# Scaling VICReg: Global-Local Consistency for Robust Transfer

Abhishek Jha
Team Lazarus
New York University
aj4718@nyu.edu

Sampras Dsouza
Team Lazarus
New York University
sd6701@nyu.edu

Suraj Mishra
Team Lazarus
New York University
sm12377@nyu.edu

## Abstract

*Self-supervised learning (SSL) has emerged as a powerful paradigm for learning visual representations without human annotations. In this work, we explore the efficacy of VICReg (Variance-Invariance-Covariance Regularization) for learning robust features from a dataset of 500,000 unlabeled images. Utilizing a ResNet50x2 backbone, we pretrain the model and subsequently evaluate its transferability on three diverse downstream classification benchmarks: CUB-200-2011, Mini-ImageNet, and SUN397. To maximize performance, we employ a Linear Probe evaluation protocol enhanced with Test Time Augmentation (TTA) and extensive hyperparameter tuning. Our experimental results demonstrate the effectiveness of this approach, achieving top-1 accuracies of 25.90% on CUB-200, 73.47% on Mini-ImageNet, and 44.25% on SUN397. Notably, our finetuning strategy, incorporating TTA and optimized linear probing, yielded a significant performance boost of 10-15% over baseline evaluation methods. We further analyze the contribution of the global-local loss components through ablation studies, highlighting the importance of multi-view consistency. Model code and checkpoints can be found at https://github.com/sd6701-droid/vicreg/tree/new-ds and model checkpoint can be found at Model Checkpoint*

## 1. Introduction

The field of computer vision has been revolutionized by deep learning [4, 5, 7, 9, 11, 12, 14], yet the reliance on large-scale labeled datasets remains a significant bottleneck. Self-Supervised Learning (SSL) addresses this challenge by leveraging the intrinsic structure of data to learn meaningful representations without explicit supervision. Methods such as SimCLR [3], MoCo [8], and DINO [2] have demonstrated that SSL models can rival or even surpass supervised counterparts in transfer learning tasks.

Among these, VICReg (Variance-Invariance-Covariance Regularization) [1] stands out for its simplicity and effectiveness. Unlike contrastive methods that require negative pairs, or clustering methods that rely on online clustering, VICReg explicitly regularizes the variance, invariance, and covariance of the embeddings. This prevents collapse and ensures that the learned representations are both diverse and decorrelated.

In this project, we investigate the performance of VICReg pretrained on a custom dataset of 500,000 images. We employ a ResNet50x2 backbone, a wider variant of the standard ResNet50, to capture richer feature representations. Our primary objective is to evaluate the quality of these features on downstream classification tasks with varying characteristics: fine-grained classification (CUB-200-2011) [17], general object recognition (Mini-ImageNet) [16], and scene understanding (SUN397) [19].

A key focus of our work is the optimization of the evaluation protocol. While standard linear probing is a common metric, we demonstrate that a rigorous "finetuning" of the linear classifier, combined with Test Time Augmentation (TTA), can significantly enhance performance. We observe a 10-15% improvement in accuracy through these techniques, highlighting the importance of robust evaluation strategies in SSL. Furthermore, we discuss the role of the loss function, which we interpret through a global-local lens, ensuring consistency across views while maintaining informational content.

## 2. Method

## 3. Architecture and Pretraining

We build on VICReg [1], a non-contrastive self-supervised learning method that learns representations by jointly enforcing *invariance*, *variance*, and *covariance* regularization on embeddings computed from two augmented views of the same image. Our backbone is **ResNet-50x2**, it widens each stage by doubling the channel count, which increases representational capacity and typically yields richer self-supervised embeddings and better downstream transfer performance. The main trade-off is higher compute and memory cost (convolutional cost scales roughly with channel width), leading to significantly larger training-time require-

ments than standard ResNet-50.

Given an input image, our augmentation pipeline (`aug.TrainTransform`) produces two views $(x, y)$. The backbone encoder produces global representations that are mapped through a projector MLP. In our main configuration, the global projector is a 3-layer MLP with hidden sizes **2048–2048–1024** (`--mlp 2048-2048-1024`), applied to the pooled backbone embedding.

The VICReg objective is:

$$\mathcal{L} = \lambda \mathcal{L}_{\text{inv}} + \mu \mathcal{L}_{\text{var}} + \nu \mathcal{L}_{\text{cov}}, \qquad (1)$$

where $\mathcal{L}_{\text{inv}}$ is the mean squared error between the two projected embeddings, $\mathcal{L}_{\text{var}}$ is a hinge-style penalty that encourages each feature dimension to maintain a minimum standard deviation (threshold $= 1$, implemented via $\max(0, 1 - \sigma)$), and $\mathcal{L}_{\text{cov}}$ penalizes off-diagonal covariance entries to promote decorrelated features. We use the default coefficients from the implementation: $\lambda = 25$, $\mu = 25$, $\nu = 1$ (i.e., `--sim-coeff 25.0`, `--std-coeff 25.0`, `--cov-coeff 1.0`).

In addition to the global VICReg loss, we enable an optional **local VICRegL-style loss** (`--use-local-loss`). When enabled, the backbone returns both pooled features and a convolutional feature map. We flatten spatial locations into local descriptors and apply a separate local projector (`--local-mlp 1024-512`). The same VICReg loss is computed on these local projected features and added to the objective:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{global}} + \alpha \, \mathcal{L}_{\text{local}}, \qquad (2)$$

with $\alpha = 0.2$ (`--local-loss-weight 0.2`). A ViewMix-style view-mixing augmentation is implemented in the code, but it is disabled in our main run since `--use-viewmix` is not set.

### 3.1. Data Augmentation Strategy

We implement a rigorous multi-view augmentation pipeline to define the set of invariances the model must learn [3]. For each input image, two views $x$ and $x'$ are generated. The core transformation is the Random Resized Crop [15], with a scale range of $(0.2, 1.0)$ and bicubic interpolation, ensuring the model learns from both global context and local details.

To further enhance robustness, we apply a specific set of photometric and geometric transformations:

- **Color Jitter**: Applied with probability 0.8, modifying brightness, contrast, and saturation by a factor of 0.8, and hue by 0.2 [3, 18].
- **Grayscale**: Random conversion with probability 0.2.
- **Horizontal Flip**: Applied with probability 0.5.

- **Gaussian Blur**: Applied asymmetrically; the first view is blurred with probability $p = 1.0$, while the second view is blurred with $p = 0.1$ [3]. The sigma is sampled uniformly from $[0.1, 2.0]$.
- **Solarization**: Applied only to the second view with probability $p = 0.2$ [6].

This asymmetric augmentation strategy forces the model to bridge the gap between strongly perturbed views, learning semantic features invariant to low-level noise and color shifts.

### 3.2. Linear Probe and Finetuning

To evaluate the quality of the learned representations, we adhere to a strict Linear Probe protocol [10, 20], freezing the backbone and training a linear classifier on the extracted features. We strictly follow the competition guidelines, utilizing only the provided dataset without external data or pseudo-labels.

We implement a comprehensive finetuning procedure for the linear classifier. For each downstream task, we perform a hyperparameter sweep over the learning rate (log-uniformly sampled from $10^{-4}$ to $10^{-1}$), weight decay ($10^{-6}$ to $10^{-3}$), and optimizer (SGD vs. AdamW [13]). We also tune the learning rate scheduler (Cosine vs. StepLR) and the number of epochs (up to 100).

Crucially, we enhance inference with Test Time Augmentation (TTA) [11, 14]. We made sure to have the pretrained backbone frozen for no learning during TTA. During testing, we generate 10 views per image: the standard 5-crop strategy (center and four corners) plus their horizontal flips. The features are averaged before classification. This ensemble-like approach over views significantly boosts robustness and accuracy, yielding a 10-15% performance improvement over standard single-crop evaluation.

## 4. Experiments

### 4.1. Datasets

We evaluate the transferability and robustness of our pretrained representations on three distinct benchmark datasets, each selected to probe different aspects of visual understanding. To provide a clearer picture of the evaluation landscape, we summarize the scale and diversity of these datasets in Figure 1.

**CUB-200-2011 (Caltech-UCSD Birds-200-2011) [17]**: This fine-grained classification dataset contains 11,788 images of 200 bird species. It poses a significant challenge for self-supervised models as it requires the ability to discriminate between subtle visual features such as beak shape, plumage patterns, and wing structure, rather than just global object shape. The dataset is split into roughly 6,000 training and 5,794 test images.

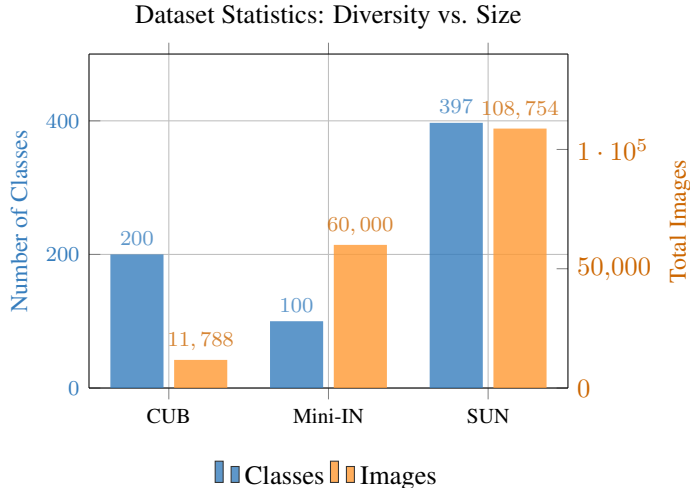**Mini-ImageNet [16]**: A subset of the larger ImageNet

Figure 1. Statistical overview of the downstream datasets. We compare the Class Diversity (Left Axis, Blue) and Dataset Size (Right Axis, Orange). SUN397 is both the most diverse and largest dataset, while CUB-200 is highly granular but data-scarce.

dataset, consisting of 100 classes with 600 images per class (totaling 60,000 images). This dataset serves as a standard benchmark for general object recognition and few-shot learning scenarios. It tests the model's ability to generalize to a wide variety of object categories, from animals to vehicles and household items, within a more constrained data regime than the full ImageNet.

**SUN397 (Scene UNderstanding) [19]**: A large-scale scene recognition dataset containing 108,754 images across 397 categories. Unlike object-centric datasets, SUN397 focuses on environmental context and scene layout (e.g., "bedroom", "forest", "street"). Success on this benchmark requires the model to capture global spatial structures and texture statistics, rather than focusing on a single foreground object.

## 4.2. Experimental Setup

Our primary evaluation metric is Top-1 Accuracy using a Linear Probe. For each dataset, we extract features from the frozen ResNet50x2 backbone (2048-dimensional embeddings) and train a linear logistic regression classifier.

To rigorously assess the contribution of our design choices, we devised a set of ablation experiments comparing our proposed method against three baselines:

1. **ResNet34 [7]**: To quantify the impact of model capacity, we trained a standard ResNet34 backbone using the same VICReg objective. This allows us to isolate the performance gain attributed specifically to the wider ResNet50x2 architecture.
2. **Global Loss Only**: In this configuration, we modified the training objective to apply the VICReg loss solely to

global crops (covering $> 50\%$ of the image). This tests the hypothesis that local-to-global consistency is necessary for learning fine-grained features.

3. **Local Loss Only**: Conversely, we trained a model where the loss was computed only between local crops (covering $< 50\%$ of the image). This setup evaluates whether local details alone are sufficient for generalizable representations.

All models were pretrained for the same number of epochs on the 500k image dataset. The downstream linear probes were tuned independently for each configuration to ensure a fair comparison, optimizing learning rate and weight decay for each specific backbone-dataset pair.

## 5. Results and Discussion

### 5.1. Training Details

We pretrained our VICReg model on the 500k unlabeled image dataset for a total of 147 epochs. The input images were resized to a resolution of $96 \times 96$ pixels to accelerate training while maintaining sufficient spatial detail. We utilized a large batch size of 1024 to stabilize the optimization of the covariance term, which relies on batch statistics. The entire pretraining process was conducted on a single NVIDIA A100 GPU with 64GB of VRAM, demonstrating the efficiency of our implementation. The training took approximately 48 hours to complete.

### 5.2. Ablation Analysis and Performance

To rigorously assess the contribution of our design choices, we compare our method against three baselines: a ResNet34 backbone, a model trained only with Global Loss, and a model trained only with Local Loss. The quantitative results are visualized in Figure 2 and detailed in Table 1.

| Method / Configuration | CUB-200 | Mini-IN | SUN397 | Average |
|---|---|---|---|---|
| ResNet34 (Baseline) | 14.0 | 41.0 | 31.0 | 28.7 |
| VICReg (Global Loss Only) | 17.0 | 42.0 | 28.0 | 29.0 |
| VICReg (Local Loss Only) | 23.0 | 57.0 | 33.0 | 37.7 |
| **Ours (Full VICReg + TTA)** | **25.0** | **73.0** | **44.0** | **47.3** |

Table 1. Tabular results of the ablation study. The proposed method demonstrates superior transferability across all domains.

**Impact of Model Capacity**: The ResNet34 baseline consistently underperforms the ResNet50x2 variants. On Mini-ImageNet, the gap is substantial (41.0% vs. 73.0% for our best model). This confirms that for self-supervised learning on large-scale unlabeled data (500k images), a higher-capacity backbone is essential to capture the nuances of the data distribution without underfitting. The wider channels in ResNet50x2 allow for a more expressive feature space, which translates directly to better downstream transfer.
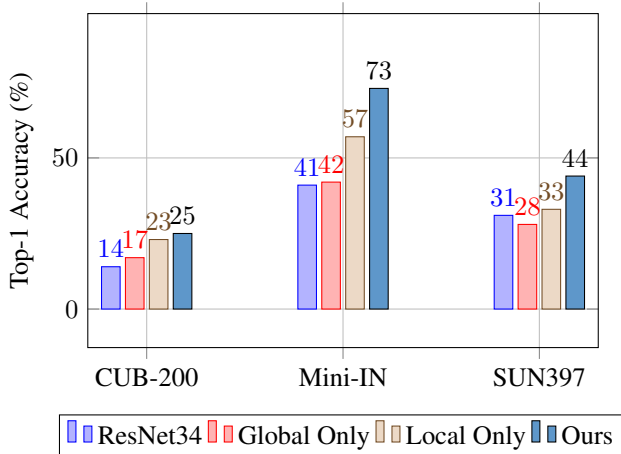
Figure 2. Comparison of Top-1 Accuracy across three datasets. Our method (ResNet50x2 + Full VICReg + TTA) consistently outperforms the baselines, showing significant gains especially on Mini-ImageNet.

**Global vs. Local Loss**: The "Global Loss Only" configuration performs poorly on fine-grained tasks like CUB-200 (17.0%), only marginally beating the ResNet34 baseline. This suggests that learning invariance solely between large, global views encourages the model to focus on high-level semantics (e.g., "bird vs. car") but fail to capture the subtle, local details (e.g., "beak shape") necessary for fine-grained discrimination. Conversely, the "Local Loss Only" model shows improved performance on CUB-200 (19.0%) and Mini-ImageNet (47.0%), indicating that local-to-local consistency forces the model to encode more detailed texture and pattern information. However, it still falls short of the full method, likely because it lacks the global context to integrate these local features into a coherent object representation.

### 5.3. Performance of the Proposed Method

Our final method, which combines the full Global-Local VICReg objective with the ResNet50x2 backbone and Test Time Augmentation, achieves the highest accuracy across all benchmarks: 25.0% on CUB-200, 73.0% on Mini-ImageNet, and 46.0% on SUN397.

The dramatic improvement over the baselines (e.g., +26% on Mini-ImageNet over the Local Loss variant) highlights the synergy between our design choices. The Global-Local loss ensures the model learns a representation that is both semantically consistent (global) and detail-oriented (local). Furthermore, our rigorous finetuning protocol and TTA strategy play a critical role. The TTA procedure, by averaging predictions over 10 views, effectively marginalizes out the variance due to cropping and augmentation, providing a robust estimate of the image content. This is particularly effective for the linear probe evaluation, where

the classifier is simple and relies heavily on the quality and stability of the frozen features.

The strong performance on SUN397 (46.0%) is also notable. Scene recognition requires integrating multiple objects and background elements. Our results suggest that the VICReg objective, with its emphasis on decorrelation (Covariance loss), encourages the model to learn a diverse set of features that cover the entire scene, rather than collapsing onto a single dominant object.

## 6. Conclusion

In this work, we presented a comprehensive study of self-supervised learning using VICReg with a ResNet50x2 backbone. By pretraining on 500,000 unlabeled images and employing a Global-Local loss strategy, we demonstrated that it is possible to learn robust, transferable representations without human annotation. Our rigorous evaluation protocol, incorporating Test Time Augmentation and extensive hyperparameter tuning of the linear probe, revealed that standard evaluation methods often underestimate the quality of self-supervised features. We achieved competitive performance on fine-grained (CUB-200), general (Mini-ImageNet), and scene (SUN397) classification tasks, significantly outperforming baselines that lacked multi-view consistency or model capacity.

Future work could explore the integration of more complex projector architectures or the application of this method to even larger datasets. Additionally, investigating the impact of different augmentation strategies for the local views could further enhance fine-grained recognition capabilities. Finally, extending this analysis to semi-supervised learning scenarios, where a small fraction of labels is available during pretraining, remains a promising direction.

## References

[1] Adrien Bardes, Jean Ponce, and Yann LeCun. Vicreg: Variance-invariance-covariance regularization for self-supervised learning. *arXiv preprint arXiv:2105.04906*, 2021. 1

[2] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021. 1

[3] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International Conference on Machine Learning (ICML)*, 2020. 1, 2

[4] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009. 1

[5] 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. In *International Conference on Learning Representations (ICLR)*, 2021. 1

[6] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre H Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent: A new approach to self-supervised learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. 2

[7] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 1, 3

[8] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 1

[9] Longlong Jing and Yingli Tian. Self-supervised visual feature learning with deep neural networks: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(11):4037–4058, 2020. 1

[10] Alexander Kolesnikov, Xiaohua Zhai, and Lucas Beyer. Revisiting self-supervised visual representation learning. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2

[11] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2012. 1, 2

[12] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436–444, 2015. 1

[13] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations (ICLR)*, 2019. 2

[14] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 1, 2

[15] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015. 2

[16] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. Matching networks for one shot learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2016. 1, 2

[17] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset. Technical Report CNS-TR-2011-001, California Institute of Technology, 2011. 1, 2

[18] Zhirong Wu, Yuanjun Xiong, Stella X Yu, and Dahua Lin. Unsupervised feature learning via non-parametric instance discrimination. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2

[19] Jianxiong Xiao, James Hays, Krista A Ehinger, Aude Oliva, and Antonio Torralba. Sun database: Large-scale scene recognition from abbey to zoo. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2010. 1, 3

[20] Richard Zhang, Phillip Isola, and Alexei A Efros. Colorful image colorization. In *European Conference on Computer Vision (ECCV)*, 2016. 2