

# 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. Recent progress in self-supervised learning (SSL) has introduced regularization techniques that bolster feature transferability. In this work, we adapt an SSL regularization technique from the VICReg method to supervised learning contexts, introducing Variance-Covariance Regularization (VCRReg). This adaptation encourages the network to learn a high-variance, low-covariance representation, promoting the learning of more diverse features. We outline best practices for implementing this regularization framework into various neural network architectures and present an optimized strategy for regularizing intermediate representations. Through extensive empirical evaluation, we demonstrate that our method significantly enhances transfer learning, achieving excellent performance across numerous tasks and datasets. VCRReg also improves performance in scenarios like long-tail learning, and hierarchical classification. Additionally, we conduct analyses to suggest that its effectiveness may stem from its success in addressing challenges like gradient starvation and neural collapse. In summary, VCRReg offers a universally applicable regularization framework that significantly advances the state of transfer learning, highlights the connection between gradient starvation, neural collapse, and feature transferability, and potentially opens a new avenue for regularization in this domain.

## 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 (VCRReg), aims to encourage the learning of representations with high variance and low covariance, thus avoiding the overemphasis on features that merely minimize supervised loss. Crucially, our detailed tests reveal that the effectiveness of VCRReg strongly depends on how well it is integrated into different neural network designs. Instead of simply applying VCRReg to the final representation of the network, we explore the most effective ways to incorporate it throughout the intermediate representations of the network.

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 neural network architectures. Subsequently, we conduct analyses on the learned representations to demonstrate VCR<sub>g</sub>'s effectiveness in mitigating common issues in transfer learning, such as neural collapse and gradient starvation. This finding suggests a promising avenue for future research in transfer learning: focusing on resolving issues like gradient starvation and neural collapse, particularly in the context of transfer learning, has the potential to significantly improve performance.

Our paper makes the following contributions:

1. We introduce a robust strategy for applying VCR<sub>g</sub> to neural networks, including integrating it into the intermediate layers.
2. We propose a computationally efficient implementation of VCR<sub>g</sub>. This implementation is optimized to ensure minimal impact from additional computational overhead, allowing for seamless integration into existing workflows while maintaining high training speed and resource efficiency.
3. Through extensive experiments on benchmark datasets, we demonstrate that using VCR<sub>g</sub> yields notable improvements 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, with preliminary results, we also find that VCR<sub>g</sub> could improve performance in scenarios like long-tail learning, and hierarchical classification.
4. We investigate the learned representation of VCR<sub>g</sub>, revealing its effectiveness in combating challenges such as gradient starvation (Pezeshki et al., 2021), neural collapse (Papayan et al., 2020), and information compression (Shwartz-Ziv, 2022).

## 2 RELATED WORK

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

VICReg (Bardes et al., 2021) is a novel SSL method. VICReg encourages learned representations to be invariant to data augmentation. However, by optimizing only the invariant criterion, the network will learn to generate a constant representation for all inputs. This means the representations will be invariant not only to data augmentation, but also to the input itself.

VICReg primarily regularizes the network by using a combination of variance loss and covariance loss. The variance loss encourages high variance in the learned representations, thereby promoting the learning of a wide range of 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 has been found to be effective in promoting diverse feature learning (Shwartz-Ziv et al., 2022). In this work, we borrow the feature collapse prevention mechanism from VICReg and propose the variance-covariance regularization method for supervised network training to improve transfer learning performance.

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  $x_i$ . With the neural network  $f_\theta(\cdot)$  and the final representations  $z'_i = f_\theta(x'_i)$  and  $z''_i = f_\theta(x''_i)$ , the VICReg minimizes the following loss (we defer the detailed formulation of the variance and covariance loss terms to the subsequent section where we introduce our methods):

$$\begin{aligned} \ell_{\text{VICReg}}(z'_1, \dots, z'_n, z''_1, \dots, z''_n) &= \alpha \ell_{\text{var}}(z'_1, \dots, z'_n) + \alpha \ell_{\text{var}}(z''_1, \dots, z''_n) \\ &+ \beta \ell_{\text{cov}}(z'_1, \dots, z'_n) + \beta \ell_{\text{cov}}(z''_1, \dots, z''_n) \\ &+ \sum_{i=1}^n \ell_{\text{inv}}(z'_i, z''_i) \end{aligned} \quad (1)$$

Notice that the only loss term that requires two augmented images is the invariance loss. We usually avoid using two augmented images for each training step in supervised learning. This is because it would approximately double the total computation, as we would need to perform two forward passes at each step. Furthermore, as discussed in some previous works (Shwartz-Ziv, 2022; Shwartz-Ziv

et al., 2023), the invariance term is not the essential factor that helps diversify the features. Therefore, in our adaptation to the supervised regime, we omit the invariance term from the regularization.

## 2.2 REPRESENTATION WHITENING AND FEATURE DIVERSITY REGULARIZERS

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

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

However, the studies cited above primarily concentrate on the benefits of optimization and generalization for the source task, frequently overlooking their implications for transfer learning. VCReg sets itself apart by explicitly targeting enhancements in transfer learning performance. Our results indicate that such regularization techniques yield only modest performance improvements in in-domain evaluations. This may be attributed to the fact that modern optimizers and regularizers have already significantly alleviated challenges related to in-domain optimization and generalization. Therefore, the most impactful domain for these types of regularization appears to be transfer learning.

## 2.3 GRADIENT STARVATION AND NEURAL COLLAPSE

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

# 3 VARIANCE-COVARIANCE REGULARIZATION

## 3.1 VANILLA VCReg: AN INTRODUCTION TO THE BASIC FORMULATION

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 Vanilla VCReg is to ensure that the  $D$ -dimensional input representation  $h_i$  to the last layer of the network exhibit both high variance and low covariance. To achieve this, we employ variance and covariance losses, same as mentioned in equation 1:

$$\ell_{\text{vcreg}}(h_1, \dots, h_N) = \alpha \ell_{\text{var}}(h_1, \dots, h_N) + \beta \ell_{\text{cov}}(h_1, \dots, h_N) \quad (2)$$

The variance and covariance loss functions are defined as:

$$\ell_{\text{var}} = \frac{1}{D} \sum_{i=1}^D \max(0, 1 - \sqrt{C_{ii}}) \quad (3)$$

$$\ell_{\text{cov}} = \frac{1}{D(D-1)} \sum_{i \neq j} C_{ij}^2 \quad (4)$$

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

Intuitively speaking, the covariance matrix captures the interdependencies among the dimensions of the feature vectors  $z_i$ . 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 then becomes:

$$L_{\text{vanilla}} = \alpha \ell_{\text{var}}(z_1, \dots, z_N) + \beta \ell_{\text{cov}}(z_1, \dots, z_N) + \frac{1}{N} \sum_{i=1}^N \ell_{\text{sup}}(\tilde{y}_i, y_i) \quad (5)$$

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:

$$L_{\text{VCReg}} = \sum_{j=1}^M \left[ \alpha \ell_{\text{var}}(h_1^{(j)}, \dots, h_N^{(j)}) + \beta \ell_{\text{cov}}(h_1^{(j)}, \dots, h_N^{(j)}) \right] + \frac{1}{N} \sum_{i=1}^N \ell_{\text{sup}}(\tilde{y}_i, y_i) \quad (6)$$

**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, a VCReg is usually applied subsequent to each block within the neural network architecture, often succeeding residual connections. This placement allows for seamless incorporation into current network architecture 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:

$$\text{SmoothL1}(x) = \begin{cases} x^2, & \text{if } |x| \leq \delta \\ 2\delta|x| - \delta^2, & \text{otherwise} \end{cases} \quad (7)$$

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 implementation speed, we take advantage of the fact that VCR<sub>eg</sub> only affects the loss function and not the forward pass. This allows us to focus on directly modifying the backward function for improvements. Specifically, we sidestep the usual process of calculating the VCR<sub>eg</sub> loss and subsequent backpropagation. Instead, we directly adjust the computed gradients, which is feasible since the VCR<sub>eg</sub> loss calculation relies solely on the current representation. Further details of this speed-optimized technique are outlined in Appendix B.

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 1. For the sake of comparison, we also include the latencies associated with adding Batch Normalization (BN) layers, revealing that our optimized VCR<sub>eg</sub> implementation exhibits similar latencies to BN layers.

**Table 1: 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 VCR<sub>eg</sub> layer in terms of computational time, compared to both naive VCR<sub>eg</sub> and Batch Normalization (BN) layers.

Network	Number of Inserted Layers	Identity	VCR <sub>eg</sub> (Naive)	VCR <sub>eg</sub> (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

## 4 EXPERIMENTS

In this section, we initially outline the experimental framework and findings to highlight the effectiveness of our proposed regularization approach, VCR<sub>eg</sub>, within the realm of transfer learning that utilizes supervised pretraining. Subsequent to that discussion, we extend our experiments beyond the scope of supervised pretraining to suggest that VCR<sub>eg</sub> could be applicable across various learning paradigms. For guidelines on reproducing these experiments, please consult Appendix C.

### 4.1 TRANSFER LEARNING WITH SUPERVISED PRETRAINING

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 pretrained 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 followed the standard PyTorch recipes (Paszke et al., 2019) for all networks and did not modify any hyperparameters other than those related to VCR<sub>eg</sub> to ensure a fair baseline comparison. Subsequently, we performed a linear probing evaluation across a variety of datasets to evaluate the transfer learning performance.

For ResNet-50, we included two other feature diversity regularizer methods, namely DeCov (Cogswell et al., 2015) and WLD-Reg (Laakom et al., 2023), for comparison. We conducted experiments solely with ResNet-50 because it is the principal architecture used in the WLD-Reg paper. To ensure a fair comparison, we sourced hyperparameters from Laakom et al. (2023) for both DeCov and WLD-Reg.

The results presented in Table 2 depict significant improvements in transfer learning performance across all downstream datasets when VCR<sub>eg</sub> is applied to the three architectures used. There is strong evidence to suggest that VCR<sub>eg</sub> can help boost overall transfer learning performance, and it is effective for both ConvNet and Transformer architectures.

### 4.2 BEYOND TRANSFER LEARNING WITH SUPERVISED LEARNING

In this section, we explore the versatility of the VCR<sub>eg</sub> regularization method by extending its application beyond transfer learning with supervised pretraining. We focus on three specialized

**Table 2: 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 VCRreg consistently improves performance and beats other feature diversity regularizer.

Architecture	ImageNet	iNat18	Places	Food	Cars	Aircraft	Pets	Flowers	DTD
ResNet-50	76.1%	42.8%	50.6%	69.1%	43.6%	54.8%	91.9%	77.1%	68.7%
ResNet-50 (DeCov)	75.9%	43.1%	50.4%	69.0%	45.7%	55.5%	90.6%	79.2%	69.1%
ResNet-50 (WLD-Reg)	<b>76.5%</b>	43.9%	<b>51.2%</b>	70.2%	43.9%	58.7%	91.4%	80.7%	69.0%
ResNet-50 (VCRreg)	76.3%	<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>
ConvNeXt-T	<b>82.5%</b>	51.6%	53.8%	78.4%	62.9%	74.7%	93.9%	91.3%	72.9%
ConvNeXt-T (VCRreg)	82.4%	<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>
ViT-Base-32	75.9%	39.1%	47.9%	70.6%	51.2%	63.8%	90.3%	84.6%	66.1%
ViT-Base-32 (VCRreg)	<b>76.3%</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>

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 VCRreg across various data distributions and learning paradigms, thereby evaluating its broader utility in machine learning applications.

**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 conducted experiments to assess how well VCRreg addresses this issue through long-tail learning. We evaluated VCRreg using the CIFAR10-LT and CIFAR100-LT Krizhevsky et al. (2009) datasets, both engineered to have an imbalance ratio of 100. These experiments were conducted using a ResNet-32 backbone architecture. The per-class sample sizes ranged from 5,000 to 50 for CIFAR10-LT and from 500 to 5 for CIFAR100-LT.

**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 3 shows that models augmented with VCRreg consistently outperformed 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 establishes VCRreg as a valuable tool for real-world applications where class imbalance is often a concern.

**Enhancing Self-Supervised Learning with VCRreg** Our subsequent investigation focuses on examining the synergy between VCRreg and existing self-supervised learning paradigms. We employed 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 conducted 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 resulted in consistent performance improvements for transfer learning. Specifically, the linear probing accuracies were 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.

**Investigating Hierarchical Classification Capabilities** To evaluate the efficacy of the learned representations across multiple levels of class granularity, we conducted 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

**Table 4: Impact of VCRreg on Self-Supervised Learning Methods:** This table presents a comparative analysis of ResNet-50 models pretrained with SimCLR and VCRreg 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.

Pretraining Methods	ImageNet	iNat18	Places	Food	Cars	Aircraft	Pets	Flowers	DTD
SimCLR	67.2%	37.2%	52.1%	66.4%	35.7%	<b>62.3%</b>	76.3%	82.6%	68.1%
SimCLR (VCRreg)	<b>67.1%</b>	<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>
VCRreg	65.2%	<b>41.7%</b>	48.2%	61.0%	27.3%	51.2%	79.1%	74.3%	65.4%
VCRreg (VCRreg)	<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>

that while samples sharing the same subclass label also share the same superclass label, the reverse does not necessarily hold true. Initially, the model was trained using only the superclass labels, i.e., the  $(x_i, y_i^{\text{sup}})$  pairs. Subsequently, linear probing was employed with the subclass labels  $(x_i, y_i^{\text{sub}})$  to assess the quality of features abstracted at the superclass level.

**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. The models were initially trained using superclass labels and subsequently probed using subclass labels. VCRreg consistently boosts performance in subclass classification tasks.

	Subsets of ImageNet					
	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>

Table 5 presents key performance metrics, highlighting the substantive 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. This attribute is particularly advantageous in real-world applications where class categorizations often exist within a hierarchical framework.

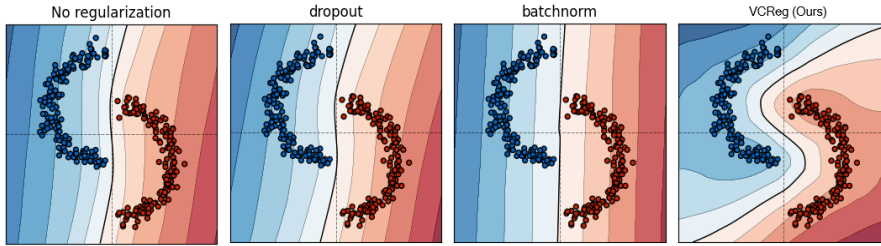
## 5 EXPLORING THE BENEFITS OF VCRREG

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), 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 replicated this setting with a three-layer network and applied our method during training. Our visualized results in Figure 1 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.

These results are significant for multiple reasons. Firstly, the encouraging outcomes with the ‘two-moon’ synthetic dataset set the stage for investigating VCRreg’s applicability in more complex, high-



**Figure 1: Comparative Evaluation of VCRreg and Traditional Regularization Techniques 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 visual clarity. The contrast between VCRreg and conventional methods underscores the latter’s limitations in forming intricate decision boundaries, while highlighting VCRreg’s effectiveness in generating meaningful ones.

dimensional tasks, thus cementing its status as a potent tool in contemporary machine learning. Second, VCRreg’s capability to mitigate gradient starvation indicates that neural networks trained using this method excel at learning complex, non-linear mappings—an essential trait for tasks that require a sophisticated understanding of data distributions. Lastly, VCRreg surpasses traditional regularization techniques by generating a feature space that is both discriminative and rich in information. This highlights its potential to boost the generalizability of neural networks, which is crucial in real-world scenarios where models need to be both robust and flexible.

## 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 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 (Li et al., 2018) for transfer learning.

Our investigation is based on 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

$$V_f(S_1, S_2) = \frac{\text{Var}_f(S_1) + \text{Var}_f(S_2)}{2\|\mu_f(S_1) - \mu_f(S_2)\|^2}, \quad (8)$$

where  $\mu_f(S)$  and  $\text{Var}_f(S)$  signify the mean and variance of the set  $\{f(x) \mid x \in S\}$ . This metric measures the degree of clustering of the features extracted from  $S_1$  and  $S_2$ , in relation to the distance between their respective features. A value approaching zero indicates perfect clustering.

**Nearest Class-Center Classifier (NCC)** This classifier is defined as

$$\hat{h}(x) = \arg \min_{c \in [C]} \|f(x) - \mu_f(S_c)\| \quad (9)$$

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** We next address the prevention of information compression during the learning process. 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 this, 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 (in terms of number of bits) it encodes.



**Table 6: VCRReg 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 in each metric for the VCRReg model indicate reduced neural collapse and richer feature representations.

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

We evaluate the learned representations of two ConvNeXt models (Liu et al., 2022), which are trained on ImageNet with supervised loss. One model was trained with VCRReg, while the other was trained without VCRReg. As demonstrated in Table 6, both types of collapse, measured by CDNV and NCC, and the mutual information estimation reveal that VCRReg representations have significantly more diverse features (lower neural collapse) and contain more information compared to regular training. This suggests that not only does VCRReg achieve superior results, but also its underlying representation contains more information.

In summary, the VCRReg method mitigates the neural collapse phenomenon 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 VCRReg as a valuable addition to the deep learning toolbox, significantly increasing the generalizability of learned representations.

## 6 CONCLUSION

In this work, we addressed prevalent challenges in supervised pretraining for transfer learning by introducing Variance-Covariance Regularization (VCRReg). Building on the regularization technique of the self-supervised VICReg method, VCRReg is designed to cultivate robust and generalizable features. Unlike conventional methods that attach regularization only to the final layer, we strategically incorporate VCRReg across intermediate layers to optimize its efficacy.

Our key contributions are threefold:

1. We present a computationally efficient VCRReg implementation that is adaptable to various network architectures.
2. We provide empirical evidence through comprehensive evaluations on multiple benchmarks, demonstrating that using VCRReg yields notable improvements in transfer learning performance across various network architectures and different learning paradigms.
3. Our in-depth analyses confirm VCRReg’s effectiveness in overcoming typical transfer learning hurdles such as neural collapse and gradient starvation.

To conclude, VCRReg stands out as a potent and adaptable regularization technique that elevates the quality and applicability of learned representations. It enhances both the performance and reliability of models in transfer learning settings, and paves the way for further research aimed at achieving highly optimized and generalizable machine learning models.

## REFERENCES

- 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, Subhransu 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.
- 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.
- Logan Engstrom, Andrew Ilyas, Shibani Santurkar, and Dimitris Tsipras. Robustness (python library), 2019. URL <https://github.com/MadryLab/robustness>.
- Jonas Geiping, Micah Goldblum, Gowthami Somepalli, Ravid Shwartz-Ziv, Tom Goldstein, and Andrew Gordon Wilson. How much data are augmentations worth? an investigation into scaling laws, invariance, and implicit regularization. *arXiv preprint arXiv:2210.06441*, 2022.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Agnan Kessy, Alex Lewin, and Korbinian Strimmer. Optimal whitening and decorrelation. *The American Statistician*, 72(4):309–314, 2018.
- Simon Kornblith, Ting Chen, Honglak Lee, and Mohammad Norouzi. Why do better loss functions lead to less transferable features? *Advances in Neural Information Processing Systems*, 34: 28648–28662, 2021.
- Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *Proceedings of the IEEE international conference on computer vision workshops*, pp. 554–561, 2013.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- Firas Laakom, Jenni Raitoharju, Alexandros Iosifidis, and Moncef Gabbouj. Wld-reg: A data-dependent within-layer diversity regularizer. *arXiv preprint arXiv:2301.01352*, 2023.
- Yann LeCun, Léon Bottou, Genevieve B Orr, and Klaus-Robert Müller. Efficient backprop. In *Neural networks: Tricks of the trade*, pp. 9–50. Springer, 2002.

- Chunyuan Li, Heerad Farkhoor, Rosanne Liu, and Jason Yosinski. Measuring the intrinsic dimension of objective landscapes. *arXiv preprint arXiv:1804.08838*, 2018.
- Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11976–11986, 2022.
- S. Maji, J. Kannala, E. Rahtu, M. Blaschko, and A. Vedaldi. Fine-grained visual classification of aircraft. Technical report, 2013.
- Ishan Misra and Laurens van der Maaten. Self-supervised learning of pretext-invariant representations. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 6707–6717, 2020.
- Behnam Neyshabur, Srinadh Bhojanapalli, David McAllester, and Nati Srebro. Exploring generalization in deep learning. *Advances in neural information processing systems*, 30, 2017.
- Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of classes. In *2008 Sixth Indian Conference on Computer Vision, Graphics & Image Processing*, pp. 722–729. IEEE, 2008.
- Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *IEEE Transactions on knowledge and data engineering*, 22(10):1345–1359, 2010.
- Vardan Papyan, XY Han, and David L Donoho. Prevalence of neural collapse during the terminal phase of deep learning training. *Proceedings of the National Academy of Sciences*, 117(40):24652–24663, 2020.
- 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.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems 32*, pp. 8024–8035. Curran Associates, Inc., 2019. URL <http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>.
- Mohammad Pezeshki, Oumar Kaba, Yoshua Bengio, Aaron C Courville, Doina Precup, and Guillaume Lajoie. Gradient starvation: A learning proclivity in neural networks. *Advances in Neural Information Processing Systems*, 34:1256–1272, 2021.
- Ravid Shwartz-Ziv. Information flow in deep neural networks. *arXiv preprint arXiv:2202.06749*, 2022.
- Ravid Shwartz-Ziv and Alexander A Alemi. Information in infinite ensembles of infinitely-wide neural networks. In *Symposium on Advances in Approximate Bayesian Inference*, pp. 1–17. PMLR, 2020.
- Ravid Shwartz-Ziv and Yann LeCun. To compress or not to compress—self-supervised learning and information theory: A review. *arXiv preprint arXiv:2304.09355*, 2023.
- Ravid Shwartz-Ziv, Amichai Painsky, and Naftali Tishby. Representation compression and generalization in deep neural networks, 2018.
- Ravid Shwartz-Ziv, Randall Balestriero, and Yann LeCun. What do we maximize in self-supervised learning? *arXiv preprint arXiv:2207.10081*, 2022.
- Ravid Shwartz-Ziv, Randall Balestriero, Kenji Kawaguchi, Tim GJ Rudner, and Yann LeCun. An information-theoretic perspective on variance-invariance-covariance regularization. *arXiv preprint arXiv:2303.00633*, 2023.

- Grant Van Horn, Oisin Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alex Shepard, Hartwig Adam, 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.
- Karl R. Weiss, Taghi M. Khoshgoftaar, and Dingding Wang. A survey of transfer learning. *Journal of Big Data*, 3, 2016.
- Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. How transferable are features in deep neural networks? *Advances in neural information processing systems*, 27, 2014.
- Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stéphane Deny. Barlow twins: Self-supervised learning via redundancy reduction. *arXiv preprint arXiv:2103.03230*, 2021.
- Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. corr abs/1611.03530 (2016). *arXiv preprint arxiv:1611.03530*, 2016.
- Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3):107–115, 2021.
- Bolei Zhou, Agata Lapedriza, Jianxiong Xiao, Antonio Torralba, and Aude Oliva. Learning deep features for scene recognition using places database. *Advances in neural information processing systems*, 27, 2014.
- Fuzhen Zhuang, Zhiyuan Qi, Keyu Duan, Dongbo Xi, Yongchun Zhu, Hengshu Zhu, Hui Xiong, and Qing He. A comprehensive survey on transfer learning. *Proceedings of the IEEE*, 109(1): 43–76, 2020.

## 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 this 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.

## 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 pre-training 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.

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

---

```

#  $\alpha$ ,  $\beta$  and  $\epsilon$  : hyperparameters
# mm: matrix-matrix multiplication

class VarianceCovarianceRegularizationFunction(Function):
    # forward pass
    # We assume the input has zero mean per channel
    # In practice, we apply a batch demean operation before call the function
    def forward(ctx, input):
        ctx.save_for_backward(input)
        return input
    # backward pass
    def backward(ctx, grad_output):
        input, = ctx.saved_tensors
        # reshape the input to have (n, d) shape
        flattened_input = input.flatten(start_dim=0, end_dim=-2)
        n, d = flattened_input.shape
        # calculate the covariance matrix
        covariance_matrix = mm(flattened_input.t(), flattened_input) / (n - 1)
        # calculate the gradient
        diagonal = F.threshold(rsqrt(covariance_matrix.diagonal()) + \epsilon, 1.0, 0.0)
        std_grad_input = diagonal * flattened_input
        cov_grad_input = torch.mm(flattened_input, covariance_matrix.fill_diagonal_(0))

        grad_input = grad_output
            -  $\alpha/(d(n - 1)) * \text{std\_grad\_input.view(grad\_output)}$ 
            +  $4\beta/(d(d - 1)) * \text{cov\_grad\_input}$ 

        return grad_input

```

---

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 2.

## C.2 SUBCLASS LINEAR PROBING RESULT WITH NETWORK PRETRAINED ON SUPERCLASS LABEL

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

## C.3 LONG-TAIL LEARNING RESULT

For our long-tail learning experiments, we use ResNet-32 as a backbone for experiments on the CIFAR10-LT and CIFAR100-LT datasets. We trained 100 epochs with batch size 256, Adam optimizer with two learning rate of {0.016, 0.008} with a 10-epoch linear warm-up 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  $\beta \in \{0.16, 0.08, 0.04, 0.02, 0.01\}$ .

#### C.4 VCReg WITH SELF-SUPERVISED LEARNING METHODS

We trained a ResNet-50 model in four different setups, using either the SimCLR loss or the VICReg loss with the ImageNet dataset. The application of the VCReg is the same as described in Appendix A.

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

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

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

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

For VCReg, we just applied the VCReg 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 followed 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.