



[bit.ly/3gu32tw](https://bit.ly/3gu32tw)

# Contrastive Learning Inverts the Data Generating Process

Roland S. Zimmermann\*, Yash Sharma\*, Steffen Schneider\*,  
Matthias Bethge<sup>†</sup>, Wieland Brendel<sup>†</sup>



EBERHARD KARLS  
UNIVERSITÄT  
TÜBINGEN



imprs-is



Bundesministerium  
für Bildung  
und Forschung



DFG

\* Equal contribution <sup>†</sup> Joint supervision

Contrastive Learning



Disentanglement  
Nonlinear ICA

**Contrastive Learning**  
empirically successful



**Disentanglement**  
**Nonlinear ICA**

Contrastive Learning



Disentanglement  
Nonlinear ICA

Image Space

Representation Space



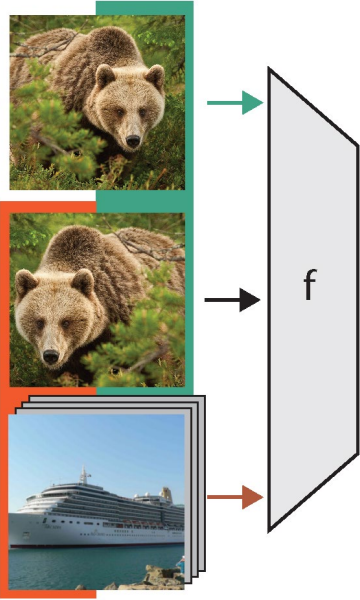
Contrastive Learning

Our Theory

Disentanglement  
Nonlinear ICA

Image Space

Representation Space



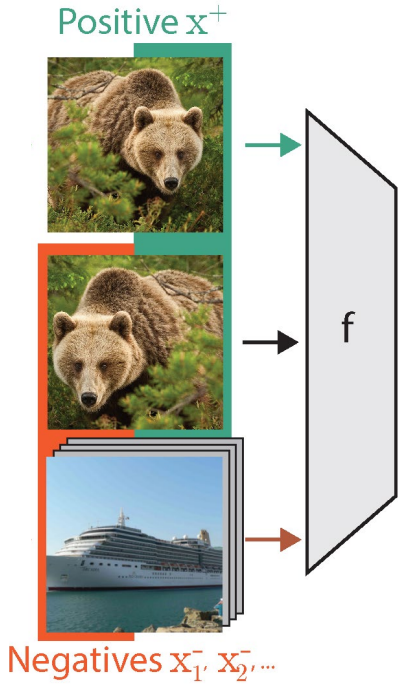
Contrastive Learning

Our Theory

Disentanglement  
Nonlinear ICA

Image Space

Representation Space



Contrastive Learning

Our Theory

Disentanglement  
Nonlinear ICA

Image Space

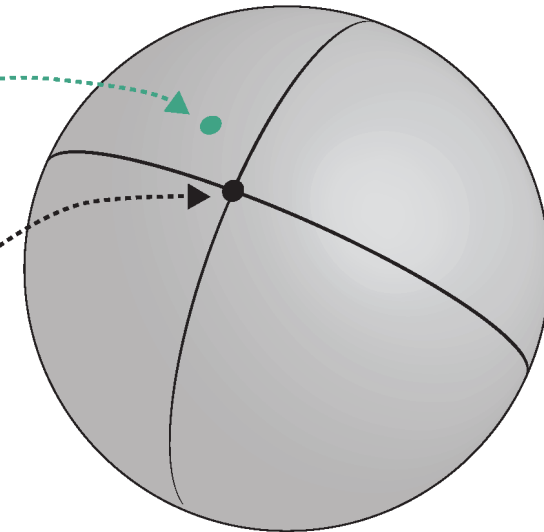
Representation Space

Positive  $x^+$



Negatives  $x_1^-, x_2^-, \dots$

f



Contrastive Learning

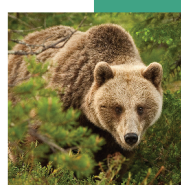
Our Theory

Disentanglement  
Nonlinear ICA

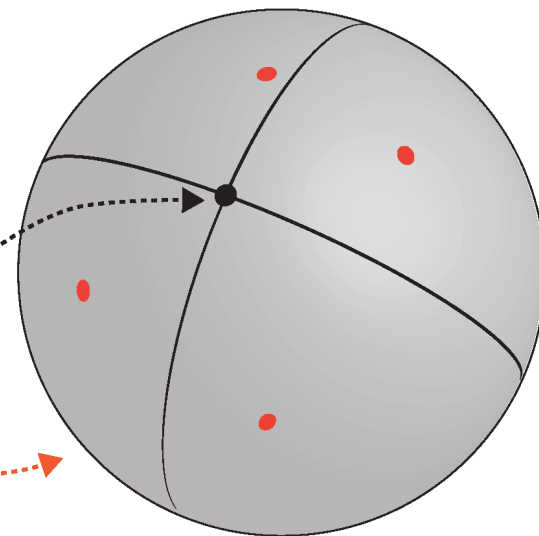
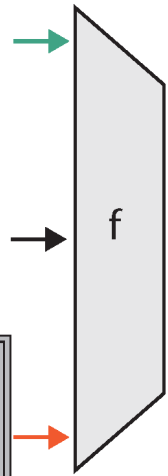
Image Space

Representation Space

Positive  $x^+$



Negatives  $x_1^-, x_2^-, \dots$





Contrastive Learning

Our Theory

Disentanglement  
Nonlinear ICA

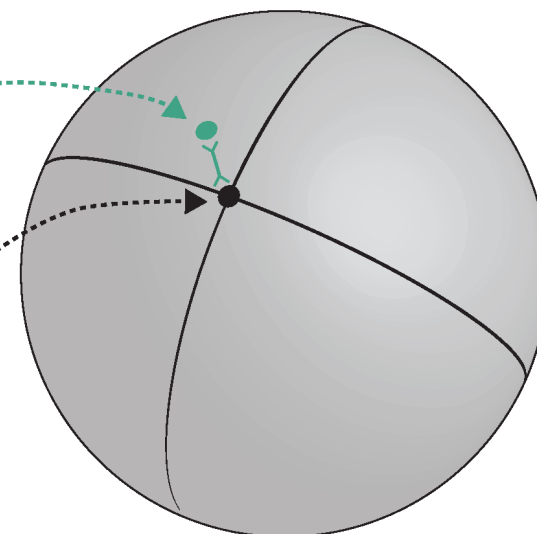
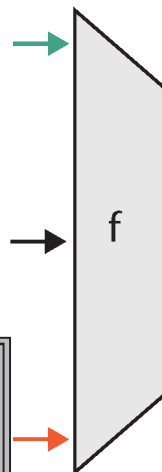
Image Space

Representation Space

Positive  $x^+$



Negatives  $x_1^-, x_2^-, \dots$



Learned by

$$\mathcal{L} = - \underset{\text{attract}}{f(x)^T f(x^+)}$$

Contrastive Learning

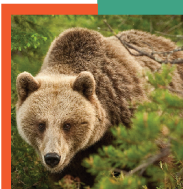
Our Theory

Disentanglement  
Nonlinear ICA

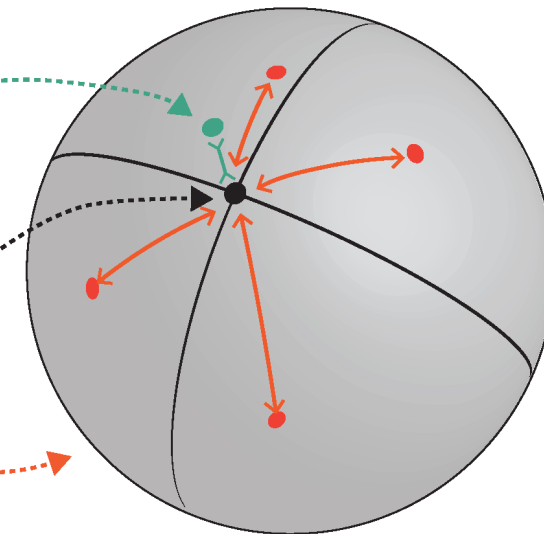
Image Space

Representation Space

Positive  $x^+$



Negatives  $x_1^-, x_2^-, \dots$

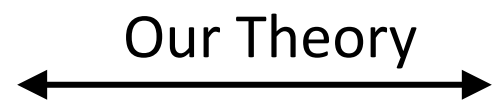


Learned by

$\mathcal{L} =$

$$\log \sum_{\{x^+, x_1^-, \dots\}} \exp(\underbrace{f(x)^T f(x')}_{\text{repel}})$$

Contrastive Learning

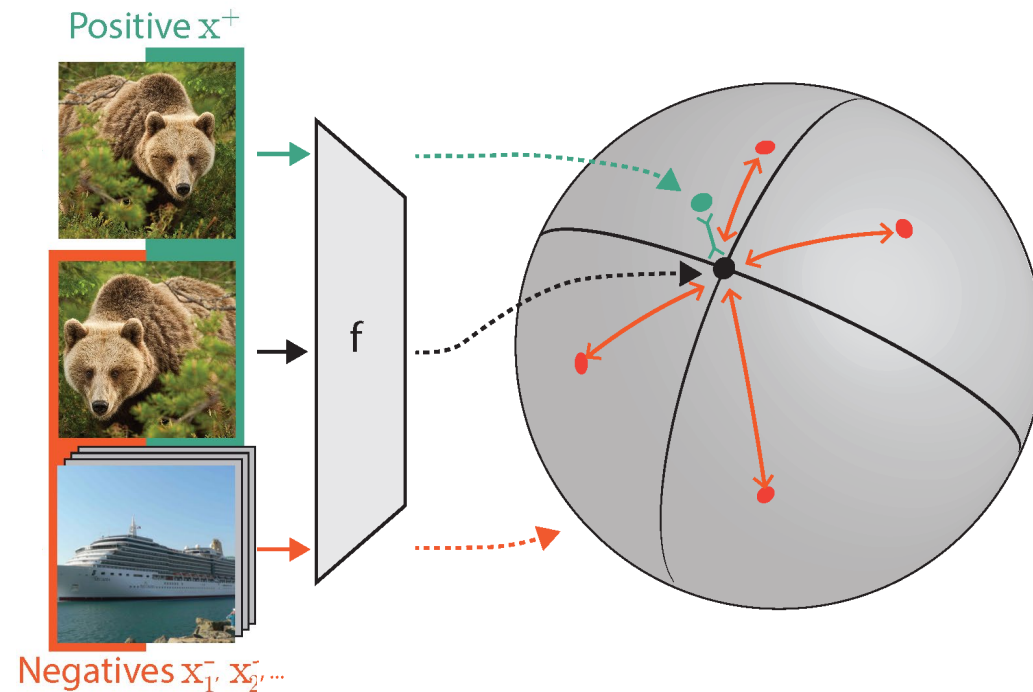


Disentanglement  
Nonlinear ICA

Contrastive Learning

Our Theory

Disentanglement  
Nonlinear ICA



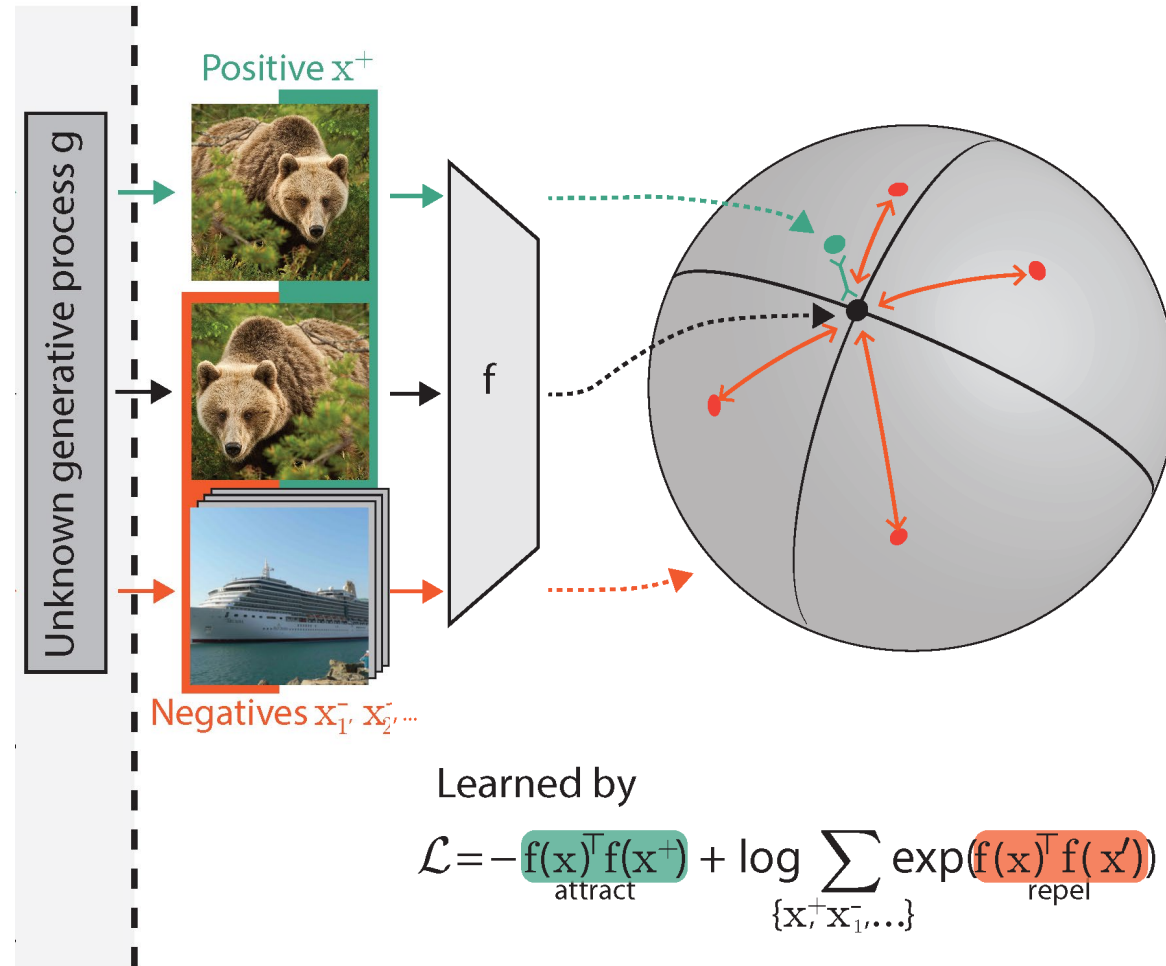
Learned by

$$\mathcal{L} = -\underset{\text{attract}}{f(x)^T f(x^+)} + \log \sum_{\{x^+, x_1^-, \dots\}} \exp(\underset{\text{repel}}{f(x)^T f(x')})$$

Contrastive Learning

Our Theory

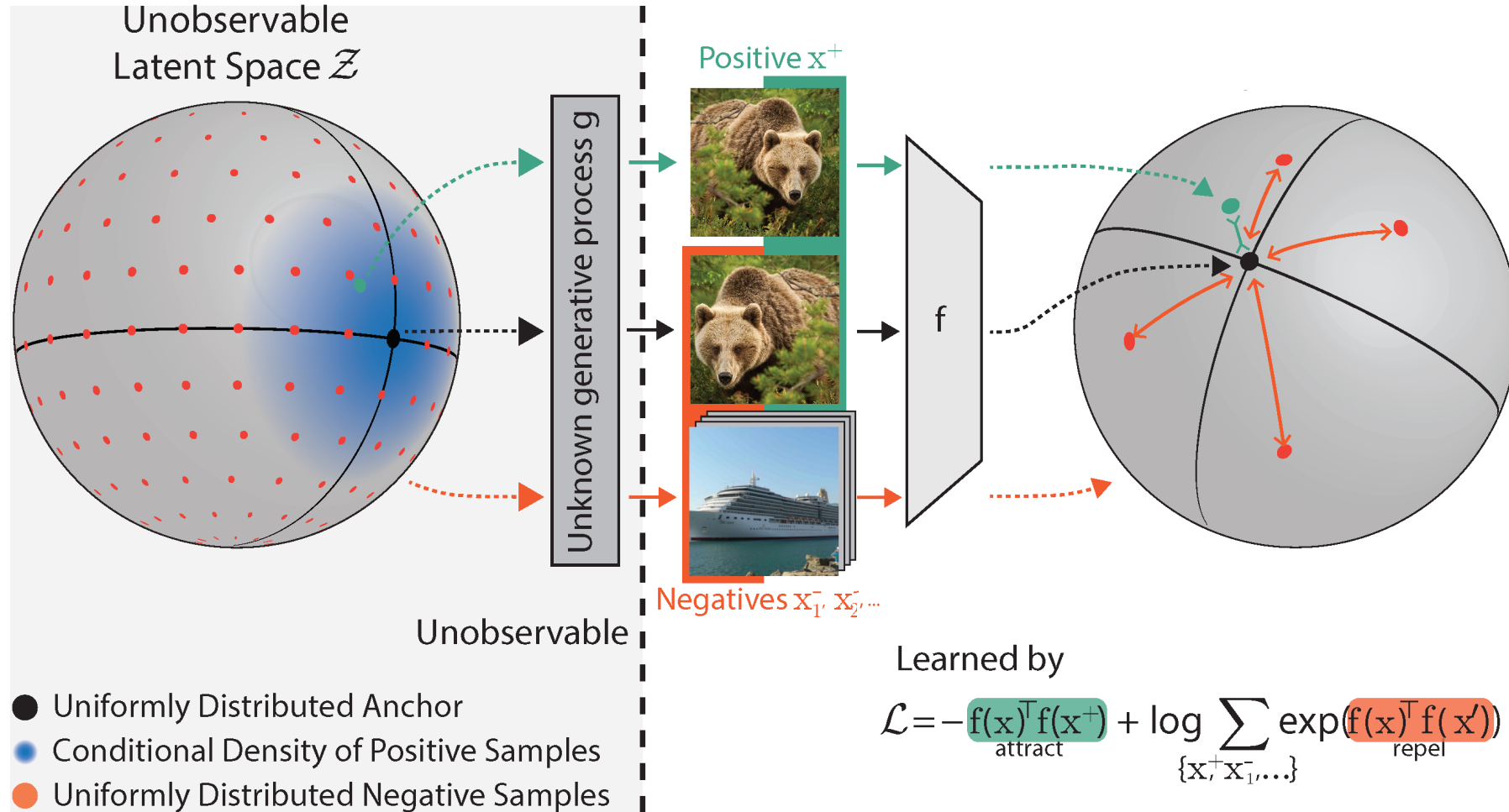
Disentanglement  
Nonlinear ICA



Contrastive Learning

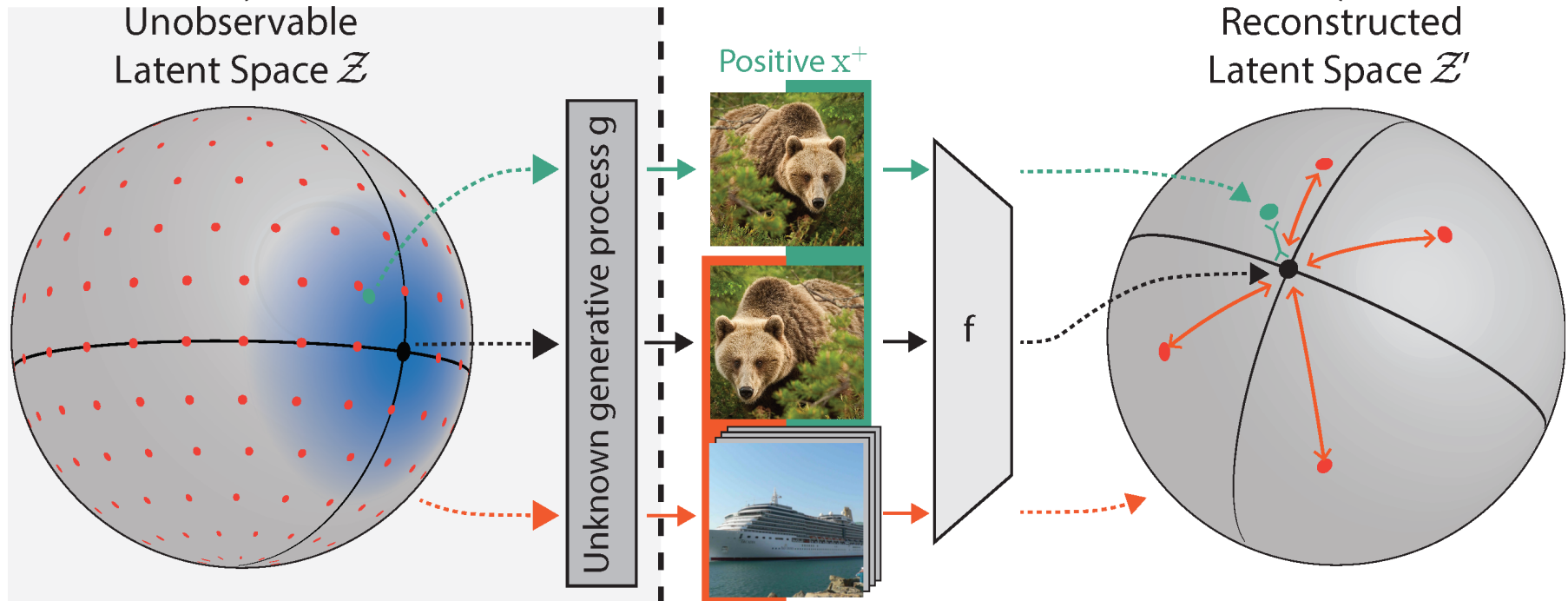
Our Theory

Disentanglement  
Nonlinear ICA



Contrastive Learning ← Our Theory → Disentanglement  
Nonlinear ICA

Goal: Recover structure



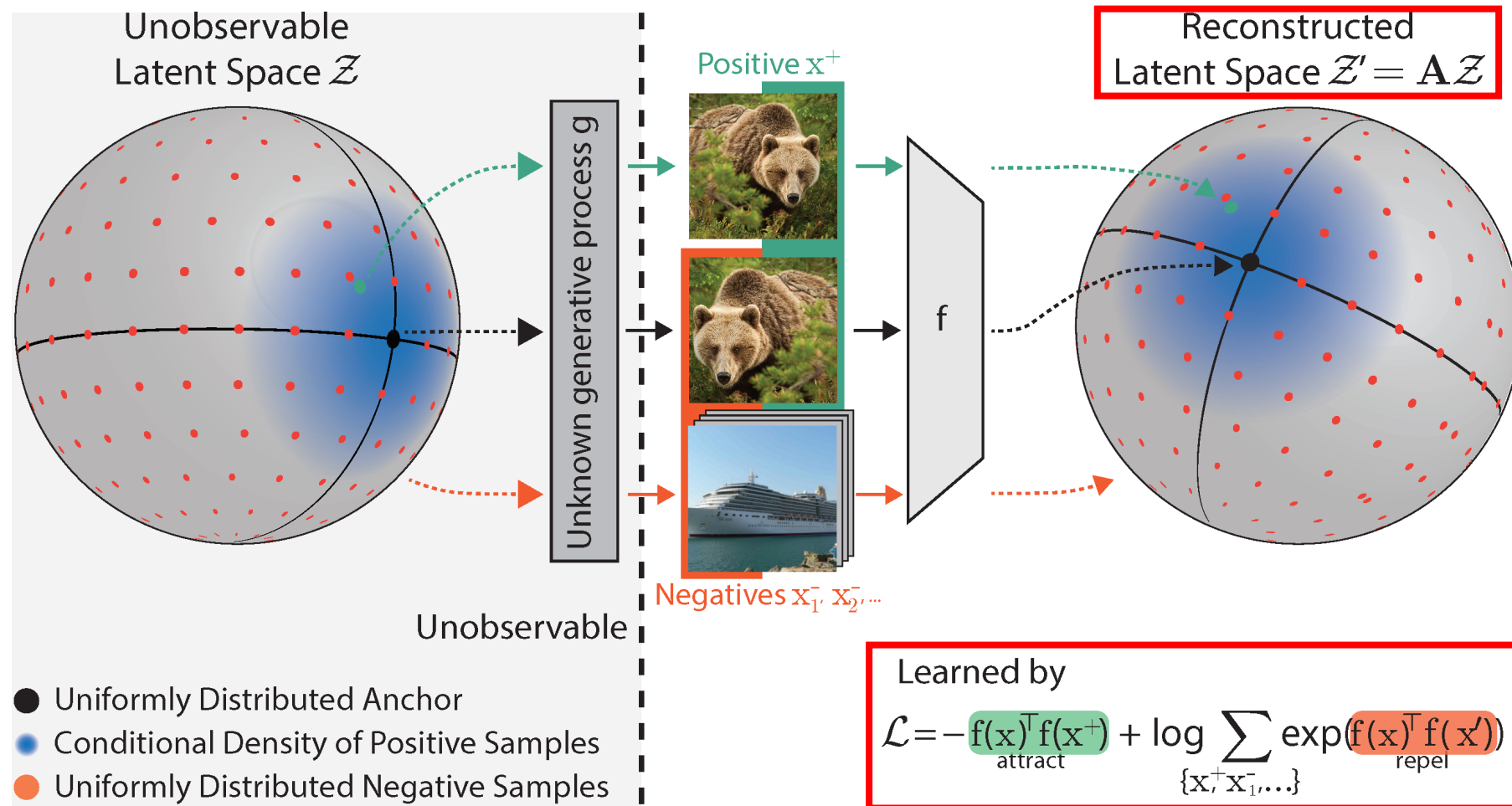
- Uniformly Distributed Anchor
- Conditional Density of Positive Samples
- Uniformly Distributed Negative Samples

Learned by

$$\mathcal{L} = -\underset{\text{attract}}{f(x)^T f(x^+)} + \log \sum_{\{x^+, x_1^-, \dots\}} \exp(\underset{\text{repel}}{f(x)^T f(x')})$$

We prove

# Contrastive Learning Inverts Data Generative Process





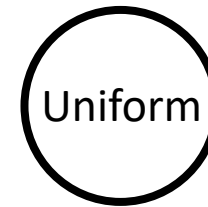
We prove

# Contrastive Learning Inverts Data Generative Process

Learned by

$$\mathcal{L} = -\underbrace{f(\mathbf{x})^\top f(\mathbf{x}^+)}_{\text{attract}} + \log \sum_{\{\mathbf{x}^+, \mathbf{x}_1^-, \dots\}} \exp(\underbrace{f(\mathbf{x})^\top f(\mathbf{x}')}_{\text{repel}})$$

Assumptions



+

von Mises-Fisher distribution

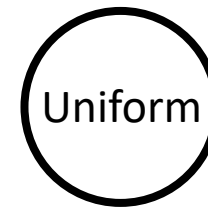
We prove

# Contrastive Learning Inverts Data Generative Process

Learned by

$$\mathcal{L} = -\underbrace{f(\mathbf{x})^\top f(\mathbf{x}^+)}_{\text{attract}} + \log \sum_{\{\mathbf{x}^+, \mathbf{x}_1^-, \dots\}} \exp(\underbrace{f(\mathbf{x})^\top f(\mathbf{x}')}_{\text{repel}})$$

Assumptions implicitly encoded



+

von Mises-Fisher distribution

We prove

# Contrastive Learning Inverts Data Generative Process

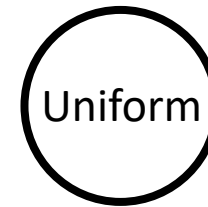
Learned by

$$\mathcal{L} = - \underbrace{f(\mathbf{x})^\top f(\mathbf{x}^+)}_{\text{attract}} + \log \sum_{\{\mathbf{x}^+, \mathbf{x}_1^-, \dots\}} \exp(\underbrace{f(\mathbf{x})^\top f(\mathbf{x}')}_{\text{repel}})$$

Learned by

$$\mathcal{L} = \underbrace{\|f(\mathbf{x}) - f(\mathbf{x}^+)\|_2^2}_{\text{attract}} + \log \sum_{\{\mathbf{x}^+, \mathbf{x}_1^-, \dots\}} \exp(-\underbrace{\|f(\mathbf{x}) - f(\mathbf{x}')\|_2^2}_{\text{repel}})$$

Assumptions encoded



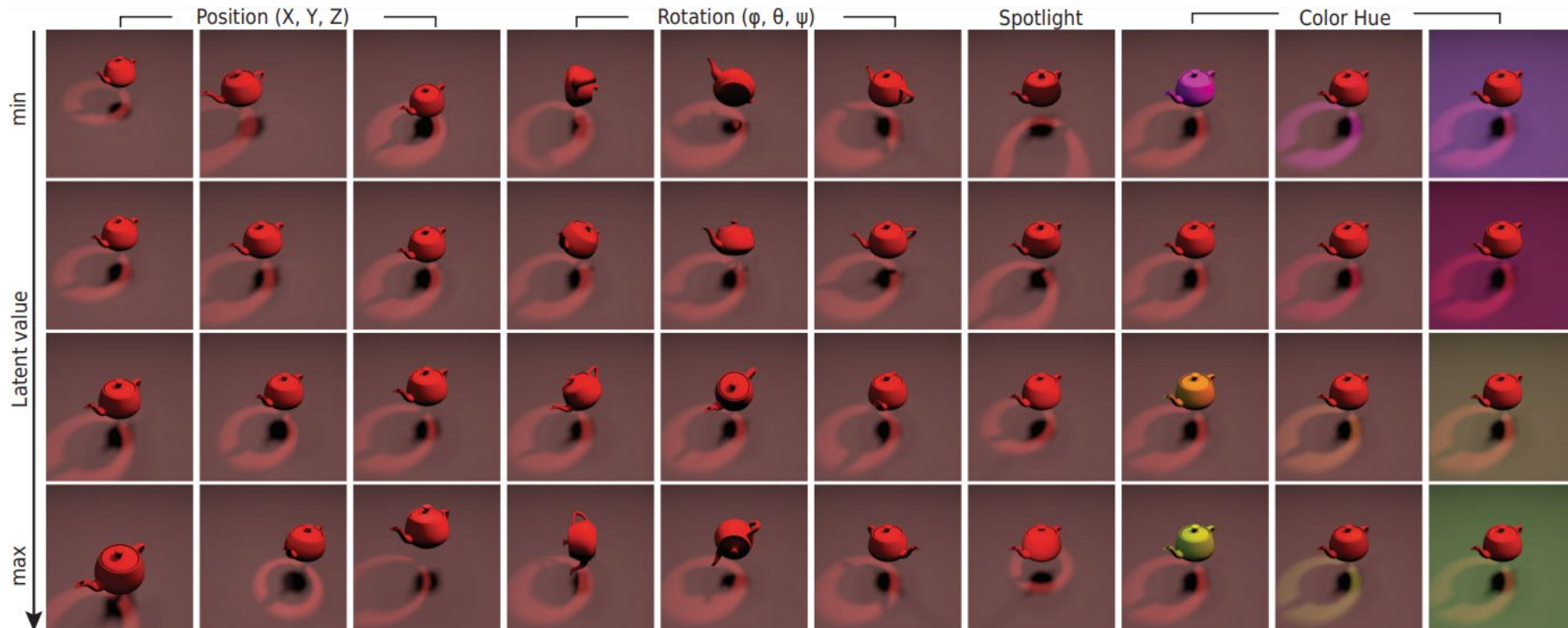
+ von Mises-Fisher distribution



+ Normal distribution

We show

# Successful Disentanglement on Visually Complex Images



3DIdent Dataset:

- High resolution
- Hallmarks of natural environments (shadows, lighting conditions, 3D object)

## Summary

**Prove that contrastive learning (InfoNCE) inverts specified data generating processes**

**Works on complex visual data and is robust to mismatches**

## Outlook

**Theory provides insights in constructing more effective contrastive losses**

Poster, Paper & Code



[bit.ly/3r4CyIX](https://bit.ly/3r4CyIX)