

# VARIANCE-COVARIANCE REGULARIZATION IMPROVES REPRESENTATION LEARNING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Transfer learning plays a key role in advancing machine learning models, yet conventional supervised pretraining often undermines feature transferability by prioritizing features that minimize the pretraining loss. In this work, we adapt a self-supervised learning regularization technique from the VICReg method to supervised learning contexts, introducing Variance-Covariance Regularization (VCReg). This adaptation encourages the network to learn high-variance, low-covariance representations, promoting learning more diverse features. We outline best practices for an efficient implementation of our framework, including applying it to the intermediate representations. Through extensive empirical evaluation, we demonstrate that our method significantly enhances transfer learning for images and videos, achieving state-of-the-art performance across numerous tasks and datasets. VCReg also improves performance in scenarios like long-tail learning and hierarchical classification. Additionally, we show its effectiveness may stem from its success in addressing challenges like gradient starvation and neural collapse. In summary, VCReg offers a universally applicable regularization framework that significantly advances transfer learning and highlights the connection between gradient starvation, neural collapse, and feature transferability.

## 1 INTRODUCTION

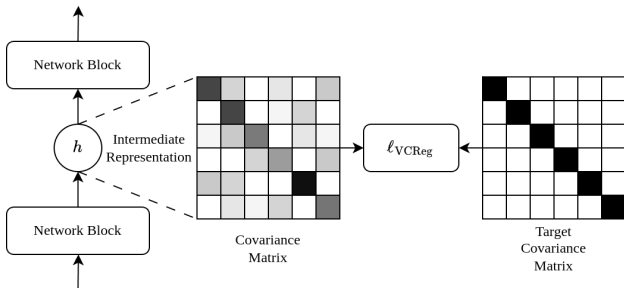
Transfer learning enables models to apply knowledge from one domain to enhance performance in another, particularly when data are scarce or costly to obtain (Pan & Yang, 2010; Weiss et al., 2016; Zhuang et al., 2020; Bommasani et al., 2021). One of the key challenges arises during the supervised pretraining phase. In this phase, models often lack detailed information about the downstream tasks to which they will be applied. Nevertheless, they must aim to capture a broad spectrum of features beneficial across various applications (Bengio, 2012; Caruana, 1997; Yosinski et al., 2014). Without proper regularization techniques, these supervised pretrained models tend to overly focus on features that minimize supervised loss, resulting in limited generalization capabilities and issues such as gradient starvation and neural collapse (Zhang et al., 2016; Neyshabur et al., 2017; Zhang et al., 2021; Pezeshki et al., 2021; Pappan et al., 2020; Shwartz-Ziv, 2022).

To tackle these challenges, we adapt the regularization techniques of the self-supervised VICReg method (Bardes et al., 2021) for the supervised learning paradigm. Our method, termed Variance-Covariance Regularization (VCReg), aims to encourage the learning of representations with high variance and low covariance, thus avoiding the overemphasis on features that merely minimize supervised loss. Instead of simply applying VCReg to the final representation of the network, we explore the most effective ways to incorporate it throughout intermediate representations.

The structure of the paper is as follows: we begin with an introduction of our method, including an outline of a fast implementation strategy designed to minimize computational overhead. Following this, we present a series of experiments aimed at validating the method’s efficacy across a wide range of tasks, datasets, and architectures. Subsequently, we conduct analyses on the learned representations to demonstrate VCReg’s effectiveness in mitigating common issues in transfer learning, such as neural collapse and gradient starvation.

Our paper makes the following contributions:

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107



**Figure 1: VCRReg regularizes the network by encouraging the intermediate representations to have high variance and low covariance.** VCRReg is applied to the output of each network block to make all the intermediate representations capture diverse features.

1. We introduce a robust strategy for applying VCRReg to neural networks, including integrating it into the intermediate layers.
2. We propose a computationally efficient implementation of VCRReg. This implementation is optimized to ensure minimal additional computational overhead, allowing for seamless integration into existing workflows.
3. Through extensive experiments on benchmark datasets both in images and videos, we demonstrate that VCRReg suppresses the prior state-of-the-art results in transfer learning performance across various network architectures, including ResNet (He et al., 2016), ConvNeXt (Liu et al., 2022), and ViT (Dosovitskiy et al., 2020). Moreover, we also show that VCRReg improves performance in diverse scenarios like long-tail learning and hierarchical classification.
4. We investigate the representations learned by VCRReg, revealing its effectiveness in combating challenges such as gradient starvation (Pezeshki et al., 2021), neural collapse (Papayan et al., 2020), information compression (Shwartz-Ziv, 2022), and sensitivity to noise.

Before delving into VCRReg’s details in the following sections, it is key to note its divergence from VICReg, namely by omitting the invariance loss and focusing on variance and covariance loss for a wider application, especially in transfer learning. This approach tackles challenges like gradient starvation and neural collapse, improving neural network training across various architectures. Our work further distinguishes itself by exploring optimal regularization strategies, moving beyond generic application to significantly enhance its effectiveness.

## 2 RELATED WORK

### 2.1 VARIANCE-INVARIANCE-COVARIANCE REGULARIZATION (VICREG)

VICReg (Bardes et al., 2021) is a novel SSL method that encourages the learned representation to be invariant to data augmentation. However, focusing solely on this invariance criterion can result in the network producing a constant representation, making it invariant to both data augmentation and the input data itself. VICReg primarily regularizes the network by combining variance loss and covariance loss. The variance loss encourages high variance in the learned representations, thereby promoting the learning of diverse features. The covariance loss, on the other hand, aims to minimize redundancy in the learned features by reducing the overlap in information captured by different dimensions of the representation. This dual-objective optimization framework effectively promotes diverse feature learning for SSL (Shwartz-Ziv et al., 2022). To improve the performance of supervised network training, we adapt the SSL feature collapse prevention mechanism from VICReg and propose a variance-covariance regularization method.

To calculate the loss function of VICReg with a batch of data  $\{x_1 \dots x_n\}$ , we first need to have a pair of inputs  $(x'_i, x''_i)$  such that  $x'_i$  and  $x''_i$  are two augmented versions of the original input  $x_i$ . Given the neural network  $f_\theta(\cdot)$  and the final representations  $z'_i = f_\theta(x'_i)$  and  $z''_i = f_\theta(x''_i)$  such that

108  $z'_i, z''_i \in \mathbb{R}^D$ , VICReg minimizes the following loss:

$$109$$

$$110 \ell_{\text{VICReg}}(z'_1 \dots z'_n, z''_1 \dots z''_n) = \alpha \ell_{\text{var}}(z'_1, \dots, z'_n) + \alpha \ell_{\text{var}}(z''_1, \dots, z''_n) \quad (1)$$

$$111$$

$$112 + \beta \ell_{\text{cov}}(z'_1, \dots, z'_n) + \beta \ell_{\text{cov}}(z''_1, \dots, z''_n) + \sum_{i=1}^n \ell_{\text{inv}}(z'_i, z''_i).$$

$$113$$

$$114$$

115 The variance and covariance loss functions are respectively defined as:

$$116$$

$$117 \ell_{\text{var}} = \frac{1}{D} \sum_{i=1}^D \max(0, 1 - \sqrt{C_{ii}}), \quad \ell_{\text{cov}} = \frac{1}{D(D-1)} \sum_{i \neq j} C_{ij}^2 \quad (2)$$

$$118$$

$$119$$

$$120$$

121 where  $C = \frac{1}{N-1} \sum_{i=1}^N (z_i - \bar{z})(z_i - \bar{z})^T$  denotes the covariance matrix, and  $\bar{z}$  represents the mean  
 122 vector, given by  $\bar{z} = \frac{1}{N} \sum_{i=1}^N z_i$ .

123 Building on insights from prior studies (Shwartz-Ziv, 2022; Shwartz-Ziv et al., 2023), it is understood  
 124 that the invariance term does not play a pivotal role in diversifying features. Consequently, in adapting  
 125 to the supervised regime, we exclude the invariance term from the regularization.

## 126 2.2 REPRESENTATION WHITENING AND FEATURE DIVERSITY REGULARIZERS

127 Representation whitening is a technique for processing inputs before they enter a network layer. It  
 128 transforms the input so that its components are uncorrelated with unit variance (Kessy et al., 2018).  
 129 This transformation achieves enhanced model optimization and generalization. It uses a whitening  
 130 matrix derived from the data’s covariance matrix and results in an identity covariance matrix, thereby  
 131 aiding gradient flow during training and acting as a lightweight regularizer to reduce overfitting and  
 132 encourage robust data representations (LeCun et al., 2002).

133 In addition to whitening as a processing step, additional regularization terms can be introduced to  
 134 enforce decorrelation in the representations. Various prior works have explored these feature diversity  
 135 regularization techniques to enhance neural network training (Cogswell et al., 2015; Ayinde et al.,  
 136 2019; Laakom et al., 2023). These methods encourage diverse features in the representation by adding  
 137 a regularization term. Recent methods like WLD-Reg (Laakom et al., 2023) and DeCov (Cogswell  
 138 et al., 2015) also employ covariance-matrix-based regularization to promote feature diversity, similarly  
 139 to our approach.

140 However, the studies above mainly focus on the benefits of optimization and generalization for the  
 141 source task, often neglecting their implications for supervised transfer learning. VICReg distinguishes  
 142 itself by explicitly targeting enhancements in transfer learning performance. Our results indicate  
 143 that such regularization techniques yield only modest performance improvements in in-domain  
 144 evaluations.

## 145 2.3 GRADIENT STARVATION AND NEURAL COLLAPSE

146 Gradient starvation and neural collapse are two recently recognized phenomena that can significantly  
 147 affect the quality of learned representations and a network’s generalization ability (Pezeshki et al.,  
 148 2021; Pappas et al., 2020; Ben-Shaul et al., 2023). Gradient starvation occurs when certain parameters  
 149 in a deep learning model receive very small gradients during the training process, thereby leading  
 150 to slower or non-existent learning for these parameters (Pezeshki et al., 2021). Neural collapse,  
 151 on the other hand, is a phenomenon observed during the late stages of training when the internal  
 152 representations of the network tend to collapse towards each other, resulting in a loss of feature  
 153 diversity (Pappas et al., 2020). Both phenomena are particularly relevant in the context of transfer  
 154 learning, where models are initially trained on a source task before being fine-tuned for a target task.  
 155 Our work, through the use of VICReg, seeks to mitigate these issues, offering a pathway to more  
 156 effective transfer learning.

### 3 VARIANCE-COVARIANCE REGULARIZATION

#### 3.1 VANILLA VCREG

Consider a labeled dataset comprising  $N$  samples, denoted as  $\{(x_1, y_1) \dots (x_N, y_N)\}$  and a neural network  $f_\theta(\cdot)$ , which takes these inputs  $x_i$  and produces final predictions  $\tilde{y}_i = f_\theta(x_i)$ . In standard supervised learning, the loss is defined as  $L_{\text{sup}} = \frac{1}{N} \sum_{i=1}^N \ell_{\text{sup}}(\tilde{y}_i, y_i)$ .

The core objective of the Vanilla VCREg is to ensure that the  $D$ -dimensional input representations  $\{h_i\}_{i=1}^N$  to the last layer of the network exhibit both high variance and low covariance. To achieve this, we employ variance and covariance regularization, same as mentioned in equation 1:

$$\ell_{\text{vcreg}}(h_1 \dots h_N) = \alpha \ell_{\text{var}}(h_1 \dots h_N) + \beta \ell_{\text{cov}}(h_1 \dots h_N) \quad (3)$$

Intuitively speaking, the covariance matrix captures the interdependencies among the dimensions of the feature vectors  $h_i$ . Maximizing  $\ell_{\text{var}}$  encourages each feature dimension to contain unique, non-redundant information, while minimizing  $\ell_{\text{cov}}$  aims to reduce the correlation between different dimensions, thus promoting feature independence. The overall training loss, which includes also the supervised loss, then becomes:

$$L_{\text{vanilla}} = \alpha \ell_{\text{var}}(h_1 \dots h_N) + \beta \ell_{\text{cov}}(h_1 \dots h_N) + \frac{1}{N} \sum_{i=1}^N \ell_{\text{sup}}(\tilde{y}_i, y_i). \quad (4)$$

Here,  $\alpha$  and  $\beta$  serve as hyperparameters to control the strength of each regularization term.

#### 3.2 EXTENDING VCREG TO INTERMEDIATE REPRESENTATIONS

While regularizing the final layer in a neural network offers certain benefits, extending this approach to intermediate layers via VCREg provides additional advantages (for empirical evidence supporting this claim, please refer to Appendix A). Regularizing intermediate layers enables the model to capture more complex, higher-level abstractions. This strategy minimizes internal covariate shifts across layers, which in turn improves both the stability of training and the model’s generalization capabilities. Furthermore, it fosters the development of feature hierarchies and enriches the latent space, leading to enhanced model interpretability and improved transfer learning performance.

To implement this extension, VCREg is applied at  $M$  strategically chosen layers throughout the neural network. For each intermediate layer  $j$ , we denote the feature representation for an input  $x_i$  as  $h_i^{(j)} \in \mathbb{R}^{D_j}$ . This culminates in a composite loss function, expressed as follows:

$$L_{\text{VCREg}} = \sum_{j=1}^M \left[ \alpha \ell_{\text{var}}(h_1^{(j)} \dots h_N^{(j)}) + \beta \ell_{\text{cov}}(h_1^{(j)} \dots h_N^{(j)}) \right] + \frac{1}{N} \sum_{i=1}^N \ell_{\text{sup}}(\tilde{y}_i, y_i). \quad (5)$$

**Spatial Dimensions** However, applying VCREg to intermediate layers of real-world neural networks presents challenges due to the spatial dimensions in these intermediate representations. Naively reshaping these representations into long vectors would lead to unmanageably large covariance matrices, thereby increasing computational costs and risking numerical instability. To address this issue, we adapt VCREg to accommodate networks with spatial dimensions. Each vector at a different spatial location is treated as an individual sample when calculating the covariance matrix. Both the variance loss and the covariance loss are then calculated based on this modified covariance matrix. In terms of practical implementation, VCREg is usually applied subsequently to each block within the neural network architecture, often succeeding residual connections. This placement allows for seamless incorporation into current network architectures and training paradigms.

**Addressing Outliers with Smooth L1 Loss** After treating spatial locations as independent samples for covariance computation, the resulting samples are no longer statistically independent. This can lead to outliers in the covariance matrix and unstable gradient updates. To address this, we introduce a smooth L1 penalty into the covariance loss term. Specifically, we replace the traditional squared covariance values  $C_{ij}$  in  $\ell_{\text{cov}}$  with a smooth L1 function:

$$\text{SmoothL1}(x) = \begin{cases} x^2, & \text{if } |x| \leq \delta \\ 2\delta|x| - \delta^2, & \text{otherwise} \end{cases} \quad (6)$$

**Table 1: Transfer Learning Performance with ImageNet Supervised Pretraining**

The table shows performance metrics for different architectures. Each model is pretrained on the full ImageNet dataset and then tested on different downstream datasets using linear probing. Application of VCReg consistently improves performance and beats other feature diversity regularizer. Averages are calculated excluding ImageNet results.

Architecture	iNat18	Places	Food	Cars	Aircraft	Pets	Flowers	DTD	Average
ResNet-50	42.8%	50.6%	69.1%	43.6%	54.8%	91.9%	77.1%	68.7%	62.33%
ResNet-50 (DeCov)	43.1%	50.4%	69.0%	45.7%	55.5%	90.6%	79.2%	69.1%	62.83%
ResNet-50 (WLD-Reg)	43.9%	<b>51.2%</b>	70.2%	43.9%	58.7%	91.4%	80.7%	69.0%	63.63%
<b>ResNet-50 (VCReg)</b>	<b>45.3%</b>	<b>51.2%</b>	<b>71.7%</b>	<b>54.1%</b>	<b>70.5%</b>	<b>92.1%</b>	<b>88.0%</b>	<b>70.8%</b>	<b>67.96%</b>
ConvNeXt-T	51.6%	53.8%	78.4%	62.9%	74.7%	93.9%	91.3%	72.9%	72.44%
<b>ConvNeXt-T (VCReg)</b>	<b>52.3%</b>	<b>54.7%</b>	<b>79.6%</b>	<b>64.2%</b>	<b>76.3%</b>	<b>94.1%</b>	<b>92.7%</b>	<b>73.3%</b>	<b>73.40%</b>
ViT-Base-32	39.1%	47.9%	70.6%	51.2%	63.8%	90.3%	84.6%	66.1%	64.20%
<b>ViT-Base-32 (VCReg)</b>	<b>40.6%</b>	<b>48.1%</b>	<b>70.9%</b>	<b>52.0%</b>	<b>65.8%</b>	<b>91.0%</b>	<b>86.6%</b>	<b>66.5%</b>	<b>65.19%</b>

By implementing this modification, we ensure that the loss function increases in a more controlled manner with respect to large covariance values. Empirically, this minimizes the impact of outliers, thereby enhancing the stability of the training process.

### 3.3 FAST IMPLEMENTATION

To optimize VCReg speed, we use the fact that VCReg only affects the loss function and not the forward pass. This allows us to focus on modifying the backward function for improvements. Specifically, we sidestep the usual process of calculating the VCReg loss and subsequent backpropagation. Instead, we directly adjust the computed gradients, which is feasible since the VCReg loss calculation relies solely on the current representation. Further details of this speed-optimized technique are outlined in Appendix B. Our optimized VCReg implementation exhibits similar latency as batch normalization layers and is more than 5 times faster than the naive VCReg implementation. The results are presented in Table 8.

## 4 EXPERIMENTS

In this section, we first outline the experimental framework and findings highlighting the effectiveness of our proposed regularization approach, VCReg, within the realm of transfer learning that utilizes supervised pretraining for both images and videos. Subsequently, we extend our experiments to three specialized learning scenarios: 1) class imbalance via long-tail learning, 2) synergizing with self-supervised learning frameworks, and 3) hierarchical classification problems. The objective is to assess the adaptability of VCReg across various data distributions and learning paradigms, thereby evaluating its broader utility in machine learning applications. For details on reproducing our experiments, please consult Appendix C.

### 4.1 TRANSFER LEARNING FOR IMAGES

In this section, we adhere to evaluation protocols established by seminal works such as (Chen et al., 2020; Kornblith et al., 2021; Misra & Maaten, 2020) for our transfer learning experiments. Initially, we pretrain models using three different architectures: ResNet-50 (He et al., 2016), ConvNeXt-Tiny (Liu et al., 2022), and ViT-Base-32 (Dosovitskiy et al., 2020), on the full ImageNet dataset. We follow the standard PyTorch recipes (Paszke et al., 2019) for all networks and do not modify any hyperparameters other than those related to VCReg to ensure a fair baseline comparison. Subsequently, we perform a linear probing evaluation across 9 different benchmark to evaluate the transfer learning performance. For ResNet-50, we include two other feature diversity regularizer methods for comparison: DeCov (Cogswell et al., 2015) and WLD-Reg (Laakom et al., 2023). We conduct experiments solely with ResNet-50 because it is the principal architecture used in the WLD-Reg paper. To ensure a fair comparison, we source hyperparameters from Laakom et al. (2023) for both DeCov and WLD-Reg.

**Table 2: Transfer Learning Performance with Kinetics-400 and Kinetics-710 pretrained models:** The table shows fine-tuning performance of Kinetics pre-trained models on HMDB51. VideoMAE-S, VideoMAE-B, and ViViT-B are pretrained on Kinetics-400 dataset while VideoMAEv2-S and VideoMAEv2-B are pre-trained on Kinetics-710. We apply VCRReg only to the networks’ output preceding the classification head. The results show that VCRReg can boost the transfer learning classification performance for networks pre-trained on video data.

Method	Backbone	HMDB51
VideoMAE-S	ViT-S	79.9%
VideoMAE-S (VCRReg)	ViT-S	<b>80.6%</b>
VideoMAE-B	ViT-B	82.2%
VideoMAE-B (VCRReg)	ViT-B	<b>83.0%</b>
VideoMAEv2-S	ViT-S	83.6%
VideoMAEv2-S (VCRReg)	ViT-S	<b>83.9%</b>
VideoMAEv2-B	ViT-B	86.5%
VideoMAEv2-B (VCRReg)	ViT-B	<b>86.9%</b>
ViViT-B	ViT-B	70.9%
ViViT-B (VCRReg)	ViT-B	<b>71.6%</b>

The results in Table 1 demonstrate that VCRReg significantly enhances performance in transfer learning for images, achieving the highest performance for 9 out of 10 datasets, and for all three architectures. Clearly, VCRReg acts as a versatile plug-in, effectively boosting transfer learning outcomes. Its effectiveness spans ConvNet and Transformer architectures, confirming its wide-ranging applicability.

## 4.2 TRANSFER LEARNING FOR VIDEOS

To extend our evaluation of VCRReg’s efficacy, we conduct experiments using networks pretrained on video datasets. Specifically, we utilize models pretrained on Kinetics-400 Kay et al. (2017) and Kinetics-710 Li et al. (2022), subsequently finetuning them for action recognition on HMDB51 Kuehne et al. (2011). We experiment with models pretrained with self-supervised learning objectives (VideoMAE Tong et al. (2022) and VideoMAEv2 Wang et al. (2023)), as well as models pretrained with conventional supervised learning objectives (ViViT Arnab et al. (2021)). We follow the finetuning protocols detailed by Tong et al. (2022) and the conventional evaluation method used in the field, where the final performance is measured by the mean classification accuracy across three provided splits Simonyan & Zisserman (2014). To pinpoint the optimal VCRReg coefficients, we conduct a grid search based on validation set accuracy. For simplicity, in this setup, VCRReg regularization is exclusively applied to the final output of each network during finetuning, just before the classification head.

Table 2 illustrates that incorporating VCRReg as a plugin regularizer improves the transfer learning performance for action recognition across various methods (VideoMAE, VideoMAE2, and ViViT-B) and backbone architectures (ViT-B and ViT-S). This solidifies VCRReg’s status as a practical and versatile regularizer, capable of substantially improving the performance of pretrained networks in transfer learning scenarios.

## 4.3 CLASS IMBALANCE WITH LONG-TAIL LEARNING

Class imbalance is a pervasive issue in many real-world datasets and poses a considerable challenge to standard neural network training algorithms. We conduct experiments to assess how well VCRReg addresses this issue through long-tail learning. We evaluate VCRReg using the CIFAR10-LT and CIFAR100-LT Krizhevsky et al. (2009) datasets, both engineered to have an imbalance ratio of 100. These experiments use a ResNet-32 backbone architecture. The per-class sample sizes ranges from 5,000 to 50 for CIFAR10-LT and from 500 to 5 for CIFAR100-LT.

Table 3 shows that models augmented with VCRReg consistently outperform the standard ResNet-32 models on imbalanced datasets. These results are noteworthy because they demonstrate that VCRReg effectively enhances the model’s ability to discriminate between classes in imbalanced settings. This

**Table 3: Performance Comparison on Class-Imbalanced Datasets Using VCRReg:** This table shows the accuracy of standard ResNet-32 with and without VCRReg when trained on class-imbalanced CIFAR10-LT and CIFAR100-LT datasets. The VCRReg-enhanced models show improved performance, demonstrating the method’s effectiveness in addressing class imbalance.

Training Methods	CIFAR10-LT	CIFAR100-LT
ResNet-32	69.6%	37.4%
ResNet-32 (VCRReg)	<b>71.2%</b>	<b>40.4%</b>

**Table 4: Impact of VCRReg on Self-Supervised Learning Methods:** This table presents a comparative analysis of ResNet-50 models pretrained with SimCLR and VICReg losses on ImageNet, both with and without the VCRReg applied. The models are evaluated using linear probing on various downstream task datasets. The VCRReg models consistently outperform the non-VCRReg models, showcasing the method’s broad utility in transfer learning for self-supervised learning scenarios. Averages are calculated excluding ImageNet results.

Pretraining	ImageNet	iNat18	Places	Food	Cars	Aircraft	Pets	Flowers	DTD	Average
SimCLR	<b>67.2%</b>	37.2%	52.1%	66.4%	35.7%	<b>62.3%</b>	76.3%	82.6%	68.1%	60.09%
<b>SimCLR+VCR</b>	67.1%	<b>41.3%</b>	<b>52.3%</b>	<b>67.7%</b>	<b>40.6%</b>	61.9%	<b>76.6%</b>	<b>83.6%</b>	<b>69.0%</b>	<b>61.63%</b>
VICReg	65.2%	<b>41.7%</b>	48.2%	61.0%	27.3%	51.2%	79.1%	74.3%	65.4%	56.03%
<b>VICReg+VCR</b>	<b>66.3%</b>	41.4%	<b>49.6%</b>	<b>61.6%</b>	<b>29.3%</b>	<b>54.2%</b>	<b>79.7%</b>	<b>74.5%</b>	<b>66.5%</b>	<b>57.10%</b>

establishes VCRReg as a valuable tool for real-world applications where class imbalance is often a concern.

#### 4.4 SELF-SUPERVISED LEARNING WITH VCRREG

Our subsequent investigation focuses on examining the synergy between VCRReg and existing self-supervised learning paradigms. As mentioned in the previous sections, we apply VCRReg not only to the final but also to intermediate representations. So in all of the following experiments for self-supervised learning with VCRReg, we apply the original loss function to the output of the network, and the VCRReg loss to all the intermediate representations. We employ a ResNet-50 architecture, training it for 100 epochs under four different configurations: using either SimCLR loss or VICReg loss, coupled with the ImageNet dataset. For evaluation, we conduct linear probing tests on multiple downstream task datasets, following the protocols prescribed by Misra & Maaten (2020); Zbontar et al. (2021).

As indicated in Table 4, integrating VCRReg into self-supervised learning paradigms such as SimCLR and VICReg results in consistent performance improvements for transfer learning. Specifically, the linear probing accuracies are enhanced across nearly all the evaluated datasets. These gains underscore the broad applicability and versatility of VCRReg, demonstrating its potential to enhance various machine learning methodologies.

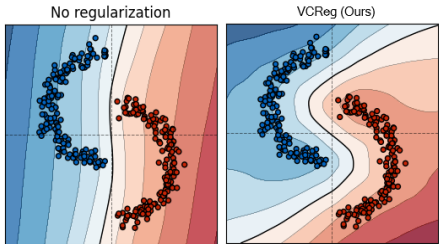
#### 4.5 HIERARCHICAL CLASSIFICATION

To evaluate the efficacy of the learned representations across multiple levels of class granularity, we conduct experiments on the CIFAR100 dataset as well as five distinct subsets of ImageNet (Engstrom et al., 2019). In each dataset, every data sample is tagged with both superclass and subclass labels, denoted as  $(x_i, y_i^{\text{sup}}, y_i^{\text{sub}})$ . Note that while samples sharing the same subclass label also share the same superclass label, the reverse does not necessarily hold true. Initially, the model is trained using only the superclass labels, i.e., the  $(x_i, y_i^{\text{sup}})$  pairs. Subsequently, linear probing is employed with the subclass labels  $(x_i, y_i^{\text{sub}})$  to assess the quality of abstract features at the superclass level.

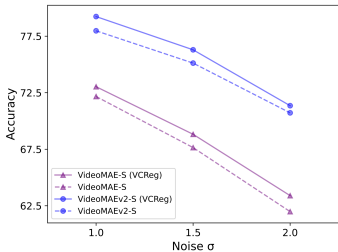
Table 5 presents key performance metrics, highlighting the substantial improvements VCRReg brings to subclass classification. The improvements are consistent across all datasets, with the CIFAR100 dataset showing the most significant gain—an increase in accuracy from 60.7% to 72.9%. These results underscore VCRReg’s capability to assist neural networks in generating feature representations that are not only discriminative at the superclass level but are also well-suited for subclass distinctions.

**Table 5: Impact of VCRreg on Hierarchical Classification in ConvNeXt Models:** This table summarizes the classification accuracies obtained with ConvNeXt models, both with and without the VCRreg regularization, across multiple datasets featuring hierarchical class structures including CIFAR100 and several subsets of ImageNet. The models were initially trained using superclass labels and subsequently probed using subclass labels. VCRreg consistently boosts performance in subclass classification tasks.

	CIFAR100	living_9	mixed_10	mixed_13	geirhos_16	big_12
Superclass Count	20	9	10	13	16	12
Subclass Count	100	72	60	78	32	240
ConvNeXt	60.7%	53.4%	60.3%	61.1%	60.5%	51.8%
ConvNeXt (VCRreg)	<b>72.9%</b>	<b>62.2%</b>	<b>67.7%</b>	<b>66.0%</b>	<b>70.1%</b>	<b>61.5%</b>



**Figure 2: Comparative evaluation between training with and without VCRreg on a “Two-Moon” Synthetic Dataset.** Decision boundaries are averaged over ten distinct runs with random data point sampling and model initialization. A single run’s data points are displayed for clarity. While “No regularization” has limitations in forming intricate decision boundaries, VCRreg is effective in generating meaningful ones.



**Figure 3: Impact of VCRreg amidst noisy data:** This figure shows the top-1 accuracy of VideoMAE-S and VideoMAEv2-S when fine-tuned for action recognition using HMDB51 corrupted with synthetic noise. We corrupt the data with Gaussian noise with standard deviation  $\sigma \in \{1, 1.5, 2\}$ . Models with VCRreg outperform their non-regularized counterparts in this setting.

This attribute is particularly advantageous in real-world applications where class categorizations often exist within a hierarchical framework.

## 5 EXPLORING THE BENEFITS OF VCRREG

This section aims to thoroughly unpack the multi-faceted benefits of VCRreg in the context of supervised neural network training. Specifically, we discuss its capability to address challenges such as gradient starvation (Pezeshki et al., 2021), neural collapse (Papayan et al., 2020), noisy data, and the preservation of information richness during model training (Shwartz-Ziv, 2022).

### 5.1 MITIGATING GRADIENT STARVATION

In line with the original study on gradient starvation (Pezeshki et al., 2021), we observe that most traditional regularization techniques fall short of capturing the vital features for the “two-moon” dataset experiment. To assess the effectiveness of VCRreg, we replicate this setting with a three-layer network and apply our method during training. Our visualized results in Figure 2 make it apparent that VCRreg has a marked advantage over traditional regularization techniques, particularly in the aspects of separation margins. Thus, it is reasonable to conclude that VCRreg can help mitigate gradient starvation. Please check section E for the detailed information about experiments related to the “two-moon” dataset.

### 5.2 PREVENTING NEURAL COLLAPSE AND INFORMATION COMPRESSION

To deepen our understanding of VCRreg and its training dynamics, we closely examine its learned representations. A recent study (Papayan et al., 2020) observed a peculiar trend in deep networks trained for classification tasks: the top-layer feature embeddings of training samples from the same



**Table 6: VCReg learns richer representation and prevents neural collapse and information compression** Metrics include Class-Distance Normalized Variance (CDNV), Nearest Class-Center Classifier (NCC), and Mutual Information (MI). Higher values indicate reduced neural collapse and richer feature representations.

Network	CDNV	NCC	MI
ConvNeXt	0.28	0.99	2.8
ConvNeXt (VCReg)	<b>0.56</b>	<b>0.81</b>	<b>4.6</b>

class tend to cluster around their respective class means, which are as distant from each other as possible. However, this phenomenon could potentially result in a loss of diversity among the learned features (Papayan et al., 2020), thus curtailing the network’s capacity to grasp the complexity of the data and leading to suboptimal performance for transfer learning (Li et al., 2018). Our neural collapse investigation includes two key metrics.

**Class-Distance Normalized Variance (CDNV)** For a feature map  $f : \mathbb{R}^d \rightarrow \mathbb{R}^p$  and two unlabeled sets of samples  $S_1, S_2 \subset \mathbb{R}^d$ , the CDNV is defined as

$$V_f(S_1, S_2) = \frac{\sigma_f^2(S_1) + \sigma_f^2(S_2)}{2\|\mu_f(S_1) - \mu_f(S_2)\|^2}, \quad (7)$$

where  $\mu_f(S)$  and  $\sigma_f^2(S)$  signify the mean and variance of the set  $\{f(x) \mid x \in S\}$ . This metric measures the degree of clustering of the features, in relation to their distance.

**Nearest Class-Center Classifier (NCC)** This classifier is defined as  $\arg \min_{c \in [C]} \|f(x) - \mu_f(S_c)\|$

According to this measure, during training, collapsed feature embeddings in the penultimate layer become separable, and the classifier converges to the “nearest class-center classifier”.

**Preventing Information Compression** Although effective compression often yields superior representations, overly aggressive compression might cause the loss of crucial information about the target task (Shwartz-Ziv et al., 2018; Shwartz-Ziv & Alemi, 2020; Shwartz-Ziv & LeCun, 2023). To investigate the compression during the learning, we use the mutual information neural estimation (MINE) (Belghazi et al., 2018), a method specifically designed to estimate the mutual information between the input and its corresponding embedded representation. This metric effectively gauges the complexity level of the representation, essentially indicating how much information it encodes.

We evaluate the learned representations of two ConvNeXt models (Liu et al., 2022), which are trained on ImageNet with supervised loss. One model is trained with VCReg, while the other is trained without. As demonstrated in Table 6, both types of collapse, measured by CDNV and NCC, and the mutual information reveal that VCReg representations have significantly more diverse features and information compared to regular training. This suggests that the VCReg mitigates the neural collapse and prevents excessive information compression, two crucial factors that often limit the effectiveness of deep learning models in transfer learning tasks. Our findings highlight the potential of VCReg as a valuable addition to the deep learning toolbox, significantly increasing the generalizability of learned representations.

### 5.3 PROVIDING ROBUSTNESS TO NOISE

In real-world scenarios, encountering noise is a common challenge, making robustness against noise a crucial feature for any effective transfer learning algorithm. Recognizing the ubiquity of noise in practical applications, we aim to evaluate the capability of VCReg to bolster transfer learning performance in noisy environments. For this purpose, we utilize video networks initially pretrained on Kinetics-400 and Kinetics-710, as mentioned in section 4.2. We then finetune these networks on the HMDB51 dataset, which is deliberately subjected to varying levels of Gaussian noise. The findings in Figure 3 reveal a clear advantage: incorporating VCReg notably improves the resilience of VideoMAE-S and VideoMAEv2-S models to noisy data. Appendix D shows that this trend of increased durability against noise is consistently seen in larger models, such as VideoMAE-B and VideoMAEv2-B.

## 6 CONCLUSION

In this work, we address prevalent challenges in supervised pretraining for transfer learning by introducing an efficient and adaptable regularization technique called Variance-Covariance Regularization (VCRReg). Our comprehensive evaluation reveals that using VCRReg yields significant improvements in transfer learning performance across various network architectures, learning paradigms, and data modalities. Moreover, our in-depth analysis confirms VCRReg’s effectiveness in overcoming typical transfer learning hurdles such as neural collapse, gradient starvation, and noisy data. Our work paves the way for further research to achieve highly optimized and generalizable machine learning models.

## REFERENCES

- Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid. Vivit: A video vision transformer. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 6836–6846, 2021.
- Babajide O Ayinde, Tamer Inanc, and Jacek M Zurada. Regularizing deep neural networks by enhancing diversity in feature extraction. *IEEE transactions on neural networks and learning systems*, 30(9):2650–2661, 2019.
- Adrien Bardes, Jean Ponce, and Yann LeCun. Vicreg: Variance-invariance-covariance regularization for self-supervised learning. *arXiv preprint arXiv:2105.04906*, 2021.
- Mohamed Ishmael Belghazi, Aristide Baratin, Sai Rajeswar, Sherjil Ozair, Yoshua Bengio, Aaron Courville, and R Devon Hjelm. Mine: mutual information neural estimation. *arXiv preprint arXiv:1801.04062*, 2018.
- Ido Ben-Shaul, Ravid Shwartz-Ziv, Tomer Galanti, Shai Dekel, and Yann LeCun. Reverse engineering self-supervised learning. *arXiv preprint arXiv:2305.15614*, 2023.
- Yoshua Bengio. Deep learning of representations for unsupervised and transfer learning. In *Proceedings of ICML workshop on unsupervised and transfer learning*, pp. 17–36. JMLR Workshop and Conference Proceedings, 2012.
- Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.
- Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101—mining discriminative components with random forests. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part VI 13*, pp. 446–461. Springer, 2014.
- Rich Caruana. Multitask learning. *Machine learning*, 28:41–75, 1997.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pp. 1597–1607. PMLR, 2020.
- Mircea Cimpoi, Subhansu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. Describing textures in the wild. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3606–3613, 2014.
- Michael Cogswell, Faruk Ahmed, Ross Girshick, Larry Zitnick, and Dhruv Batra. Reducing overfitting in deep networks by decorrelating representations. *arXiv preprint arXiv:1511.06068*, 2015.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.

- 540 Logan Engstrom, Andrew Ilyas, Shibani Santurkar, and Dimitris Tsipras. Robustness (python library),  
541 2019. URL <https://github.com/MadryLab/robustness>.  
542
- 543 Jonas Geiping, Micah Goldblum, Gowthami Somepalli, Ravid Shwartz-Ziv, Tom Goldstein, and  
544 Andrew Gordon Wilson. How much data are augmentations worth? an investigation into scaling  
545 laws, invariance, and implicit regularization. *arXiv preprint arXiv:2210.06441*, 2022.
- 546 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
547 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,  
548 pp. 770–778, 2016.
- 549 Geoffrey E Hinton, Nitish Srivastava, Alex Krizhevsky, Ilya Sutskever, and Ruslan R Salakhutdinov.  
550 Improving neural networks by preventing co-adaptation of feature detectors. *arXiv preprint*  
551 *arXiv:1207.0580*, 2012.  
552
- 553 Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by  
554 reducing internal covariate shift. In *International conference on machine learning*, pp. 448–456.  
555 pmlr, 2015.
- 556 Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan,  
557 Fabio Viola, Tim Green, Trevor Back, Paul Natsev, et al. The kinetics human action video dataset.  
558 *arXiv preprint arXiv:1705.06950*, 2017.
- 559 Agnan Kessy, Alex Lewin, and Korbinian Strimmer. Optimal whitening and decorrelation. *The*  
560 *American Statistician*, 72(4):309–314, 2018.  
561
- 562 Simon Kornblith, Ting Chen, Honglak Lee, and Mohammad Norouzi. Why do better loss functions  
563 lead to less transferable features? *Advances in Neural Information Processing Systems*, 34:  
564 28648–28662, 2021.
- 565 Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained  
566 categorization. In *Proceedings of the IEEE international conference on computer vision workshops*,  
567 pp. 554–561, 2013.  
568
- 569 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- 570 Hildegard Kuehne, Hueihan Jhuang, Estíbaliz Garrote, Tomaso Poggio, and Thomas Serre. Hmdb: a  
571 large video database for human motion recognition. In *2011 International conference on computer*  
572 *vision*, pp. 2556–2563. IEEE, 2011.  
573
- 574 Firas Laakom, Jenni Raitoharju, Alexandros Iosifidis, and Moncef Gabbouj. Wld-reg: A data-  
575 dependent within-layer diversity regularizer. *arXiv preprint arXiv:2301.01352*, 2023.
- 576 Yann LeCun, Léon Bottou, Genevieve B Orr, and Klaus-Robert Müller. Efficient backprop. In *Neural*  
577 *networks: Tricks of the trade*, pp. 9–50. Springer, 2002.  
578
- 579 Chunyuan Li, Heerad Farkhoor, Rosanne Liu, and Jason Yosinski. Measuring the intrinsic dimension  
580 of objective landscapes. *arXiv preprint arXiv:1804.08838*, 2018.
- 581 Kunchang Li, Yali Wang, Yinan He, Yizhuo Li, Yi Wang, Limin Wang, and Yu Qiao. Uniformerv2:  
582 Spatiotemporal learning by arming image vits with video uniformer, 2022.  
583
- 584 Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie.  
585 A convnet for the 2020s. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*  
586 *Pattern Recognition*, pp. 11976–11986, 2022.
- 587 S. Maji, J. Kannala, E. Rahtu, M. Blaschko, and A. Vedaldi. Fine-grained visual classification of  
588 aircraft. Technical report, 2013.
- 589 Ishan Misra and Laurens van der Maaten. Self-supervised learning of pretext-invariant representations.  
590 In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 6707–  
591 6717, 2020.  
592
- 593 Behnam Neyshabur, Srinadh Bhojanapalli, David McAllester, and Nati Srebro. Exploring generaliza-  
tion in deep learning. *Advances in neural information processing systems*, 30, 2017.

- 594 Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number  
595 of classes. In *2008 Sixth Indian Conference on Computer Vision, Graphics & Image Processing*,  
596 pp. 722–729. IEEE, 2008.
- 597  
598 Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *IEEE Transactions on knowledge*  
599 *and data engineering*, 22(10):1345–1359, 2010.
- 600  
601 Vardan Papyan, XY Han, and David L Donoho. Prevalence of neural collapse during the terminal  
602 phase of deep learning training. *Proceedings of the National Academy of Sciences*, 117(40):  
603 24652–24663, 2020.
- 604  
605 Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, and CV Jawahar. Cats and dogs. In *2012*  
*IEEE conference on computer vision and pattern recognition*, pp. 3498–3505. IEEE, 2012.
- 606  
607 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,  
608 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas  
609 Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy,  
610 Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-  
611 performance deep learning library. In *Advances in Neural Information Processing Systems 32*, pp.  
612 8024–8035. Curran Associates, Inc., 2019. URL [http://papers.neurips.cc/paper/](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf)  
613 [9015-pytorch-an-imperative-style-high-performance-deep-learning-library.](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf)  
614 [pdf](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf).
- 615  
616 Mohammad Pezeshki, Oumar Kaba, Yoshua Bengio, Aaron C Courville, Doina Precup, and Guil-  
617 laume Lajoie. Gradient starvation: A learning proclivity in neural networks. *Advances in Neural*  
*Information Processing Systems*, 34:1256–1272, 2021.
- 618  
619 Ravid Shwartz-Ziv. Information flow in deep neural networks. *arXiv preprint arXiv:2202.06749*,  
620 2022.
- 621  
622 Ravid Shwartz-Ziv and Alexander A Alemi. Information in infinite ensembles of infinitely-wide  
623 neural networks. In *Symposium on Advances in Approximate Bayesian Inference*, pp. 1–17. PMLR,  
624 2020.
- 625  
626 Ravid Shwartz-Ziv and Yann LeCun. To compress or not to compress—self-supervised learning and  
627 information theory: A review. *arXiv preprint arXiv:2304.09355*, 2023.
- 628  
629 Ravid Shwartz-Ziv, Amichai Painsky, and Naftali Tishby. Representation compression and general-  
630 ization in deep neural networks, 2018.
- 631  
632 Ravid Shwartz-Ziv, Randall Balestrieri, and Yann LeCun. What do we maximize in self-supervised  
633 learning? *arXiv preprint arXiv:2207.10081*, 2022.
- 634  
635 Ravid Shwartz-Ziv, Randall Balestrieri, Kenji Kawaguchi, Tim GJ Rudner, and Yann LeCun. An  
636 information-theoretic perspective on variance-invariance-covariance regularization. *arXiv preprint*  
637 *arXiv:2303.00633*, 2023.
- 638  
639 Karen Simonyan and Andrew Zisserman. Two-stream convolutional networks for action recognition  
640 in videos. *Advances in neural information processing systems*, 27, 2014.
- 641  
642 Zhan Tong, Yibing Song, Jue Wang, and Limin Wang. Videomae: Masked autoencoders are data-  
643 efficient learners for self-supervised video pre-training. *Advances in neural information processing*  
644 *systems*, 35:10078–10093, 2022.
- 645  
646 Grant Van Horn, Oisín Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alex Shepard, Hartwig Adam,  
647 Pietro Perona, and Serge Belongie. The inaturalist species classification and detection dataset. In  
*Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 8769–8778,  
2018.
- 648  
649 Limin Wang, Bingkun Huang, Zhiyu Zhao, Zhan Tong, Yinan He, Yi Wang, Yali Wang, and Yu Qiao.  
650 Videomae v2: Scaling video masked autoencoders with dual masking. In *Proceedings of the*  
*IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14549–14560, 2023.

648 Karl R. Weiss, Taghi M. Khoshgoftaar, and Dingding Wang. A survey of transfer learning. *Journal*  
649 *of Big Data*, 3, 2016.  
650

651 Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. How transferable are features in deep  
652 neural networks? *Advances in neural information processing systems*, 27, 2014.

653 Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stéphane Deny. Barlow twins: Self-supervised  
654 learning via redundancy reduction. *arXiv preprint arXiv:2103.03230*, 2021.  
655

656 Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding  
657 deep learning requires rethinking generalization. corr abs/1611.03530 (2016). *arXiv preprint*  
658 *arxiv:1611.03530*, 2016.

659 Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep  
660 learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3):107–115,  
661 2021.  
662

663 Bolei Zhou, Agata Lapedriza, Jianxiong Xiao, Antonio Torralba, and Aude Oliva. Learning deep  
664 features for scene recognition using places database. *Advances in neural information processing*  
665 *systems*, 27, 2014.

666 Fuzhen Zhuang, Zhiyuan Qi, Keyu Duan, Dongbo Xi, Yongchun Zhu, Hengshu Zhu, Hui Xiong, and  
667 Qing He. A comprehensive survey on transfer learning. *Proceedings of the IEEE*, 109(1):43–76,  
668 2020.  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701

## A EXPERIMENTAL INVESTIGATION ON EFFECTIVE APPLICATION OF VCREG TO STANDARD NETWORKS

To determine the optimal manner of integrating the VCREg into a standard network, we conducted several experiments utilizing the ConvNeXt-Atto architecture, trained on ImageNet following the torchvision (Paszke et al., 2019) training recipe. To reduce the training time, we limited the network training to 90 epochs with a batch size of 4096. The complete configuration comprised 90 epochs, a batch size of 4096, two learning rate of  $\{0.016, 0.008\}$  with a 5 epochs linear warmup followed by a cosine annealing decay. The weight decay was set at 0.05 and the norm layers were excluded from the weight decay. we experimented with  $\alpha \in \{1.28, 0.64, 0.32, 0.16\}$  and  $\beta \in \{0.16, 0.08, 0.04, 0.02, 0.01\}$ .

We experimented with incorporating the VCREg layers in four different locations:

1. Applying the VCREg exclusively to the second last representation (the input of the classification layer).
2. Applying VCREg to the output of each ConvNeXt block.
3. Applying VCREg to the output of each downsample layer.
4. Applying VCREg to the output of both, each ConvNeXt block and each downsample layer.

The VCREg layer was implemented as detailed in 1, with the addition of a mean removal layer along the batch preceding the VCREg layer to ensure that the VCREg input exhibited a zero mean.

**Table 7: Transfer Learning Experiments with Different VCREg Configurations**

Architecture	Food	Cars	Aircraft	Pets	Flowers	DTD
ConvNeXt-Atto (VCREg1)	63.2%	39.6%	55.9%	89.1%	85.3%	65.1%
ConvNeXt-Atto (VCREg2)	<b>66.8%</b>	48.1%	<b>60.4%</b>	<b>91.1%</b>	<b>86.4%</b>	<b>66.4%</b>
ConvNeXt-Atto (VCREg3)	64.0%	40.9%	56.5%	89.4%	85.9%	65.1%
ConvNeXt-Atto (VCREg4)	66.7%	<b>48.3%</b>	59.6%	90.6%	85.6%	66.1%

The results in Table 7 indicate superior performance when the VCREg layer is applied to the output of each block (second setup) or applied to the output of blocks and downsample layers (fourth setup) compared to the other setups. Considering architectures like ViT lack downsample layers, for consistency across different architectures, we decided to use the second configuration for further experiments.

## B THE FAST IMPLEMENTATION OF THE VCREG

The VCREg does not affect the forward pass in any way, allowing us to substantially speed up the implementation by modifying the backward function directly. Instead of computing the VCREg loss and backpropagating it, we can directly alter the calculated gradient. This is possible since the VCREg loss calculation only requires the current representation. The specifics of this speed-optimized implementation are outlined in Algorithm 1.

We quantify the computational overhead by measuring the average time required for one NVIDIA A100 GPU to execute both the forward and backward passes on the entire network for a batch size of 128 using the ImageNet dataset. These results are summarized in Table 8. For the sake of comparison, we also include the latencies associated with adding Batch Normalization (BN) layers, revealing that our optimized VCREg implementation exhibits similar latencies to BN layers and is almost 5 times faster than the naive implementation.

**Algorithm 1: PyTorch-Style Pseudocode for Fast VCREg Implementation**

```

756
757
758 #  $\alpha$ ,  $\beta$  and  $\epsilon$  : hyperparameters
759 # mm: matrix-matrix multiplication
760 class VarianceCovarianceRegularizationFunction(Function):
761     # forward pass
762     # We assume the input has zero mean per channel
763     # In practice, we apply a batch demean operation before calling the function
764     def forward(ctx, input):
765         ctx.save_for_backward(input)
766         return input
767     # backward pass
768     def backward(ctx, grad_output):
769         input, = ctx.saved_tensors
770         # reshape the input to have (n, d) shape
771         flattened_input = input.flatten(start_dim=0, end_dim=-2)
772         n, d = flattened_input.shape
773         # calculate the covariance matrix
774         covariance_matrix = mm(flattened_input.t(), flattened_input) / (n - 1)
775         # calculate the gradient
776         diagonal = F.threshold(rsqrt(covariance_matrix.diagonal()) + \epsilon, 1.0, 0.0)
777         std_grad_input = diagonal * flattened_input
778         cov_grad_input = torch.mm(flattened_input, covariance_matrix.fill_diagonal_(0))
779
780         grad_input = grad_output
781             -  $\alpha/(d(n-1)) * \text{std\_grad\_input.view(grad\_output)}$ 
782             +  $4\beta/(d(d-1)) * \text{cov\_grad\_input}$ 
783
784         return grad_input

```

**Table 8: Average Time Required for One Forward and Backward Pass with Various Layers Inserted**  
Comparison of computational latencies across different configurations of ViT and ConvNeXt networks. The table demonstrates the efficacy of the optimized VCREg layer in terms of computational time, compared to both naive VCREg and Batch Normalization (BN) layers.

Network	Number of Inserted Layers	Identity	VCREg (Naive)	VCREg (Fast)	BN
ViT-Base-32	12	0.223s	1.427s	0.245s	0.247s
ConvNeXt-T	18	0.442s	2.951s	0.471s	0.468s

## C IMPLEMENTATION DETAILS

### C.1 TRANSFER LEARNING EXPERIMENTS WITH IMAGENET PRETRAINING

In conducting the transfer learning experiments, we adhered primarily to the training recipe specified by PyTorch Paszke et al. (2019) for each respective architecture during the supervised pretraining phase. We abstained from pretraining any of the baseline models, instead opting to directly download the weights from PyTorch’s own repository. The only modifications applied were to the parameters associated with VCREg loss, and we experimented with  $\alpha \in \{1.28, 0.64, 0.32, 0.16\}$  and  $\beta \in \{0.16, 0.08, 0.04, 0.02, 0.01\}$ .

For iNaturalist 18 Van Horn et al. (2018) and Place205 Zhou et al. (2014), we relied on the experimental settings detailed in Zbontar et al. (2021) for the linear probe evaluation.

Regarding Food-101 Bossard et al. (2014), Stanford Cars Krause et al. (2013), FGVC Aircraft Maji et al. (2013), Oxford-IIIT Pets Parkhi et al. (2012), Oxford 102 Flowers Nilsback & Zisserman (2008), and the Describable Textures Dataset (DTD) Cimpoi et al. (2014), we complied with the evaluation protocol provided by Chen et al. (2020); Kornblith et al. (2021). An  $L_2$ -regularized multinomial logistic regression classifier was trained on features extracted from the frozen pretrained network. Optimization of the softmax cross-entropy objective was conducted using L-BFGS, without the application of data augmentation. All images were resized to 224 pixels along the shorter side through bicubic resampling, followed by a 224 x 224 center crop. The  $L_2$ -regularization parameter was selected from a range of 45 logarithmically spaced values between 0.00001 and 100000.

All experiments were run three times, with the average results presented in Table 1.

## 810 C.2 TRANSFER LEARNING EXPERIMENTS WITH KINETICS PRE-TRAINED MODELS

811  
812 In conducting experiments with video-pretrained models, we utilize the publicly available code bases  
813 and model checkpoints provided for VideoMAE and VideoMAEv2 (<https://github.com/MCG-NJU/VideoMAE> and <https://github.com/OpenGVLab/VideoMAEv2>). For both  
814 VideoMAE and VideoMAEv2 we use ViT-Small and ViT-Base checkpoints. VideoMAE models are  
815 pre-trained on Kinetics-400 while VideoMAEv2 on Kinetics-710. We use the pre-trained checkpoint  
816 for ViViT-B (ViT-Base backbone) pre-trained on Kinetics-400 from HuggingFace. For evaluation,  
817 we adopt the inference protocol of 10 clips  $\times$  3 crops. For VCRReg hyperparameters experiments with  
818 values for  $\alpha \in 1, 3, 5$  and  $\beta \in \{0.1, 0.3, 0.5\}$ . For the rest of the finetuning hyperparameters as well  
819 as the data pre-processing and evaluation protocol, we use the configuration for HMDB51 available  
820 in VideoMAE Tong et al. (2022) and its corresponding code base (linked above).  
821

## 823 C.3 SUBCLASS LINEAR PROBING RESULT WITH NETWORK PRETRAINED ON SUPERCLASS LABEL

824 LABEL  
825  
826 For our subclass linear probing experiments, we employed a ConvNeXt-Atto network. Each model  
827 was pretrained for 200 epochs using the superclasses, adhering to the same procedure detailed in the  
828 Appendix A. Subsequent to this pretraining phase, we initiated a linear probing process using the  
829 subclass labels. This linear classifier was trained for 100 epochs, using a base learning rate of 0.016  
830 in conjunction with a cosine learning rate schedule. The optimizer used was AdamW, which worked  
831 to minimize cross-entropy loss with a weight decay set at 0.05. We processed our training data in  
832 batches of 256.  
833

## 834 C.4 LONG-TAIL LEARNING RESULT

835  
836 For our long-tail learning experiments, we use ResNet-32 as a backbone for experiments on the  
837 CIFAR10-LT and CIFAR100-LT datasets. We trained 100 epochs with batch size 256, Adam optimizer  
838 with two learning rate of  $\{0.016, 0.008\}$  with a 10-epoch linear warm-up followed by a cosine  
839 annealing decay. The weight decay was set at 0.05 and the norm layers were excluded from the weight  
840 decay. we experimented with  $\alpha \in \{1.28, 0.64, 0.32, 0.16\}$  and  $\beta \in \{0.16, 0.08, 0.04, 0.02, 0.01\}$ .  
841

## 842 C.5 VCREG WITH SELF-SUPERVISED LEARNING METHODS

843  
844 We train a ResNet-50 model in four different setups, using either the SimCLR loss or the VICReg  
845 loss with the ImageNet dataset. The application of the VCRReg is the same as described in Appendix  
846 A.

847 We closely follow the original setting in Chen et al. (2020) for SimCLR pretraining and Bardes et al.  
848 (2021) for VICReg pretraining.

849 **Augmentation** For both methods, we use the same augmentation methods. Each augmented view  
850 is generated from a random set of augmentations of the same input image. We apply a series of  
851 standard augmentations for each view, including random cropping, resizing to 224x224, random  
852 horizontal flipping, random color-jittering, randomly converting to grayscale, and a random Gaussian  
853 blur. These augmentations are applied symmetrically on two branches Geiping et al. (2022)

854 **Architecture** For SimCLR, the encoder is a ResNet-50 network without the final classification  
855 layer followed by a projector. The projector is a two-layer MLP with input dimension 2048, hidden  
856 dimension 2048, and output dimension 256. The projector has ReLU between the two layers and  
857 batch normalization after every layer. This 256-dimensional embedding is fed to the infoNCE loss.  
858

859 For VICReg, the online encoder is a ResNet-50 network without the final classification layer. The  
860 online projector is a two-layer MLP with input dimension 2048, hidden dimension 8192, and output  
861 dimension 8192. The projector has ReLU between the two layers and batch normalization after every  
862 layer. This 8192-dimensional embedding is fed to the infoNCE loss.

863 For VCRReg, we just applied the VCRReg layers to the ResNet-50 network as described in the Appendix  
A.



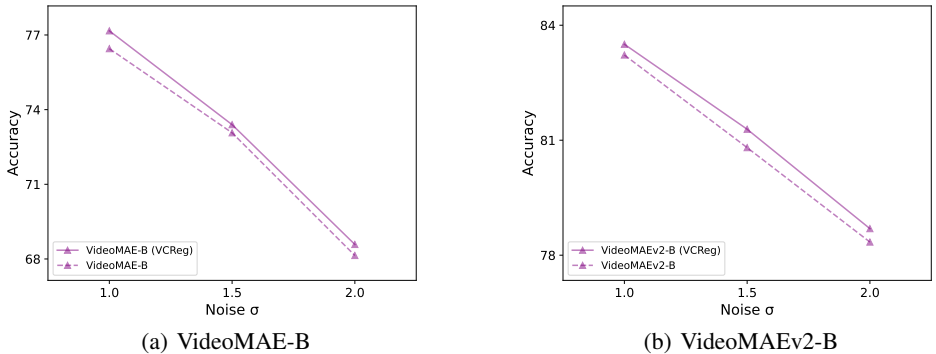
**Optimization** We follow the training protocol in Zbontar et al. (2021). For SimCLR experiments, we used a LARS optimizer and a base learning rate 0.3 with cosine learning rate decay schedule. We pretrain the model for 100 epochs with 5 epochs warm-up with batch size 4096.

For VICReg, we use a LARS optimizer and a base learning rate 0.2 using cosine learning rate decay schedule. We pretrain the model for 100 epochs with 5 epochs warm-up with batch size 4096.

**Evaluation** We follow the standard evaluation protocol as prescribed by Misra & Maaten (2020); Zbontar et al. (2021), performing linear probing evaluations, on iNaturalist 18 Van Horn et al. (2018) and Place205 Zhou et al. (2014) datasets.

## D ROBUSTNESS TO NOISE

This section provides additional results on measuring VCRreg’s ability to enhance transfer learning performance in the presence of noise. In these experiments we start with VideoMAE-B and VideoMAEv2-B networks (from section 4.2) pre-trained on Kinetics-400 and Kinetics-710, respectively, then fine-tune them on HMDB51 corrupted with varying levels of Gaussian noise. During fine-tuning, we compare the transfer learning performance of VideoMAE-B and VideoMAEv2-B networks with and without the addition of VCRreg. When VCRreg is added, it is only applied to the final layer of these networks preceding the classification head. Figure 4 shows that VCRreg models outperform their non-regularized counterparts in this setting.



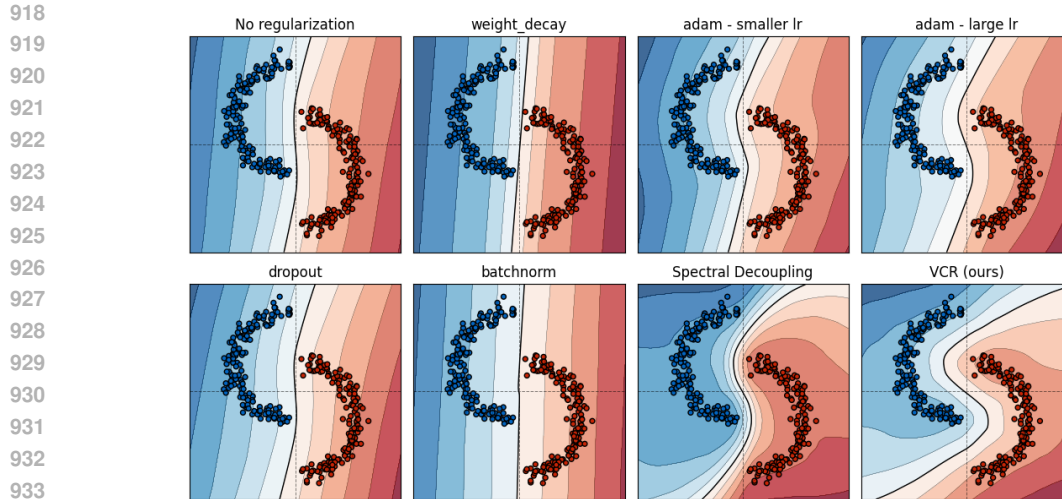
**Figure 4: Impact of VCRreg amidst noisy data:** This figure shows the top-1 accuracy of VideoMAE-B and VideoMAEv2-B when fine-tuned for action recognition using HMDB51 with synthetic noise. We corrupt the data with Gaussian noise with standard deviation  $\sigma \in \{1, 1.5, 2\}$ . Models with VCRreg outperform their non-regularized counterparts in this setting.

## E TWO-MOON DATASET

In alignment with the original gradient starvation study Pezeshki et al. (2021), we notice that most regular routine regularization techniques do not sufficiently capture the necessary features for the “two-moon” dataset experiment. To evaluate our approach, we mirrored this setting and applied the VCRreg during the training.

The synthetic “two-moon” dataset comprises two classes of points, each forming a moon-like shape. The gradient starvation study highlighted an issue where if the gap between the two moons is wide enough for a straight line to separate the two classes, the network stops learning additional features and focuses solely on a single feature. We duplicated this situation using a three-layer network and applied all the initially tested methods in the original study. The resulting decision boundary after training with the “two-moon” dataset is visualized in Figure 5.

From the visualization, it becomes apparent that not only does VCRreg outperform other conventional regularization techniques in separation margins, but also it shows superior performance compared to spectral decoupling, a method specifically designed for this task. VCRreg is effective in maximizing



**Figure 5: The effect of conventional regularization methods and the VCReg on a simple task of two-moon classification.** Shown decision boundaries are the average over 10 runs in which data points and the model initialization parameters are sampled randomly. Here, only the data points of one particular seed are plotted for visual clarity. It can be seen that conventional regularizations of deep learning seem not to help with learning a curved decision boundary.

the variance while minimizing the covariance in the feature space, an achievement that is not obtained by other techniques such as L2, dropout Hinton et al. (2012), and batch normalization Ioffe & Szegedy (2015). Consequently, these other techniques yield features that are less discriminative and informative.

## F COMPUTE RESOURCES

The majority of our experiments were run using AMD MI50 GPUs. The longest pretraining for ConvNeXt-Tiny takes about 48 hours on 2 nodes, where each node has 8 MI50 GPUs attached. We estimate that the total amount of compute resources used for all the experiments can be roughly approximated by  $60 \text{ (days)} \times 24 \text{ (hours per day)} \times 8 \text{ (nodes)} \times 8 \text{ (GPUs per nodes)} = 92,160 \text{ (GPU hours)}$ .

We are aware of potential environmental impact of consuming a lot of compute resources needed for this work, such as atmospheric CO<sub>2</sub> emissions due to the electricity used by the servers. However, we also believe that advancements in representation learning and transfer learning can potentially help mitigate these effects by reducing the need for data and compute resources in the future.

## G LIMITATIONS

Due to a lack of compute resources, we were unable to conduct a large number of experiments with the goal of tuning hyperparameters and searching for the best configurations. Therefore, the majority of hyperparameters and network configurations used in this work are the same as provided by PyTorch Paszke et al. (2019). The only hyperparameters that were tuned were  $\alpha$  and  $\beta$ , the coefficients for VCR. All the other hyperparameters may not be optimal.

In addition, all models were pretrained on the ImageNet Deng et al. (2009) and Krizhevsky et al. (2009) dataset, so their performances might differ if pretrained with other datasets containing different data distributions or different types of images (e.g., x-rays). We encourage further exploration in this direction for current and future self-supervised learning frameworks.