# An Identifiability Perspective on Representation Learning

Yash Sharma

# Outline

1. Towards Nonlinear Disentanglement in Natural Data with Temporal Sparse Coding
   a. Identifiability w/ assumptions derived from natural video

# Outline

1. Towards Nonlinear Disentanglement in Natural Data with Temporal Sparse Coding
   a. Identifiability w/ assumptions derived from natural video
2. Contrastive Learning Inverts the Data Generating Process
   a. Identifiability & InfoNCE

# Outline

1. Towards Nonlinear Disentanglement in Natural Data with Temporal Sparse Coding
   a. Identifiability w/ assumptions derived from natural video
2. Contrastive Learning Inverts the Data Generating Process
   a. Identifiability & InfoNCE
3. Self-Supervised Learning with Data Augmentations Provably Isolates Content from Style
   a. Identifiability when augmentations leave factors invariant

# Towards Nonlinear Disentanglement in Natural Data with Temporal Sparse Coding

# Overview

1. Problem Statement
   a. What is disentanglement?

# Overview

1. Problem Statement
   a. What is disentanglement?

# Overview

1. Problem Statement
   a. What is disentanglement?
   b. What do we need to solve disentanglement?

# Overview

1. Problem Statement
    a. What is disentanglement?
    b. What do we need to solve disentanglement?
    c. Can we find what we need in natural video?

# Overview

1. Problem Statement
   a. What is disentanglement?
   b. What do we need to solve disentanglement?
   c. Can we find what we need in natural video?
2. Theoretical Contributions
   a. A prior based on natural statistics provably enables disentanglement

# Overview

1. Problem Statement
   a. What is disentanglement?
   b. What do we need to solve disentanglement?
   c. Can we find what we need in natural video?
2. Theoretical Contributions
   a. A prior based on natural statistics provably enables disentanglement
3. Empirical Contributions
   a. Qualitative and quantitative results on existing and contributed datasets demonstrate outperformance in aggregate

# What is Disentanglement?

State of the world, **z**:

- shape
- scale
- orientation
- x position
- y position

# What is Disentanglement?

State of the world, **z**:

- shape
- scale
- orientation
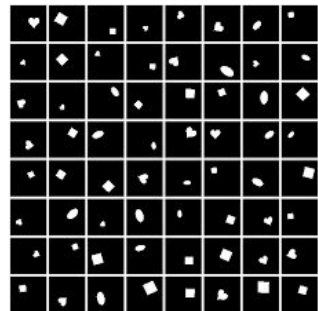- x position
- y position



data generating process

# What is Disentanglement?

State of the world, **z**:

- shape
- scale
- orientation
- x position
- y position

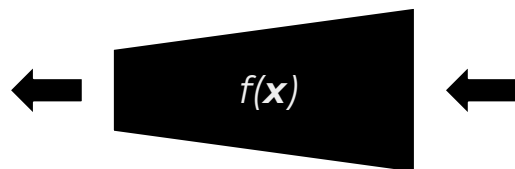*g(**z**)*

data generating process

Observations, **x**:

# What is Disentanglement?

State of the world, **z**:

- shape
- scale
- orientation
- x position
- y position

*g(**z**)*

data generating process

*f(**x**)*

representation learning

Observations, **x**:

# What is Disentanglement?

State of the world, **z**:

- shape
- scale
- orientation
- x position
- y position

$g(\mathbf{z})$

data generating process

$f(\mathbf{x})$

representation learning

Observations, **x**:

**Too simplistic**

# What is Disentanglement?

State of the world, **z**:

- shape
- scale
- orientation
- x position
- y position

$g(\boldsymbol{z})$

data generating process

Observations, **x**:

$f(\boldsymbol{x})$
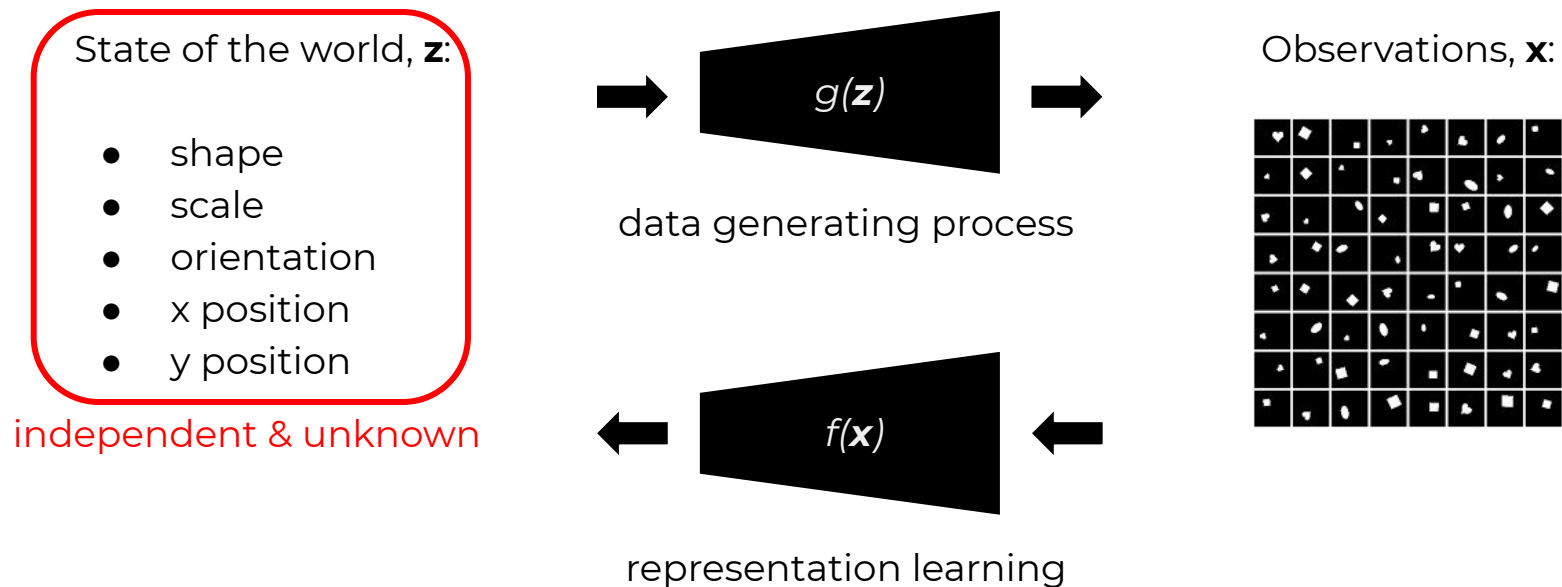
representation learning

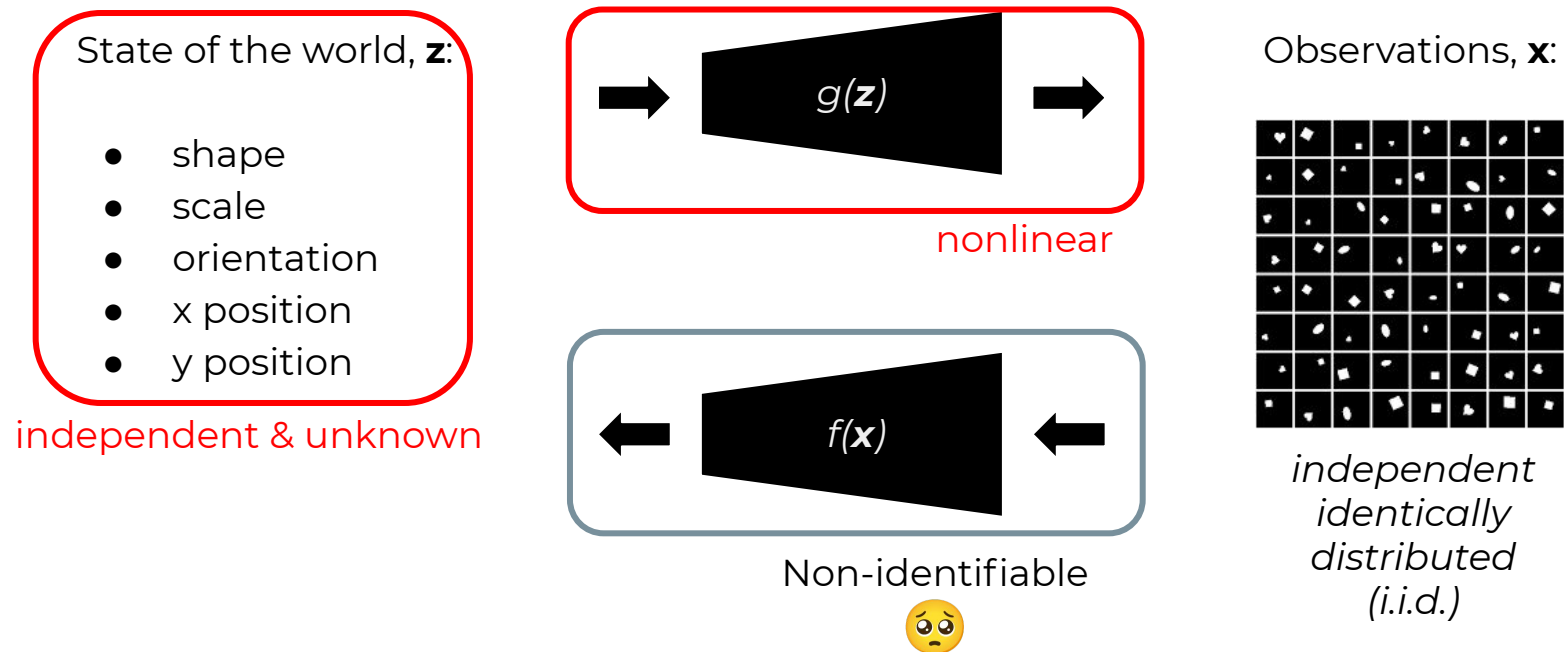**prev solutions do not account for observed natural scene statistics**
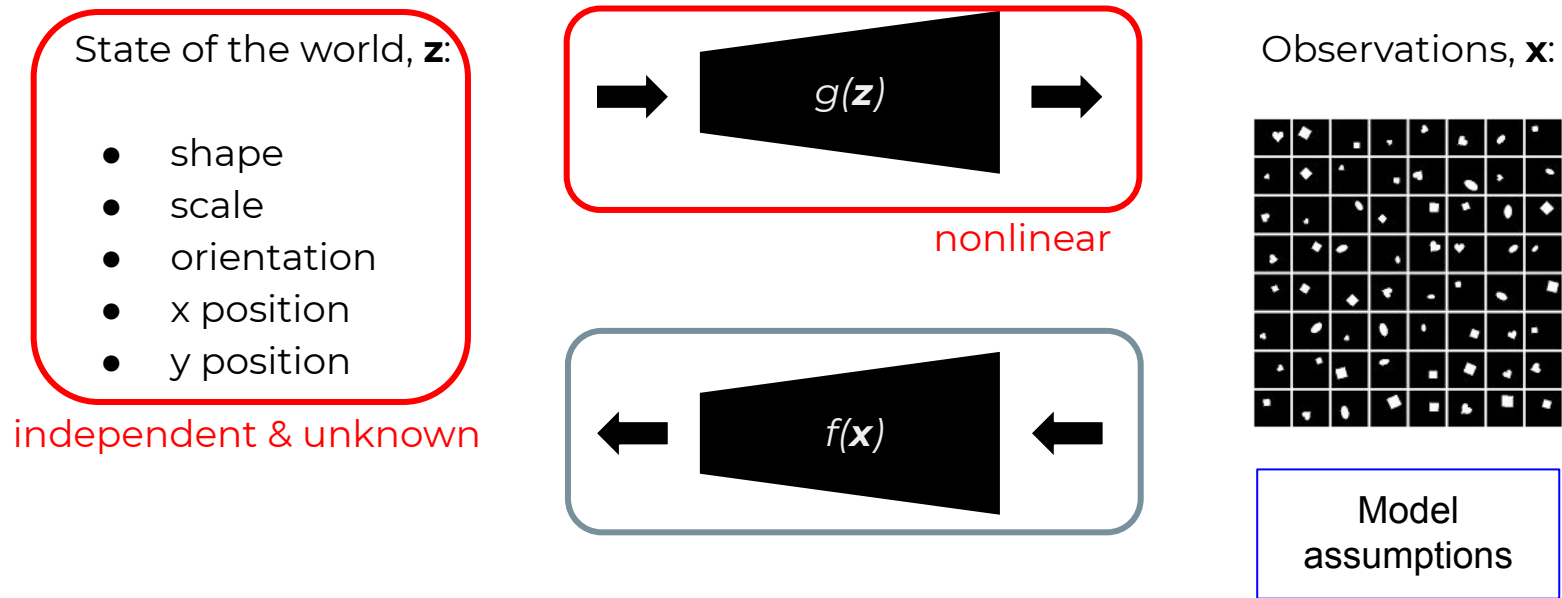
**Too simplistic**

# Non-identifiability



State of the world, **z**:

- shape
- scale
- orientation
- x position
- y position

independent & unknown

$g(\boldsymbol{z})$

data generating process

$f(\boldsymbol{x})$

representation learning

Observations, **x**:

Hyvärinen & Pajunen (1999) *Nonlinear independent component analysis: Existence and uniqueness results*

# Non-identifiability

State of the world, **z**:

- shape
- scale
- orientation
- x position
- y position

independent & unknown

$g(\boldsymbol{z})$

nonlinear

$f(\boldsymbol{x})$

representation learning

Observations, **x**:

Hyvärinen & Pajunen (1999) *Nonlinear independent component analysis: Existence and uniqueness results*

# Non-identifiability



State of the world, **z**:

- shape
- scale
- orientation
- x position
- y position

**independent & unknown**

$g(\textbf{z})$

**nonlinear**

$f(\textbf{x})$

Non-identifiable
🥺

Observations, **x**:

*independent
identically
distributed
(i.i.d.)*

Locatello et al. (2018) *Challenging Common Assumptions in the Unsupervised Learning of Disentangled Representations*

Hyvärinen & Pajunen (1999) *Nonlinear independent component analysis: Existence and uniqueness results*

# Nonlinear Disentanglement

State of the world, **z**:

- shape
- scale
- orientation
- x position
- y position

independent & unknown

$g(\boldsymbol{z})$

nonlinear

$f(\boldsymbol{x})$

Observations, **x**:



Model
assumptions

Locatello et al. (2020) *Weakly-Supervised Disentanglement Without Compromises*

Hyvärinen & Morioka (2017) *Nonlinear ICA of Temporally Dependent Stationary Sources*

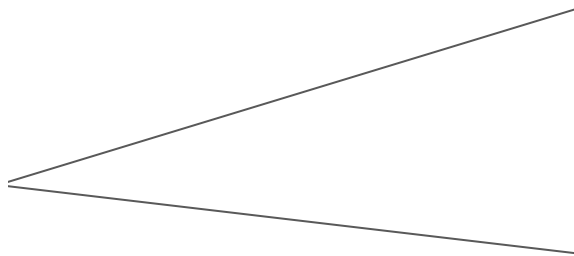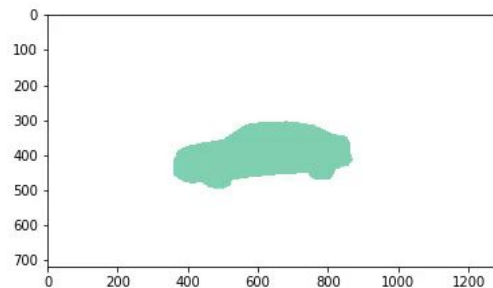Hyvärinen & Morioka (2016) *Unsupervised Feature Extraction by Time-Contrastive Learning and Nonlinear ICA*

# Natural Video

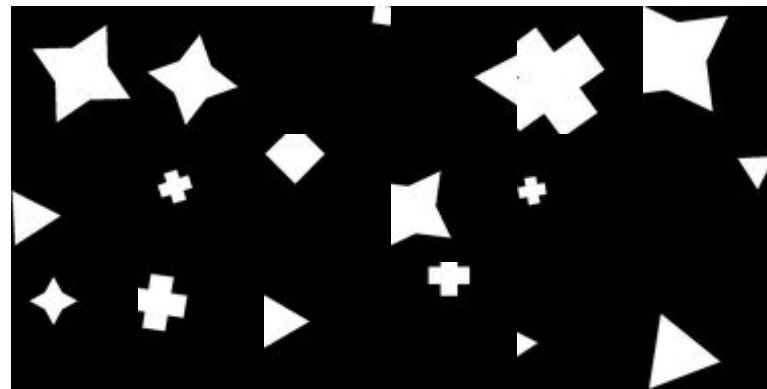# Natural Video

# Natural Video



- scale
- x position
- y position
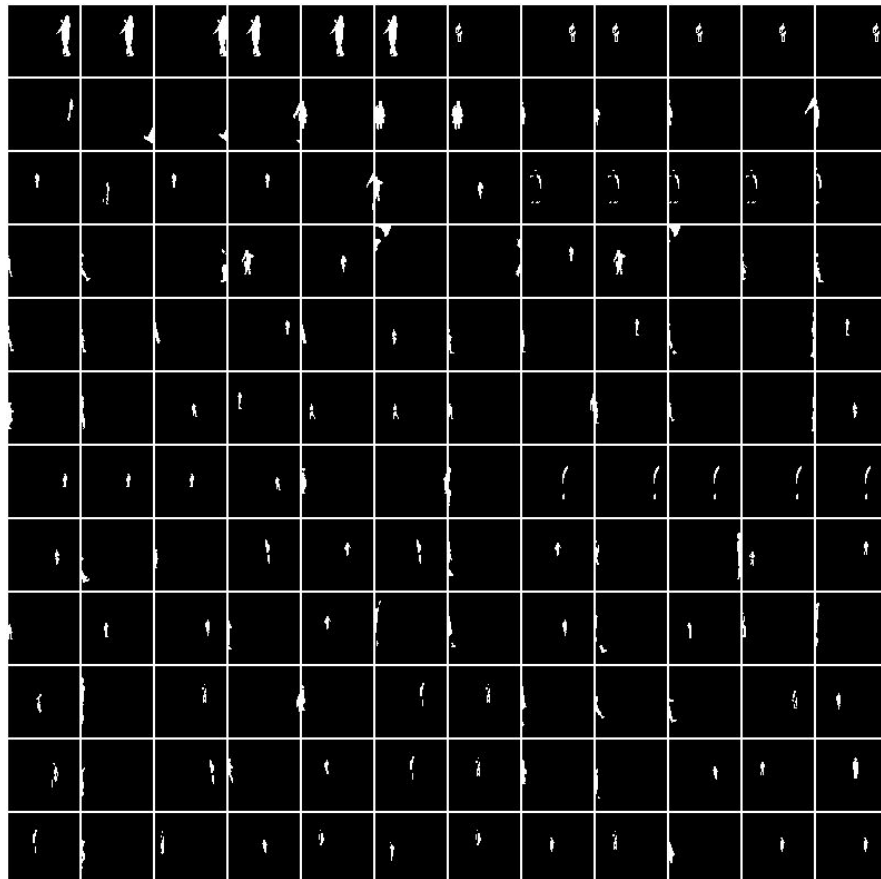
# Natural Data Analysis

# Natural Sprites

- Images generated online using renderer

- Simple, well-controlled objects

- Transitions sampled from YouTube-VOS

# KITTI Masks

- Pedestrian masks extracted directly from autonomous recorded videos

- Realistic objects & transitions

# The world is not *i.i.d.*

# The world is not *i.i.d.*
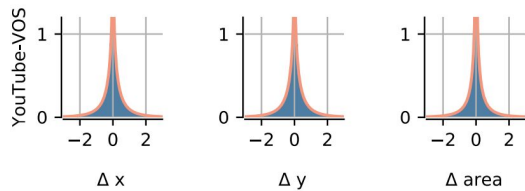
State of the world, **z**:

{shape, scale, orientation,
x position, y position}

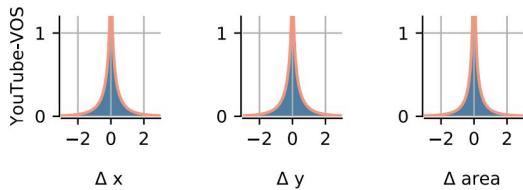# The world is not *i.i.d.*

State of the world, **z**:

{shape, scale, orientation,
x position, y position}

And dynamics:

# The world is not *i.i.d.*

State of the world, **z**:

{shape, scale, orientation,
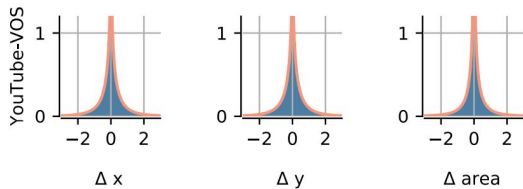x position, y position}

And dynamics:



*g(**z**)*

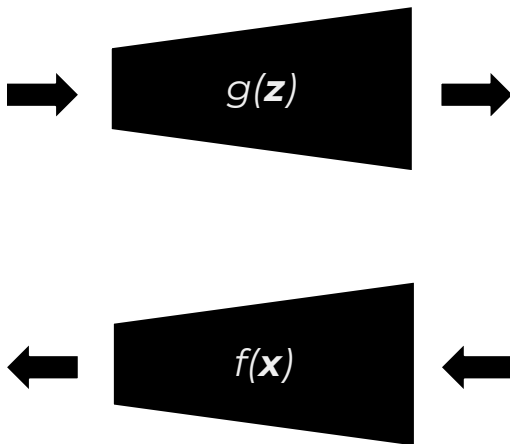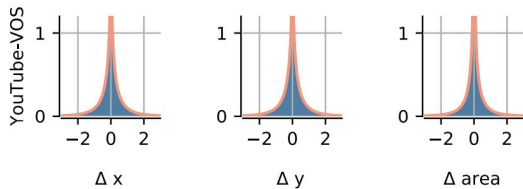# The world is not *i.i.d.*

State of the world, **z**:

{shape, scale, orientation,
x position, y position}

And dynamics:



Observations:

# The world is not *i.i.d.*

State of the world, **z**:

{shape, scale, orientation,
x position, y position}

And dynamics:



$g(\boldsymbol{z})$
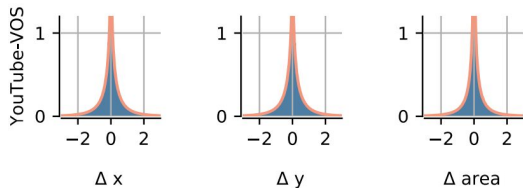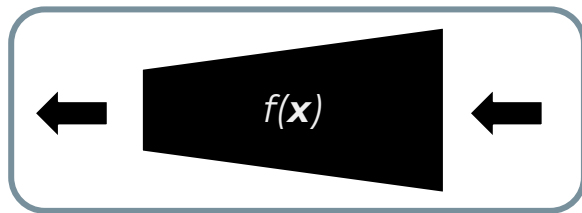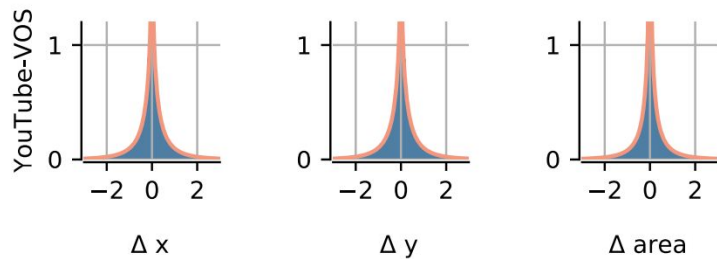
$f(\boldsymbol{x})$

Observations:



$\boldsymbol{x}_{t-1}$ $\quad$ $\boldsymbol{x}_{t}$

# The world is not *i.i.d.*

State of the world, **z**:

{shape, scale, orientation, x position, y position}

And dynamics:



$g(\boldsymbol{z})$

$f(\boldsymbol{x})$

identifiable 🤩

Observations:

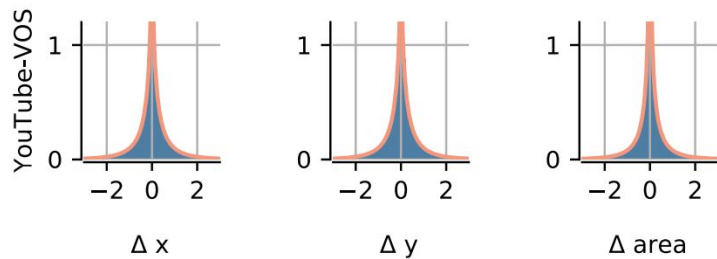$\boldsymbol{x}_{t\text{-}1}$   $\boldsymbol{x}_t$

# Identifiability Proof Intuition

# Identifiability Proof Intuition



**Prior:**
*objects in nature*
*change sparsely*

# Identifiability Proof Intuition



YouTube-VOS

Δ x          Δ y          Δ area

**Prior:**
*objects in nature change sparsely*

y position

x position

Probability of next position given previous position

# Identifiability Proof Intuition


YouTube-VOS

Δ x    Δ y    Δ area

**Prior:**
*objects in nature change sparsely*

y position

x position

Probability of next position given previous position

# Identifiability Proof Intuition



YouTube-VOS

Δ x   Δ y   Δ area

**Prior:**

*objects in nature change sparsely*

y position

x position

Probability of next position given previous position

# Identifiability Proof Intuition



True Model

$y^*$

$x^*$

**Prior:**
*objects in nature change sparsely*

x position

y position

# Identifiability Proof Intuition



**Prior:**
*objects in nature change sparsely*

# Identifiability Proof Intuition

YouTube-VOS

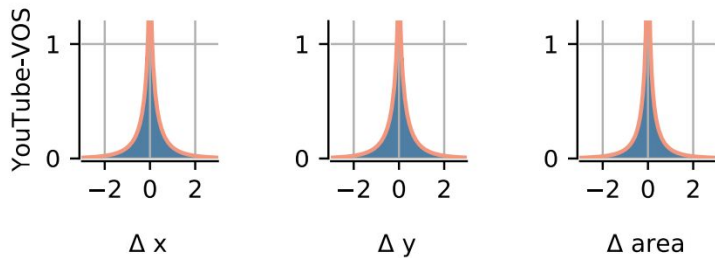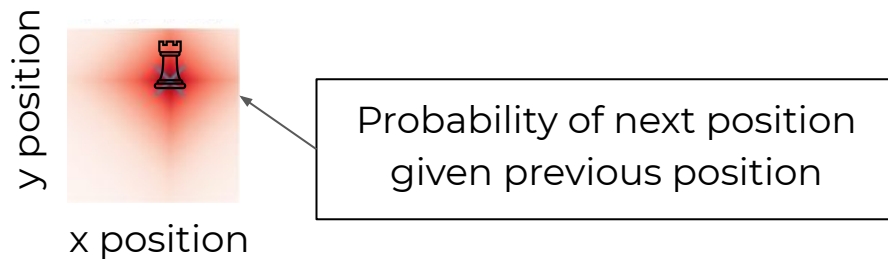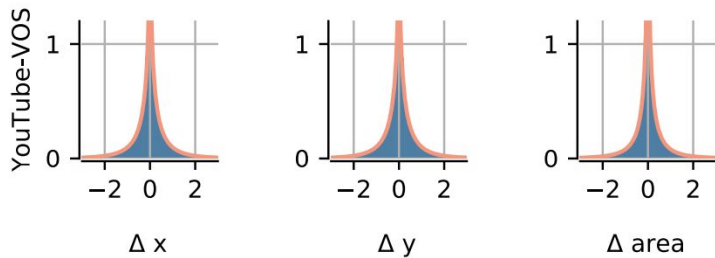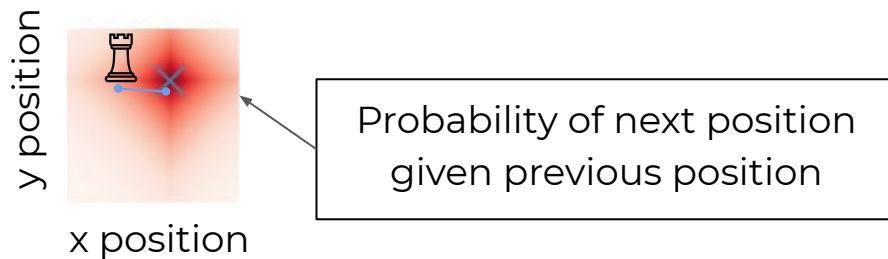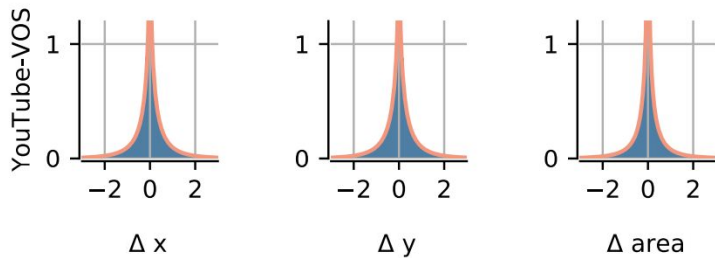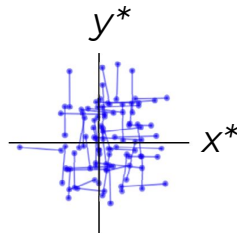$\Delta x$    $\Delta y$    $\Delta$ area

True Model

$y^*$    $x^*$

$g^*(z)$

**Prior:**
*objects in nature change sparsely*

y position

x position

# Identifiability Proof Intuition



YouTube-VOS

Δ x   Δ y   Δ area

**Prior:**
*objects in nature change sparsely*

y position

x position

True Model

$y^*$

$x^*$

≠

Learned Model

$y^*$

$y$

$x^*$

$x$

model latents

$g^*(\mathbf{z})$

# Identifiability Proof Intuition



True Model    Learned Model

$y^*$    $y^*$    $y$

$x^*$    ≠    $x^*$

$x$

$g^*(z)$    $g(z)$

YouTube-VOS

1    1    1

0    0    0

−2  0  2    −2  0  2    −2  0  2

Δ x    Δ y    Δ area

**Prior:**
*objects in nature change sparsely*

y position

x position

# Identifiability Proof Intuition



True Model    Learned Model

$y^*$    $y^*$    $y$

$x^*$    $\neq$    $x^*$

$x$

$g^*(\mathbf{z})$    $g(\mathbf{z})$

YouTube-VOS

1    1    1

0    0    0

−2  0  2    −2  0  2    −2  0  2

Δ x    Δ y    Δ area

**Prior:**
*objects in nature change sparsely*

y position

x position

# Identifiability Proof Intuition



**Prior:**
*objects in nature change sparsely*

True Model ≠ Learned Model

$g^*(z)$   $g(z)$

sparse transitions

# Identifiability Proof Intuition

True Model    Learned Model

$y^*$    $y^*$    $y$

$x^*$    $\neq$    $x^*$

$x$

Prior:
*objects in nature
change sparsely*

y position

x position

$g^*(\boldsymbol{z})$    $g(\boldsymbol{z})$

sparse transitions    $\neq$    dense transitions

YouTube-VOS

1

0

−2  0  2

Δ x

1

0

−2  0  2

Δ y

1

0

−2  0  2

Δ area

# Slow Variational Autoencoder at time *t-1*

# Slow Variational Autoencoder at time *t-1*

# Quantitative Results - dSprites

# Quantitative Results - dSprites

**Data**

$X_{t-1}$ $\quad$ $X_t$

# Quantitative Results - dSprites

**Data**

$X_{t-1}$   $X_t$



Ada-GVAE [Locatello et al. 19]

PCL [Hyvärinen et al. 17]

SlowVAE (ours)

# Quantitative Results - dSprites

**Data**

$X_{t-1}$    $X_t$

**Disentanglement Performance**



Ada-GVAE [Locatello et al. 19]

PCL [Hyvärinen et al. 17]

SlowVAE (ours)

# Quantitative Results - dSprites

**Data**

$X_{t-1}$   $X_t$

**Disentanglement Performance**



Average Model Rank

1st

2nd

3rd

dSprites
(Laplace)

Ada-GVAE [Locatello et al. 19]

PCL [Hyvärinen et al. 17]

SlowVAE (ours)

# Quantitative Results - dSprites

**Data**



**Disentanglement Performance**

Ada-GVAE [Locatello et al. 19]

PCL [Hyvärinen et al. 17]

SlowVAE (ours)

# Qualitative Results - dSprites



**SlowVAE (Latent Walk)**

# Qualitative Results - dSprites



**SlowVAE (Latent Walk)**

# Qualitative Results - dSprites

**SlowVAE (Latent Walk)**

# Failure Cases - dSprites

# Failure Cases - dSprites

# Failure Cases - dSprites

# Failure Cases - dSprites

# Results on Natural Data

**KITTI-Masks**

# Results on Natural Data

**KITTI-Masks**



**Disentanglement Performance**



MCC-Metric

75

50

25

0

KITTI-Masks
(Δt=0.15s)

Ada-VAE [Locatello et al. 19]

PCL [Hyvärinen et al. 17]

SlowVAE (ours)

# Results - KITTI-Masks

**SlowVAE Latent Walk** 🕺



Scale          X-Position          Y-Position

# Paper Contributions

Objects in natural scenes have **sparse** marginal transition statistics

# Paper Contributions

Objects in natural scenes have **sparse** marginal transition statistics

Intuitive proof for **identifiability** in Nonlinear ICA & Disentanglement

# Paper Contributions

Objects in natural scenes have **sparse** marginal transition statistics



Intuitive proof for **identifiability** in Nonlinear ICA & Disentanglement



Empirical results using a **Flow** and a **VAE** based implementation of the theoretical model

# Paper Contributions

Objects in natural scenes have **sparse** marginal transition statistics



Intuitive proof for **identifiability** in Nonlinear ICA & Disentanglement



Two challenging new **datasets** to push disentanglement towards **natural** video



Empirical results using a **Flow** and a **VAE** based implementation of the theoretical model

# Contrastive Learning Inverts the Data Generating Process

# Overview

1. Theoretical Connection between InfoNCE & Nonlinear ICA

# Overview

1. Theoretical Connection between InfoNCE & Nonlinear ICA
2. Empirical Test on robustness to mismatch (in assumptions)

# Overview

1. Theoretical Connection between InfoNCE & Nonlinear ICA
2. Empirical Test on robustness to mismatch (in assumptions)
3. Identifiability on 3DIdent

# Overview

1. Theoretical Connection between InfoNCE & Nonlinear ICA
2. Empirical Test on robustness to mismatch (in assumptions)
3. Identifiability on 3DIdent
   a. complex, high-resolution images

# Nonlinear ICA

# Theoretical Framework



Unobservable Latent Space $\mathcal{Z}$

Positive $\mathbf{x}^+$

Unknown Generative Process g

Unobservable

f

Negatives $\mathbf{x}_1^-, \mathbf{x}_2^-, ...$

Reconstructed Latent Space $\mathcal{Z}' = \mathbf{A}\mathcal{Z}$

● Uniformly Distributed Anchor
● Conditional Density of Positive Samples
● Uniformly Distributed Negative Samples

Learned by $\mathcal{L} = -f(\mathbf{x})^\top f(\mathbf{x}^+) + \log \sum_{\{\mathbf{x}^+, \mathbf{x}_1^-, ...\}} \exp(f(\mathbf{x})^\top f(\mathbf{x}'))$
          attract                                    repel

# Theoretical Framework

**Theorem 2.** *Let $\mathcal{Z} = \mathbb{S}^{N-1}$, the ground-truth marginal be uniform, and the conditional a vMF distribution (cf. Eq. 2). Let the mixing function $g$ be differentiable and injective. If the assumed form of $q_\mathrm{h}$, as defined above, matches that of $p$, and if $f$ is differentiable and minimizes the CL loss (1), then for fixed $\tau > 0$ and $M \to \infty$, $h = f \circ g$ is linear, i.e., $f$ recovers the latent sources up to orthogonal linear transformations.*

$$\mathcal{L}_{\text{contr}}(f; \tau, M) := \tag{1}$$

$$\mathop{\mathbb{E}}_{\substack{(\mathbf{x},\tilde{\mathbf{x}})\sim p_{\text{pos}} \\ \{\mathbf{x}_i^-\}_{i=1}^M \overset{\text{i.i.d.}}{\sim} p_{\text{data}}}} \left[ -\log \frac{e^{f(\mathbf{x})^\mathsf{T} f(\tilde{\mathbf{x}})/\tau}}{e^{f(\mathbf{x})^\mathsf{T} f(\tilde{\mathbf{x}})/\tau} + \sum_{i=1}^M e^{f(\mathbf{x}_i^-)^\mathsf{T} f(\tilde{\mathbf{x}})/\tau}} \right].$$

$$p(\mathbf{z}) = |\mathcal{Z}|^{-1}, \qquad p(\mathbf{z}|\tilde{\mathbf{z}}) = C_p^{-1} e^{\kappa \mathbf{z}^\mathsf{T} \tilde{\mathbf{z}}}$$

$$\text{with} \quad C_p := \int e^{\kappa \mathbf{z}^\mathsf{T} \tilde{\mathbf{z}}} \, \mathrm{d}\tilde{\mathbf{z}} = \text{const.}, \quad \mathbf{x} = g(\mathbf{z}). \tag{2}$$

$$q_\mathrm{h}(\tilde{\mathbf{z}}|\mathbf{z}) = C_h(\tilde{\mathbf{z}})^{-1} e^{h(\tilde{\mathbf{z}})^\mathsf{T} h(\mathbf{z})/\tau}$$

$$\text{with} \quad C_h(\mathbf{z}) := \int e^{h(\tilde{\mathbf{z}})^\mathsf{T} h(\mathbf{z})/\tau} \, \mathrm{d}\mathbf{z},$$

# Theoretical Framework

**Theorem 1** ($\mathcal{L}_{\text{contr}}$ converges to the cross-entropy between latent distributions). *If the ground-truth marginal distribution $p$ is uniform, then for fixed $\tau > 0$, as the number of negative samples $M \to \infty$, the (normalized) contrastive loss converges to*

$$\lim_{M \to \infty} \mathcal{L}_{\text{contr}}(f; \tau, M) - \log M + \log |\mathcal{Z}| = $$
$$\mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} \left[ H(p(\cdot|\mathbf{z}), q_h(\cdot|\mathbf{z})) \right] \tag{14}$$

*where $H$ is the cross-entropy between the ground-truth conditional distribution $p$ over positive pairs and a conditional distribution $q_{\text{h}}$ parameterized by the model $f$, and $C_h(\mathbf{z}) \in \mathbb{R}^+$ is the partition function of $q_{\text{h}}$ (see Appendix A.1.1):*

$$q_{\text{h}}(\tilde{\mathbf{z}}|\mathbf{z}) = C_h(\mathbf{z})^{-1} e^{h(\tilde{\mathbf{z}})^{\top} h(\mathbf{z})/\tau}$$
$$\textit{with} \quad C_h(\mathbf{z}) := \int e^{h(\tilde{\mathbf{z}})^{\top} h(\mathbf{z})/\tau} \, d\tilde{\mathbf{z}}. \tag{15}$$

# Theoretical Framework

**Proposition 1** (Minimizers of the cross-entropy maintain the dot product). *Let* $\mathcal{Z} = \mathbb{S}^{N-1}$, $\tau > 0$ *and consider the ground-truth conditional distribution of the form* $p(\tilde{\mathbf{z}}|\mathbf{z}) = C_p^{-1}\exp(\kappa\tilde{\mathbf{z}}^\top\mathbf{z})$. *Let* $h$ *map onto a hypersphere with radius* $\sqrt{\tau\kappa}$.[3] *Consider the conditional distribution* $q_h$ *parameterized by the model, as defined above in Theorem 1, where the hypothesis class for* $h$ *is assumed to be sufficiently flexible such that* $p(\tilde{\mathbf{z}}|\mathbf{z})$ *and* $q_h(\tilde{\mathbf{z}}|\mathbf{z})$ *can match. If* $h$ *is a minimizer of the cross-entropy* $\mathbb{E}_{p(\tilde{\mathbf{z}}|\mathbf{z})}[-\log q_h(\tilde{\mathbf{z}}|\mathbf{z})]$, *then* $p(\tilde{\mathbf{z}}|\mathbf{z}) = q_h(\tilde{\mathbf{z}}|\mathbf{z})$ *and* $\forall\mathbf{z},\tilde{\mathbf{z}} : \kappa\mathbf{z}^\top\tilde{\mathbf{z}} = h(\mathbf{z})^\top h(\tilde{\mathbf{z}})$.

**Proposition 2** (Extension of the Mazur-Ulam theorem to hyperspheres and the dot product). *Let* $\mathcal{Z} = \mathbb{S}^{N-1}$. *If* $h : \mathcal{Z} \to \mathcal{Z}$ *maintains the dot product up to a constant factor, i.e.,* $\forall\mathbf{z},\tilde{\mathbf{z}} \in \mathcal{Z} : \kappa\mathbf{z}^\top\tilde{\mathbf{z}} = h(\mathbf{z})^\top h(\tilde{\mathbf{z}})$, *then* $h$ *is an orthogonal linear transformation.*

# Theoretical Framework

**Theorem 5.** *Let $\mathcal{Z}$ be a convex body in $\mathbb{R}^N$, $h = f \circ g : \mathcal{Z} \to \mathcal{Z}$, and $\delta$ be a metric. Further, let the ground-truth marginal distribution be uniform and the conditional distribution be as (5). Let the mixing function $g$ be differentiable and injective. If the assumed form of $q_h$ matches that of $p$, i.e.,*

$$q_{\mathrm{h}}(\tilde{\mathbf{z}}|\mathbf{z}) = C_q^{-1}(\mathbf{z})e^{-\delta(h(\tilde{\mathbf{z}}),h(\mathbf{z}))/\tau}$$

$$\text{with} \quad C_q(\mathbf{z}) := \int e^{-\delta(h(\tilde{\mathbf{z}}),h(\mathbf{z}))/\tau}\,\mathrm{d}\tilde{\mathbf{z}}, \tag{7}$$

*and if $f$ is differentiable and minimizes the $\mathcal{L}_{\delta\text{-contr}}$ objective in (6) for $M \to \infty$, we find that $h = f \circ g$ is invertible and affine, i.e., we recover the latent sources up to affine transformations.*

$$p(\mathbf{z}) = |\mathcal{Z}|^{-1}, \qquad p(\mathbf{z}|\tilde{\mathbf{z}}) = C_p^{-1}e^{-\delta(\mathbf{z},\tilde{\mathbf{z}})}$$

$$\text{with} \quad C_p(\mathbf{z}) := \int e^{-\delta(\mathbf{z},\tilde{\mathbf{z}})}\,\mathrm{d}\tilde{\mathbf{z}}, \quad \mathbf{x} = g(\mathbf{z}), \tag{5}$$

$$\mathcal{L}_{\delta\text{-contr}}(f;\tau,M) := \tag{6}$$

$$\mathop{\mathbb{E}}_{\substack{(\mathbf{x},\tilde{\mathbf{x}})\sim p_{\mathrm{pos}} \\ \{\mathbf{x}_i^-\}_{i=1}^M \overset{\text{i.i.d.}}{\sim} p_{\mathrm{data}}}}\left[-\log \frac{e^{-\delta(f(\mathbf{x}),f(\tilde{\mathbf{x}}))/\tau}}{e^{-\delta(f(\mathbf{x}),f(\tilde{\mathbf{x}}))/\tau}+\sum_{i=1}^M e^{-\delta(f(\mathbf{x}_i^-),f(\tilde{\mathbf{x}}))/\tau}}\right].$$

# Theoretical Framework

**Theorem 3.** *Let $\delta$ be a semi-metric and $\tau, \lambda > 0$ and let the ground-truth marginal distribution $p$ be uniform. Consider a ground-truth conditional distribution $p(\tilde{\mathbf{z}}|\mathbf{z}) = C_p^{-1}(\mathbf{z}) \exp(-\lambda \delta(\tilde{\mathbf{z}}, \mathbf{z}))$ and the model conditional distribution*

$$q_{\mathrm{h}}(\tilde{\mathbf{z}}|\mathbf{z}) = C_h^{-1}(\mathbf{z}) e^{-\delta(h(\tilde{\mathbf{z}}), h(\mathbf{z}))/\tau}$$

$$\text{with} \quad C_h(\mathbf{z}) := \int_{\mathcal{Z}} e^{-\delta(h(\tilde{\mathbf{z}}), h(\mathbf{z}))/\tau} d\tilde{\mathbf{z}}. \tag{61}$$

*Then the cross-entropy between $p$ and $q_{\mathrm{h}}$ is given by*

$$\lim_{M \to \infty} \mathcal{L}_{\delta\text{-contr}}(f; \tau, M) - \log M + \log |\mathcal{Z}| =$$
$$\mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} \left[ H(p(\cdot|\mathbf{z}), q_{\mathrm{h}}(\cdot|\mathbf{z}) \right], \tag{62}$$

*which can be implemented by sampling data from the accessible distributions.*

**Theorem 4.** *Let $\mathcal{Z} = \mathcal{Z}'$ be a convex body in $\mathbb{R}^N$. Let the mixing function $g$ be differentiable and invertible. If the assumed form of $q_{\mathrm{h}}$ as defined in (4) matches that of $p$, and if $f$ is differentiable and minimizes the cross-entropy between $p$ and $q_{\mathrm{h}}$, then we find that $h = f \circ g$ is affine, i.e., we recover the latent sources up to affine transformations.*

# Theoretical Framework

**Theorem 6.** *Let $\mathcal{Z}$ be a convex body in $\mathbb{R}^N$, $h : \mathcal{Z} \to \mathcal{Z}$, and $\delta$ be an $L^\alpha$ metric for $\alpha \geq 1, \alpha \neq 2$. Further, let the ground-truth marginal distribution be uniform and the conditional distribution be as (5), and let the mixing function $g$ be differentiable and invertible. If the assumed form of $q_{\mathrm{h}}(\cdot|\mathbf{z})$ matches that of $p(\cdot|\mathbf{z})$, i.e., both use the same metric $\delta$ up to a constant scaling factor, and if $f$ is differentiable and minimizes the $\mathcal{L}_{\delta\text{-contr}}$ objective in (6) for $M \to \infty$, we find that $h = f \circ g$ is a composition of input independent permutations, sign flips and rescaling.*

**Theorem D.** *Suppose $1 \leq \alpha \leq \infty$ and $\alpha \neq 2$. An $n \times n$ matrix $\mathbf{A}$ is an isometry of $L^\alpha$-norm if and only if $\mathbf{A}$ is a generalized permutation matrix, i.e., $\forall \mathbf{z} : (\mathbf{A}\mathbf{z})_{\mathbf{i}} = \alpha_{\mathbf{i}}\mathbf{z}_{\sigma(\mathbf{i})}$, with $\alpha_i = \pm 2$ and $\sigma$ being a permutation.*

*Proof.* See Li & So (1994). Note that this can also be concluded from the Banach-Lamperti Theorem (Lamperti et al., 1958). □

$$p(\mathbf{z}) = |\mathcal{Z}|^{-1}, \qquad p(\mathbf{z}|\tilde{\mathbf{z}}) = C_p^{-1} e^{-\delta(\mathbf{z},\tilde{\mathbf{z}})}$$

$$\text{with} \quad C_p(\mathbf{z}) := \int e^{-\delta(\mathbf{z},\tilde{\mathbf{z}})} \, \mathrm{d}\tilde{\mathbf{z}}, \quad \mathbf{x} = g(\mathbf{z}), \tag{5}$$

$$\mathcal{L}_{\delta\text{-contr}}(f; \tau, M) := \tag{6}$$

$$\mathbb{E}_{\substack{(\mathbf{x},\tilde{\mathbf{x}}) \sim p_{\mathrm{pos}} \\ \{\mathbf{x}_i^-\}_{i=1}^M \stackrel{\mathrm{i.i.d.}}{\sim} p_{\mathrm{data}}}} \left[ -\log \frac{e^{-\delta(f(\mathbf{x}),f(\tilde{\mathbf{x}}))/\tau}}{e^{-\delta(f(\mathbf{x}),f(\tilde{\mathbf{x}}))/\tau} + \sum_{i=1}^M e^{-\delta(f(\mathbf{x}_i^-),f(\tilde{\mathbf{x}}))/\tau}} \right].$$

# Different Assumptions, Different Losses



$\bigcirc$ + vMF $\longrightarrow$ $\exp(f(x)^\top f(x'))$

$\square$ + Normal $\longrightarrow$ $\exp(-\|f(x) - f(x')\|_2)$

$\square$ + Laplace $\longrightarrow$ $\exp(-\|f(x) - f(x')\|_1)$

# Empirical Results

| | Generative process $g$ | | | Model $f$ | | $R^2$ Score [%] | | |
|---|---|---|---|---|---|---|---|---|
| Space | $p(\cdot)$ | $p(\cdot\|\cdot)$ | Space | $q_{\mathrm{h}}(\cdot\|\cdot)$ | M. | Identity | Supervised | Unsupervised |
| Sphere | Uniform | vMF($\kappa$=1) | Sphere | vMF($\kappa$=1) | ✓ | $66.98 \pm 2.79$ | $99.71 \pm 0.05$ | $99.42 \pm 0.05$ |
| Sphere | Uniform | vMF($\kappa$=10) | Sphere | vMF($\kappa$=1) | ✗ | —"— | —"— | $99.86 \pm 0.01$ |
| Sphere | Uniform | Laplace($\lambda$=0.05) | Sphere | vMF($\kappa$=1) | ✗ | —"— | —"— | $99.91 \pm 0.01$ |
| Sphere | Uniform | Normal($\sigma$=0.05) | Sphere | vMF($\kappa$=1) | ✗ | —"— | —"— | $99.86 \pm 0.00$ |
| Box | Uniform | Normal($\sigma$=0.05) | Unbounded | Normal | ✗ | $67.93 \pm 7.40$ | $99.78 \pm 0.06$ | $99.60 \pm 0.02$ |
| Box | Uniform | Laplace($\lambda$=0.05) | Unbounded | Normal | ✗ | —"— | —"— | $99.64 \pm 0.02$ |
| Box | Uniform | Laplace($\lambda$=0.05) | Unbounded | GenNorm($\beta$=3) | ✗ | —"— | —"— | $99.70 \pm 0.02$ |
| Box | Uniform | Normal($\sigma$=0.05) | Unbounded | GenNorm($\beta$=3) | ✗ | —"— | —"— | $99.69 \pm 0.02$ |
| Sphere | Normal($\sigma$=1) | Laplace($\lambda$=0.05) | Sphere | vMF($\kappa$=1) | ✗ | $63.37 \pm 2.41$ | $99.70 \pm 0.07$ | $99.02 \pm 0.01$ |
| Sphere | Normal($\sigma$=1) | Normal($\sigma$=0.05) | Sphere | vMF($\kappa$=1) | ✗ | —"— | —"— | $99.02 \pm 0.02$ |
| Unbounded | Laplace($\lambda$=1) | Normal($\sigma$=1) | Unbounded | Normal | ✗ | $62.49 \pm 1.65$ | $99.65 \pm 0.04$ | $98.13 \pm 0.14$ |
| Unbounded | Normal($\sigma$=1) | Normal($\sigma$=1) | Unbounded | Normal | ✗ | $63.57 \pm 2.30$ | $99.61 \pm 0.17$ | $98.76 \pm 0.03$ |

# Empirical Results

# Empirical Results

| Generative process $g$ | | | Model $f$ | | | MCC Score [%] | | |
|---|---|---|---|---|---|---|---|---|
| Space | $p(\cdot)$ | $p(\cdot|\cdot)$ | Space | $q_{\mathrm{h}}(\cdot|\cdot)$ | M. | Identity | Supervised | Unsupervised |
| Box | Uniform | Laplace($\lambda$=0.05) | Box | Laplace | ✓ | $46.55 \pm 1.34$ | $99.93 \pm 0.03$ | $98.62 \pm 0.05$ |
| Box | Uniform | GenNorm($\beta$=3; $\lambda$=0.05) | Box | GenNorm($\beta$=3) | ✓ | ——"—— | ——"—— | $99.90 \pm 0.06$ |
| Box | Uniform | Normal($\sigma$=0.05) | Box | Normal | ✗ | ——"—— | ——"—— | $99.77 \pm 0.01$ |
| Box | Uniform | Laplace($\lambda$=0.05) | Box | Normal | ✗ | ——"—— | ——"—— | $99.76 \pm 0.02$ |
| Box | Uniform | GenNorm($\beta$=3; $\lambda$=0.05) | Box | Laplace | ✗ | ——"—— | ——"—— | $98.80 \pm 0.02$ |
| Box | Uniform | Laplace($\lambda$=0.05) | Unbounded | Laplace | ✗ | ——"—— | $99.97 \pm 0.03$ | $98.57 \pm 0.02$ |
| Box | Uniform | GenNorm($\beta$=3; $\lambda$=0.05) | Unbounded | GenNorm($\beta$=3) | ✗ | ——"—— | ——"—— | $99.85 \pm 0.01$ |
| Box | Uniform | Normal($\sigma$=0.05) | Unbounded | Normal | ✗ | ——"—— | ——"—— | $58.26 \pm 3.00$ |
| Box | Uniform | Laplace($\lambda$=0.05) | Unbounded | Normal | ✗ | ——"—— | ——"—— | $59.67 \pm 2.33$ |
| Box | Uniform | Normal($\sigma$=0.05) | Unbounded | GenNorm($\beta$=3) | ✗ | ——"—— | ——"—— | $43.80 \pm 2.15$ |

# KITTI Masks

Table 3. **KITTI Masks**. Mean $\pm$ standard deviation over 10 random seeds. $\overline{\Delta t}$ indicates the average temporal distance of frames used.

|  | Model | Model Space | MCC [%] |
|---|---|---|---|
| $\overline{\Delta t} = 0.05s$ | SlowVAE | Unbounded | $66.1 \pm 4.5$ |
|  | Laplace | Unbounded | $77.1 \pm 1.0$ |
|  | Laplace | Box | $74.1 \pm 4.4$ |
|  | Normal | Unbounded | $58.3 \pm 5.4$ |
|  | Normal | Box | $59.9 \pm 5.5$ |
| $\overline{\Delta t} = 0.15s$ | SlowVAE | Unbounded | $79.6 \pm 5.8$ |
|  | Laplace | Unbounded | $79.4 \pm 1.9$ |
|  | Laplace | Box | $80.9 \pm 3.8$ |
|  | Normal | Unbounded | $60.2 \pm 8.7$ |
|  | Normal | Box | $68.4 \pm 6.7$ |

# 3DIdent

# 3DIdent

| Dataset | Model $f$ | | | Identity [%] | Unsupervised [%] | |
| $p(\cdot\|\cdot)$ | Space | $q_{\mathrm{h}}(\cdot\|\cdot)$ | M. | $R^2$ | $R^2$ | MCC |
|---|---|---|---|---|---|---|
| Normal | Box | Normal | ✓ | $5.25 \pm 1.20$ | $96.73 \pm 0.10$ | $98.31 \pm 0.04$ |
| Normal | Unbounded | Normal | ✗ | —‖— | $96.43 \pm 0.03$ | $54.94 \pm 0.02$ |
| Laplace | Box | Normal | ✗ | —‖— | $96.87 \pm 0.08$ | $98.38 \pm 0.03$ |
| Normal | Sphere | vMF | ✗ | —‖— | $65.74 \pm 0.01$ | $42.44 \pm 3.27$ |

# Ongoing Work

# Ongoing Work

1. Extend framework to object-centric methods
   a. MONet, IODINE, Slot Attention etc.

# Ongoing Work

1. Extend framework to object-centric methods
   a. MONet, IODINE, Slot Attention etc.
2. Extend framework to data augmentations
   a. Content & Style Disambiguation
   b. Invariant factors == delta conditional

# Self-supervised learning with data augmentations provably isolates content from style

with **Julius von Kügelgen***, **Yash Sharma***, **Luigi Gresele***, Wieland Brendel, Michel Besserve, Francesco Locatello

Formalise generation $x = f(z)$ and augmentation $\tilde{x} = f(\tilde{z})$ processes as latent variable model with a content-style partition $z = (c, s)$:

- *invariant content $c$:*   always shared between pairs $(x, \tilde{x})$ of views;
- *varying style $s$:*   may change across pairs $(x, \tilde{x})$ of views.

Allow causal dependence of style on content (*Causal3DIdent* dataset):

augmented view $\tilde{x}$ = counterfactual under soft style intervention on $x$.

**Theory:** Can identify* invariant content partition in generative and discriminative learning with entropy maximisation (e.g., SimCLR).
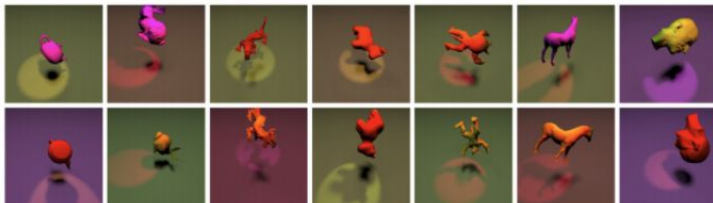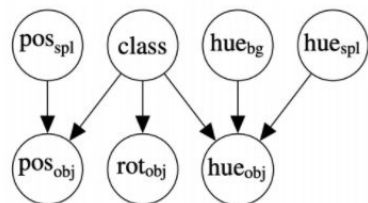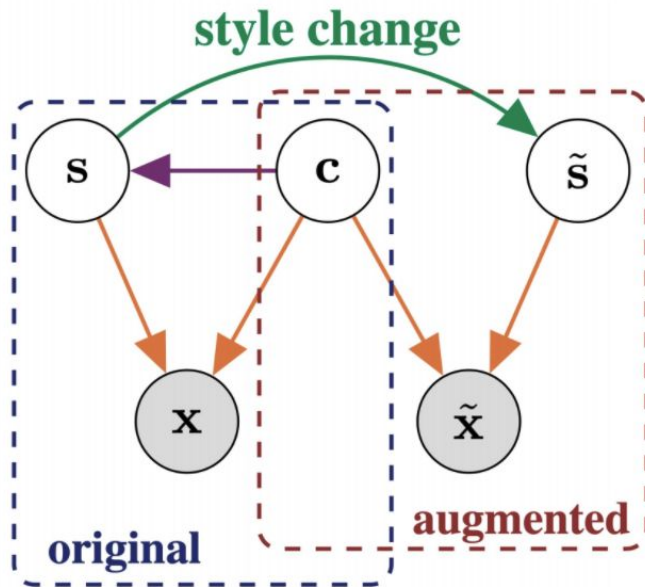


Figure 2: *(Left)* Causal graph for the *Causal3DIdent* dataset. *(Right)* Two samples from each object class.

*up to invertible transformation

# Ongoing Work

1.  Extend framework to object-centric methods

    a.  MONet, IODINE, Slot Attention etc.

2.  Extend framework to data augmentations

    a.  Content & Style Disambiguation

    b.  Invariant factors == delta conditional

3.  Extend framework for causal discovery

    a.  Robustness in downstream tasks?

# Thank you for your attention!



David Klindt

Lukas Schott

Roland Zimmermann

Steffen Schneider

Ivan Ustyuzhaninov

Wieland Brendel

Matthias Bethge

Dylan Paiton

If you are interested in this research, feel free to reach out!
**Mail**: ysharma1126@gmail.com
**Twitter**: @yash_j_sharma