

Formalisms for Understanding Progress in Representation Learning

Yash Sharma^{1,2,3}

¹ IMPRS for Intelligent Systems, ³ University of Tübingen, ³ MPI for Intelligent Systems Tübingen,

imprs-is

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

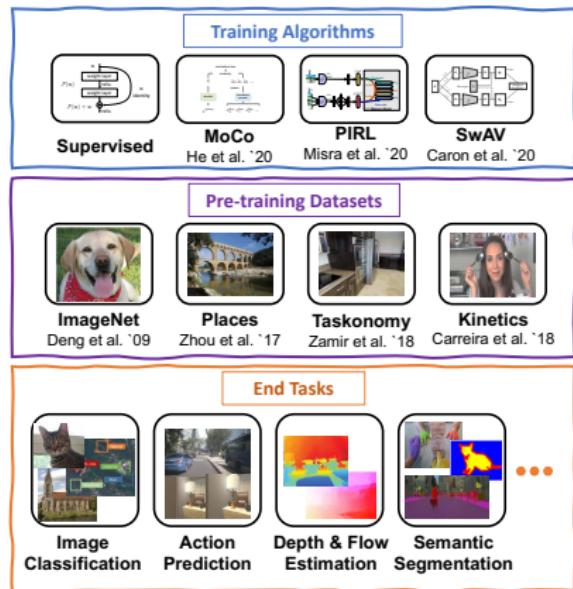
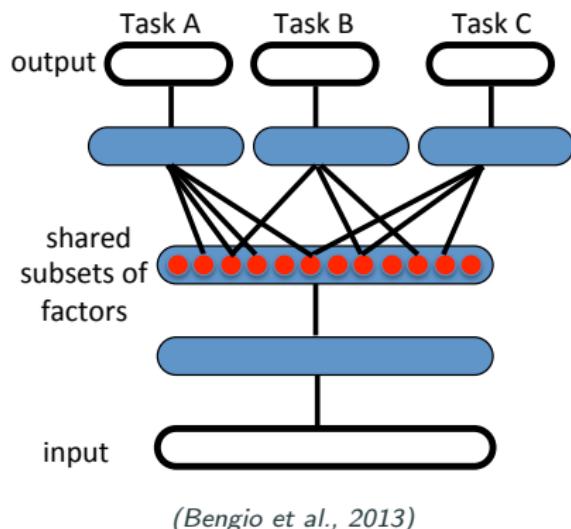


MAX PLANCK INSTITUTE
FOR INTELLIGENT SYSTEMS



Representation Learning

Learning representations of the data that facilitate the extraction of useful information for solving downstream tasks.

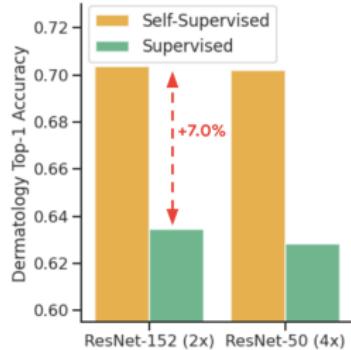
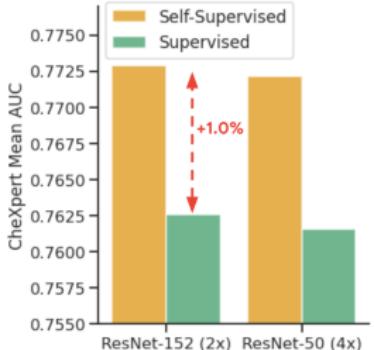


Recent Progress

		COCO keypoint detection		
pre-train		Ap ^{kp}	Ap ^{kp} ₅₀	Ap ^{kp} ₇₅
random init.		65.9	86.5	71.7
super. IN-1M		65.8	86.9	71.9
MoCo IN-1M		66.8 (+1.0)	87.4 (+0.5)	72.5 (+0.6)
MoCo IG-1B		66.9 (+1.1)	87.8 (+0.9)	73.0 (+1.1)
		COCO dense pose estimation		
pre-train		AP ^{dip}	AP ^{dip} ₅₀	AP ^{dip} ₇₅
random init.		39.4	78.5	35.1
super. IN-1M		48.3	85.6	50.6
MoCo IN-1M		50.1 (+1.8)	86.8 (+1.2)	53.9 (+3.3)
MoCo IG-1B		50.6 (+2.3)	87.0 (+1.4)	54.3 (+3.7)
		LVIS v0.5 instance segmentation		
pre-train		AP ^{mk}	AP ^{mk} ₅₀	AP ^{mk} ₇₅
random init.		22.5	34.8	23.8
super. IN-1M [†]		24.4	37.8	25.8
MoCo IN-1M		24.1 (-0.3)	37.4 (-0.4)	25.5 (-0.3)
MoCo IG-1B		24.9 (+0.5)	38.2 (+0.4)	26.4 (+0.6)
		Cityscapes instance seg.		Semantic seg. (mIoU)
pre-train		AP ^{mk}	AP ^{mk} ₅₀	
random init.		25.4	51.1	65.3
super. IN-1M		32.9	59.6	74.6
MoCo IN-1M		32.3 (-0.6)	59.3 (-0.3)	75.3 (+0.7)
MoCo IG-1B		32.9 (-0.0)	60.3 (+0.7)	75.5 (+0.9)
				39.5
				74.4
				72.5 (-1.9)
				73.6 (-0.8)
		Cityscapes	VOC	

Table 6. **MoCo** vs. ImageNet supervised pre-training, fine-tuned on various tasks. For each task, the same architecture and

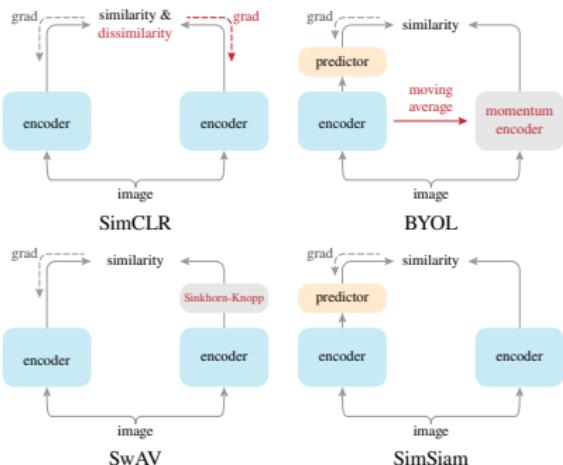
(He et al., 2020)



(Azizi et al., 2021)

Architecture Comparison

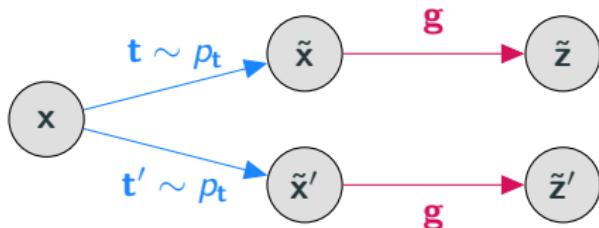
All methods maximize similarity between augmentations, subject to different conditions for avoiding collapse.



(Chen & He, 2020)

(Chen et al., 2020)

Preliminaries: SSL with Data Augmentation



1. Specify a set \mathcal{T} of transformations with distribution p_t .
2. For each observation x , sample pair of transformations $t, t' \sim p_t$ and apply to x to form positive pair $(\tilde{x}, \tilde{x}') = (t(x), t'(x))$.
3. Learn encoder g s.t. representations $(\tilde{z}, \tilde{z}') = (g(\tilde{x}), g'(\tilde{x}'))$ are similar while avoiding a collapsed (trivial) representation.
 - Contrastive learning (CL): e.g., SIMCLR (*Chen et al., 2020*)

$$\mathcal{L}_{\text{InfoNCE}}(g) = \mathbb{E}_{\{\mathbf{x}_i\}_{i=1}^K \sim p_x} \left[- \sum_{i=1}^K \log \frac{\exp\{\text{sim}(\tilde{\mathbf{z}}_i, \tilde{\mathbf{z}}'_i)\}}{\sum_{j=1}^K \exp\{\text{sim}(\tilde{\mathbf{z}}_i, \tilde{\mathbf{z}}'_j)\}} \right]$$

(*Gutmann & Hyvärinen, 2010; van den Oord et al., 2018*)

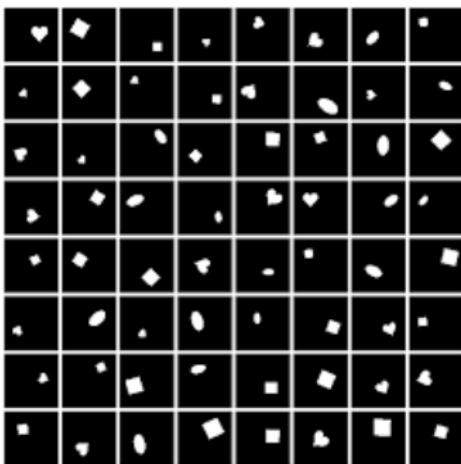
Problem Setting I: Latent Variable Model

Formalise problem setting from latent variable model perspective:

Data generation: $\mathbf{z} \sim p_{\mathbf{z}}, \quad \mathbf{x} = \mathbf{f}(\mathbf{z}),$

Augmentation: $\tilde{\mathbf{z}} \sim p_{\tilde{\mathbf{z}}|\mathbf{z}}, \quad \tilde{\mathbf{x}} = \mathbf{f}(\tilde{\mathbf{z}}),$

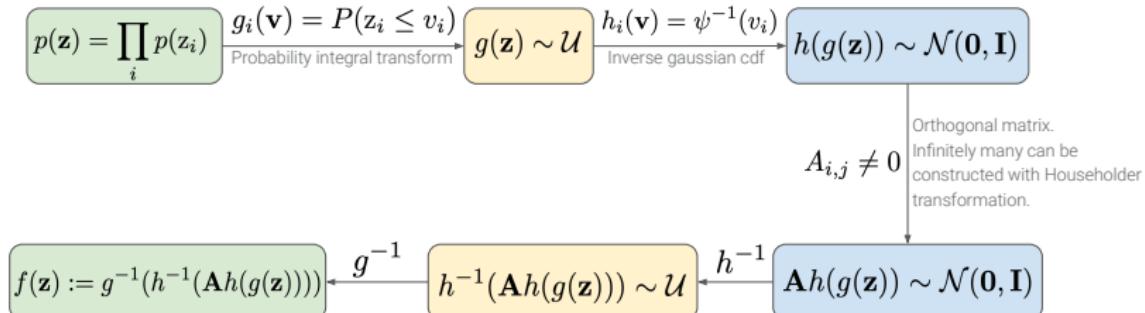
with invertible decoder or mixing function $\mathbf{f} : \mathcal{Z} \rightarrow \mathcal{X}.$



(Matthey et al., 2017)

Identifiability: Impossible?

Identifiability: recover properties of the true generative process.



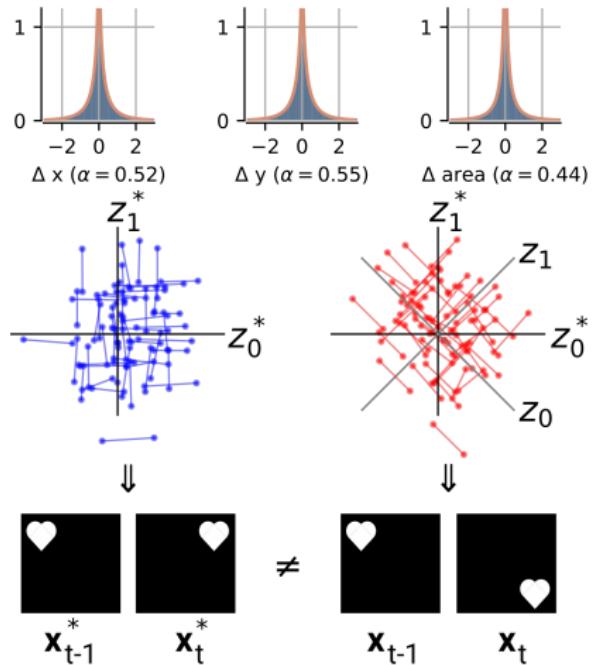
Issue: $P(\mathbf{z} \leq \mathbf{u}) = P(f(\mathbf{z}) \leq \mathbf{u}), \quad \frac{\partial f_i(\mathbf{z})}{\partial \mathbf{z}_j} \neq 0$

(Locatello, 2020)

Earlier work proved similar results (Hyvarinen & Pajunen, 1999).

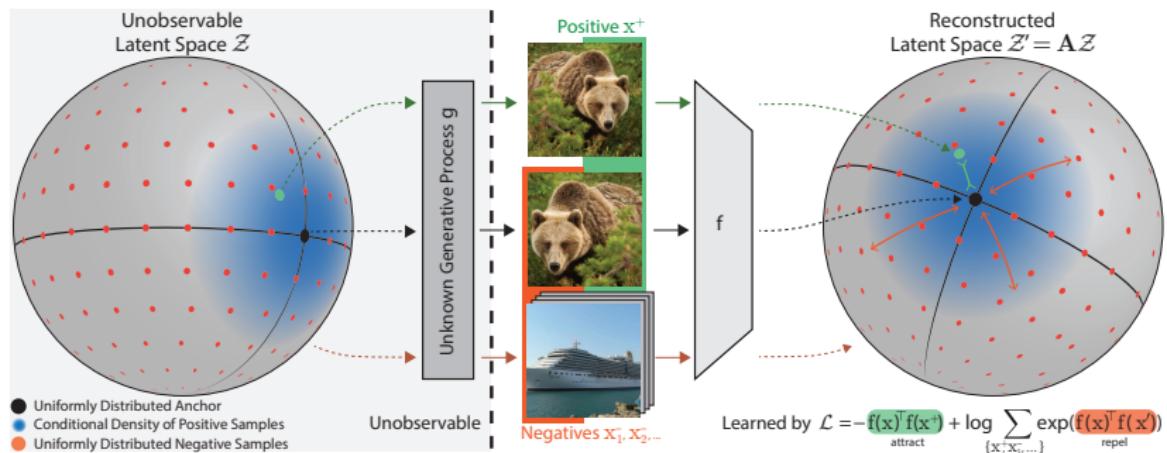
Identifiability: Additional Assumptions

Learn a generative model for which both the marginal p_z **and** the conditional $p_{\bar{z}|z}$ distributions match ground-truth.



Identifiability: Additional Assumptions

Can we obtain a similar result if we solely learn an encoder via contrastive learning?



(Zimmermann*, Sharma*, Schneider* et al., 2021)

Identifiability: Additional Assumptions

Can we obtain a similar result if we solely learn an encoder via contrastive learning?

⇒ Yes, if the distribution encoded in $\mathcal{L}_{\text{InfoNCE}}(\mathbf{g})$ matches that of ground-truth

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 $M \rightarrow \infty$, the (normalized) contrastive loss converges to

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

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

$$q_h(\bar{\mathbf{z}}|\mathbf{z}) = C_h(\mathbf{z})^{-1} e^{h(\bar{\mathbf{z}})^T h(\mathbf{z})/\tau}$$

with $C_h(\mathbf{z}) := \int e^{h(\bar{\mathbf{z}})^T h(\mathbf{z})/\tau} d\bar{\mathbf{z}}$. (15)

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 restriction of the mixing function g to \mathcal{Z} be injective and h be differentiable in a vicinity of \mathcal{Z} . If the assumed form of q_h , as defined above, matches that of p , and if f is differentiable and minimizes the CL loss as defined in Eq. (1), then for fixed $\tau > 0$ and $M \rightarrow \infty$, $h = f \circ g$ is linear, i.e., f recovers the latent sources up to an orthogonal linear transformation and a constant scaling factor.

Theorem 6. Let \mathcal{Z} be a convex body in \mathbb{R}^N , $h : \mathcal{Z} \rightarrow \mathcal{Z}$, and δ be an L^α metric or semi-metric (cf. Lemma 1 in Appx. A.2.4) for $\alpha \geq 1, \alpha \neq 2$. Further, let the ground-truth marginal distribution be uniform and the conditional distribution be as Eq. (5), and let the mixing function g be differentiable and invertible. If the assumed form of $q_h(\cdot|\mathbf{z})$ matches that of $p(\cdot|\mathbf{z})$, i.e., both use the same metric δ up to a constant scaling factor, and if f is differentiable and minimizes the $\mathcal{L}_{\delta, \text{contr}}$ objective in Eq. (6) for $M \rightarrow \infty$, we find that $h = f \circ g$ is a composition of input independent permutations, sign flips and rescaling.

How well do our assumptions match practice? (e.g., SimCLR)

Prior work on identifiable representation learning and multi-view nonlinear ICA (e.g., Hyvarinen & Morioka, 2016, 2017; Gresele et al., 2019; Locatello et al., 2020; Klindt et al., 2020; Zimmermann et al., 2021)

1. **Assume factorized p_z , i.e. independent latents**
(unknown statistical/causal dependence between latents of interest)
2. **Assume that all latents z_i change across views (x, \tilde{x})**
(augmentations chosen to leave certain aspects invariant)

Can we perform identifiability analysis without these assumptions?

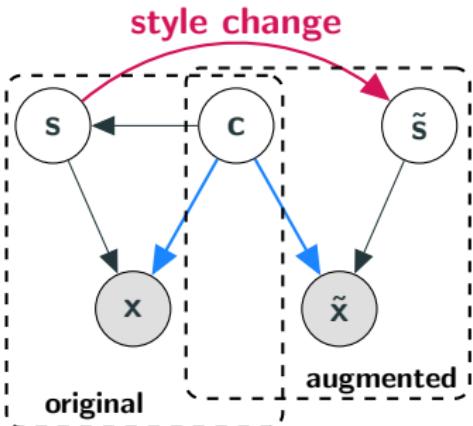
⇒ **Yes!** (*von Kügelgen*, Sharma*, Gresele* et al., 2021*)

Problem Setting II: Content-Style Partition

Partition the latent representation \mathbf{z} into two disjoint blocks:

$$\mathbf{z} = (\mathbf{c}, \mathbf{s}), \quad \mathbf{c} = \mathbf{z}_{1:n_c}, \quad \mathbf{s} = \mathbf{z}_{n_c+1:n}$$

- **invariant content block \mathbf{c} is always shared across $(\mathbf{x}, \tilde{\mathbf{x}})$**
- **varying style block \mathbf{s} may change across $(\mathbf{x}, \tilde{\mathbf{x}})$**



Note:

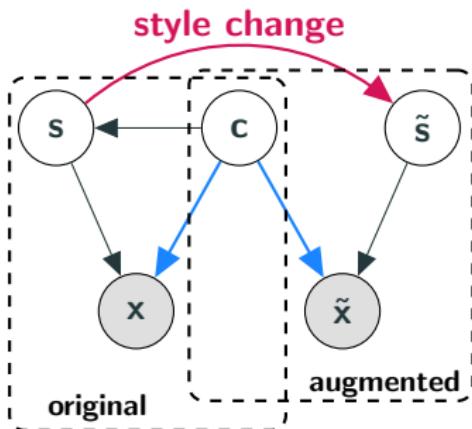
In practice, the content-style partition is implicitly defined by the set \mathcal{T} of transformations $\mathbf{t} \sim p_{\mathbf{t}}$ used for data augmentation.

Problem Setting III: Main Assumptions

Assumption 1: Content Invariance

The conditional $p_{\tilde{z}|z}$ takes the form

$$p_{\tilde{z}|z}(\tilde{z}|z) = \delta(\tilde{\mathbf{c}} - \mathbf{c}) p_{\tilde{s}|s}(\tilde{s}|s).$$



Assumption 2: Random Style Subsets Change

Let p_A be a distribution over style subsets $A \subseteq \{1, \dots, n_s\}$.

Then, the style conditional $p_{\tilde{s}|s}$ is given by

$$A \sim p_A, \quad p_{\tilde{s}|s,A}(\tilde{s}|s, A) = \delta(\tilde{s}_{A^c} - s_{A^c}) p_{\tilde{s}_A|s_A}(\tilde{s}_A|s_A),$$

Theorem

Theorem: Identifying Content with a Non-Invertible Encoder

Assume the same data generating process and conditions.

Let $\mathbf{g} : \mathcal{X} \rightarrow (0, 1)^{n_c}$ be any smooth function that minimises

$$\mathcal{L}_{\text{AlignMaxEnt}}(\mathbf{g}) := \mathbb{E}_{(\mathbf{x}, \tilde{\mathbf{x}}) \sim p_{\mathbf{x}, \tilde{\mathbf{x}}}} \left[(\mathbf{g}(\mathbf{x}) - \mathbf{g}(\tilde{\mathbf{x}}))^2 \right] - H(\mathbf{g}(\mathbf{x}))$$

where $H(\cdot)$ denotes differential entropy.

Then \mathbf{g} block-identifies the true content variables.

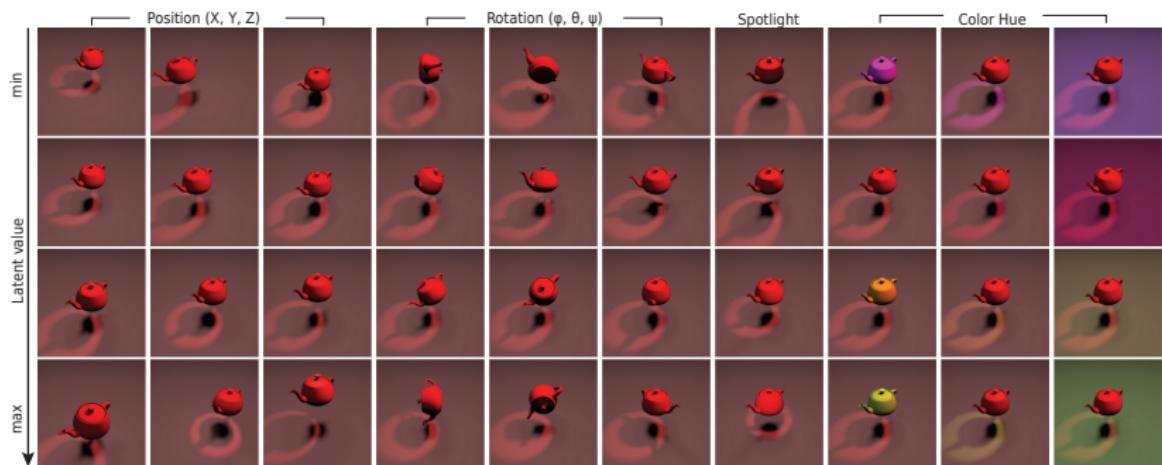
Definition: Block-Identifiability

The true content $\mathbf{c} = \mathbf{f}^{-1}(\mathbf{x})_{1:n_c}$ is *block-identified* by $\mathbf{g} : \mathcal{X} \rightarrow \mathcal{Z}$ if the inferred content $\hat{\mathbf{c}} = \mathbf{g}(\mathbf{x})_{1:n_c}$ contains *all* and *only* information about \mathbf{c} , i.e., if \exists *invertible* $\mathbf{h} : \mathbb{R}^{n_c} \rightarrow \mathbb{R}^{n_c}$ s.t. $\hat{\mathbf{c}} = \mathbf{h}(\mathbf{c})$.

Main Result

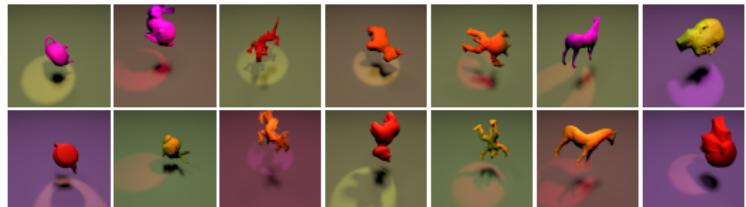
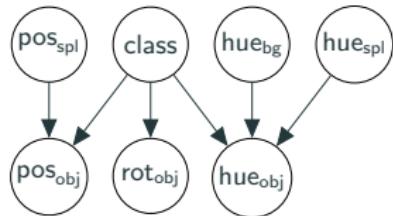
CL with InfoNCE (e.g., SimCLR) asymptotically isolates content.

3DIdent (Zimmermann et al., 2021)

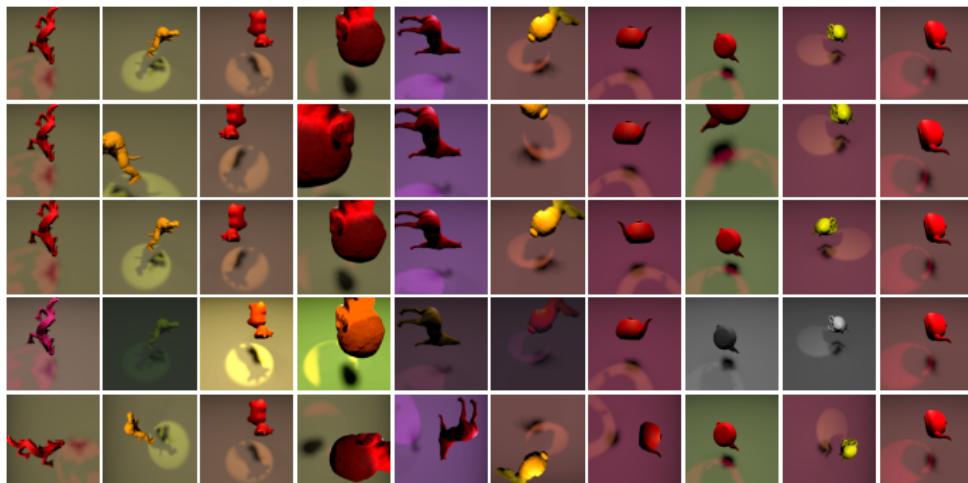


Causal3DIdent

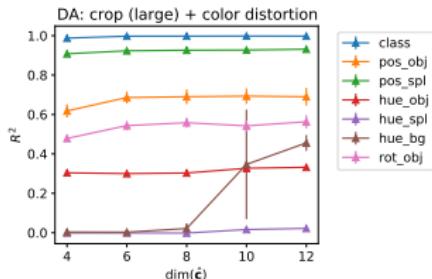
Dataset:



Augmentations:

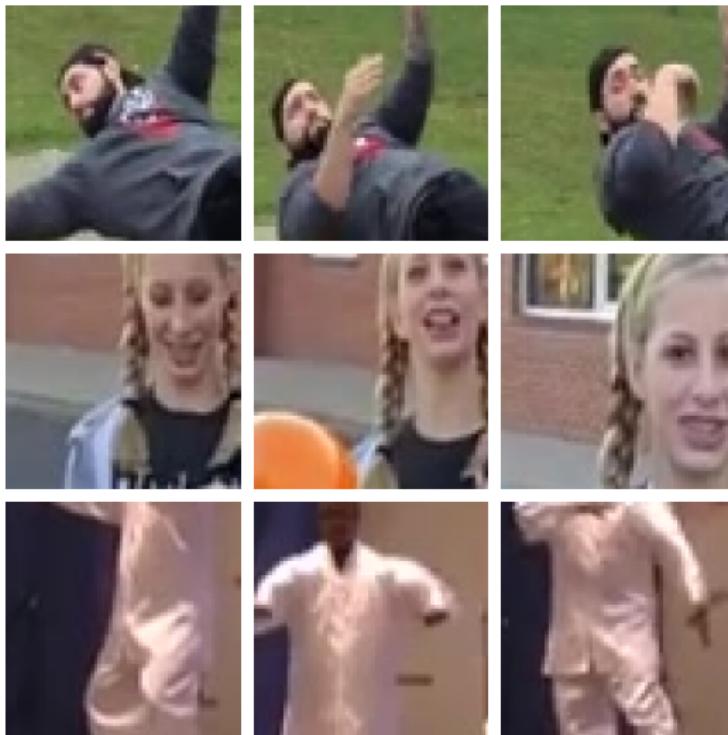


SimCLR Results



Views generated by	Class	Positions		Hues			Rotations
		object	spotlight	object	spotlight	background	
DA: colour distortion	0.42 ± 0.01	0.61 ± 0.10	0.17 ± 0.00	0.10 ± 0.01	0.01 ± 0.00	0.01 ± 0.00	0.33 ± 0.02
LT: change hues	1.00 ± 0.00	0.59 ± 0.33	0.91 ± 0.00	0.30 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.30 ± 0.01
DA: crop (large)	0.28 ± 0.04	0.09 ± 0.08	0.21 ± 0.13	0.87 ± 0.00	0.09 ± 0.02	1.00 ± 0.00	0.02 ± 0.02
DA: crop (small)	0.14 ± 0.00	0.00 ± 0.01	0.00 ± 0.01	0.00 ± 0.00	0.00 ± 0.00	1.00 ± 0.00	0.00 ± 0.00
LT: change positions	1.00 ± 0.00	0.16 ± 0.23	0.00 ± 0.01	0.46 ± 0.02	0.00 ± 0.00	0.97 ± 0.00	0.29 ± 0.01
DA: crop (large) + colour distortion	0.97 ± 0.00	0.59 ± 0.07	0.59 ± 0.05	0.28 ± 0.00	0.01 ± 0.01	0.01 ± 0.00	0.74 ± 0.03
DA: crop (small) + colour distortion	1.00 ± 0.00	0.69 ± 0.04	0.93 ± 0.00	0.30 ± 0.01	0.00 ± 0.00	0.02 ± 0.03	0.56 ± 0.03
LT: change positions + hues	1.00 ± 0.00	0.22 ± 0.22	0.07 ± 0.08	0.32 ± 0.02	0.00 ± 0.01	0.02 ± 0.03	0.34 ± 0.06
DA: rotation	0.33 ± 0.06	0.17 ± 0.09	0.23 ± 0.12	0.83 ± 0.01	0.30 ± 0.12	0.99 ± 0.00	0.05 ± 0.03
LT: change rotations	1.00 ± 0.00	0.53 ± 0.33	0.90 ± 0.00	0.41 ± 0.00	0.00 ± 0.00	0.97 ± 0.00	0.28 ± 0.00
DA: rotation + colour distortion	0.59 ± 0.01	0.58 ± 0.06	0.21 ± 0.01	0.12 ± 0.02	0.01 ± 0.00	0.01 ± 0.00	0.33 ± 0.04
LT: change rotations + hues	1.00 ± 0.00	0.57 ± 0.34	0.91 ± 0.00	0.30 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.28 ± 0.00

Video?

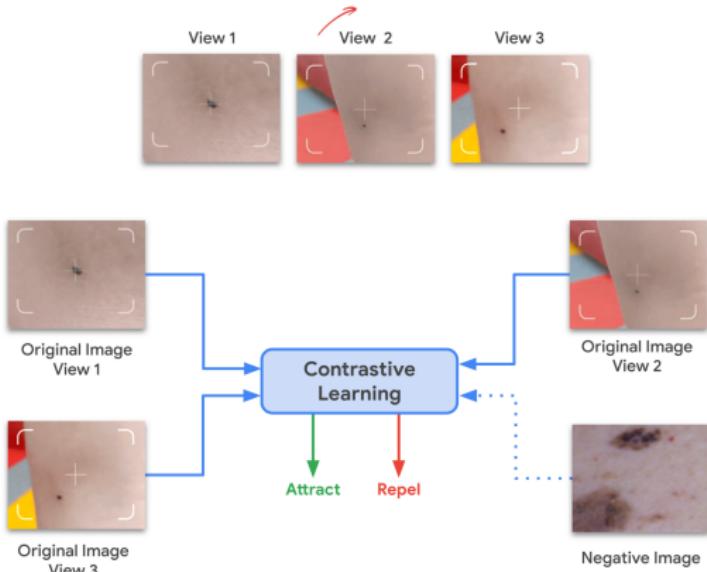


(Sharma et al., 2022)

Take-home Message

Given the practical considerations of medical applications, can we move beyond augmentations?

Images with the same pathology but captured from different view is used to create positive pairs.



(Azizi et al., 2021)