligned}\text{Minimum risk } (R(f) = \mathbb{E}[\|N\|^2]) &\implies f = \Gamma Z \\ &\implies \hat{\Gamma} \hat{Z} = \Gamma Z\end{aligned}$$

Since $\hat{\Gamma} \in R^{k,d}$ is full column rank ($k > d$), we can find a set S of d rows such that $\hat{W} = [\hat{\Gamma}]_{S,:}$ is invertible

Linear Identifiability of ERM

$$\begin{aligned}R(f) &= \mathbb{E}[\|Y - f(X)\|^2] \\&= \mathbb{E}[\|\Gamma Z + N - f \circ g(Z)\|^2] \\&= \mathbb{E}[\|\Gamma Z - f \circ g(Z)\|^2] + \mathbb{E}[\|N\|^2] - 2 * \mathbb{E}[(\Gamma Z - f \circ g(Z))^T N] \\&= \mathbb{E}[\|\Gamma Z - f \circ g(Z)\|^2] + \mathbb{E}[\|N\|^2] \quad (\text{since } Z \perp N \text{ and } \mathbb{E}[N] = 0)\end{aligned}$$

$$\begin{aligned}\text{Minimum risk } (R(f) = \mathbb{E}[\|N\|^2]) &\implies f = \Gamma Z \\&\implies \hat{\Gamma} \hat{Z} = \Gamma Z \\&\implies \hat{W} \hat{Z} = WZ \\&\implies \hat{Z} = (\hat{W})^{-1} WZ \\&\implies \hat{Z} = AZ\end{aligned}$$

Linear Identifiability of IC-ERM

$$\begin{aligned} R(f) &= \mathbb{E}[\|Y - f(X)\|^2] \\ &= \mathbb{E}[\|\Gamma Z + N - f \circ g(Z)\|^2] \\ &= \mathbb{E}[\|\Gamma Z - f \circ g(Z)\|^2] + \mathbb{E}[\|N\|^2] - 2 * \mathbb{E}[(\Gamma Z - f \circ g(Z))^T N] \\ &= \mathbb{E}[\|\Gamma Z - f \circ g(Z)\|^2] + \mathbb{E}[\|N\|^2] \quad (\text{since } Z \perp N \text{ and } \mathbb{E}[N] = 0) \end{aligned}$$

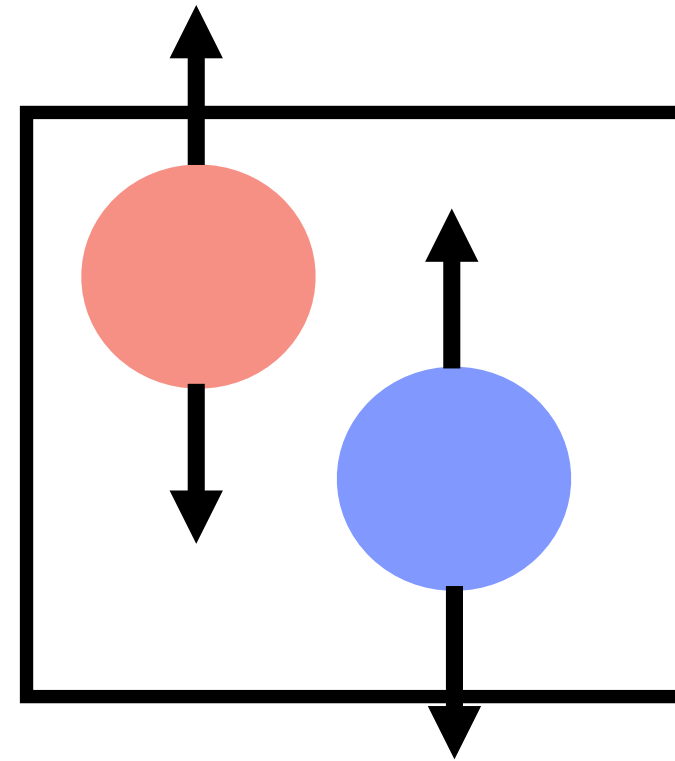
$$\begin{aligned} \text{Minimum risk } (R(f) = \mathbb{E}[\|N\|^2]) &\implies f = \Gamma Z \\ &\implies \hat{Z} = AZ \end{aligned}$$

\hat{Z} has mutually independent components; similar to the **Linear ICA** problem!
Hence, A must be **permutation & scaling** matrix.

Datasets

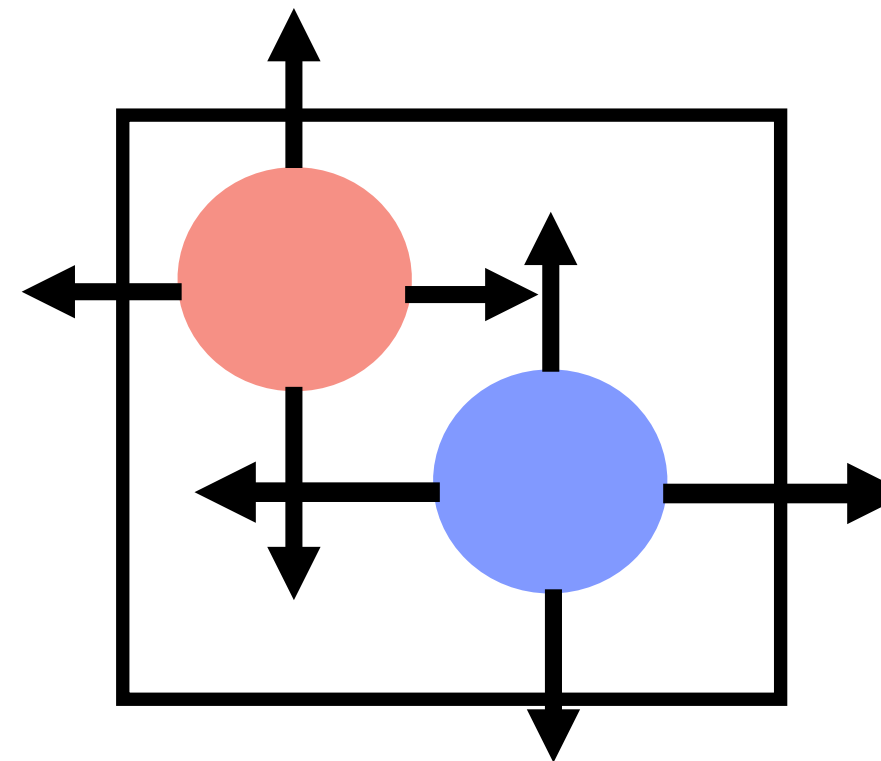
Scalar Latent Dataset: Balls move only along y-axis

$$\dim(Z) = 2$$
$$\mathcal{B} = \{\{1\}, \{2\}\}$$



Block Latent Dataset: Balls move along both x and y-axis

$$\dim(Z) = 4$$
$$\mathcal{B} = \{\{1,2\}, \{3,4\}\}$$



Independent Case: $z_{B_1} \perp z_{B_2}$

Dependent Case: $z_{B_1} \not\perp z_{B_2}$

Disentanglement

Scalar Latent Block Latent Independent Block Latent Dependent

Non-Additive Decoder

70.6 (5.2)

53.9 (7.6)

78.1 (2.9)

Additive Decoder

91.5 (3.6)

92.2 (4.9)

99.9 (0.0)

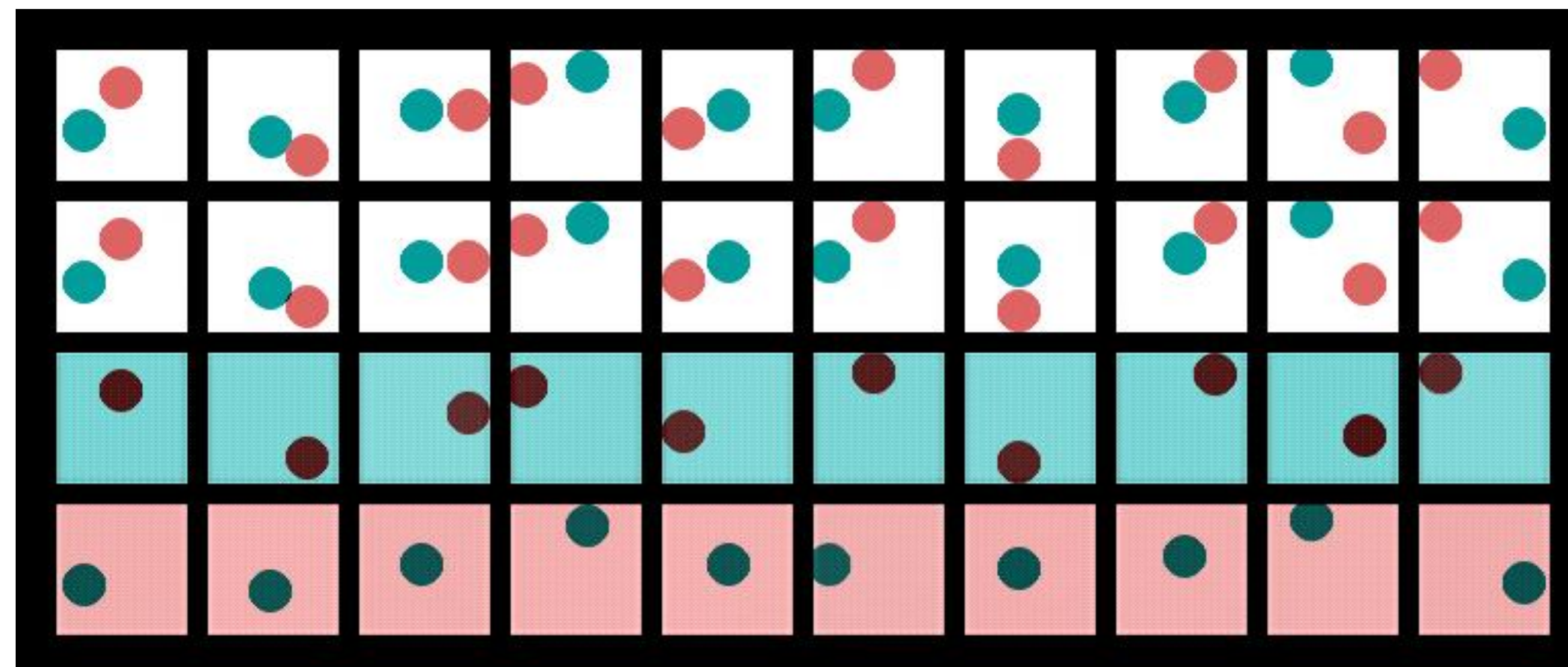
Modified MCC score (Higher implies more disentangled)

Original Images

Reconstructions

Block-specific decoder #1

Block-specific decoder #2



Proof Sketch

Samples from the CPE still have the support for marginal distributions!

$$\text{Supp}(Z_B^{tr}) = \text{Supp}(Z_B^{te})$$

Precisely, for all $B \in \mathcal{B}$ we have

$$\hat{f}^{(B)}(z_B) = f^{(\pi(B))}(v_{\pi(B)}(z_B)) + c^{(B)}$$

Block decoders equality holds true for all blocks on samples from CPE
Additivity enables you to get final image via addition of block decoders!

Global Disentanglement is necessary for compositionality!

We need the learned block decoders to correspond in unique manner to true block decoders!

Proof Sketch

$$\mathbb{E}^{\text{train}} \|\mathbf{x} - \hat{\mathbf{f}}(\hat{\mathbf{g}}(\mathbf{x}))\|^2 = \mathbb{E}^{\text{train}} \|\mathbf{f}(\mathbf{z}) - \hat{\mathbf{f}}(\hat{\mathbf{g}}(\mathbf{f}(\mathbf{z})))\|^2 = 0,$$



$\hat{\mathbf{f}}$: Diffeomorphism
 $\hat{\mathbf{g}}$: continuous

$$\mathbf{f} \circ \mathbf{v}(\mathbf{z}) = \hat{\mathbf{f}}(\mathbf{z}) \quad \forall \mathbf{z} \in \hat{\mathcal{Z}}^{\text{train}} \quad \mathbf{v} := \mathbf{f}^{-1} \circ \hat{\mathbf{f}} \text{ is a } C^2\text{-diffeomorphism}$$

$$\sum_{B \in \mathcal{B}} \mathbf{f}^{(B)}(\mathbf{v}_B(\mathbf{z})) = \sum_{B \in \mathcal{B}} \hat{\mathbf{f}}^{(B)}(\mathbf{z}_B) \quad \forall \mathbf{z} \in \hat{\mathcal{Z}}^{\text{train}}.$$



Derivative w.r.t $z_j \in J$ for some $J \in \mathcal{B}$

$$\sum_{B \in \mathcal{B}} \sum_{i \in B} D_i \mathbf{f}^{(B)}(\mathbf{v}_B(\mathbf{z})) D_j \mathbf{v}_i(\mathbf{z}) = D_j \hat{\mathbf{f}}^{(J)}(\mathbf{z}_J).$$

Proof Sketch

$$f \circ v(z) = \hat{f}(z) \quad \forall z \in \hat{\mathcal{Z}}^{\text{train}} \quad v := f^{-1} \circ \hat{f} \text{ is a } C^2\text{-diffeomorphism}$$

$$\sum_{B \in \mathcal{B}} f^{(B)}(v_B(z)) = \sum_{B \in \mathcal{B}} \hat{f}^{(B)}(z_B) \quad \forall z \in \hat{\mathcal{Z}}^{\text{train}}.$$

Derivative w.r.t $z_j \in J$ for some $J \in \mathcal{B}$

$$\sum_{B \in \mathcal{B}} \sum_{i \in B} D_i f^{(B)}(v_B(z)) D_j v_i(z) = D_j \hat{f}^{(J)}(z_J).$$

Derivative w.r.t $z_{j'} \in J'$ for some $J' \in \mathcal{B} / \{J\}$

$$\sum_{B \in \mathcal{B}} \sum_{i \in B} \left[D_i f^{(B)}(v_B(z)) D_{j,j'}^2 v_i(z) + \sum_{i' \in B} D_{i,i'}^2 f^{(B)}(v_B(z)) D_{j'} v_{i'}(z) D_j v_i(z) \right] = 0$$

Proof Sketch

$$\sum_{B \in \mathcal{B}} \sum_{i \in B} D_i \mathbf{f}^{(B)}(\mathbf{v}_B(\mathbf{z})) D_j \mathbf{v}_i(\mathbf{z}) = D_j \hat{\mathbf{f}}^{(J)}(\mathbf{z}_J).$$

Derivative w.r.t $z_{j'} \in J'$ for some $J' \in \mathcal{B} / \{J\}$

$$\sum_{B \in \mathcal{B}} \sum_{i \in B} \left[D_i \mathbf{f}^{(B)}(\mathbf{v}_B(\mathbf{z})) D_{j,j'}^2 \mathbf{v}_i(\mathbf{z}) + \sum_{i' \in B} D_{i,i'}^2 \mathbf{f}^{(B)}(\mathbf{v}_B(\mathbf{z})) D_{j'} \mathbf{v}_{i'}(\mathbf{z}) D_j \mathbf{v}_i(\mathbf{z}) \right] = 0$$

Simplifying the expression

$$\mathbf{M}(\mathbf{z}) \mathbf{w}(\mathbf{v}(\mathbf{z}), k) = 0.$$

$$\mathbf{w}(\mathbf{z}, k) := ((D_i \mathbf{f}_k^{(B)}(\mathbf{z}_B))_{i \in B}, (D_{i,i}^2 \mathbf{f}_k^{(B)}(\mathbf{z}_B))_{i \in B}, (D_{i,i'}^2 \mathbf{f}_k^{(B)}(\mathbf{z}_B))_{(i,i') \in B^2})_{B \in \mathcal{B}}$$

$$\mathbf{M}(\mathbf{z}) := [[\vec{a}_i(\mathbf{z})]_{i \in B}, [\vec{b}_i(\mathbf{z})]_{i \in B}, [\vec{c}_{i,i'}(\mathbf{z})]_{(i,i') \in B^2}]_{B \in \mathcal{B}},$$

Proof Sketch

$$\mathbf{M}(\mathbf{z})\mathbf{w}(\mathbf{v}(\mathbf{z}), k) = 0.$$

$$\mathbf{w}(\mathbf{z}, k) := ((D_i \mathbf{f}_k^{(B)}(\mathbf{z}_B))_{i \in B}, (D_{i,i}^2 \mathbf{f}_k^{(B)}(\mathbf{z}_B))_{i \in B}, (D_{i,i'}^2 \mathbf{f}_k^{(B)}(\mathbf{z}_B))_{(i,i') \in B_{\neq}^2})_{B \in \mathcal{B}}$$

$$\mathbf{M}(\mathbf{z}) := [[\vec{a}_i(\mathbf{z})]_{i \in B}, [\vec{b}_i(\mathbf{z})]_{i \in B}, [\vec{c}_{i,i'}(\mathbf{z})]_{(i,i') \in B_{\neq}^2}]_{B \in \mathcal{B}},$$



$$\mathbf{W}(\mathbf{v}(\mathbf{z}))^T = [w(\mathbf{v}(\mathbf{z}), 1), \dots, w(\mathbf{v}(\mathbf{z}), d_x)]$$

$$\mathbf{W}(\mathbf{v}(\mathbf{z}))\mathbf{M}(\mathbf{z})^T = 0$$

Assumption of sufficient non-linearity on f implies $\mathbf{W}(\mathbf{v}(\mathbf{z}))$ has full-column rank



$$\mathbf{M}(\mathbf{z})^T = 0$$

Proof Sketch

$$W(\mathbf{v}(\mathbf{z}))\mathbf{M}(\mathbf{z})^\top = 0$$

Assumption of sufficient non-linearity on f implies $W(\mathbf{v}(\mathbf{z}))$ has full-column rank

$$\begin{array}{c} \downarrow \\ \mathbf{M}(\mathbf{z}) := [[\vec{a}_i(\mathbf{z})]_{i \in B}, [\vec{b}_i(\mathbf{z})]_{i \in B}, [\vec{c}_{i,i'}(\mathbf{z})]_{(i,i') \in B^2_\zeta}]_{B \in \mathcal{B}} \\ \downarrow \\ \mathbf{M}(\mathbf{z})^\top = 0 \end{array}$$

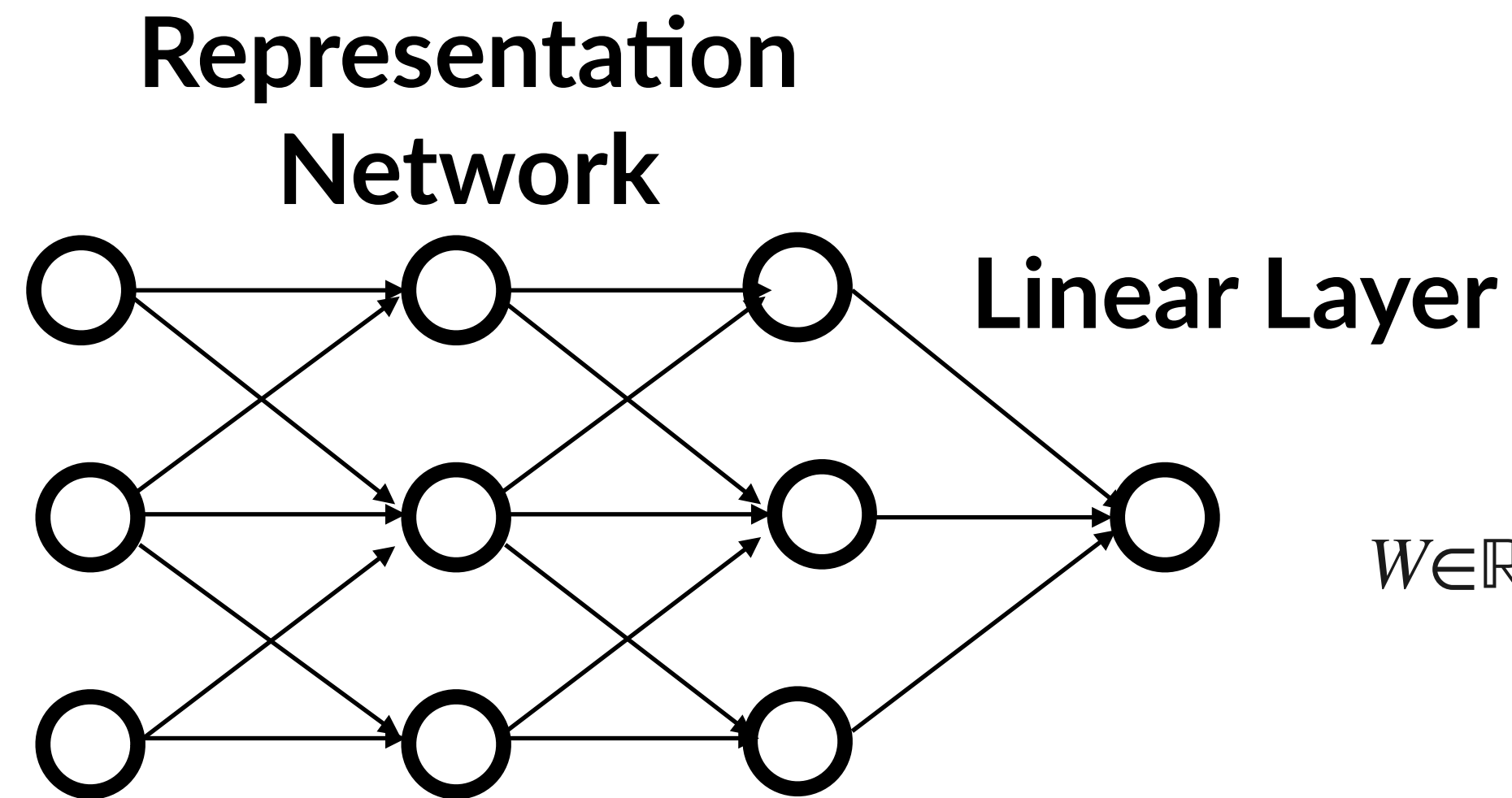
$$\forall i \in \{1, \dots, d_z\}, \vec{b}_i(\mathbf{z}) = 0,$$

$$\forall i \in \{1, \dots, d_z\}, \forall (j, j') \in S^c, D_j \mathbf{v}_i(\mathbf{z}) D_{j'} \mathbf{v}_i(\mathbf{z}) = 0$$

The final equality can be implied further to show $D\mathbf{v}(\mathbf{z})$ is \mathcal{B} -block permutation matrix

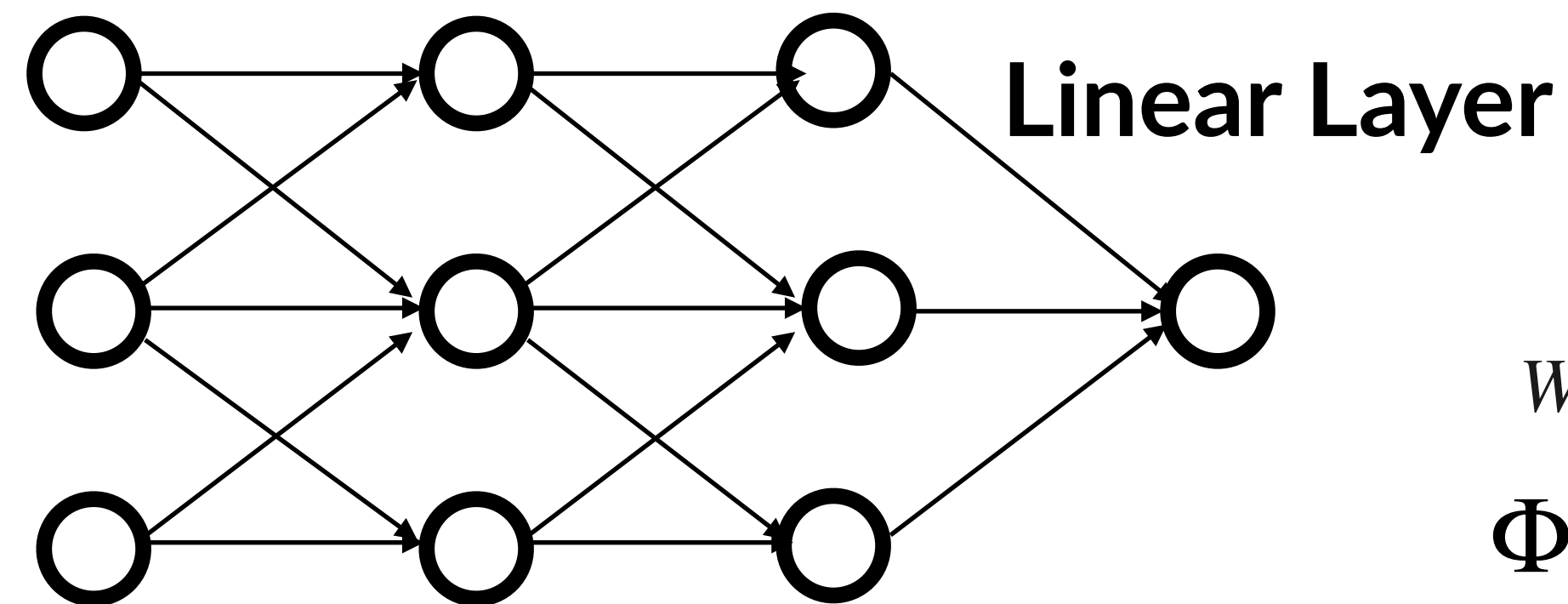
Independence Constrained ERM

ERM



$$\min_{W \in \mathbb{R}^{d \times k}, \Phi \in \mathcal{H}_\Phi} \sum_{i=1}^N \ell(W \circ \Phi(x_i), y_i)$$

IC-ERM



$$\min_{W \in \mathbb{R}^{d \times k}, \Phi} \sum_{i=1}^N \ell(W \circ \Phi(x_i), y_i)$$

$\Phi(X)$ is i.i.d.

Representation
Network