

---

# On Improving the Sample Efficiency of Non-Contrastive SSL

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 In this work, we provide theoretical insights on the implicit bias of the BarlowTwins  
2 and VICReg loss that can explain these heuristics and guide the development of  
3 more principled recommendations. Our first insight is that the orthogonality of the  
4 features is more important than projector dimensionality for learning good represen-  
5 tations. Based on this, we empirically demonstrate that low-dimensional projector  
6 heads are sufficient with appropriate regularization, contrary to the existing heuristic.  
7 Our second theoretical insight suggests that using multiple data augmentations  
8 better represents the desiderata of the SSL objective. Based on this, we demonstrate  
9 that leveraging more augmentations per sample improves representation quality  
10 and trainability. In particular, it improves optimization convergence, leading to  
11 better features emerging earlier in the training. Remarkably, we demonstrate that  
12 we can reduce the pretraining dataset size by up to 4x while maintaining accuracy  
13 and improving convergence simply by using more data augmentations. Combining  
14 these insights, we present pretraining recommendations that improve wall-clock  
15 time by 2x and downstream performance on CIFAR-10/STL-10 datasets.

## 16 1 Introduction

17 A prominent subgroup among non-contrastive SSL methods is the family of Canonical Correlation  
18 Analysis (CCA) algorithms, which includes BarlowTwins [Zbontar et al., 2021] and VICReg [Bardes  
19 et al., 2021]. These methods aim to enforce orthogonality among the learned features in addition to  
20 learning to map similar images to nearby points in feature space and have been shown to achieve  
21 competitive performance on benchmark computer vision datasets. These methods have become the  
22 preferred strategy for representation learning in several domains due to the lack of need for negative  
23 samples and their simple formulation. However, despite the apparent simplicity of their loss functions,  
24 the behavior of this family of algorithms is not well understood. Therefore, researchers often use  
25 empirically driven heuristics to design successful applications, such as using (i) a high-dimensional  
26 projector head or (ii) two augmentations per image.

27 Alongside relying on heuristics and researchers’ intuition for design, existing SSL algorithms are  
28 extremely data-hungry. In particular, state-of-the-art algorithms often rely on large-scale datasets  
29 [Russakovsky et al., 2015] or data engines [Oquab et al., 2023] to achieve good representations.  
30 While this strategy works exceptionally well in natural-image settings, its application is limited in  
31 other critical domains, such as medical imaging, where the number of samples is scarce.

32 With these challenges in mind, the primary focus of this work is making progress toward establishing  
33 theoretical foundations underlying the family of non-contrastive SSL algorithms (NC-SSL) with an  
34 eye toward sample efficiency. In particular, we analyse the BarlowTwins and VICReg losses and  
35 show that they implicitly learn the data similarity kernel that is defined by the chosen augmentations.  
36 We find that learning the data similarity kernel is helped by greater orthogonality in the projector

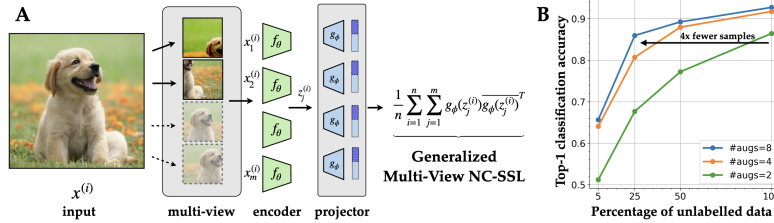


Figure 1: Existing SSL algorithms make design choices often driven by heuristics. (A) We investigate the theoretical underpinnings of two choices (i) the number of augmentations and (ii) the dimensionality of the projector. (B) We show that the generalized NC-SSL algorithm with multiple augmentations and low-dimensional projectors outperforms existing heuristics, using  $\sim 4\times$  fewer samples.

37 outputs and more data augmentations. As such, increasing the orthogonality of the projector output  
 38 eliminates the requirement for a high-dimensional projector head, and increasing the number of data  
 39 augmentations decreases the number of unique samples required.

40 We empirically verify our theoretical insights using the popular ResNet-50 backbone on benchmark  
 41 datasets, CIFAR-10 and STL-10. Strikingly, we show that our multi-augmentation approach can learn  
 42 good features even with a quarter of the number of samples in the pretraining dataset. In summary,  
 43 our core contributions are:

- 44 • **Eigenfunction interpretation:** We demonstrate that the loss functions of the CCA family  
 45 of non-contrastive SSL algorithms are equivalent to the objective of learning eigenfunctions  
 46 of the augmentation-defined data kernel.
- 47 • **Role of heuristics:** We provide a mechanistic explanation for the role of projector di-  
 48 mensionality and the number of data augmentations, and empirically demonstrate that  
 49 low-dimensional projector heads are sufficient and using more augmentations leads to  
 50 learning better representations.
- 51 • **Data efficient NC-SSL:** Leveraging the convergence benefits of the multi-augmentation  
 52 framework, we demonstrate that we can learn good features with significantly smaller  
 53 datasets (upto 25%) without harming downstream performance.

## 54 2 Data augmentation kernel perspective of non-contrastive SSL

55 We will define two notions of the data augmentation kernel. Given two images,  $x, z$ , the first kernel,  
 56 which we call the forward data augmentation covariance kernel, is given by

$$k^{DAF}(x, z) = \mathbb{E}_{x_0 \sim \rho_X} [p(x | x_0)p(z | x_0)] \quad (1)$$

57 This covariance kernel measures the similarity between  $x, z$  in terms of how likely they are to be  
 58 reached from  $x_0$ , weighted by the distribution of  $x_0$ . Note that this is indeed the edge strength  
 59 between nodes  $x, z$  in the augmentation graph. We can also define a (backward) data augmentation  
 60 covariance kernel  $k^{DAB}(x, z)$ , which reverses the roles of  $(x, z)$  and  $x_0$ .

61 SSL aims to learn features that preserve the covariance kernel structure (imposed by this choice of  
 62 mapping  $M$ ) [Dubois et al., 2022]. Therefore, we want to define a loss which determines *vector*  
 63 *features*,  $F : X \rightarrow \mathbb{R}^d$ , which factor a data augmentation kernel  $k^{DA}(x, z) = F(x)^\top F(z)$ . Doing  
 64 this directly is prohibitively data intensive at scale, since it involves a search over data augmented  
 65 images. However, since the covariance kernels are PSD, they define a Reproducing Kernel Hilbert  
 66 space (RKHS). This allows us to apply Mercer’s theorem to find vector features as in Deng et al.  
 67 [2022a,b], Pfau et al. [2018].

68 **Theorem 2.1.** *Let  $G(x)$  be the infinite Mercer features of the backward data augmentation covariance*  
 69 *kernels,  $k^{DAB}$ . Let  $F(x) = (f_1(x), f_2(x), \dots, f_k(x))$  be the features given by minimizing the*  
 70 *following data augmentation invariance loss*

$$L(F) = \sum_{i=1}^{N_k} \|T_M f_i - f_i\|_{L^2(\rho_X)}^2, \quad \text{subject to} \quad (f_i, f_j)_{\rho_X} = \delta_{ij} \quad (2)$$

71 *which includes the orthogonality constraint. Then,  $V(F) \subset V(G)$ ,  $V(F) \rightarrow V(G)$  as  $N_k \rightarrow \infty$ .*

## 72 3 Experiments

73 In our experiments, we seek to serve two purposes (i) provide empirical support for our theoretical  
 74 insights and (ii) present practical primitives for designing efficient self-supervised learning routines.  
 75 In summary, with extensive experiments across learning algorithms (BarlowTwins, VICReg) and  
 76 training datasets (CIFAR-10/STL-10), we establish that

- 77 • **low-dimensional projectors** as sufficient for learning *good representations*.
- 78 • multi-Augmentation **improves sample efficiency** in SSL pretraining, i.e. recovering similar  
 79 performance with significantly fewer unlabelled samples.

80 **Experiment Setup:** We evaluate the effectiveness of different pretraining approaches for non-  
 81 contrastive SSL algorithms using image classification as the downstream task. Across all experiments,  
 82 we use linear probing with Resnet-50 as the feature encoder backbone. On CIFAR-10, all models are  
 83 pretrained for 100 epochs, and STL-10 models are pretrained for 50 epochs (averaged over 3 seeds).

### 84 3.1 Sufficiency of Low-dimensional projectors

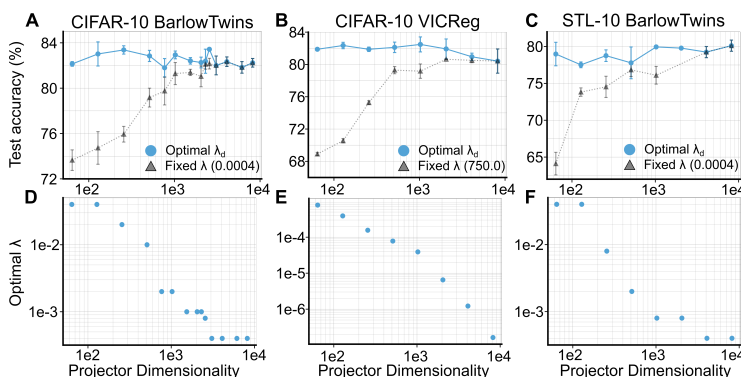


Figure 2: Low-dimensional projectors are sufficient for good feature learning. We demonstrate that using a higher orthogonality constraint ( $\lambda$  for D, F and  $\lambda_{eff} = \frac{1}{d\lambda}$  for E) for lower projector dimensionality can achieve similar performance over a wide range of projector dimensions ( $d$ ).

85 Existing works recommend using high-dimensional MLPs as projectors (e.g.,  $d=8192$  for Imagenet  
 86 in Zbontar et al. [2021], Bardes et al. [2021]), and show significant degradation in performance for a  
 87 fixed redundancy coefficient ( $\lambda$ ). To reproduce this result, we run a grid search to find the optimal  
 88 coefficient ( $\lambda_{8192}^*$ ) for  $d = 8192$  and show that performance progressively degrades for lower  $d$  if the  
 89 same coefficient  $\lambda_{8192}^*$  is reused for  $d \in \{64, 128, 256, 512, 1024, 2048, 4096, 8192\}$ .

90 Our insights in Appendix B.2 suggest low-dimensional projectors should recover similar performance  
 91 with appropriate orthogonalization. To test this, we find the best  $\lambda$  by performing a grid search  
 92 independently for each  $d \in \{64, 128, 256, 512, 1024, 2048, 4096, 8192\}$ . As illustrated in Figure 2,  
 93 low-dimensional projectors are indeed sufficient. Strikingly, we also observe that the optimal  
 94  $\lambda_d \propto 1/d$ , is in alignment with our theoretical insights.

### 95 3.2 Sample Efficient Multi-View Learning

96 Although some SSL pretraining approaches, like SWaV, incorporate more than two views, the most  
 97 widely used heuristic in non-contrastive SSL algorithms involve using two views jointly encoded by  
 98 a shared backbone. In line with this observation, our baselines for examining the role of multiple  
 99 augmentations use two views for computing the cross-correlation matrix.

100 To understand the role of multiple augmentations in pretraining in light of the augmentation-kernel  
 101 interpretation, we propose Equation (10), which generalizes Barlow-Twins and VICReg to the  
 102 multi-augmentation setting. In particular, for  $\#aug \in \{2, 4, 8\}$ , we pretrain Resnet-50 with the  
 103 generalized NC-SSL loss for 100 epochs on CIFAR-10 and 50-epochs for STL-10. Building on the  
 104 insight from the previous section, we use a 256-dimensional projector head for all experiments. Here,

105 we use the linear evaluation protocol as outlined by Chen et al. [2022]. In line with previous work,  
 106 we observe that pretraining with multiple augmentations outperforms the 2-augmentation baseline  
 107 (see Appendix). Although using more augmentations increases the per-epoch time during pretraining,  
 108 we observe that the four-augmentation pre-trained models achieve the same accuracy faster (both  
 109 in terms of the number of epochs and wall-clock time) than their two-augmentation counterparts.  
 110 Data Augmentation can be viewed as a form of data-inflation, where the number of training samples  
 111 is increased by a factor of  $k$  (for  $k$  augmentations). Therefore, we seek to investigate if multiple  
 112 augmentations in SSL pretraining pipeline can compensate for less unique samples in the dataset.

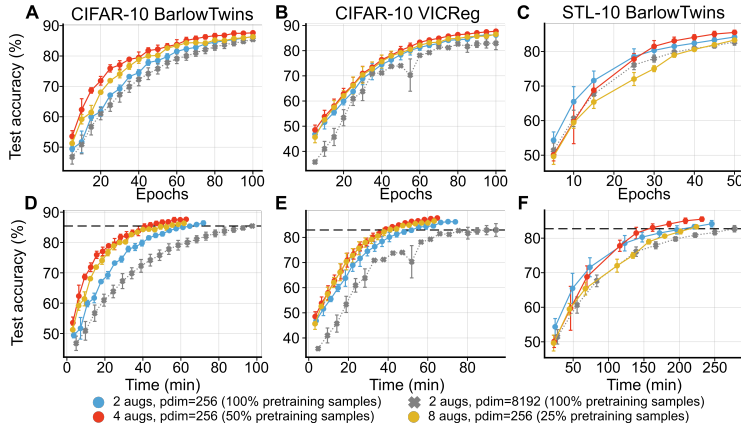


Figure 3: Multi-augmentation improves sample efficiency, recovering similar performance with significantly less number of unique samples in the pretraining dataset. Across BarlowTwins and VICReg pretraining on CIFAR-10 and STL-10, for the same effective dataset size ( $\#aug \times \#unique\_samples$ ), using more patches improves performance at the same epoch (A-C) or wall clock time (D-F). However, there exists a tradeoff wherein doing more data augmentations fails to improve performance in the very low data regime.

113 To this effect, we fixed the effective size of the inflated dataset by varying the fraction of the unique  
 114 samples in the pretraining dataset depending on the number of augmentations  $k \in \{2, 4, 8\}$ , e.g.  
 115 we use 1/2 the dataset for 4 views. We then evaluate the performance of the pre-trained models  
 116 on the downstream task, where the linear classifier is trained on the same set of labeled samples.  
 117 Strikingly, Figure 3 shows that using multiple augmentations can achieve similar (sometimes even  
 118 better) performance with lesser pretraining samples, thereby indicating that more data augmentations  
 119 can be used to compensate for smaller pretraining datasets.

## 120 4 Discussion

**Pareto Optimal SSL** In the context of sample efficiency, training a model using two augmentations with different fractions of the dataset leads to a natural Pareto frontier, i.e. training on the full dataset achieves the best error but takes the most time (**Baseline (2-Aug)**). Our extensive experiments demonstrate that using more than two augmentations improves the overall Pareto frontier, i.e. achieves better convergence while maintaining accuracy (**Multi-Aug**). Strikingly, as shown in Figure 4, we observe that for a target error level, we can either use a larger pretraining dataset or more augmentations. Therefore, the number of augmentations can be used as a knob to control the sample efficiency of the pretraining routine.

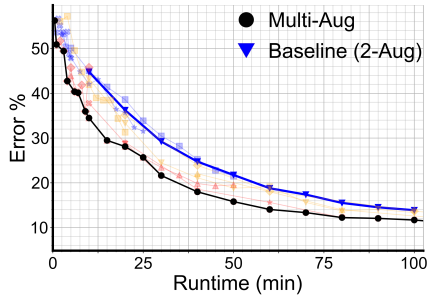


Figure 4: Using  $> 2$  augmentations with a fraction of dataset improves Pareto frontier, with runtime boost by  $\sim 2\times$ .

121 **Limitations** Our algorithm relies on multiple views of the same image to improve the estimation of  
 122 the data-augmentation kernel. Although this approach does add some extra computational overhead,  
 123 it significantly speeds up the learning process.

124 **References**

- 125 Adrien Bardes, Jean Ponce, and Yann LeCun. Vicreg: Variance-invariance-covariance regularization  
126 for self-supervised learning. *arXiv preprint arXiv:2105.04906*, 2021.
- 127 Yubei Chen, Adrien Bardes, Zengyi Li, and Yann LeCun. Intra-instance vicreg: Bag of self-supervised  
128 image patch embedding. *arXiv preprint arXiv:2206.08954*, 2022.
- 129 Zhijie Deng, Jiaxin Shi, Hao Zhang, Peng Cui, Cewu Lu, and Jun Zhu. Neural eigenfunctions are  
130 structured representation learners. *arXiv preprint arXiv:2210.12637*, 2022a.
- 131 Zhijie Deng, Jiaxin Shi, and Jun Zhu. Neuralef: Deconstructing kernels by deep neural networks. In  
132 *International Conference on Machine Learning*, pages 4976–4992. PMLR, 2022b.
- 133 Yann Dubois, Stefano Ermon, Tatsunori B Hashimoto, and Percy S Liang. Improving self-supervised  
134 learning by characterizing idealized representations. *Advances in Neural Information Processing*  
135 *Systems*, 35:11279–11296, 2022.
- 136 Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov,  
137 Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning  
138 robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023.
- 139 David Pfau, Stig Petersen, Ashish Agarwal, David GT Barrett, and Kimberly L Stachenfeld. Spectral  
140 inference networks: Unifying deep and spectral learning. *arXiv preprint arXiv:1806.02215*, 2018.
- 141 Phil Pope, Chen Zhu, Ahmed Abdelkader, Micah Goldblum, and Tom Goldstein. The intrinsic dimen-  
142 sion of images and its impact on learning. In *International Conference on Learning Representations*,  
143 2020.
- 144 Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang,  
145 Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition  
146 challenge. *International journal of computer vision*, 115:211–252, 2015.
- 147 James B Simon, Maksis Knutins, Liu Ziyin, Daniel Geisz, Abraham J Fetterman, and Joshua Albrecht.  
148 On the stepwise nature of self-supervised learning. *arXiv preprint arXiv:2303.15438*, 2023.
- 149 Roman Vershynin. Introduction to the non-asymptotic analysis of random matrices. *arXiv preprint*  
150 *arXiv:1011.3027*, 2010.
- 151 Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stéphane Deny. Barlow twins: Self-supervised  
152 learning via redundancy reduction. In *International Conference on Machine Learning*, pages  
153 12310–12320. PMLR, 2021.

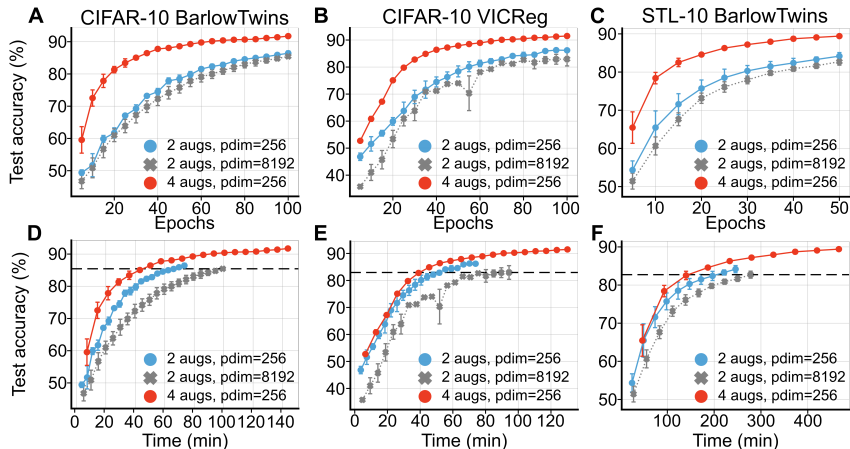


Figure 5: Using multiple augmentations improves representation learning performance and convergence. (A-C) Across BarlowTwins and VICReg for CIFAR-10 and STL-10 pretraining, using 4 augmentations instead of 2 helps improve performance. (D-F) Although the 4-augmentations take longer for each epoch, its performance still trumps the 2-augmentation version of the algorithm at the same wall clock time.

## 155 **B Data augmentation kernel perspective of non-contrastive SSL**

156 Following the previous section, we will now present an augmentation kernel perspective of BarlowTwins and VICReg losses. Specifically, we show that these losses are equivalent to the optimization problem of learning eigenfunctions of the augmentation-defined data covariance kernel. 157  
 158 Subsequently, we argue that using a high-dimensional projector yields better overlap with the top eigenvectors of the data augmentation kernel at initialization as compared to a low-dimensional projector. 159  
 160 Therefore, our analysis suggests using a stronger orthogonalization constraint during optimization for lower-dimensional projectors to ensure that features learned are equivalent to those learned with high-dimensional projectors. 161  
 162 Furthermore, we also argue that using more number of augmentations improves our estimate of the augmentation-defined data covariance kernel, thereby 163  
 164 aiding the eigenfunction optimization problem. Therefore, our analysis suggests using an averaging operator with more data augmentations to better estimate the true augmentation kernel. 165  
 166

### 167 **B.1 Features in terms of data augmentation kernels**

168 We will define two notions of the data augmentation kernel. Given two images,  $x, z$ , the first kernel, which we call the forward data augmentation covariance kernel, is given by 169

$$k^{DAF}(x, z) = \mathbb{E}_{x_0 \sim \rho_X} [p(x | x_0)p(z | x_0)] \quad (3)$$

170 This covariance kernel measures the similarity between  $x, z$  in terms of how likely they are to be reached from  $x_0$ , weighted by the distribution of  $x_0$ . Note that this is indeed the edge strength 171  
 172 between nodes  $x, z$  in the augmentation graph. We can also define a (backwards) data augmentation covariance kernel which reverses the roles of  $(x, z)$  and  $x_0$ : 173

$$k^{DAB}(x, z) = \mathbb{E}_{x_0 \sim \rho_X} [p(x_0 | x)p(x_0 | z)] \quad (4)$$

174 The goal of SSL is to learn features that preserve the covariance kernel structure (imposed by this choice of mapping  $M$ ) [Dubois et al., 2022]. Therefore, we want to define a loss which determines 175  
 176 vector features,  $F : X \rightarrow \mathbb{R}^d$ , which factor a data augmentation kernel  $k^{DA}(x, z) = F(x)^\top F(z)$ . Doing this directly is prohibitively data intensive at scale, since it involves a search over data 177  
 178 augmented images. However, since the covariance kernels are PSD, they define a Reproducing Kernel Hilbert space (RKHS). This allows us to apply Mercer’s theorem to find vector features as in Deng 179  
 180 et al. [2022a,b], Pfau et al. [2018].

181 The construction of features using Mercer’s theorem goes as follows. Given a PSD data augmentation  
 182 kernel,  $k^{DA}$ , define the  $T_k$  operator, which takes a function  $f$  and returns its convolution with the  
 183 data augmentation kernel.

$$T_k f(x) = \mathbb{E}_{z \sim \rho_X} [k(z, x) f(z)] \quad (5)$$

184 We will also make use of the the following operator,

$$T_M f(x) = \mathbb{E}_{x_0} [p(x_0 | x) f(x_0)] \quad (6)$$

185 which averages the values of the function,  $f$ , over the augmented images  $x_0 = M(x)$  of the data,  $x$ .

186 Since the operator  $T_k$  is compact and positive, it has a spectral decomposition consisting of eigen-  
 187 functions  $\phi_i$  and corresponding eigenvalues  $\lambda_i$ . Using these eigenpairs, we can define the (infinite  
 188 sequence of square summable) spectral features,  $G : X \rightarrow \ell_2$ , (where  $\ell_2$  represents square summable  
 189 sequences), by

$$G(x) = (\sqrt{\lambda_1} \phi_1(x), \dots, \sqrt{\lambda_d} \phi_d(x), \dots) \quad (7)$$

190 Then, Mercer’s theorem gives

$$k^{DA}(x, z) = G(x) \cdot G(z) \quad (\text{Mercer})$$

191 and ensures that the inner product is finite. These are the desired features, which factor the kernel.  
 192 However, computing the eigenfunctions of  $T_k$  is costly. Instead we propose an alternative using the  
 193 more efficient operator  $T_M$ . Both operators lead to equivalent features, according to Definition B.1.

194 **Definition B.1.** Let  $F(x) = (f_1(x), \dots, f_d(x))$  be a  $d$ -dimensional feature vector (a vector of  
 195 functions). Define the subspace

$$V = V(F) = \{h : X \rightarrow \mathbb{R} \mid h(x) = w \cdot F(x), \quad w \in \mathbb{R}^d\} \quad (8)$$

196 to be the span of the components of  $F$ . Given an  $n$ -dimensional feature vector,  $G(x) =$   
 197  $(g_1(x), \dots, g_n(x))$  we say the features  $G$  and  $F$  are equivalent, if  $V(F) = V(G)$ .

198 **Theorem B.2.** Let  $G(x)$  be the infinite Mercer features of the backward data augmentation covariance  
 199 kernels,  $k^{DAB}$ . Let  $F(x) = (f_1(x), f_2(x), \dots, f_k(x))$  be the features given by minimizing the  
 200 following data augmentation invariance loss

$$L(F) = \sum_{i=1}^{N_k} \|T_M f_i - f_i\|_{L^2(\rho_X)}^2, \quad \text{subject to} \quad (f_i, f_j)_{\rho_X} = \delta_{ij} \quad (9)$$

201 which includes the orthogonality constraint. Then,  $V(F) \subset V(G)$ ,  $V(F) \rightarrow V(G)$  as  $N_k \rightarrow \infty$ .

202 The idea of the proof uses the fact that, as linear operators,  $T_{k^{DAB}} = T_M^\top T_M$  and that  $T_{k^{DAF}} =$   
 203  $T_M T_M^\top$ . Then we use spectral theory of compact operators, which is analogue of the Singular Value  
 204 Decomposition in Hilbert Space, to show that eigenfunctions of  $T_M^\top T_M$  operator are the same as  
 205 those obtained from optimizing  $L(F)$ . A similar result can be obtained using  $k^{DAF}$  and  $T_M^\top$ .

206 Note that  $L(F)$  is the constrained optimization formulation of the BarlowTwins loss. Furthermore,  
 207  $L(F)$  with the additional constraint that  $(f_i, f_i) \geq \gamma \forall i \in \{1, 2 \dots N_k\}$  is the constrained optimiza-  
 208 tion formulation of the VICReg loss.

## 209 B.2 Corollary 1: Low-dimensional projectors are sufficient

210 While BarlowTwins and VICReg frameworks have advocated the use of high-dimensional projectors  
 211 to facilitate good feature learning on Imagenet, our kernel perspective challenges this notion. Since the  
 212 intrinsic dimensionality of Imagenet is estimated to be  $\sim 40$  [Pope et al., 2020], it is not unreasonable  
 213 to expect that the span of desired features would be of similar dimensionality. It is, thus, intriguing  
 214 that these frameworks mandate the use of an  $\sim 8192 - d$  projector head to capture the intricacies  
 215 of corresponding data augmentation kernel. This discrepancy can be explained by observing the  
 216 learning dynamics of a linearized model under the BarlowTwins loss optimization [Simon et al.,  
 217 2023]. These dynamics reveal that initializing the projection weight matrix in alignment with the  
 218 eigenfunctions of the data kernel retains this alignment throughout the learning process. Notably,  
 219 a high-dimensional projector is more likely to have a greater span at initialization compared to its  
 220 low-dimensional counterpart, increasing the likelihood of overlap with the relevant eigenfunctions.  
 221 We hypothesize that it is possible to rectify this issue by using a stronger orthogonalization constraint  
 222 for low-dimensional projectors, thereby rendering them sufficient for good feature learning.

223 **B.3 Corollary 2: Multiple augmentations improve optimization**

224 Theorem B.2 implies that the invariance loss optimization would ideally entail using the  $T_M$  operator,  
 225 thereby requiring many augmentations for each sample  $x$ . Using only two augmentations per sample  
 226 yields a noisy estimate of  $T_M$ , yielding spurious eigenpairs [Vershynin, 2010] (see Appendix). These  
 227 spurious eigenpairs add stochasticity to the learning dynamics, and hinder the alignment of the  
 228 learned features with the eigenfunctions of the data kernel [Simon et al., 2023]. We hypothesize that  
 229 improving this estimation error by increasing the number of augmentations could ameliorate this  
 230 issue and improve the speed and quality of feature learning.

231 Increasing the number of augmentations (say  $m$ ) in BarlowTwins and VICReg comes with added  
 232 compute costs. A straightforward approach would involve computing the invariance loss for every  
 233 pair of augmentations, resulting in  $\mathcal{O}(m^2)$  operations. However, Theorem B.2 proposes an alternative  
 234 method that uses the sample estimate of  $T_M$ , thereby requiring only  $\mathcal{O}(m)$  operations. Both these  
 235 strategies are functionally equivalent (see Appendix), but the latter is computationally more efficient.  
 236 In summary, Theorem B.2 establishes a mechanistic role for the number of data augmentations,  
 237 paving the way for a computationally efficient multi-augmentation framework:

$$\widehat{L}(F) = \mathbb{E}_{x \sim \rho_X} \left[ \sum_{i=1}^{N_k} \sum_{j=1}^m \| \overline{f_i(x)} - f_i(x_j) \|_{L^2(\rho_X)}^2 \right], \quad \text{subject to} \quad (f_i, f_j)_{\rho_X} = \delta_{ij} \quad (10)$$

238 where  $\overline{f_i(x)} = \frac{1}{m} \sum_{j=1}^m f_i(x_j)$  is the sample estimate of  $T_M f_i(x)$ .

239 **C Data augmentation kernel perspective of non-contrastive SSL**

240 **Theorem C.1.** *Let  $G(x)$  be the infinite Mercer features of the backward data augmentation covari-*  
 241 *ance kernels,  $k^{DAB}$ . Let  $F(x) = (f_1(x), f_2(x), \dots, f_k(x))$  be the features given by minimizing the*  
 242 *following data augmentation invariance loss*

$$L(F) = \sum_{i=1}^{N_k} \| T_M f_i - f_i \|_{L^2(\rho_X)}^2, \quad \text{subject to} \quad (f_i, f_j)_{\rho_X} = \delta_{ij} \quad (11)$$

243 *which includes the orthogonality constraint. Then,  $V(F) \subset V(G)$ ,  $V(F) \rightarrow V(G)$  as  $N_k \rightarrow \infty$ .*

244 The idea of the proof uses the fact that, as linear operators,  $T_{k^{DAB}} = T_M^\top T_M$  and that  $T_{k^{DAF}} =$   
 245  $T_M T_M^\top$ . Then we use spectral theory of compact operators, which is analogue of the Singular Value  
 246 Decomposition in Hilbert Space, to show that eigenfunctions of  $T_M^\top T_M$  operator are the same as  
 247 those obtained from optimizing  $L(F)$ . A similar result can be obtained using  $k^{DAF}$  and  $T_M^\top$ .

248 Note that  $L(F)$  is the constrained optimization formulation of the BarlowTwins loss. Furthermore,  
 249  $L(F)$  with the additional constraint that  $(f_i, f_i) \geq \gamma \forall i \in \{1, 2 \dots N_k\}$  is the constrained optimiza-  
 250 tion formulation of the VICReg loss.

251 **C.1 Proof of theorem 3.2**

252 We show we can factor the linear operator, leading to a practical algorithm. Here, we show that we  
 253 can capture the backward data augmentation kernel with the forward data augmentation averaging  
 254 operator

255 **Lemma C.2.** *Using the definitions above, and with  $k$  in equation 5 given by  $k^{DAB}$ ,*

$$T_k = T_M^\top T_M$$

256 *Proof.* First, define the non-negative definite bilinear form

$$B^{VAR}(f, g) = (T_M f, T_M g)_{\rho_X} \quad (12)$$

257 Given the backwards data augmentation covariance kernel,  $k^{DAB}$ , define

$$B^{DAB}(f, g) = (T_k f, g)_{\rho_X}$$



258 We claim, that

$$B^{VAR} = B^{DA,B} \quad (13)$$

259 This follows from the following calculation,

$$B^{DA,B}(f, g) = (T_k f, g)_{\rho_X} \quad (14)$$

$$= \mathbb{E}_x [T_k f(x), g(x)] = \mathbb{E}_x \mathbb{E}_z [k_{DA,B}(z, x) f(z) g(x)] \quad (15)$$

$$= \mathbb{E}_x \mathbb{E}_z \mathbb{E}_{x_0} [p(x | x_0) p(z | x_0) f(z) g(x)] \quad (16)$$

$$= \mathbb{E}_{x_0} [\mathbb{E}_x [p(x | x_0) g(x)], \mathbb{E}_z [p(z | x_0) f(z)], ] = \mathbb{E}_{x_0} T_M f(x_0) T_M g(x_0) \quad (17)$$

$$= (T_M f, T_M g)_{\rho_X} = B^{VAR}(f, g) \quad (18)$$

260

□

261 For implementations, it is more natural to consider *invariance* to data augmentations.

262 **Theorem C.3** (equivalent eigenfunctions). *Assume that  $T_M$  is a compact operator. Define the*  
 263 *invariance bilinear form*

$$B^{INV}(f, g) = (T_M f - f, T_M g - g) \quad (19)$$

264 *Then  $B^{INV}$ ,  $B^{VAR}$  share the same set of eigenfunctions. Moreover, these are the same as the*  
 265 *eigenfunctions of  $B^{DA,B}$ . In particular, for any eigenfunction  $f_j$  of  $B^{VAR}$ , with eigenvalue  $\lambda_j$ , then*  
 266  *$f_j$  is also an eigenfunction of  $B^{INV}$ , with the corresponding eigenvalue given by  $(\sqrt{\lambda_j} - 1)^2$ .*

267 *Proof.* Define  $T_{MM}$  by,

$$T_{MM} f = T_M^\top T_M f \quad (20)$$

268 Define

$$T_{MS} = (T_M - I)^\top (T_M - I) \quad (21)$$

269 Note, by the assumption of compactness,  $T_M$  has the Singular Value Decomposition, (see the Hilbert  
 270 Space section for equation SVD),

$$T_M(h) = \sum_{j=1}^{\infty} \lambda_j(h, g_j) f_j \quad (\text{SVD})$$

271 Let  $f_j$  be any right eigenvector of  $T_M$ , with eigenvalue  $\mu_j$ . Then  $f_j$  is also a right eigenvector  $T_M - I$ ,  
 272 with eigenvalue  $\mu_j - 1$ . So we see that  $T_{MM}$  has  $f_j$  as an eigenvector, with eigenvalue  $\lambda_j = \mu_j^2$  and  
 273  $T_{MS}$  has  $f_j$  as an eigenvector, with eigenvalue  $(\sqrt{\lambda_j} - 1)^2$ . Finally, the fact that there are no other  
 274 eigenfunctions also follows from equation SVD.

275 The final part follows from the previous lemma. □