
Self-Supervised Disentanglement by Leveraging Structure in Data Augmentations

Cian Eastwood^{*1,2,3} Julius von Kügelgen^{1,4} Linus Ericsson²
Diane Bouchacourt³ Pascal Vincent³ Bernhard Schölkopf¹ Mark Ibrahim³

¹ Max Planck Institute for Intelligent Systems, Tübingen
² University of Edinburgh ³ Meta AI ⁴ University of Cambridge

Abstract

Self-supervised representation learning often uses data augmentations to induce some invariance to “style” attributes of the data. However, with downstream tasks generally unknown at training time, it is difficult to deduce *a priori* which attributes of the data are indeed “style” and can be safely discarded. To address this, we introduce a more principled approach that seeks to *disentangle* style features rather than discard them. The key idea is to add multiple style embedding spaces where: (i) each is invariant to *all-but-one* augmentation; and (ii) *joint* entropy is maximized. We empirically demonstrate the benefits of our approach on synthetic datasets and then present promising but limited results on ImageNet.

1 Introduction

Learning useful representations from unlabelled data is widely recognized as an important step towards more capable machine-learning systems (Bengio et al., 2013). In recent years, *self-supervised learning* (SSL) has made significant progress towards this goal, approaching the performance of supervised methods on many downstream tasks (Ericsson et al., 2021). The main idea is to leverage known data structures to construct proxy tasks or objectives that act as a form of (self-)supervision. This could involve predicting one part of an observation from another (Brown et al., 2020), or, as we focus on in this work, leveraging **data augmentations/transformations** to perturb different attributes of the data.

Most current approaches are based on the joint-embedding framework and use data augmentations as weak supervision to determine what information to retain (termed “content”) and what information to discard (termed “style”) (Bromley et al., 1994; Chen et al., 2020a; Zbontar et al., 2021; Bardes et al., 2022). In particular, they do so by optimizing for representation similarity or **invariance** across transformations of the same observation, subject to some form of **entropy** regularization, with this invariance-entropy trade-off tuned for some particular task (e.g., ImageNet object classification).

However, at pre-training time, it is unclear what information should be discarded as **one task’s style may be another’s content**. Ericsson et al. (2021) illustrated this point, finding ImageNet object-classification accuracy (the task optimized for in pre-training) to be poorly correlated with downstream object-detection and dense-prediction tasks, concluding that “*universal pre-training is still unsolved*”.

Example 1.1 (Color and Rotation). Suppose we want to make use of **color** and **rotation** transformations. While some invariance to (or discarding of) an image’s color and orientation features can be *beneficial* for ImageNet object classification (Chen et al., 2020a), it can also be *detrimental* for other tasks like segmentation or fine-grained species classification (Cole et al., 2022).

To address this shortcoming and ultimately learn more universal/transferrable representations, we introduce a new joint-embedding framework for SSL which uses data augmentations to **disentangle style attributes of the data rather than discard them**. In particular, as illustrated in Fig. 1, we leverage M transformations to learn $M+1$ *disentangled* embedding spaces capturing both content and style information—with one style space per (group of) transformation(s).

*Work completed during an internship at Meta AI (FAIR), New York.

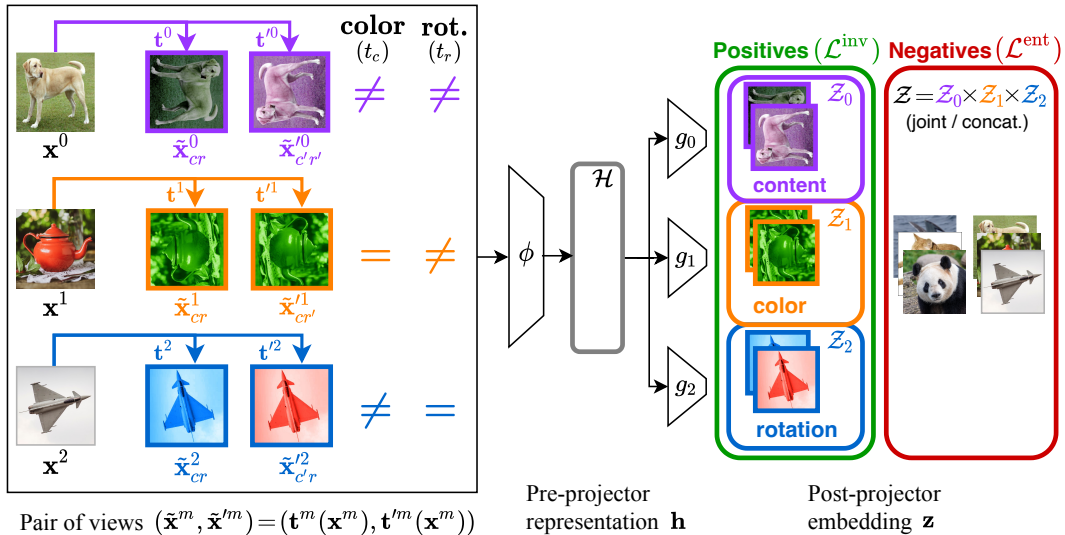


Figure 1: **Framework overview.** Given M atomic transformations like **color distortion** or **rotation** (here, $M=2$), we learn a “content” embedding space (\mathcal{Z}_0) that is invariant to *all* transformations and M “style” embedding spaces ($\mathcal{Z}_1, \mathcal{Z}_2$) that are invariant to *all-but-the- m^{th}* transformation. To do so, we construct $M+1$ transformation pairs $(\mathbf{t}^m, \mathbf{t}'^m)$ sharing different transformation parameters and use these to create $M+1$ transformed image pairs $(\tilde{\mathbf{x}}^m, \tilde{\mathbf{x}}'^m)$ sharing different features. After routing each pair to a different space, we: (i) enforce *invariance within* each space; and (ii) maximize *entropy across* the joint spaces. The result is $M+1$ *disentangled* embedding spaces.

2 Background: Using data augmentations to discard

Joint-embedding methods are often categorized as contrastive or non-contrastive; while both employ some invariance criterion \mathcal{L}^{inv} to encourage the same embedding across different views of the same image (e.g., cosine similarity or mean squared error), they differ in how they regularize this invariance criterion to avoid collapsed or trivial solutions. In particular, contrastive methods (Chen et al., 2020a,b, 2021; He et al., 2020) do so by pushing apart the embeddings of different images, while non-contrastive methods do so by architectural design (Grill et al., 2020; Chen and He, 2021) or by regularizing the covariance of embeddings (Zbontar et al., 2021; Bardes et al., 2022; Ermolov et al., 2021). We focus on **contrastive** and **covariance-based non-contrastive** methods which can both be expressed as a combination of **invariance** \mathcal{L}^{inv} and **entropy** \mathcal{L}^{ent} terms (Garrido et al., 2023),

$$\mathcal{L}^{\text{SSL}} = \mathcal{L}^{\text{inv}} + \mathcal{L}^{\text{ent}}. \quad (2.1)$$

Note these terms have also been called alignment and uniformity (Wang and Isola, 2020), respectively. For concreteness, Table 3 of App. C.2 specifies \mathcal{L}^{inv} and \mathcal{L}^{ent} for some common SSL methods.

In general, the joint-embedding framework involves an unlabelled dataset of observations or images \mathbf{x} and M transformation distributions $\mathcal{T}_1, \dots, \mathcal{T}_M$ from which to sample M atomic transformations t_1, \dots, t_M , with $t_m \sim \mathcal{T}_m$, composed together to form $\mathbf{t} = t_1 \circ \dots \circ t_M$. Critically, **each atomic transformation t_m is designed to perturb a different “style” attribute of the data** deemed nuisance for the task at hand. Returning to Example 1.1, this could mean sampling parameters for a **color distortion** $t_c \sim \mathcal{T}_c$ and **rotation** $t_r \sim \mathcal{T}_r$, and then composing them as $\mathbf{t} = t_c \circ t_r$. For brevity, this sample-and-compose operation is often written as $\mathbf{t} \sim \mathcal{T}$.

For each image \mathbf{x} , a pair of transformations $\mathbf{t}, \mathbf{t}' \sim p_{\mathbf{t}}$ is sampled and applied to form a pair of views $(\tilde{\mathbf{x}}, \tilde{\mathbf{x}}') = (\mathbf{t}(\mathbf{x}), \mathbf{t}'(\mathbf{x}))$. The views are then passed through a shared backbone network ϕ to form a pair of representations $(\mathbf{h}, \mathbf{h}')$, with $\mathbf{h} = \phi(\tilde{\mathbf{x}})$, and then through a smaller “projector” network g to form a pair of embeddings $(\mathbf{z}, \mathbf{z}')$, with $\mathbf{z} = g(\mathbf{h}) = g(\phi(\tilde{\mathbf{x}})) \in \mathcal{Z}$. Critically, the single embedding space \mathcal{Z} seeks invariance to all transformations, thereby **discarding each of the “style” attributes**.

3 Framework: Using data augmentations to disentangle

We now describe our framework for using data augmentations to *disentangle* style attributes of the data, rather than discard them—see Fig. 1 for an illustration. Given M transformations, we learn $M+1$ embedding spaces $\{\mathcal{Z}_m\}_{m=0}^M$ capturing both content (\mathcal{Z}_0) and style ($\{\mathcal{Z}_m\}_{m=1}^M$) information—with one style space per (group of) atomic transformation(s).

Views. We start by constructing $M + 1$ transformation pairs $\{(\mathbf{t}^m, \mathbf{t}'^m)\}_{m=0}^M$ which *share different transformation parameters*. For $m=0$, we independently sample two transformations $\mathbf{t}^0, \mathbf{t}'^0 \sim \mathcal{T}$, which will generally not share any transformation parameters (i.e., $t_k^0 \neq t_k'^0 \forall k$). For $1 \leq m \leq M$, we also independently sample two transformations $\mathbf{t}^m, \mathbf{t}'^m \sim \mathcal{T}$, but then enforce that **the parameters of the m^{th} transformation are shared** by setting $t_m'^m := t_m^m$. Finally, we apply each of these transformation pairs to a different image to form a pair of views $(\tilde{\mathbf{x}}^m, \tilde{\mathbf{x}}'^m) = (\mathbf{t}^m(\mathbf{x}^m), \mathbf{t}'^m(\mathbf{x}^m))$.

Example 1.1 (continued). Suppose we can sample parameters for two transformations: **color distortion** $t_c \sim \mathcal{T}_c$ and **rotation** $t_r \sim \mathcal{T}_r$. As depicted in Fig. 1, we can then construct three transformation pairs *sharing different parameters*: $(\mathbf{t}^0, \mathbf{t}'^0) = (t_c^0 \circ t_r^0, t_c'^0 \circ t_r'^0)$ with **no shared parameters**; $(\mathbf{t}^1, \mathbf{t}'^1) = (t_c^1 \circ t_r^1, t_c^1 \circ t_r'^1)$ with **shared color parameters** t_c^1 ; and $(\mathbf{t}^2, \mathbf{t}'^2) = (t_c^2 \circ t_r^2, t_c'^2 \circ t_r^2)$ with **shared rotation parameters** t_r^2 . Applying each transformation pair to a different image, we get three pairs of views: $(\tilde{\mathbf{x}}_{c,r}^0, \tilde{\mathbf{x}}_{c',r'}^0)$ for which only “**content**” information is shared as both color and rotation differ; $(\tilde{\mathbf{x}}_{c,r}^1, \tilde{\mathbf{x}}_{c',r'}^1)$ for which “**content**” and **color information is shared**, but rotation differs; and $(\tilde{\mathbf{x}}_{c,r}^2, \tilde{\mathbf{x}}_{c',r'}^2)$ for which “**content**” and **rotation information is shared**, but color differs.

Embedding spaces. As depicted in Fig. 1, the pairs of views $(\tilde{\mathbf{x}}^m, \tilde{\mathbf{x}}'^m)$ are passed through a shared backbone network ϕ to form pairs of representations $(\mathbf{h}^m, \mathbf{h}'^m)$ and subsequently through *separate* projectors g_l to form pairs of embeddings $(\mathbf{z}_l^m, \mathbf{z}_l'^m)$, with

$$\mathbf{z}_l^m = g_l(\mathbf{h}^m) = g_l \circ \phi(\tilde{\mathbf{x}}^m) = g_l \circ \phi \circ \mathbf{t}^m(\mathbf{x}^m) \in \mathcal{Z}_l \quad (3.1)$$

the embedding of view $\tilde{\mathbf{x}}^m$ in embedding space \mathcal{Z}_l . We call \mathcal{Z}_0 “**content**” space as it seeks invariance to *all* transformations, thereby discarding all style attributes and leaving only content. We call the other M spaces $\{\mathcal{Z}_m\}_{m=1}^M$ “**style**” spaces as they seek invariance to *all-but-one* transformation t_m , thereby discarding *all-but-one* style attribute (that which is perturbed by t_m).

Loss. Given $M+1$ pairs of views $\{(\tilde{\mathbf{x}}^m, \tilde{\mathbf{x}}'^m)\}_{m=0}^M$ sharing different transformation parameters, we learn $M+1$ disentangled embedding spaces by minimizing the following objective:

$$\begin{aligned} \mathcal{L}^{\text{ours}}(f, \{g_m\}_{m=0}^M; \{(\tilde{\mathbf{x}}^m, \tilde{\mathbf{x}}'^m)\}_{m=0}^M) &= \underbrace{\mathcal{L}_{\mathcal{Z}_0}^{\text{inv}} + \mathcal{L}_{\mathcal{Z}_0}^{\text{ent}}}_{\text{standard loss (content} \rightarrow \mathcal{Z}_0)} + \underbrace{\left(\sum_{m=1}^M \mathcal{L}_{\mathcal{Z}_m}^{\text{inv}}\right) + \mathcal{L}_{\mathcal{Z}}^{\text{ent}}}_{\text{additional terms (style} \rightarrow \mathcal{Z}_m\text{'s)}}, \quad (3.2) \\ &= \underbrace{\left(\sum_{m=0}^M \mathcal{L}_{\mathcal{Z}_m}^{\text{inv}}\right)}_{M+1 \text{ inv. terms}} + \underbrace{\mathcal{L}_{\mathcal{Z}}^{\text{ent}}}_{\text{joint entropy}} + \underbrace{\mathcal{L}_{\mathcal{Z}_0}^{\text{ent}}}_{\text{content entropy}}, \quad (3.3) \end{aligned}$$

where the individual invariance ($\mathcal{L}_{\mathcal{Z}_m}^{\text{inv}}$) and (content / joint) entropy ($\mathcal{L}_{\mathcal{Z}_0}^{\text{ent}}$ / $\mathcal{L}_{\mathcal{Z}}^{\text{ent}}$) terms are given by $\mathcal{L}_{\mathcal{Z}_m}^{\text{inv}} = \lambda_m \mathcal{L}^{\text{inv}}(\mathbf{z}_m^m, \mathbf{z}_m'^m)$, $\mathcal{L}_{\mathcal{Z}_0}^{\text{ent}} = \mathcal{L}^{\text{ent}}(\{\mathbf{z}_0^m, \mathbf{z}_0'^m\}_{m=0}^M)$, $\mathcal{L}_{\mathcal{Z}}^{\text{ent}} = \mathcal{L}^{\text{ent}}(\{\mathbf{z}^m, \mathbf{z}'^m\}_{m=0}^M)$, with $\mathbf{z}^m = [\mathbf{z}_0^m, \dots, \mathbf{z}_M^m] \in \mathcal{Z}_0 \times \dots \times \mathcal{Z}_M$ the concatenated embeddings of $\tilde{\mathbf{x}}^m$ across all spaces. Eq. (3.2) highlights the additional terms we add to the standard contrastive loss. In particular, note that we require **two different entropy terms to ensure disentangled embedding spaces**. Since “content” is invariant to all transformations (by definition), we require $\mathcal{L}_{\mathcal{Z}}^{\text{ent}}$ to prevent redundancy (M additional copies of content, one per style space) and $\mathcal{L}_{\mathcal{Z}_0}^{\text{ent}}$ to ensure content is indeed encoded in \mathcal{Z}_0 (otherwise it could be spread across all $M+1$ spaces). As detailed in App. D.2, this is a key difference compared to Xiao et al. (2021), who learn multiple embedding spaces but do not achieve disentanglement.

4 Experiments

We now present our experimental results which: (i) use a numerical dataset and a synthetic-image dataset to illustrate how *adapting our λ hyperparameter helps to fully disentangle content* (see App. C.3); and (ii) use ImageNet to illustrate the *downstream performance benefits of retaining more style information*. App. C gives full implementation details.

Numerical dataset: Recovering only content. Following von Kùgelgen et al. (2021, Sec. 5.1), we generate synthetic data pairs $(\mathbf{x}, \tilde{\mathbf{x}}) = (f(\mathbf{c}, \mathbf{s}), f(\mathbf{c}, \tilde{\mathbf{s}}))$ with shared content \mathbf{c} and perturbed style $\mathbf{s}, \tilde{\mathbf{s}}$ (see App. C.4 for details.). We then train a simple encoder ϕ (3-layer MLP) with SimCLR using (i) fixed λ (ii) our adaptive λ (see App. C.3) to get learned embeddings $\mathbf{z} = \phi(\mathbf{x})$. We then report the r^2 coefficient of determination in predicting the ground-truth \mathbf{c} and \mathbf{s} from \mathbf{z} . Fig. 2 shows how

Table 1: **SimCLR’s sensitivity to augmentation strengths with fixed and adaptive λ on ColorDSprites.** r^2 in predicting the ground-truth factor values from the post-projector embedding \mathbf{z} with a linear classifier. Adapting λ ensures that \mathbf{z} captures all of content ($C = 1$) and almost no style ($\bar{S} = 0$), regardless of the augmentation strengths.

λ	Augm. Strength	Content (C)		Style (S)				$\bar{S} (\downarrow)$
		Shape	Color	Scale	Orient.	PosX	PosY	
fixed	weak	1.0	0.93	0.89	0.30	0.82	0.83	0.75
	medium	1.0	0.73	1.00	0.19	0.89	0.89	0.74
	strong	1.0	0.31	1.00	0.05	0.23	0.30	0.38
adaptive	weak	1.0	0.21	0.17	0.00	0.01	0.01	0.08
	medium	1.0	0.10	0.16	0.00	0.00	0.00	0.05
	strong	1.0	0.10	0.11	0.00	0.00	0.00	0.04

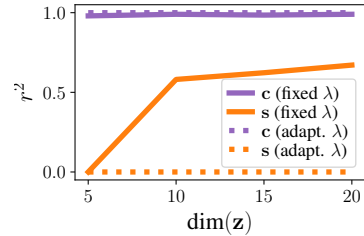


Figure 2: **Recovering *only* content.** r^2 in predicting the ground-truth **content \mathbf{c}** and **style \mathbf{s}** from the learned embedding \mathbf{z} .

Table 2: **Linear evaluation on ImageNet and a broad range of downstream tasks.** We show top-1 accuracies (%) for all but CUB_{bbox} (r^2), CUB_{kpt} (r^2), and VOC (AP_{50}). We use frozen representations \mathbf{h} and embeddings \mathbf{z} (post-projector). FT: our framework *fine-tunes* a base SimCLR model. Ct101: Caltech101. Cf10: CIFAR10.

Alg.	Feat.	ImNt	AcfT	Ct101	Cars	Cf10	Cf100	CUB_{bbox}	CUB_{cls}	CUB_{kpt}	DTD	Flwrs	Pets	SUN	VOC	Avg.
SimCLR	\mathbf{z}	56.5	14.6	70.9	13.0	76.7	50.5	35.6	22.5	12.0	66.4	66.8	70.3	48.9	74.6	47.9
SimCLR-Ours	\mathbf{z}	49.0	25.9	77.6	14.4	81.8	56.2	60.5	15.1	17.6	64.4	63.4	60.7	45.9	73.6	50.6
SimCLR-Ours-FT	\mathbf{z}	57.8	15.9	72.4	14.6	77.8	53.6	36.2	22.5	12.6	67.0	67.3	70.6	49.5	74.9	48.8
SimCLR	\mathbf{h}	68.1	50.9	88.2	50.7	89.3	73.0	71.3	48.6	32.5	75.2	93.5	82.4	60.3	79.7	68.9
SimCLR-Ours	\mathbf{h}	61.7	46.2	86.5	37.5	87.2	67.4	70.6	29.6	23.6	72.6	86.8	73.3	55.1	77.2	62.6
SimCLR-Ours-FT	\mathbf{h}	67.9	51.0	88.1	51.0	89.4	72.9	71.4	48.5	32.9	75.9	93.5	82.6	60.2	79.7	69.0

varying the dimensionality of \mathbf{z} affects the recovery of content \mathbf{c} and style \mathbf{s} , focusing on the scenario where we have sufficient capacity to encode (all of) content \mathbf{c} (i.e., $\dim(\mathbf{z}) \geq \dim(\mathbf{c})$). Similar to von Kügelgen et al. (2021, Fig. 10), we find that with standard SimCLR (fixed λ), excess capacity is used to encode some style information (since that increases entropy). However, by adapting λ using the procedure of App. C.3, we prevent style “leaking in”, allowing us to recover *only content* in \mathbf{z} .

ColorDSprites: Sensitivity to augmentation strengths. We now make use of a colored version of the DSprites dataset (Locatello et al., 2019) which contains images of 2D shapes generated from 6 independent ground-truth factors (# values): color (10), shape (3), scale (6), orientation (40), x-position (32) and y-position (32). We first train SimCLR and VICReg models on the (unlabelled) dataset using **different augmentation strengths** (see Fig. 3 of App. C.1). We then train linear classifiers on top of frozen embeddings \mathbf{z} to predict the ground-truth factor values. Table 1 shows that: (i) **for fixed λ , the augmentation strengths severely affect SimCLR’s invariance-entropy trade-off** and, as a result, the amount (and type) of style information captured in \mathbf{z} ; and (ii) **adapting λ** (see App. C.3) makes SimCLR’s invariance-entropy trade-off much more robust to the augmentation strengths, ensuring that \mathbf{z} captures all of content ($C = 1$) and almost no style ($\bar{S} = 0$)—regardless of the augmentation strengths. Table 5 of App. E.1 gives the corresponding results for VICReg.

ImageNet: Downstream performance. We train all models for 100 epochs on a blurred-face (for legal reasons) ImageNet1k (Russakovsky et al., 2014) dataset using the standard transformations (random crop, horizontal flip, color jitter, grayscale and blur). For our method, we group these into spatial (crop, flip) and appearance (color jitter, grayscale, blur) transformations and thus learn $M = 2$ style spaces. We then follow the setup of Ericsson et al. (2021) to evaluate models on a **broad range of downstream tasks**, covering object/texture/scene classification, localization, and keypoint estimation. With the post-projector \mathbf{z} , Table 2 shows that: (i) using our framework from scratch improves downstream performance at the cost of ImageNet performance; and (ii) using our framework to fine-tune a SimCLR model (i.e., add in style spaces) leads to improved performance both downstream *and* on ImageNet. Unfortunately, this improved performance with \mathbf{z} did not translate into improved performance with the pre-projector \mathbf{h} . This highlights the importance of the projector, but also our poor understanding of its role and impact on the retention of style information.

5 Discussion

Related work. Our framework is closely related to Xiao et al. (2021) who also learn multiple embeddings by applying augmentations in a structured way. However, while the idea of additional style embedding spaces is shared, our framework goes further by enforcing that these spaces are *disentangled* (see App. D.2 for further details). In addition, we provide a theoretical analysis of when this disentanglement is possible (see App. A). Further related work is discussed in App. D.1.

Outlook. We have presented a framework for learning disentangled representations in SSL. We hope future work strives for disentanglement in SSL to enable broader transferability, and explores new data augmentations designed for disentangling (rather than discarding).

References

- Ahuja, K., Hartford, J. S., and Bengio, Y. (2022). Weakly supervised representation learning with sparse perturbations. In *Advances in Neural Information Processing Systems*, volume 35, pages 15516–15528. [Cited on p. 12 and 16.]
- Bardes, A., Ponce, J., and LeCun, Y. (2022). VICReg: variance-invariance-covariance regularization for self-supervised learning. In *International Conference on Learning Representations*. [Cited on p. 1, 2, 13, and 15.]
- Bengio, Y., Courville, A., and Vincent, P. (2013). Representation learning: A review and new perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(8):1798–1828. [Cited on p. 1 and 16.]
- Bouchacourt, D., Tomioka, R., and Nowozin, S. (2018). Multi-level variational autoencoder: Learning disentangled representations from grouped observations. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32. [Cited on p. 16.]
- Brehmer, J., De Haan, P., Lippe, P., and Cohen, T. (2022). Weakly supervised causal representation learning. In *Advances in Neural Information Processing Systems*, volume 35, pages 38319–38331. [Cited on p. 12 and 16.]
- Bromley, J., Guyon, I., LeCun, Y., Säckinger, E., and Shah, R. (1994). Signature verification using a "siamese" time delay neural network. *Advances in Neural Information Processing Systems*, 6. [Cited on p. 1.]
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. (2020). Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. [Cited on p. 1.]
- Chen, T., Kornblith, S., Norouzi, M., and Hinton, G. (2020a). A simple framework for contrastive learning of visual representations. In *International Conference on Machine Learning*, pages 1597–1607. [Cited on p. 1, 2, and 13.]
- Chen, T., Luo, C., and Li, L. (2021). Intriguing properties of contrastive losses. *Advances in Neural Information Processing Systems*, 34:11834–11845. [Cited on p. 2.]
- Chen, X., Fan, H., Girshick, R., and He, K. (2020b). Improved baselines with momentum contrastive learning. *arXiv:2003.04297*. [Cited on p. 2.]
- Chen, X. and He, K. (2021). Exploring simple siamese representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15750–15758. [Cited on p. 2.]
- Cimpoi, M., Maji, S., Kokkinos, I., Mohamed, S., and Vedaldi, A. (2014). Describing textures in the wild. In *Proceedings of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. [Cited on p. 13.]
- Cole, E., Yang, X., Wilber, K., Mac Aodha, O., and Belongie, S. (2022). When does contrastive visual representation learning work? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14755–14764. [Cited on p. 1.]
- Dangovski, R., Jing, L., Loh, C., Han, S., Srivastava, A., Cheung, B., Agrawal, P., and Soljagic, M. (2022). Equivariant self-supervised learning: Encouraging equivariance in representations. In *International Conference on Learning Representations*. [Cited on p. 15.]
- Darmois, G. (1951). Analyse des liaisons de probabilité. In *Proc. Int. Stat. Conferences 1947*, page 231. [Cited on p. 12.]
- Daunhawer, I., Bizeul, A., Palumbo, E., Marx, A., and Vogt, J. E. (2023). Identifiability results for multimodal contrastive learning. In *The Eleventh International Conference on Learning Representations*. [Cited on p. 10 and 16.]
- Desjardins, G., Courville, A., and Bengio, Y. (2012). Disentangling factors of variation via generative entangling. *arXiv preprint arXiv:1210.5474*. [Cited on p. 16.]

- Eastwood, C. and Williams, C. K. I. (2018). A framework for the quantitative evaluation of disentangled representations. In *International Conference on Learning Representations*. [Cited on p. 16.]
- Ericsson, L., Gouk, H., and Hospedales, T. M. (2021). How well do self-supervised models transfer? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5414–5423. [Cited on p. 1, 4, and 15.]
- Ermolov, A., Siarohin, A., Sangineto, E., and Sebe, N. (2021). Whitening for self-supervised representation learning. In *International Conference on Machine Learning*, pages 3015–3024. [Cited on p. 2.]
- Everingham, M., Van Gool, L., Williams, C. K. I., Winn, J., and Zisserman, A. (2007). The PASCAL Visual Object Classes Challenge 2007 (VOC2007) Results. <http://www.pascal-network.org/challenges/VOC/voc2007/workshop/index.html>. [Cited on p. 13.]
- Garrido, Q., Chen, Y., Bardes, A., Najman, L., and LeCun, Y. (2023). On the duality between contrastive and non-contrastive self-supervised learning. In *International Conference on Learning Representations*. [Cited on p. 2.]
- Goyal, P., Dollár, P., Girshick, R. B., Noordhuis, P., Wesolowski, L., Kyrola, A., Tulloch, A., Jia, Y., and He, K. (2017). Accurate, large minibatch sgd: Training imagenet in 1 hour. *ArXiv*. [Cited on p. 15.]
- Gresele, L., Rubenstein, P. K., Mehrjou, A., Locatello, F., and Schölkopf, B. (2019). The incomplete rosetta stone problem: Identifiability results for multi-view nonlinear ICA. In *Proceedings of the 35th Conference on Uncertainty in Artificial Intelligence*, volume 115, pages 217–227. PMLR. [Cited on p. 12 and 16.]
- Gresele, L., Von Kügelgen, J., Stimper, V., Schölkopf, B., and Besserve, M. (2021). Independent mechanism analysis, a new concept? In *Advances in Neural Information Processing Systems*, volume 34, pages 28233–28248. [Cited on p. 12.]
- Grill, J.-B., Strub, F., Altché, F., Tallec, C., Richemond, P. H., Buchatskaya, E., Doersch, C., Pires, B. A., Guo, Z. D., Azar, M. G., Piot, B., Kavukcuoglu, K., Munos, R., and Valko, M. (2020). Bootstrap Your Own Latent: A new approach to self-supervised learning. In *Advances in Neural Information Processing Systems*. [Cited on p. 2.]
- He, K., Fan, H., Wu, Y., Xie, S., and Girshick, R. (2020). Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9729–9738. [Cited on p. 2.]
- He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning for image recognition. In *CVPR*. [Cited on p. 15.]
- Higgins, I., Matthey, L., Pal, A., Burgess, C., Glorot, X., Botvinick, M., Mohamed, S., and Lerchner, A. (2017). β -VAE: Learning basic visual concepts with a constrained variational framework. In *International Conference on Learning Representations*. [Cited on p. 13 and 16.]
- Hyvärinen, A. and Pajunen, P. (1999). Nonlinear independent component analysis: Existence and uniqueness results. *Neural networks*, 12(3):429–439. [Cited on p. 12.]
- Ilse, M., Tomczak, J. M., and Forré, P. (2021). Selecting data augmentation for simulating interventions. In *International Conference on Machine Learning*, pages 4555–4562. PMLR. [Cited on p. 11.]
- Jing, L., Vincent, P., LeCun, Y., and Tian, Y. (2022). Understanding dimensional collapse in contrastive self-supervised learning. In *International Conference on Learning Representations*. [Cited on p. 16.]
- Krause, J., Stark, M., Deng, J., and Fei-Fei, L. (2013). 3d object representations for fine-grained categorization. In *2013 IEEE International Conference on Computer Vision Workshops*. [Cited on p. 13.]
- Krizhevsky, A. (2009). Learning multiple layers of features from tiny images. [Cited on p. 13.]

- Kulkarni, T. D., Whitney, W. F., Kohli, P., and Tenenbaum, J. (2015). Deep convolutional inverse graphics network. *Advances in neural information processing systems*, 28. [Cited on p. 16.]
- Lee, H., Hwang, S. J., and Shin, J. (2020). Self-supervised label augmentation via input transformations. In *International Conference on Machine Learning*, pages 5714–5724. [Cited on p. 15.]
- Lee, H., Lee, K., Lee, K., Lee, H., and Shin, J. (2021). Improving transferability of representations via augmentation-aware self-supervision. In *Advances in Neural Information Processing Systems*. [Cited on p. 15.]
- Li, F.-F., Andreeto, M., Ranzato, M., and Perona, P. (2022). Caltech 101. [Cited on p. 13.]
- Locatello, F., Bauer, S., Lucic, M., Raetsch, G., Gelly, S., Schölkopf, B., and Bachem, O. (2019). Challenging common assumptions in the unsupervised learning of disentangled representations. In *36th International Conference on Machine Learning*, pages 7247–7283. Curran Associates, Inc. [Cited on p. 4, 13, and 16.]
- Locatello, F., Poole, B., Rätsch, G., Schölkopf, B., Bachem, O., and Tschannen, M. (2020). Weakly-supervised disentanglement without compromises. In *International Conference on Machine Learning*, pages 6348–6359. PMLR. [Cited on p. 12 and 16.]
- Loshchilov, I. and Hutter, F. (2017). Sgdr: Stochastic gradient descent with restarts. In *ICLR*. [Cited on p. 15.]
- Lyu, Q., Fu, X., Wang, W., and Lu, S. (2021). Understanding latent correlation-based multiview learning and self-supervision: An identifiability perspective. In *International Conference on Learning Representations*. [Cited on p. 12 and 16.]
- Maji, S., Rahtu, E., Kannala, J., Blaschko, M., and Vedaldi, A. (2013). Fine-grained visual classification of aircraft. *arXiv preprint arXiv:1306.5151*. [Cited on p. 13.]
- Mitrovic, J., McWilliams, B., Walker, J. C., Buesing, L. H., and Blundell, C. (2021). Representation learning via invariant causal mechanisms. In *International Conference on Learning Representations*. [Cited on p. 11.]
- Nilsback, M.-E. and Zisserman, A. (2008). Automated flower classification over a large number of classes. In *Proceedings of the Indian Conference on Computer Vision, Graphics and Image Processing*. [Cited on p. 13.]
- Papamakarios, G., Nalisnick, E., Rezende, D. J., Mohamed, S., and Lakshminarayanan, B. (2021). Normalizing flows for probabilistic modeling and inference. *The Journal of Machine Learning Research*, 22(1):2617–2680. [Cited on p. 12.]
- Parkhi, O., Vedaldi, A. ., Zisserman, A., and Jawahar, C. V. (2012). Cats and dogs. In *IEEE Conference on Computer Vision and Pattern Recognition*. [Cited on p. 13.]
- Pearl, J. (2009). *Causality: Models, Reasoning, and Inference*. Cambridge University Press. [Cited on p. 11.]
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M. S., Berg, A. C., and Fei-Fei, L. (2014). Imagenet large scale visual recognition challenge. *CoRR*. [Cited on p. 4 and 13.]
- Schölkopf, B., Locatello, F., Bauer, S., Ke, N. R., Kalchbrenner, N., Goyal, A., and Bengio, Y. (2021). Toward causal representation learning. *Proceedings of the IEEE*, 109(5):612–634. [Cited on p. 10.]
- Squires, C., Seigal, A., Bhate, S., and Uhler, C. (2023). Linear causal disentanglement via interventions. In *40th International Conference on Machine Learning*. [Cited on p. 12.]
- Suter, R., Miladinovic, D., Schölkopf, B., and Bauer, S. (2019). Robustly disentangled causal mechanisms: Validating deep representations for interventional robustness. In *International Conference on Machine Learning*, pages 6056–6065. [Cited on p. 11.]
- Tenenbaum, J. and Freeman, W. (1996). Separating style and content. In *Advances in Neural Information Processing Systems*, volume 9. MIT Press. [Cited on p. 16.]

- von Kügelgen, J., Sharma, Y., Gresele, L., Brendel, W., Schölkopf, B., Besserve, M., and Locatello, F. (2021). Self-supervised learning with data augmentations provably isolates content from style. In *Advances in Neural Information Processing Systems*, volume 34, pages 16451–16467. [Cited on p. 3, 4, 10, 11, 12, 13, 15, and 16.]
- Wah, C., Branson, S., Welinder, P., Perona, P., and Belongie, S. (2011). The caltech-ucsd birds-200-2011 dataset. Technical Report CNS-TR-2011-001, California Institute of Technology. [Cited on p. 13.]
- Wang, T. and Isola, P. (2020). Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In *International Conference on Machine Learning*, pages 9929–9939. PMLR. [Cited on p. 2 and 15.]
- Wang, Y. and Jordan, M. I. (2021). Desiderata for representation learning: A causal perspective. *arXiv preprint arXiv:2109.03795*. [Cited on p. 11.]
- Xiao, J., Hays, J., Ehinger, K., Oliva, A., and Torralba, A. (2010). Sun database: Large-scale scene recognition from abbey to zoo. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. [Cited on p. 13.]
- Xiao, T., Wang, X., Efros, A. A., and Darrell, T. (2021). What should not be contrastive in contrastive learning. In *International Conference on Learning Representations*. [Cited on p. 3, 4, 9, 15, 16, and 17.]
- You, Y., Gitman, I., and Ginsburg, B. (2017). Large batch training of convolutional networks. *arXiv: Computer Vision and Pattern Recognition*. [Cited on p. 15.]
- Zbontar, J., Jing, L., Misra, I., LeCun, Y., and Deny, S. (2021). Barlow twins: Self-supervised learning via redundancy reduction. In *International Conference on Machine Learning*, pages 12310–12320. [Cited on p. 1, 2, 13, and 15.]
- Zimmermann, R. S., Sharma, Y., Schneider, S., Bethge, M., and Brendel, W. (2021). Contrastive learning inverts the data generating process. In *International Conference on Machine Learning*. [Cited on p. 13.]

Appendices

Table of Contents

A Causal Representation Learning Perspective and Identifiability Analysis	10
B Proof of Thm. A.2	12
C Implementation Details	13
C.1 Datasets	13
C.2 Invariance and entropy terms	13
C.3 Adaptive λ_m	15
C.4 Numerical dataset	15
C.5 ImageNet	15
D Related Work	15
D.1 Self-supervised learning and disentanglement	15
D.2 Detailed comparison with Xiao et al. (2021)	16
E Further Results	17
E.1 ColorDSprites	17
E.2 ImageNet	18

A Causal Representation Learning Perspective and Identifiability Analysis

In this section, we investigate *what* is actually learned by the structured use of data augmentations in § 3, through the lens of causal representation learning (Schölkopf et al., 2021). To this end, we first formalize the data generation and augmentation processes as a (causal) latent variable model, and then study the identifiability of different components of the latent representation. Our analysis strongly builds on and extends the work of von Kügelgen et al. (2021) by showing that the structure inherent to different augmentation transformations can be leveraged to identify not only the block of shared content variables, but also *individual style components* (subject to suitable assumptions).

Data-generation and augmentation processes. We assume that the observations $\mathbf{x} \in \mathcal{X}$ result from *underlying latent vectors* $\mathbf{z} \in \mathcal{Z}$ via an invertible *nonlinear mixing function* $f : \mathcal{Z} \rightarrow \mathcal{X}$,

$$\mathbf{z} \sim p_{\mathbf{z}}, \quad \mathbf{x} = f(\mathbf{z}). \quad (\text{A.1})$$

Here, $\mathcal{Z} \subseteq \mathbb{R}^d$ is a *latent space* capturing object properties such **color** or **rotation**; $p_{\mathbf{z}}$ is a distribution over latents; and \mathcal{X} denotes the d -dimensional data manifold, which is typically embedded in a higher dimensional pixel space. In the same spirit, we model the way in which augmented views $(\tilde{\mathbf{x}}, \tilde{\mathbf{x}}')$ are generated from \mathbf{x} through perturbations in the latent space:

$$\tilde{\mathbf{z}}, \tilde{\mathbf{z}}' \sim p_{\tilde{\mathbf{z}}|\mathbf{z}}, \quad \tilde{\mathbf{x}} = f(\tilde{\mathbf{z}}), \quad \tilde{\mathbf{x}}' = f(\tilde{\mathbf{z}}'). \quad (\text{A.2})$$

The conditional $p_{\tilde{\mathbf{z}}|\mathbf{z}}$, from which the pair of *augmented latents* $(\tilde{\mathbf{z}}, \tilde{\mathbf{z}}')$ is drawn given the original latent \mathbf{z} , constitutes the latent-space analogue of the image-level transformations $(\mathbf{t}, \mathbf{t}') \sim \mathcal{T}$ in § 3. More specifically, $\tilde{\mathbf{z}} \sim p_{\tilde{\mathbf{z}}|\mathbf{z}}$ captures the behavior of $f^{-1} \circ \mathbf{t} \circ f$ with $\mathbf{t} \sim \mathcal{T}$ acting on $\mathbf{x} = f(\mathbf{z})$.

Content-style partition. Typically, augmentations are designed to affect some semantic aspects of the data (e.g., color and rotation) and not others (e.g., object identity). We therefore partition the latents into *style* latents \mathbf{s} , which *are* affected by the augmentations, and *shared content latents* \mathbf{c} , which *are not* affected by the augmentations. Further, $p_{\tilde{\mathbf{z}}|\mathbf{z}}$ in (A.2) takes the form

$$p_{\tilde{\mathbf{z}}|\mathbf{z}}(\tilde{\mathbf{z}} | \mathbf{z}) = \delta(\tilde{\mathbf{c}} - \mathbf{c}) p_{\tilde{\mathbf{s}}|\mathbf{s}}(\tilde{\mathbf{s}} | \mathbf{s}) \quad (\text{A.3})$$

for some style conditional $p_{\tilde{\mathbf{s}}|\mathbf{s}}$, such that \mathbf{z} , $\tilde{\mathbf{z}}$, and $\tilde{\mathbf{z}}'$ in (A.2) are given by

$$\mathbf{z} = (\mathbf{c}, \mathbf{s}), \quad \tilde{\mathbf{z}} = (\mathbf{c}, \tilde{\mathbf{s}}), \quad \tilde{\mathbf{z}}' = (\mathbf{c}, \tilde{\mathbf{s}}'). \quad (\text{A.4})$$

For this setting, it has been shown that—under suitable additional assumptions—contrastive SSL recovers the shared content latents \mathbf{c} up to an invertible function (von Kügelgen et al., 2021, Thm. 4.4).

Beyond content identifiability: separating and recovering individual style latents. Previous analyses of SSL with data augmentations considered style latents \mathbf{s} as nuisance variables that should be discarded, thus seeking a pure content-based representation that is invariant to all augmentations (von Kügelgen et al., 2021; Daunhawer et al., 2023). The focus of our study, and its key difference to these previous analyses, is that *we seek to also identify and disentangle different style variables*, by leveraging available structure in data augmentations that has not been exploited thus far.

First, note that each class of atomic transformation \mathcal{T}_m (e.g., color distortion or rotation) typically affects a different property, meaning that it should only affect a subset of style variables. Hence, we partition the style block into more fine-grained *individual style components* \mathbf{s}_m ,

$$\mathbf{s} = (\mathbf{s}_1, \dots, \mathbf{s}_M), \quad \tilde{\mathbf{s}} = (\tilde{\mathbf{s}}_1, \dots, \tilde{\mathbf{s}}_M), \quad \tilde{\mathbf{s}}' = (\tilde{\mathbf{s}}'_1, \dots, \tilde{\mathbf{s}}'_M), \quad (\text{A.5})$$

and assume that the style conditional $p_{\tilde{\mathbf{s}}|\mathbf{s}}$ in (A.3) factorizes as follows:

$$p_{\tilde{\mathbf{s}}|\mathbf{s}}(\tilde{\mathbf{s}} | \mathbf{s}) = \prod_{m=1}^M p_{\tilde{\mathbf{s}}_m|\mathbf{s}_m}(\tilde{\mathbf{s}}_m | \mathbf{s}_m), \quad (\text{A.6})$$

where each term $p_{\tilde{\mathbf{s}}_m|\mathbf{s}_m}$ on the RHS is the latent-space analogue of $t_m \sim \mathcal{T}_m$.

Next, we wish for our latent variable model to capture the *structured* use of data augmentation through transformation pairs with *shared parameters*, as described in § 3. Specifically, note that—unlike most

prior approaches to SSL with data augmentations—we do *not* create a single dataset of (“positive”) pairs $(\tilde{\mathbf{x}}, \tilde{\mathbf{x}}')$. Instead, we construct transformation pairs $(\mathbf{t}^m, \mathbf{t}'^m)$ in $M+1$ different ways, giving rise to $M+1$ datasets of pairs $(\tilde{\mathbf{x}}^m, \tilde{\mathbf{x}}'^m)$, each differing in the shared (style) properties. In particular, the m^{th} atomic transformation is shared across $(\mathbf{t}^m, \mathbf{t}'^m)$ by construction, such that $(\tilde{\mathbf{x}}^m, \tilde{\mathbf{x}}'^m)$ should share the same perturbed m^{th} style components $\tilde{\mathbf{s}}_m = \tilde{\mathbf{s}}_m'$ —*regardless of its original value* \mathbf{s}_m . To model this procedure, we define $M+1$ different ways of jointly perturbing the style variables as follows:

$$p^{(m)}(\tilde{\mathbf{s}}, \tilde{\mathbf{s}}' | \mathbf{s}) = \prod_{l=1}^M p^{(m)}(\tilde{\mathbf{s}}_l, \tilde{\mathbf{s}}'_l | \mathbf{s}_l) \quad \text{for } m = 0, \dots, M, \quad (\text{A.7})$$

where

$$p^{(m)}(\tilde{\mathbf{s}}_l, \tilde{\mathbf{s}}'_l | \mathbf{s}_l) = \begin{cases} p_{\tilde{\mathbf{s}}_l | \mathbf{s}_l}(\tilde{\mathbf{s}}_l | \mathbf{s}_l) \delta(\tilde{\mathbf{s}}'_l - \tilde{\mathbf{s}}_l) & \text{if } l = m \\ p_{\tilde{\mathbf{s}}_l | \mathbf{s}_l}(\tilde{\mathbf{s}}_l | \mathbf{s}_l) p_{\tilde{\mathbf{s}}'_l | \mathbf{s}_l}(\tilde{\mathbf{s}}'_l | \mathbf{s}_l) & \text{otherwise} \end{cases}$$

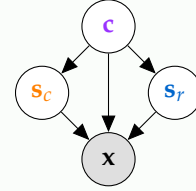
Together with $\mathbf{z} = (\mathbf{c}, \mathbf{s}) \sim p_{\mathbf{z}}$ as in (A.1), the conditionals in (A.7) induce $M+1$ different joint distributions $p_{\tilde{\mathbf{x}}, \tilde{\mathbf{x}}'}^{(m)}$ over observation pairs $(\tilde{\mathbf{x}}^m, \tilde{\mathbf{x}}'^m)$: analogous to (A.2), we have for $m = 0, \dots, M$,

$$\tilde{\mathbf{s}}^m, \tilde{\mathbf{s}}'^m \sim p_{\tilde{\mathbf{s}}, \tilde{\mathbf{s}}' | \mathbf{s}}^{(m)}, \quad \tilde{\mathbf{x}}^m = f([\mathbf{c}, \tilde{\mathbf{s}}^m]), \quad \tilde{\mathbf{x}}'^m = f([\mathbf{c}, \tilde{\mathbf{s}}'^m]). \quad (\text{A.8})$$

Remark A.1. In practice, we do not generate $M+1$ augmented pairs for each $\mathbf{x} = f(\mathbf{z})$ as described above. Instead, each pair is constructed from a different observation with $\mathbf{x}^l = f(\mathbf{z}^l)$ transformed according to $m := l \bmod M+1$. In the limit of infinite data, these two options have the same effect.

Example 1.1 (continued). Denote the style component capturing **color** by \mathbf{s}_c and that capturing **rotation** by \mathbf{s}_r . For $m = 0, 1, 2$, let $\mathbf{z}^m = (\mathbf{c}^m, \mathbf{s}_c^m, \mathbf{s}_r^m)$ be the latents underlying separate images \mathbf{x}^m . Then the augmentations shown in Fig. 1 (left) are captured by the following changes to the latents:

m	\mathbf{z}^m	$\tilde{\mathbf{z}}^m$	$\tilde{\mathbf{z}}'^m$	Shared Latents
0	$(\mathbf{c}^0, \mathbf{s}_c^0, \mathbf{s}_r^0)$	$(\mathbf{c}^0, \tilde{\mathbf{s}}_c^0, \tilde{\mathbf{s}}_r^0)$	$(\mathbf{c}^0, \tilde{\mathbf{s}}_c'^0, \tilde{\mathbf{s}}_r'^0)$	only content
1	$(\mathbf{c}^1, \mathbf{s}_c^1, \mathbf{s}_r^1)$	$(\mathbf{c}^1, \tilde{\mathbf{s}}_c^1, \tilde{\mathbf{s}}_r^1)$	$(\mathbf{c}^1, \tilde{\mathbf{s}}_c'^1, \tilde{\mathbf{s}}_r'^1)$	content & color
2	$(\mathbf{c}^2, \mathbf{s}_c^2, \mathbf{s}_r^2)$	$(\mathbf{c}^2, \tilde{\mathbf{s}}_c^2, \tilde{\mathbf{s}}_r^2)$	$(\mathbf{c}^2, \tilde{\mathbf{s}}_c'^2, \tilde{\mathbf{s}}_r'^2)$	content & rotation



Causal interpretation. The described augmentation procedure can also be interpreted in causal terms (Ilse et al., 2021; Mitrovic et al., 2021; von Kügelgen et al., 2021). Given a factual observation \mathbf{x} , the augmented views $(\tilde{\mathbf{x}}^m, \tilde{\mathbf{x}}'^m)$ constitute pairs of *counterfactuals under joint interventions on all style variables*, provided that (i) \mathbf{c} is a root node in the causal graph, to ensure content invariance in (A.3); and (ii) the style components \mathbf{s}_m do not causally influence each other, to justify the factorization in (A.6) and (A.7).² A causal graph compatible with these constraints is shown for Example 1.1 above on the right. As a structural causal model (SCM; Pearl, 2009), this can be written as

$$\mathbf{c} := \mathbf{u}_c, \quad \mathbf{s}_m := f_m(\mathbf{c}, \mathbf{u}_m), \quad \text{for } m = 1, \dots, M, \quad (\text{A.9})$$

with jointly independent exogenous variables $\mathbf{u}_c, \{\mathbf{u}_m\}_{m=0}^M$. The style conditionals $p_{\tilde{\mathbf{s}}_m | \tilde{\mathbf{s}}_m}$ in (A.6) can then arise, e.g., from *shift do* $(\mathbf{s}_m = f_m(\mathbf{c}, \mathbf{u}_m) + \tilde{\mathbf{u}}_m)$ or *perfect do* $(\mathbf{s}_m = \tilde{\mathbf{u}}_m)$ interventions with independent augmentation noise $\tilde{\mathbf{u}}_m$. Note that the latter renders $\tilde{\mathbf{s}}_m$ independent of all other variables.

Style identifiability and disentanglement. By construction, $\{\mathbf{c}, \tilde{\mathbf{s}}_m\}$ is shared across $(\tilde{\mathbf{x}}^m, \tilde{\mathbf{x}}'^m)$ and can thus be identified *up to nonlinear mixing* by contrastive SSL on the m^{th} dataset (von Kügelgen et al., 2021, Thm. 4.4). However, it remains unclear how to disentangle the two and recover only $\tilde{\mathbf{s}}_m$, i.e., how to “remove” \mathbf{c} , which can separately be recovered as the only shared latent for $m=0$. The following result, proven in App. B, shows that our approach from § 3 with $M+1$ alignment terms and joint entropy regularization indeed disentangles and recovers the individual style components.

²The allowed structure is similar to Suter et al. (2019, Fig. 1); Wang and Jordan (2021, Fig. 9). However, ours is more general as content does not only confound different \mathbf{s}_m , but also directly influences the observed \mathbf{x} .

Theorem A.2 (Identifiability). *For the data generating process in (A.1), (A.7), (A.8), assume that*

- A₁. \mathcal{Z} is open and simply connected; f is diffeomorphic onto its image; $p_{\mathbf{z}}$ is smooth and fully supported on \mathcal{Z} ; each $p_{\tilde{\mathbf{s}}_m|\mathbf{s}_m}$ is smooth and supported on an open, non-empty set around any \mathbf{s}_m ;
- A₂. $p_{\mathbf{z}}$ and $\{p_{\tilde{\mathbf{s}}_m|\mathbf{s}_m}\}_{m=1}^M$ are such that $\{\mathbf{c}\} \cup \{\tilde{\mathbf{s}}_m\}_{m=1}^M$ are jointly independent;
- A₃. the latent dims. $\{d_m\}_{m=0}^M$ are known and $\{\phi_m : \mathcal{X} \rightarrow (0, 1)^{d_m}\}_{m=0}^M$ are smooth minimizers of

$$\sum_{m=0}^M \mathbb{E}_{p_{\tilde{\mathbf{x}}, \tilde{\mathbf{x}}'}^{(m)}}} [\|\phi_m(\tilde{\mathbf{x}}) - \phi_m(\tilde{\mathbf{x}}')\|_2] - H_{p_{\tilde{\mathbf{x}}}^{(0)}}}([\phi_0(\tilde{\mathbf{x}}), \dots, \phi_M(\tilde{\mathbf{x}})]). \quad (\text{A.10})$$

Then ϕ_0 block-identifies (von Kügelgen et al., 2021, Defn. 4.1) the content \mathbf{c} , and ϕ_m block-identifies \mathbf{s}_m in the sense that $\hat{\mathbf{s}}_m = \phi_m(\mathbf{x}) = \psi_m(\mathbf{s}_m)$ for some invertible ψ_m for each $m = 1, \dots, M$.

Discussion of Thm. A.2. The technical assumption A₁ is also needed to prove content identifiability (von Kügelgen et al., 2021). Assumption A₂, which requires that the augmentation process renders \mathbf{c} and $\{\tilde{\mathbf{s}}_m\}_{m=1}^M$ independent, is specific to our extended analysis. It holds, e.g., if (a) $p_{\mathbf{z}}$ is such that \mathbf{c} and $\{\mathbf{s}_m\}_{m=1}^M$ are independent to begin with; or if (b) $p_{\tilde{\mathbf{s}}|\mathbf{s}} = p_{\tilde{\mathbf{s}}}$ does not depend on \mathbf{s} , as would be the case for *perfect* interventions. As discussed in more detail in App. D, (a) relates to work on multi-view latent correlation maximization (Lyu et al., 2021), nonlinear ICA (Gresele et al., 2019), and disentanglement (Locatello et al., 2020; Ahuja et al., 2022), whereas (b) relates to work in weakly supervised causal representation learning (Brehmer et al., 2022). In case (b), we could actually also allow for causal relations among individual style components $\mathbf{s}_m \rightarrow \mathbf{s}_{m'}$, as such links are broken by perfect interventions. When A₂ does not hold (e.g., for content-dependent style and *imperfect* interventions—arguably the most realistic setting), block-identifiability of the style components seems infeasible, consistent with existing negative results (Brehmer et al., 2022; Squires et al., 2023). However, in this case we posit that the exogenous style variables \mathbf{u}_m in (A.9), which capture any style information not due to \mathbf{c} and are jointly independent by assumption, are recovered in place of \mathbf{s}_m .

B Proof of Thm. A.2

Theorem A.2 (Identifiability). *For the data generating process in (A.1), (A.7), (A.8), assume that*

- A₁. \mathcal{Z} is open and simply connected; f is diffeomorphic onto its image; $p_{\mathbf{z}}$ is smooth and fully supported on \mathcal{Z} ; each $p_{\tilde{\mathbf{s}}_m|\mathbf{s}_m}$ is smooth and supported on an open, non-empty set around any \mathbf{s}_m ;
- A₂. $p_{\mathbf{z}}$ and $\{p_{\tilde{\mathbf{s}}_m|\mathbf{s}_m}\}_{m=1}^M$ are such that $\{\mathbf{c}\} \cup \{\tilde{\mathbf{s}}_m\}_{m=1}^M$ are jointly independent;
- A₃. the latent dims. $\{d_m\}_{m=0}^M$ are known and $\{\phi_m : \mathcal{X} \rightarrow (0, 1)^{d_m}\}_{m=0}^M$ are smooth minimizers of

$$\sum_{m=0}^M \mathbb{E}_{p_{\tilde{\mathbf{x}}, \tilde{\mathbf{x}}'}^{(m)}}} [\|\phi_m(\tilde{\mathbf{x}}) - \phi_m(\tilde{\mathbf{x}}')\|_2] - H_{p_{\tilde{\mathbf{x}}}^{(0)}}}([\phi_0(\tilde{\mathbf{x}}), \dots, \phi_M(\tilde{\mathbf{x}})]). \quad (\text{A.10})$$

Then ϕ_0 block-identifies (von Kügelgen et al., 2021, Defn. 4.1) the content \mathbf{c} , and ϕ_m block-identifies \mathbf{s}_m in the sense that $\hat{\mathbf{s}}_m = \phi_m(\mathbf{x}) = \psi_m(\mathbf{s}_m)$ for some invertible ψ_m for each $m = 1, \dots, M$.

Proof. The proof follows a similar argument as that of von Kügelgen et al. (2021, Thm. 4.4), extended to our setting with $M+1$ alignment terms instead of a single one, and with *joint* entropy regularization.

Step 1. First, we show the existence of a solution $\{\phi_m^*\}_{m=0}^M$ attaining the global minimum of zero of the objective in (A.10). To this end, we construct each ϕ_m^* by composing the inverse of the true mixing function with the cumulative distribution function (CDF) transform³ to map each latent block to a uniform version of itself. Specifically, let $\phi_0^* := F_{\mathbf{c}} \circ f_{1:d_0}^{-1}$, and for $m = 1, \dots, M$, let $\phi_m^* := F_{\mathbf{s}_m} \circ f_{a_m:b_m}^{-1}$ with $a_m = 1 + \sum_{l=0}^{m-1} d_l$ and $b_m = \sum_{l=0}^m d_l$, where F_v denotes the CDF of v . By construction, $\phi_0^*(\tilde{\mathbf{x}})$ is a function of \mathbf{c} only, and uniformly distributed on $(0, 1)^{d_0}$; similarly, $\phi_m^*(\tilde{\mathbf{x}})$ is a function of $\tilde{\mathbf{s}}_m$ only and uniform on $(0, 1)^{d_m}$ for $m = 1, \dots, M$. Recall that, with probability one,

³Sometimes also referred to as ‘‘Darmois construction’’ (Darmois, 1951; Hyvärinen and Pajunen, 1999; Gresele et al., 2021; Papamakarios et al., 2021).

\mathbf{c} is shared across $(\tilde{\mathbf{x}}, \tilde{\mathbf{x}}') \sim p_{\tilde{\mathbf{x}}, \tilde{\mathbf{x}}'}^{(0)}$ and $\tilde{\mathbf{s}}_m$ is shared across $(\tilde{\mathbf{x}}, \tilde{\mathbf{x}}') \sim p_{\tilde{\mathbf{x}}, \tilde{\mathbf{x}}'}^{(m)}$. Hence, all the alignment (expectation) terms in (A.10) are zero. Finally, since $\{\mathbf{c}\} \cup \{\tilde{\mathbf{s}}_m\}_{m=1}^M$ are mutually independent by assumption A_2 , and since each ϕ_m^* for $m = 0, \dots, M$ is uniform on $(0, 1)^{d_m}$, it follows that $[\phi_0^*(\tilde{\mathbf{x}}), \dots, \phi_M^*(\tilde{\mathbf{x}})]$ is jointly uniform on $(0, 1)^d$. Hence, the entropy term in (A.10) is also zero.

Step 2. Next, let $\{\phi_m\}_{m=0}^M$ be any other solution attaining the global minimum of (A.10). By the above existence argument, this implies that (i) $\phi_m(\tilde{\mathbf{x}}) = \phi_m(\tilde{\mathbf{x}}')$ almost surely w.r.t. $p_{\tilde{\mathbf{x}}, \tilde{\mathbf{x}}'}^{(m)}$ for $m = 0, \dots, M$; and (ii) $[\phi_0(\tilde{\mathbf{x}}), \dots, \phi_M(\tilde{\mathbf{x}})]$ is jointly uniform on $(0, 1)^d$. As shown by von Kügelgen et al. (2021), the invariance constraint (i) together with the postulated data generating process and assumption A_1 implies that each $\phi_m \circ f$ can only be a function of the latents that are shared almost surely across $(\tilde{\mathbf{x}}, \tilde{\mathbf{x}}') \sim p_{\tilde{\mathbf{x}}, \tilde{\mathbf{x}}'}^{(m)}$. That is, $\phi_0(\mathbf{x}) = \psi_0(\mathbf{c})$ and $\phi_m(\mathbf{x}) = \psi_m(\mathbf{c}, \mathbf{s}_m)$ for $m = 1, \dots, M$. By A_1 and constraint (ii), ψ_0 maps a regular density to another regular density and thus must be invertible (Zimmermann et al., 2021, Prop. 5).

Step 3. It remains to show that ψ_m is invertible and actually cannot depend on \mathbf{c} for $m = 1, \dots, M$, for this would otherwise violate the maximum entropy (uniformity) constraint (ii). Suppose for a contraction that ψ_k depends on \mathbf{c} for some $k \in \{1, \dots, M\}$. By constraint (ii), $[\phi_0(\tilde{\mathbf{x}}), \phi_k(\tilde{\mathbf{x}})] = [\psi_0(\mathbf{c}), \psi_k(\mathbf{c}, \tilde{\mathbf{s}}_k)]$ is jointly uniform on $(0, 1)^{d_0+d_k}$. Hence, $\psi_0(\mathbf{c})$ and $\psi_k(\mathbf{c}, \tilde{\mathbf{s}}_k)$ are independent. Since ψ_0 is invertible, this implies that \mathbf{c} and $\psi_k(\mathbf{c}, \tilde{\mathbf{s}}_k)$ are independent, which (by smoothness of $\psi_k = \phi_k \circ f$ and independence of \mathbf{c} and $\tilde{\mathbf{s}}_k$) contradicts the assumption that ψ_k depends on \mathbf{c} .

Thus, by contradiction, we have that $\phi_m(\tilde{\mathbf{x}}) = \psi_m(\mathbf{c}, \tilde{\mathbf{s}}_m) = \psi_m(\tilde{\mathbf{s}}_m)$ for $m = 1, \dots, M$. Finally, invertibility of $\psi_m(\tilde{\mathbf{s}}_m)$ for $m = 1, \dots, M$ follows from A_1 and Prop. 5 of Zimmermann et al. (2021). Together with $\phi_0(\tilde{\mathbf{x}}) = \psi_0(\mathbf{c})$ (established above) concludes the proof of block-identifiability. \square

C Implementation Details

C.1 Datasets

The numerical data is based on the experiments of von Kügelgen et al. (2021), and the data samples are generated programmatically. ColorDSprites is a synthetic image dataset based on DSprites (Higgins et al., 2017), and extended by Locatello et al. (2019). The rest of the experiments are based on models pretrained on ImageNet1k (Russakovsky et al., 2014), which are then evaluated on the downstream datasets FGVC Aircraft (Maji et al., 2013), Caltech-101 (Li et al., 2022), Stanford Cars (Krause et al., 2013), CIFAR10 (Krizhevsky, 2009), CIFAR100 (Krizhevsky, 2009), CUB (Wah et al., 2011), DTD (Cimpoi et al., 2014), Oxford Flowers (Nilsback and Zisserman, 2008), Oxford-IIIT Pets (Parkhi et al., 2012), SUN397 (Xiao et al., 2010) and VOC2007 (Everingham et al., 2007).

ColorDSprites samples. Fig. 3 depicts samples from the ColorDprites dataset when transformed with transformations/augmentations of different strengths.

C.2 Invariance and entropy terms

To remain general and apply to any contrastive (e.g., SimCLR, Chen et al. 2020a) or non-contrastive covariance-based SSL method (e.g., VICReg, Bardes et al. 2022), both Eq. (2.1) and Eq. (3.2) are expressed as a combination of **invariance** \mathcal{L}^{inv} and **entropy** \mathcal{L}^{ent} terms. For concreteness, Table 3 specifies these terms for some common SSL methods, namely SimCLR (Chen et al., 2020a), BarlowTwins (BTs, Zbontar et al. 2021), and VICReg (Bardes et al., 2022). Note a slight misalignment between $\mathcal{L}^{\text{ent}}(Z, Z')$ in Table 3 and our usage of it in Eq. (3.2). In particular, we write $\mathcal{L}^{\text{ent}}(\{\mathbf{z}^m, \mathbf{z}'^m\}_{m=0}^M)$ for brevity, but should write $\mathcal{L}^{\text{ent}}(Z, Z')$ with $Z = [\mathbf{z}^0, \mathbf{z}^1, \dots, \mathbf{z}^M]$ to align with Table 3.

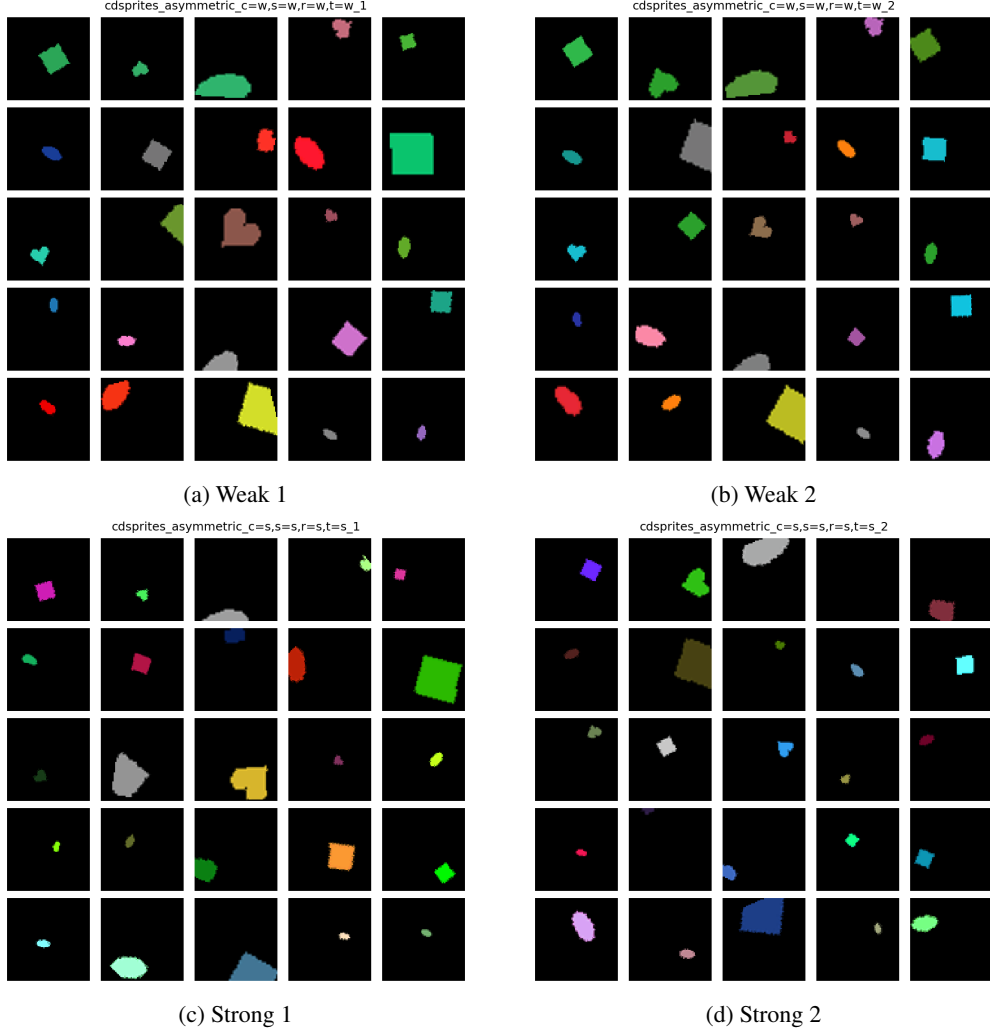


Figure 3: **Augmentation strengths on ColorDSprites.** Columns show augmentation pairs of the same strength. Note that images are more similar across (a) & (b) than across (c) & (d), in terms of color, orientation, scale, translation and X-Y position.

Table 3: **Unified Perspective on SSL Objectives Through Invariance and Entropy.** Many SSL methods can be expressed as a (weighted) combination of **invariance** \mathcal{L}^{inv} and **entropy** \mathcal{L}^{ent} terms. Here, $Z = [\mathbf{z}^1, \mathbf{z}^2, \dots, \mathbf{z}^n]$ and $Z' = [\mathbf{z}'^1, \mathbf{z}'^2, \dots, \mathbf{z}'^n]$ are two batches of n vectors of d -dimensional representations with $Z, Z' \in \mathbb{R}^{n \times d}$; $Z_j \in \mathbb{R}^n$ is a vector composed of the values at dimension j for all n vectors in Z ; $C(Z) = 1/(n-1) \sum_i (\mathbf{z}^i - \bar{\mathbf{z}})(\mathbf{z}^i - \bar{\mathbf{z}})^T$ is the (sample) covariance matrix of Z with $\bar{\mathbf{z}} = 1/n \sum_{i=1}^n \mathbf{z}^i$; and λ_v, λ_c are hyperparameters for weighting the variance and covariance terms, respectively.

Algorithm	$\mathcal{L}^{\text{inv}}(Z, Z')$	$\mathcal{L}^{\text{ent}}(Z, Z')$
SimCLR	$-\frac{1}{n} \sum_{i=1}^n \frac{\mathbf{z}_i^T \mathbf{z}'_i}{\ \mathbf{z}_i\ \ \mathbf{z}'_i\ }$	$\frac{1}{n} \sum_{i=1}^n \log \sum_{j=1, j \neq i}^n \exp \left(\frac{\mathbf{z}_i^T \mathbf{z}'_j}{\ \mathbf{z}_i\ \ \mathbf{z}'_j\ } \right)$
BTs	$\sum_{j=1}^d \left(1 - \frac{(Z_j)^T Z'_j}{\ Z_j\ \ Z'_j\ } \right)^2$	$\sum_{j=1}^d \sum_{k=1, k \neq j}^d \left(\frac{(Z_j)^T Z'_k}{\ Z_j\ \ Z'_k\ } \right)^2$
VICReg	$-\frac{1}{n} \sum_{i=1}^n \ \mathbf{z}_i - \mathbf{z}'_i\ _2^2$	$\frac{\lambda_v}{d} \left(\sum_{j=1}^d \max(0, 1 - \sqrt{\text{Var}(Z_j) + \epsilon}) + \max(0, 1 - \sqrt{\text{Var}(Z'_j) + \epsilon}) \right) + \frac{\lambda_c}{d} \left(\sum_{i=1}^n \sum_{j=1, j \neq i}^n [C(Z)]_{i,j}^2 + [C(Z')]_{i,j}^2 \right)$

C.3 Adaptive λ_m

We now describe our procedure for adaptively updating our λ_m hyperparameters in Eq. (3.2) using a dual-ascent approach. To motivate this approach, first note that our placement of λ_m on the invariance terms \mathcal{L}^{inv} differs from the standard approach of placing it on the entropy term \mathcal{L}^{ent} (Wang and Isola, 2020; Zbontar et al., 2021). Doing so allows us to:

- **View \mathcal{L}^{inv} as a constraint that should be satisfied.** We view the goal of Eq. (2.1) as the soft/unconstrained version of the following constrained problem: maximize \mathcal{L}^{ent} , subject to $\mathcal{L}^{\text{inv}} = 0$. As a result, we then view λ_m as a Lagrange multiplier which should be set such that the invariance constraint is satisfied to within some acceptable tolerance ϵ , i.e., $\mathcal{L}^{\text{inv}} < \epsilon$. This way of choosing λ diverges from the standard approach to choosing the invariance-entropy trade-off in SSL (implicitly or explicitly), where it is chosen to maximize performance on some downstream task (e.g., ImageNet object classification accuracy).
- **Iteratively updating λ_m during training using a dual-ascent approach.** While we could take a standard grid-search approach to choose λ_m such that this invariance constraint is satisfied at the end of training, we instead iteratively adapt λ_m during training using a dual-ascent approach. In particular, given a step size or learning rate η and tolerance level ϵ , we perform iterative gradient-based updates of both the model parameters θ (inner loop) and λ_m (outer loop) with $\lambda_m^t \leftarrow \lambda_m^{t-1} + \eta \cdot \text{relu}(\mathcal{L}^{\text{inv}}(\theta^t) - \epsilon)$.

C.4 Numerical dataset

Following von Kügelgen et al. (2021, Sec. 5.1), we generate synthetic data pairs $(\mathbf{x}, \tilde{\mathbf{x}}) = (f(\mathbf{c}, \mathbf{s}), f(\mathbf{c}, \tilde{\mathbf{s}}))$ with content $\mathbf{c} \sim \mathcal{N}(0, \Sigma_c)$, style $\mathbf{s} | \mathbf{c} \sim \mathcal{N}(\mathbf{a} + B\mathbf{c}, \Sigma_s)$, and perturbed style $\tilde{\mathbf{s}} \sim \mathcal{N}(\mathbf{s}, \Sigma_{\tilde{\mathbf{s}}})$. We choose the simplest setup with Σ_c , Σ_s and $\Sigma_{\tilde{\mathbf{s}}}$ set to the identity. See von Kügelgen et al. (2021, App. D) for further details on the data-generation process.

C.5 ImageNet

Pretraining. Our ImageNet1k pretraining setup is based on the settings in (Bardes et al., 2022), which can be consulted for full details. We train ResNet50 (He et al., 2016) models for only 100 epochs, with 3-layer projectors of dimension 8196. The optimizer is LARS (You et al., 2017; Goyal et al., 2017), the batch size is 2048 and the learning rate follows a cosine decay schedule (Loshchilov and Hutter, 2017).

The data augmentation also follows Bardes et al. (2022) and is applied asymmetrically to the two views. It includes crops, flips, color jitter, grayscale, solarize and blur. These atomic augmentations are split into two groups: *spatial* (crops and flips) and *appearance* (color jitter, grayscale, solarize and blur). Thus, the number of “style” attributes in this setting are $M = 2$.

While we aim for fair experiments that use default hyperparameters, projectors, and augmentation settings, we note that these are optimized for existing SSL methods that prioritize information removal. Perhaps other settings, such as different augmentations explored in Xiao et al. (2021) and Lee et al. (2021), can be beneficial in our framework which instead aims to retain and disentangle information.

Downstream evaluation. Our downstream evaluation follows that of Ericsson et al. (2021). We train linear models (logistic or ridge regression) on frozen pre-projector representations \mathbf{h} and post-projector embeddings \mathbf{z} . Images are cropped to 224×224 , with L_2 regularization searched using 5-fold cross-validation over 45 logarithmically spaced values in the range 10^{-6} to 10^5 .

D Related Work

D.1 Self-supervised learning and disentanglement

Self-supervised learning. Xiao et al. (2021) also learn multiple embedding spaces in order to capture style information. In our work, we further develop these ideas towards a fully disentangled embedding space through a different use of augmentations and embedding spaces, as well as a different objective function—see App. D.2 for a detailed comparison with Xiao et al. (2021). Other prior work sought to retain some style information by predicting the augmentation parameters (Lee et al., 2020, 2021), seeking transformation equivariance (Dangovski et al., 2022), or employing techniques that improve

Table 4: **High-level comparison with Xiao et al. (2021)**. While both use structured augmentations and multiple embedding spaces to capture style attributes of the data, only ours seeks disentangled embedding spaces and provides theoretical grounding/analyses.

Method	Structured augmentations	Multiple embeddings	Disentangled embeddings	Theoretical underpinning
Xiao et al. (2021)	✓	✓	✗	✗
Ours	✓	✓	✓	✓

performance when using a linear projector (Jing et al., 2022). Importantly, we seek to both retain and separate/disentangle style information using a theoretically-grounded framework.

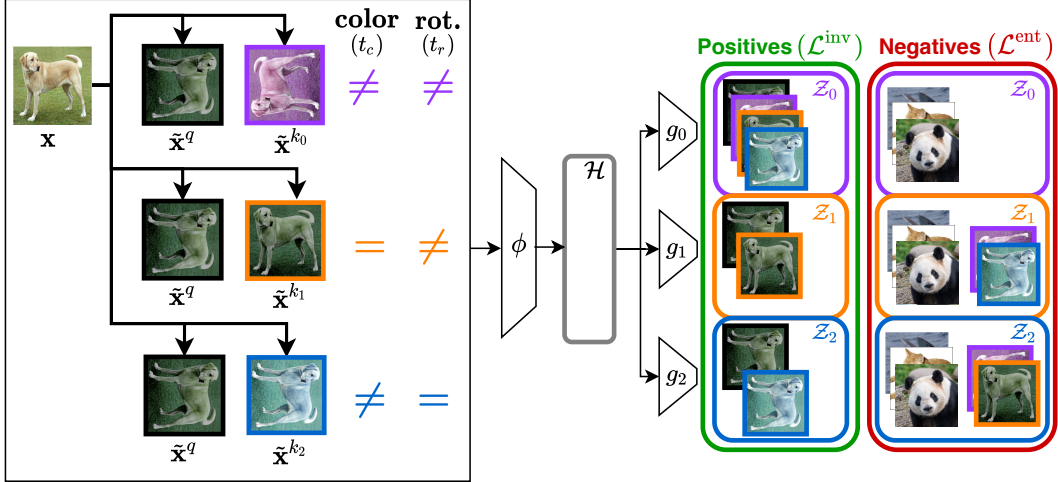
Generative disentanglement. In a generative setting, disentangled representations are commonly sought (Desjardins et al., 2012; Bengio et al., 2013; Higgins et al., 2017; Eastwood and Williams, 2018; Locatello et al., 2019), with the separation of “content” and “style” long sought in the vision-as-inverse-graphics paradigm (Tenenbaum and Freeman, 1996; Kulkarni et al., 2015). More recently, this generative disentanglement has been provably achieved with weak supervision in the form of paired data (Bouchacourt et al., 2018; Locatello et al., 2020), which is perhaps the setting that is most related to our work.

Identifiability in disentangled and causal representation learning. Our Thm. A.2 can be viewed as an extension of the content block-identifiability result of von Kügelgen et al. (2021, Thm. 4.4), which was generalized to a multi-modal setting with distinct mixing functions $f_1 \neq f_2$ and additional modality-specific latents by Daunhawer et al. (2023, Thm. 1). The two options discussed at the end of App. A for satisfying assumption A_2 —(a) independent style variables, and (b) perfect interventions—can be used to draw additional links to existing identifiability results. Option (a) relates to a result of Lyu et al. (2021, Thm. 2) showing that the entire style blocks $\tilde{\mathbf{s}}$ and $\tilde{\mathbf{s}}'$ can be block-identified through latent correlation maximization with invertible encoders, provided that \mathbf{c} , $\tilde{\mathbf{s}}$, and $\tilde{\mathbf{s}}'$ are mutually independent. Thm. A.2 establishes a more fine-grained disentanglement into individual style components. On the other hand, Gresele et al. (2019) and Locatello et al. (2020) prove identifiability of individual latents for the setting in which all latents are mutually independent and subject to change (with probability > 0), i.e., without an invariant block of content latents. Option (b) relates to a result of Brehmer et al. (2022, Thm. 1) showing that all variables (and the graph) in a causal representation learning setup can be identified through weak supervision in the form of pairs $(\mathbf{x}, \tilde{\mathbf{x}})$ arising from single-node perfect interventions by fitting a generative model via maximum likelihood. Perhaps most closely related is the work of Ahuja et al. (2022) who do not assume independence of latents, and also consider learning from M views arising from sparse perturbations, but require perturbations on all latent blocks for full identifiability.

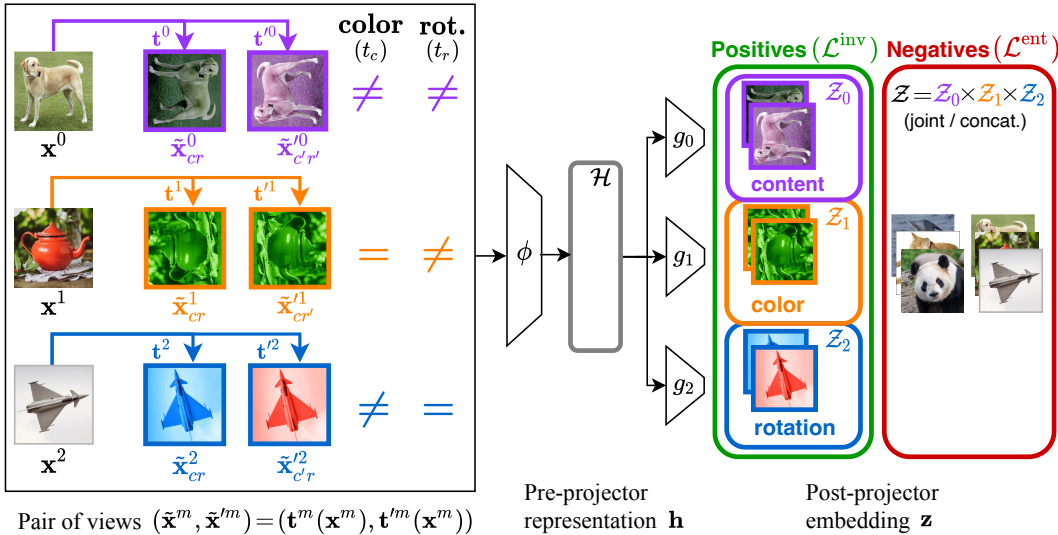
D.2 Detailed comparison with Xiao et al. (2021)

Table 4 presents the key differences between the framework of Xiao et al. (2021) and ours. Both rely on **structured augmentations**, by which views are constructed using augmentations that either share or do not share the same parameters. Both frameworks also learn **multiple embedding spaces** to capture style attributes of the data. However, our goal is not only to learn multiple embedding spaces but to **fully disentangle them**. This is achieved in our framework by the careful combination of invariance and entropy terms in Eq. (3.2), including the removal of redundant information with an entropy term across the joint embedding space. Furthermore, in Apps. A and B, we provide a **theoretical analysis** with the conditions under which our framework *identifies* the underlying style attributes or features.

In addition to these key, high-level differences, there are several smaller differences at the implementation level. In particular, we use an optimization procedure that adaptively sets the λ_m hyperparameters to guarantee the disentanglement of content and style (see App. C.3). We also adopt a more general construction of our framework, instantiating it with multiple different SSL methods, SimCLR, BTs and VICReg. Finally, our construction of image views allows more negative samples and in a given batch, compared to the query-key construction of Xiao et al. (2021)—see Fig. 4.



(a) Xiao et al. (2021)



(b) Ours

Figure 4: **Comparison with Xiao et al. (2021)**. Note the differences in data augmentation modules, as well as the embedding spaces in which positives and negatives are compared. See Xiao et al. (2021, Sec. 3) for details on their query-key notation. See App. D.2 for further details on this comparison.

E Further Results

We now present additional results.

E.1 ColorDSprites

Table 5 shows VICReg’s sensitivity to augmentation strengths with fixed and adaptive λ on ColorD-Sprites, complementing the results for SimCLR in Table 1.

Table 5: VICReg’s sensitivity to augmentation strengths with fixed and adaptive λ on ColorDSprites. r^2 in predicting the ground-truth factor values from the post-projector embedding \mathbf{z} with a linear classifier. Adapting λ ensures that \mathbf{z} captures all of content ($C = 1$) and almost no style ($\bar{S} = 0$), regardless of the augmentation strengths. SimCLR results in Table 1.

Algorithm	λ	Augm. Str.	Content (C)		Style (S)					\bar{S} (\downarrow)
			Shape	Color	Scale	Orient.	PosX	PosY		
VICReg	fixed	Weak	1.0	0.87	0.71	0.29	0.45	0.45	0.55	
	fixed	Medium	1.0	0.40	1.00	0.05	0.56	0.56	0.51	
	fixed	Strong	1.0	0.12	0.99	0.08	0.62	0.62	0.49	
	adaptive	Weak	1.0	0.20	0.17	0.00	0.00	0.00	0.07	
	adaptive	Medium	1.0	0.10	0.52	0.00	0.00	0.01	0.13	
	adaptive	Strong	1.0	0.10	0.53	0.00	0.00	0.00	0.13	

E.2 ImageNet

Table 6 give VICReg’s linear evaluation results on ImageNet for a broad range of downstream tasks, complementing the results for SimCLR in Table 2.

Table 6: **Linear evaluation on ImageNet and a broad range of downstream tasks.** We show top-1 accuracies (%) for all but CUB_{bbox} (r^2), CUB_{kpt} (r^2), and VOC (AP₅₀). We use frozen representations \mathbf{h} and embeddings \mathbf{z} (post-projector). FT: our framework *fine-tunes* a base VICReg model. Ct101: CalTech101. Cf10: CIFAR10.

Alg.	Feat.	ImNt	Acft	Ct101	Cars	Cf10	Cf100	CUB _{bbox}	CUB _{cls}	CUB _{kpt}	DTD	Flwrs	Pets	SUN	VOC	Avg.
VICReg	z	55.7	10.6	69.5	9.5	75.0	48.3	27.6	17.1	10.8	64.6	61.3	68.4	46.1	75.6	45.0
VICReg-Ours-FT	z	55.3	11.7	72.6	11.1	78.5	54.4	32.5	17.8	11.4	66.5	66.8	68.3	48.0	75.5	47.3
VICReg	h	67.2	51.1	87.6	52.6	88.3	70.1	69.1	47.4	31.9	75.3	93.4	83.1	59.7	79.6	68.4
VICReg-Ours-FT	h	66.9	50.1	87.5	52.4	88.4	70.3	69.8	47.2	32.8	75.7	93.6	82.7	59.8	79.6	68.5