
DIETing: Self-Supervised Learning with Instance Discrimination Learns Identifiable Features

Anonymous Author(s)

Affiliation

Address

email

Abstract

Self-Supervised Learning (SSL) methods often consist of elaborate pipelines with hand-crafted data augmentations and computational tricks. However, it is unclear what is the provably minimal set of building blocks that ensures good downstream performance. The recently proposed instance discrimination method, coined DIET, stripped down the SSL pipeline and demonstrated how a simple SSL algorithm can work by predicting the sample index. Our work proves that DIET recovers cluster-based latent representations, while successfully identifying the correct cluster centroids in its classification head. We demonstrate the identifiability of DIET on synthetic data adhering to and violating our assumptions, revealing that the recovery of the cluster centroids is even more robust than the feature recovery.

1 Introduction

Self-Supervised Learning (SSL) methods use unlabeled datasets to learn representations by solving an auxiliary task, thus bypassing time-consuming labelling efforts. Importantly, co-occurrence-based SSL relies on positive data pairs (similar samples, e.g., an original sample and a transformed/augmented one) and negative data pairs (dissimilar samples, often randomly drawn from the dataset). Contrastive and non-contrastive learning, the two prominent families of SSL methods, utilize positives and negatives differently, though they are theoretically connected [Balestriero and LeCun, 2022]. Contrastive Learning (CL) [Chen et al., 2020, Zimmermann et al., 2021, von Kügelgen et al., 2021, Lyu et al., 2021, Eastwood et al., 2023] attracts positive pairs’ and repels negative pairs’ representations. Non-contrastive learning [Bardes et al., 2021, Zbontar et al., 2021, Mialon et al., 2022] only uses positive pairs, and avoids representation collapse with strategies such as momentum encoders or covariance regularization. Unfortunately, the many actively developed Self-Supervised Learning methods with such computational tricks potentially hinder selecting the best performing and simplest SSL method for a given task. Recently, Ibrahim et al. [2024] proposed DIET, a SSL method that strips away unnecessary details by reducing the auxiliary task to a simple instance classification paradigm, and showed competitive performance on small datasets.

Identifiability theory, particularly Independent Component Analysis (ICA) [Comon, 1994, Hyvarinen et al., 2001] studies guarantees of probabilistic models to recover the ground-truth latent variables in a probabilistic latent variable model (LVM). Recent advances in nonlinear ICA theory proposed multiple self-supervised/weakly supervised models with identifiability guarantees [Hyvarinen et al., 2019, Gresele et al., 2019, Khemakhem et al., 2020a, Hälvä et al., 2021, Hyvarinen and Morioka, 2016, Khemakhem et al., 2020b, Locatello et al., 2020, Morioka and Hyvarinen, 2023, Morioka et al., 2021]. Several papers study a contrastive scenario, [Hyvarinen and Morioka, 2016, Hyvarinen et al., 2019, Zimmermann et al., 2021, von Kügelgen et al., 2021, Rusak et al., 2024], providing a possible theoretical explanation for CL’s practical success.

Our paper investigates whether DIET’s competitive performance can be explained by identifiability theory. We model the data generating process (DGP) in a new, cluster-based way, and show that

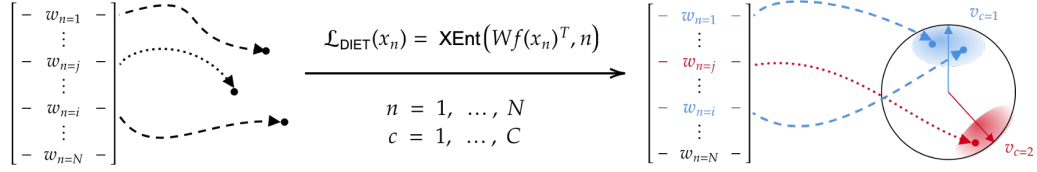


Figure 1: **DIET [Ibrahim et al., 2024] learns identifiable features:** DIET learns a linear $(N \times d)$ -dimensional classification head \mathbf{W} on top of a nonlinear encoder \mathbf{f} through an instance discrimination objective (1). For unit-normalized $\mathbf{f}(x_n)$, DIET maps samples and their augmentations close to the cluster vector \mathbf{v}_c corresponding to the class as if sampled from a von Mises-Fisher (vMF) distribution, centered around the cluster vector. In case of duplicate samples, i.e., matching class labels, the corresponding rows of \mathbf{W} will be the same, as shown for x_1 and x_i with $w_1 = w_i$

DIET’s learned representation is linearly related to the ground truth representation. We also show how DIET’s classification head recovers the cluster centroids, a connection to clustering that is absent from prior identifiability works for Self-Supervised Learning. Unlike other SSL solutions such as SimCLR [Chen et al., 2020], BYOL [Grill et al., 2020], BarlowTwins [Zbontar et al., 2021], or VICReg [Bardes et al., 2021], DIET’s training objective applies to the same representation that is used post-training for solving downstream tasks. More precisely, no projector network is removed post-training. This implies that our theoretical guarantees directly apply to the SSL representation being used post-training, as opposed to other identifiability results in SSL [Zimmermann et al., 2021, von Kügelgen et al., 2021, Daunhawer et al., 2023, Rusak et al., 2024]. We corroborate our theoretical claims on synthetic data adhering to our assumptions—we even show that good performance is possible when the assumptions are violated. Notably, we observe that cluster centroids recovery from DIET’s classification head is more robust than ground-truth representation prediction from the learned representation.

2 Identifiability guarantees for DIET

This section presents our main theoretical contribution. After summarizing DIET, we introduce a mildly constrained theoretical setup, in which DIET provably recovers the correct latents. The setup is followed by the main result and a discussion on the intuition for our theoretical model.

DIET [Ibrahim et al., 2024]. DIET solves an instance classification problem, where each sample x in the training dataset has a unique instance label i . Augmentations do not affect this label. We have a composite model $\mathbf{W} \circ \mathbf{f}$, where the backbone \mathbf{f} produces d -dimensional representations, and a linear, bias-free classification head \mathbf{W} that maps these representations to a logit vector equal in size to the cardinality of the training dataset. If the parameter vector corresponding to logit i is denoted as w_i , then \mathbf{W} effectively computes similarity scores (scalar products) between the w_i ’s and embeddings $\mathbf{f}(x)$. DIET trains this architecture to predict the correct instance label using multinomial regression (with \mathbf{f} , \mathbf{W} and temperature β as variables):

$$\mathcal{L}(\mathbf{f}, \mathbf{W}, \beta) = \mathbb{E}_{(x,i)} \left[-\ln \frac{e^{\beta \langle w_i, \mathbf{f}(x) \rangle}}{\sum_j e^{\beta \langle w_j, \mathbf{f}(x) \rangle}} \right]. \quad (1)$$

Setup. For our theory, we need to formally define an latent variable model (LVM) for the data generating process (DGP) to assess the identifiability of latent factors. For this, we take a cluster-centric approach, representing semantic classes by cluster vectors, similar to proxy-based metric learning [Kirchhof et al., 2022]. Then, we model the samples of a class with a von Mises-Fisher (vMF) distribution, centered around the class’s cluster vector. This conditional distribution jointly models intra-class sample selection and *augmentations* of samples, together called *intra-class variances*. We provide an overview of our assumptions, and defer additional details to Assums. 1C in Appx. A:

Assumptions 1 (DGP with vMF samples around cluster vectors. *Details omitted.*).

- (i) There is a finite set of semantic classes \mathcal{C} , represented by a set of unit-norm d -dimensional cluster-vectors $\{\mathbf{v}_c | c \in \mathcal{C}\} \subseteq \mathbb{S}^{d-1}$. The system $\{\mathbf{v}_c\}$ is sufficiently large and spread out.
- (ii) Any sample i belongs to exactly one class $c = \mathcal{C}(i)$.

74 (iii) The latent $\mathbf{z} \in \mathbb{S}^{d-1}$ of our data sample with instance label i is drawn from a vMF distribution
 75 around the cluster vector \mathbf{v}_c of class $c = \mathcal{C}(i)$:

$$\mathbf{z} \sim p(\mathbf{z}|c) \propto e^{\alpha \langle \mathbf{v}_c, \mathbf{z} \rangle}. \quad (2)$$

76 (iv) Sample \mathbf{x} is generated by passing latent \mathbf{z} through an injective generator function: $\mathbf{x} = \mathbf{g}(\mathbf{z})$.

77 **Main result.** Under Assums. 1, we prove the identifiability of both the latent representations and
 78 the cluster vectors, \mathbf{v}_c , in all four combinations of unit-normalized (i.e., when the latent space is the
 79 hypersphere, commonly used, e.g., in InfoNCE [Chen et al., 2020]); and non-normalized (as in the
 80 original DIET paper [Ibrahim et al., 2024]) latents, \mathbf{z} , and weight vectors, \mathbf{w}_i . We state a concise
 81 version of our result and defer the full treatment and the proof to Thm. 1C in Appx. A:

82 **Theorem 1** (Identifiability of latents drawn from vMF around cluster vectors. *Details omitted.*). *Let*
 83 *$(\mathbf{f}, \mathbf{W}, \beta)$ globally minimize the DIET objective (1) under the following additional constraints:*

84 *C3. the embeddings $\mathbf{f}(\mathbf{x})$ are unnormalized, while the \mathbf{w}_i ’s are unit-normalized. Then \mathbf{w}_i identifies*
 85 *the cluster vector $\mathbf{v}_{\mathcal{C}(i)}$ up to an orthogonal linear transformation \mathcal{O} : $\mathbf{w}_i = \mathcal{O}\mathbf{v}_{\mathcal{C}(i)}$, for any i .*
 86 *Furthermore, the inferred latents $\tilde{\mathbf{z}} = \mathbf{f}(\mathbf{x})$ identify the ground-truth latents \mathbf{z} up to the same*
 87 *orthogonal transformation, but scaled.*

88 *C4. neither the embeddings $\mathbf{f}(\mathbf{x})$ nor the \mathbf{w}_i ’s are unit-normalized. Then the cluster vectors \mathbf{v}_c and*
 89 *the latent \mathbf{z} are identified up to an affine linear and linear transformation, respectively.*

90 *In all cases, the weight vectors belonging to samples of the same class are equal, i.e., for any i, j ,*
 91 *$\mathcal{C}(i) = \mathcal{C}(j)$ implies $\mathbf{w}_i = \mathbf{w}_j$.*

92 **Intuition.** DIET assigns a different (instance) label and a unique weight vector \mathbf{w}_i to each training
 93 sample. The cross-entropy objective is optimized if the trained neural network can distinguish
 94 between the samples. Thus, the learned representation $\tilde{\mathbf{z}} = \mathbf{f}(\mathbf{x})$ should capture enough information
 95 to distinguish different samples, even from the same class.

96 However, the weight vectors \mathbf{w}_i ’s cannot be sensitive to the intra-class sample variance or the sample’s
 97 instance label i (because multiple instances will usually belong to the same class). This leads to the
 98 weight vectors taking the values of the cluster vectors. As cluster vectors only capture some statistics
 99 of the conditional, feature recovery is more fine-grained than cluster identifiability. The interaction
 100 between the two is dictated by the cross-entropy loss, which is minimized if the representation $\tilde{\mathbf{z}}$
 101 is most similar to its own assigned weight vector \mathbf{w}_i . Fig. 1 provides a visualization conveying the
 102 intuition behind Thm. 1.

103 3 Experiments

104 In the following section, we empirically verify the claims made in Thm. 1 in the synthetic setting.
 105 We generate data samples according to Assums. 1: ground-truth latents are sampled around cluster
 106 centroids \mathbf{v}_c following a vMF distribution. Data augmentations, which share the same instance label
 107 i , are sampled from the same vMF distribution around \mathbf{v}_c .

108 **Synthetic Setup.** We consider N data samples of dimensionality d generated from $\mathbf{z} \sim p(\mathbf{z}|\mathbf{v}_c)$,
 109 sampled around a set of $|\mathcal{C}|$ class vectors, \mathbf{v}_c uniformly distributed across the unit hyper-sphere. We
 110 use an invertible multi-layer perceptron (MLP) to map ground truth latents to data samples. We
 111 train a classification head $\mathbf{W} = [\mathbf{w}_i^\top]_{i=1}^N$ and an MLP encoder that maps samples to representations
 112 $\tilde{\mathbf{z}} \in \mathbb{R}^d$ using the DIET objective (1). While to verify Thm. 1 case C4., we do not normalize \mathbf{W} , we
 113 do unit-normalize the weight vectors to validate Thm. 1 case C3. We verify our theoretical claims by
 114 measuring the predictability of the ground-truth \mathbf{z} from $\tilde{\mathbf{z}}$ and \mathbf{v}_c from \mathbf{w}_i using the R^2 score on a
 115 held-out dataset. For identifiability up to orthogonal linear transformations, we train linear mappings
 116 with no intercept, assess the R^2 score and verify that the singular values of this transformation
 117 converge to one, while for identifiability up to affine linear transformations, we simply assess the
 118 predictive accuracy of a linear predictor with intercept.

119 **Results.** Tab. 1 depicts our results for synthetic experiments. For both cases, when \mathbf{W} is and
 120 is not unit-normalized, the R^2 score for both the latents and the cluster vectors is close to 100%,
 121 except when the latent dimensionality is 20—such scalability problems are a common artifact in
 122 SSL [Zimmermann et al., 2021, Rusak et al., 2024]. For unit-normalized \mathbf{W} , the MAE is close to
 123 zero even in such cases. We also observe that for a higher concentration of samples around \mathbf{v}_c (i.e.

Table 1: Identifiability in the synthetic setup. Mean \pm standard deviation across 5 random seeds. Settings that match and violate our theoretical assumptions are \checkmark and \times respectively. We report the R^2 score for linear mappings, $\tilde{z} \rightarrow z$ and $w_i \rightarrow v_c$ for cases with normalized (o) and not normalized (a) w_i . For normalized w_i , we verify that mappings $\tilde{z} \rightarrow z$ are orthogonal by reporting the mean absolute error between their singular values and those of an orthogonal transformation.

N	d	$ \mathcal{C} $	$p(z v_c)$	M.	normalized w_i cases				unnormalized w_i	
					$R_o^2(\uparrow)$	$\tilde{z} \rightarrow z$	$w_i \rightarrow v_c$	MAE _o (\downarrow)	$\tilde{z} \rightarrow z$	$w_i \rightarrow v_c$
10^3	5	100	vMF($\kappa=10$)	\checkmark	98.6 ± 0.01	99.9 ± 0.00	0.01 ± 0.00	0.00 ± 0.00	99.0 ± 0.00	99.9 ± 0.00
10^5	5	100	vMF($\kappa=10$)	\checkmark	98.2 ± 0.01	99.5 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	99.7 ± 0.00	99.8 ± 0.00
10^3	5	100	vMF($\kappa=10$)	\checkmark	98.6 ± 0.01	99.9 ± 0.00	0.01 ± 0.00	0.00 ± 0.00	99.0 ± 0.00	99.9 ± 0.00
10^3	10	100	vMF($\kappa=10$)	\checkmark	92.5 ± 0.01	99.6 ± 0.00	0.01 ± 0.00	0.00 ± 0.00	93.0 ± 0.03	99.6 ± 0.00
10^3	20	100	vMF($\kappa=10$)	\checkmark	70.8 ± 0.02	97.1 ± 0.01	0.03 ± 0.00	0.00 ± 0.00	81.9 ± 0.01	99.7 ± 0.00
10^3	5	10	vMF($\kappa=10$)	\checkmark	88.6 ± 0.05	85.7 ± 0.15	0.02 ± 0.00	0.00 ± 0.00	90.0 ± 0.05	99.0 ± 0.03
10^3	5	100	vMF($\kappa=10$)	\checkmark	98.6 ± 0.01	99.9 ± 0.01	0.01 ± 0.00	0.00 ± 0.00	99.0 ± 0.00	99.9 ± 0.00
10^3	5	1000	vMF($\kappa=10$)	\checkmark	99.3 ± 0.00	99.9 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	99.2 ± 0.00	99.9 ± 0.00
10^3	5	100	vMF($\kappa=5$)	\checkmark	98.6 ± 0.01	99.9 ± 0.01	0.01 ± 0.00	0.00 ± 0.00	0.01 ± 0.00	0.00 ± 0.00
10^3	5	100	vMF($\kappa=10$)	\checkmark	99.0 ± 0.00	99.9 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00
10^3	5	100	vMF($\kappa=50$)	\checkmark	45.0 ± 0.06	49.7 ± 0.06	0.30 ± 0.00	0.00 ± 0.00	0.30 ± 0.00	0.00 ± 0.00
10^3	5	100	vMF($\kappa=10$)	\checkmark	98.6 ± 0.01	99.9 ± 0.01	0.01 ± 0.00	0.00 ± 0.00	99.0 ± 0.00	99.9 ± 0.00
10^3	5	100	Laplace ($b=1.0$)	\times	85.2 ± 0.01	99.7 ± 0.01	0.01 ± 0.00	0.00 ± 0.00	85.4 ± 0.00	99.5 ± 0.00
10^3	5	100	Normal ($\sigma^2=1.0$)	\times	98.7 ± 0.00	99.8 ± 0.00	0.01 ± 0.00	0.00 ± 0.00	98.6 ± 0.00	99.6 ± 0.00

$\kappa=50$) as well as lower number of clusters (i.e. $|\mathcal{C}|=10$), identifiability suffers (i.e., the R^2 score decreases), which is also a common phenomenon, and is possibly explained by the content-style partitioning of latents [von Kügelgen et al., 2021] and insufficient augmentation overlap [Wang et al., 2022, Rusak et al., 2024]. Our results also suggest that even under model misspecification (last two rows with non-vMF latent distributions), identifiability still holds. We provide an additional ablation study for the concentration of v_c across the unit hyper-sphere in Appx. B.

4 Discussion

Limitations. Our analysis proves the identifiability of DIET [Ibrahim et al., 2024] with a cluster-based DGP, thus providing the first such result for self-supervised parametric instance classification methods. However, our theory cannot yet explain the importance of label smoothing in DIET, noted by Ibrahim et al. [2024], and it also remains to be seen whether such identifiability results scale for larger datasets, for which the large-dimensional classifier head in DIET in the original form is prohibitive. It also remains an issue that the vMF conditional distribution around cluster centroids jointly models intra-class sample selection and augmentations of samples, as we suspect that the supports of augmentation spaces of different samples do not overlap as much as it would be suggested by the choice of conditional. Also, we leave it for future work to investigate a formal connection to nonlinear ICA methods such as InfoNCE [Zimmermann et al., 2021] or the Generalized Contrastive Learning framework [Hyvarinen et al., 2019].

Conclusion. By modeling the DGP in DIET [Ibrahim et al., 2024] with a cluster-based latent variable model, we provide identifiability results for both the latent representation and the cluster vectors, which is the first of its kind for self-supervised instance discrimination methods. We also showcase this in synthetic settings, where we recover both the latents and cluster vectors even under model misspecification. We hope that our work inspires further research into investigating the theoretical guarantees of simplified but effective SSL methods like DIET.

References

- Randall Balestriero and Yann LeCun. Contrastive and Non-Contrastive Self-Supervised Learning Recover Global and Local Spectral Embedding Methods, June 2022. URL <http://arxiv.org/abs/2205.11508>. arXiv:2205.11508 [cs, math, stat]. 1
- Adrien Bardes, Jean Ponce, and Yann LeCun. VICReg: Variance-Invariance-Covariance Regularization for Self-Supervised Learning. *arXiv:2105.04906 [cs]*, May 2021. URL <http://arxiv.org/abs/2105.04906>. arXiv: 2105.04906. 1, 2
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A Simple Framework for Contrastive Learning of Visual Representations. *arXiv:2002.05709 [cs, stat]*, June 2020. URL <http://arxiv.org/abs/2002.05709>. arXiv: 2002.05709. 1, 2, 3
- Pierre Comon. Independent component analysis, a new concept? *Signal processing*, 36(3):287–314, 1994. 1
- Imant Daunhawer, Alice Bizeul, Emanuele Palumbo, Alexander Marx, and Julia E. Vogt. Identifiability Results for Multimodal Contrastive Learning, March 2023. URL <http://arxiv.org/abs/2303.09166>. arXiv:2303.09166 [cs, stat] version: 1. 2
- Cian Eastwood, Julius von Kügelgen, Linus Ericsson, Diane Bouchacourt, Pascal Vincent, Bernhard Schölkopf, and Mark Ibrahim. Self-Supervised Disentanglement by Leveraging Structure in Data Augmentations, November 2023. URL <http://arxiv.org/abs/2311.08815>. arXiv:2311.08815 [cs, stat]. 1
- Luigi Gresele, Paul K. Rubenstein, Arash Mehrjou, Francesco Locatello, and Bernhard Schölkopf. The Incomplete Rosetta Stone Problem: Identifiability Results for Multi-View Nonlinear ICA. *arXiv:1905.06642 [cs, stat]*, August 2019. URL <http://arxiv.org/abs/1905.06642>. arXiv: 1905.06642. 1
- Jean-Bastien Grill, Florian Strub, Florent Alché, Corentin Tallec, Pierre H. Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi Azar, Bilal Piot, Koray Kavukcuoglu, Rémi Munos, and Michal Valko. Bootstrap your own latent: A new approach to self-supervised Learning. *arXiv:2006.07733 [cs, stat]*, September 2020. URL <http://arxiv.org/abs/2006.07733>. arXiv: 2006.07733. 2
- Aapo Hyvarinen and Hiroshi Morioka. Unsupervised Feature Extraction by Time-Contrastive Learning and Nonlinear ICA. *arXiv:1605.06336 [cs, stat]*, May 2016. URL <http://arxiv.org/abs/1605.06336>. arXiv: 1605.06336. 1
- Aapo Hyvarinen, Juha Karhunen, and Erkki Oja. *Independent component analysis*. J. Wiley, New York, 2001. ISBN 978-0-471-40540-5. 1
- Aapo Hyvarinen, Hiroaki Sasaki, and Richard E. Turner. Nonlinear ICA Using Auxiliary Variables and Generalized Contrastive Learning. *arXiv:1805.08651 [cs, stat]*, February 2019. URL <http://arxiv.org/abs/1805.08651>. arXiv: 1805.08651. 1, 4
- Hermanni Hälvä, Sylvain Le Corff, Luc Lehéricy, Jonathan So, Yongjie Zhu, Elisabeth Gassiat, and Aapo Hyvarinen. Disentangling Identifiable Features from Noisy Data with Structured Nonlinear ICA. *arXiv:2106.09620 [cs, stat]*, June 2021. URL <http://arxiv.org/abs/2106.09620>. arXiv: 2106.09620. 1
- Mark Ibrahim, David Klindt, and Randall Balestriero. Occam’s Razor for Self Supervised Learning: What is Sufficient to Learn Good Representations?, June 2024. URL <http://arxiv.org/abs/2406.10743>. arXiv:2406.10743 [cs]. 1, 2, 3, 4
- Ilyes Khemakhem, Diederik Kingma, Ricardo Monti, and Aapo Hyvarinen. Variational Autoencoders and Nonlinear ICA: A Unifying Framework. In *International Conference on Artificial Intelligence and Statistics*, pages 2207–2217. PMLR, June 2020a. URL <http://proceedings.mlr.press/v108/khemakhem20a.html>. ISSN: 2640-3498. 1

195 Ilyes Khemakhem, Ricardo Pio Monti, Diederik P. Kingma, and Aapo Hyvärinen. ICE-BeeM:
196 Identifiable Conditional Energy-Based Deep Models Based on Nonlinear ICA. *arXiv:2002.11537*
197 [*cs, stat*], October 2020b. URL <http://arxiv.org/abs/2002.11537>. arXiv: 2002.11537. 1

198 Michael Kirchhof, Karsten Roth, Zeynep Akata, and Enkelejda Kasneci. A Non-isotropic Probabilistic
199 Take on Proxy-based Deep Metric Learning, July 2022. URL <http://arxiv.org/abs/2207.03784>. arXiv:2207.03784 [*cs, stat*]. 2

201 Francesco Locatello, Ben Poole, Gunnar Rätsch, Bernhard Schölkopf, Olivier Bachem, and Michael
202 Tschannen. Weakly-Supervised Disentanglement Without Compromises. *arXiv:2002.02886* [*cs,*
203 *stat*], October 2020. URL <http://arxiv.org/abs/2002.02886>. arXiv: 2002.02886. 1

204 Qi Lyu, Xiao Fu, Weiran Wang, and Songtao Lu. Latent Correlation-Based Multiview Learning
205 and Self-Supervision: A Unifying Perspective. *arXiv:2106.07115* [*cs, stat*], June 2021. URL
206 <http://arxiv.org/abs/2106.07115>. arXiv: 2106.07115. 1

207 Grégoire Mialon, Randall Balestriero, and Yann LeCun. Variance Covariance Regularization Enforces
208 Pairwise Independence in Self-Supervised Representations, September 2022. URL <http://arxiv.org/abs/2209.14905>. arXiv:2209.14905 [*cs*]. 1

210 Hiroshi Morioka and Aapo Hyvärinen. Connectivity-contrastive learning: Combining causal discov-
211 ery and representation learning for multimodal data. In *Proceedings of The 26th International*
212 *Conference on Artificial Intelligence and Statistics*, pages 3399–3426. PMLR, April 2023. URL
213 <https://proceedings.mlr.press/v206/morioka23a.html>. ISSN: 2640-3498. 1

214 Hiroshi Morioka, Hermanni Hälvä, and Aapo Hyvärinen. Independent Innovation Analysis for
215 Nonlinear Vector Autoregressive Process. *arXiv:2006.10944* [*cs, stat*], February 2021. URL
216 <https://arxiv.org/abs/2006.10944>. arXiv: 2006.10944. 1

217 Evgenia Rusak, Patrik Reizinger, Attila Juhos, Oliver Bringmann, Roland S. Zimmermann, and
218 Wieland Brendel. InfoNCE: Identifying the Gap Between Theory and Practice, June 2024. URL
219 <http://arxiv.org/abs/2407.00143>. arXiv:2407.00143 [*cs, stat*]. 1, 2, 3, 4

220 Julius von Kügelgen, Yash Sharma, Luigi Gresele, Wieland Brendel, Bernhard Schölkopf, Michel
221 Besserve, and Francesco Locatello. Self-Supervised Learning with Data Augmentations Provably
222 Isolates Content from Style, June 2021. URL <http://arxiv.org/abs/2106.04619>. arXiv:
223 2106.04619. 1, 2, 4

224 Yifei Wang, Qi Zhang, Yisen Wang, Jiansheng Yang, and Zhouchen Lin. Chaos is a Ladder: A New
225 Theoretical Understanding of Contrastive Learning via Augmentation Overlap, May 2022. URL
226 <http://arxiv.org/abs/2203.13457>. arXiv:2203.13457 [*cs, stat*]. 4

227 Wikipedia. Gibbs’ inequality, 2024a. URL [https://en.wikipedia.org/w/index.php?title=](https://en.wikipedia.org/w/index.php?title=Gibbs%27_inequality&oldid=1231436245)
228 [Gibbs%27_inequality&oldid=1231436245](https://en.wikipedia.org/w/index.php?title=Gibbs%27_inequality&oldid=1231436245). Online; accessed 10-September-2024. 9

229 Wikipedia. Tietze extension theorem, 2024b. URL [https://en.wikipedia.org/w/index.php?](https://en.wikipedia.org/w/index.php?title=Tietze_extension_theorem&oldid=1237682676)
230 [title=Tietze_extension_theorem&oldid=1237682676](https://en.wikipedia.org/w/index.php?title=Tietze_extension_theorem&oldid=1237682676). Online; accessed 10-September-
231 2024. 8

232 Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stéphane Deny. Barlow Twins: Self-Supervised
233 Learning via Redundancy Reduction. *arXiv:2103.03230* [*cs, q-bio*], June 2021. URL <http://arxiv.org/abs/2103.03230>. arXiv: 2103.03230. 1, 2

235 Roland S. Zimmermann, Yash Sharma, Steffen Schneider, Matthias Bethge, and Wieland Brendel.
236 Contrastive Learning Inverts the Data Generating Process. *arXiv:2102.08850* [*cs*], February 2021.
237 URL <http://arxiv.org/abs/2102.08850>. arXiv: 2102.08850. 1, 2, 3, 4