# FroSSL: Frobenius Norm Minimization for Efficient Multiview Self-Supervised Learning

Oscar Skean[1], Aayush Dhakal[2],
Nathan Jacobs[2], and Luis Gonzalo Sanchez Giraldo[1]
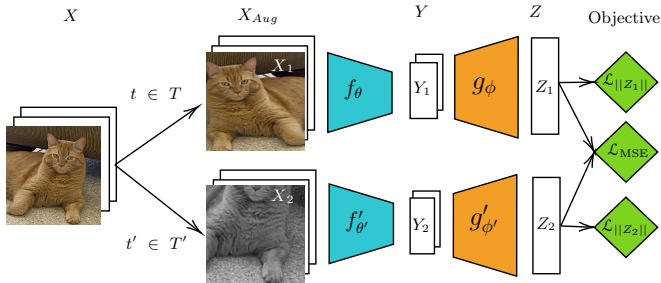
[1] University of Kentucky
[2] Washington University in St. Louis
oscar.skean@uky.edu

**Abstract.** Self-supervised learning (SSL) is a popular paradigm for representation learning. Recent multiview methods can be classified as sample-contrastive, dimension-contrastive, or asymmetric network-based, with each family having its own approach to avoiding informational collapse. While these families converge to solutions of similar quality, it can be empirically shown that some methods are epoch-inefficient and require longer training to reach a target performance. Two main approaches to improving efficiency are covariance eigenvalue regularization and using more views. However, these two approaches are difficult to combine due to the computational complexity of computing eigenvalues. We present the objective function FroSSL which reconciles both approaches while avoiding eigendecomposition entirely. FroSSL works by minimizing covariance Frobenius norms to avoid collapse and minimizing mean-squared error to achieve augmentation invariance. We show that FroSSL reaches competitive accuracies more quickly than any other SSL method and provide theoretical and empirical support that this faster convergence is due to how FroSSL affects the eigenvalues of the embedding covariance matrices. We also show that FroSSL learns competitive representations on linear probe evaluation when used to train a ResNet-18 on several datasets, including STL-10, Tiny ImageNet, and ImageNet-100. Github

## 1 Introduction

The problem of learning representations without human supervision is fundamental in machine learning. Unsupervised representation learning is particularly useful when label information is difficult to obtain or noisy. It requires the identification of structure in data with limited knowledge about what the structure is. One common way of learning structure without labels is joint embedding self-supervised learning (SSL) [1, 2, 3, 4, 5, 6, 7, 8]. The basic goal of SSL is to train neural networks to capture *semantic* input features that are *augmentation-invariant*. This goal is appealing for representation learning because the inference set often has similar semantic content to the training set.

**Fig. 1:** The SSL pipeline used in this work. In general, the encoder and projector may be asymmetric. We use symmetric encoders with shared weights and the same augmentation set for each view. We refer to $X_1$ as view 1 of $X$, and $X_2$ as view 2. Only two views are shown here, though more may be used in practice.

A trivial solution to learning augmentation-invariant features is to encode all images to the same point, often called trivial or informational collapse. The resulting networks are essentially useless for downstream tasks. Different mechanisms have been proposed to handle collapse in SSL. These can be grouped into three families: sample-contrastive, dimension-contrastive, and asymmetric network methods.

A less studied problem in all current SSL methods is their speed of convergence. When compared to traditional supervised learning, SSL methods must be trained for large numbers of iterations to reach competitive performance on downstream tasks. For example, a typical experiment in the literature is to train for 1000 epochs on ImageNet which can take several weeks even with many GPUs. An imperative direction of research is to investigate how to reduce SSL training time. An observation that is often hidden by only reporting the final epoch accuracy is that, empirically, certain SSL methods require more training time to reach competitive accuracies. This phenomenon has been observed for many dimension-contrastive methods by Simon *et al.* [9] but not discussed in detail. We provide additional support for this claim in Section 4.1. Our work attempts to answer the following research question: Does there exist an SSL method with dimension-contrastive advantages, namely simplicity via avoidance of both negative sampling and architectural restrictions, while achieving competitive accuracies more quickly than other existing SSL methods?

We propose an SSL objective which we call FroSSL. Similar to many dimension-contrastive methods, FroSSL consists of a variance and invariance term. The invariance term is simply a mean-squared error between the views and is identical to VICReg's invariance term [10]. The variance term is the logarithm of the squared Frobenius norm of the normalized covariance embedding matrices. Using the Frobenius norm of covariance matrices for improving learned representations has not been explored in SSL.

Our contribution can be summarized as:

- We introduce the FroSSL objective function and show that it is *both* dimension-contrastive and sample-contrastive up to a normalization of the embeddings.
- We introduce a theoretical framework that unifies dimension-contrastive methods that scale linearly in the number of views.
- We show that FroSSL combines two techniques to reduce training time: using more views and improving eigenvalue dynamics. We examine covariance eigenvalue trajectories during training on STL-10 to show that FroSSL learns useful, high-rank representations more quickly than other dimension-contrastive methods.
- We evaluate FroSSL on the standard setup of SSL pretraining and linear probe evaluation on CIFAR-10, CIFAR-100, STL-10, Tiny ImageNet, and ImageNet-100. We find that FroSSL achieves strong performance, especially when models are trained for fewer epochs.

## 2    Background and Notation

Consider a matrix $A \in \mathbb{R}^{m \times n}$. Let $A_{ij} \in \mathbb{R}$ denote the element at the $i$th row and $j$th column of $A$, and $A_{i,:} \in \mathbb{R}^m$ denote the $i$th column vector representing the $i$th row of $A$, and $A_{:,j}$ the $j$th column of $A$. Let $\sigma_k(A)$ denote the $k$th singular value of $A$ ordered non-increasingly. The entry-wise product (also known as Hadamard product) is denoted as $A \odot B$. The Ky Fan $p$-norm of $A$ is defined as [11]:

$$\|A\|_p = \left( \sum_{k}^{\min(m,n)} \sigma_k^p(A) \right)^{1/p}, \qquad (1)$$

which is a unitarily invariant norm. For $p = 2$, we have the Frobenius norm $\|A\|_2 = \|A\|_F = \sqrt{\sum_i \sum_j A_{ij}^2}$.

### 2.1    The Joint Embedding Self-Supervised Learning Problem

The goal of self-supervised learning is to learn useful representations without external supervision. Many visual joint embedding SSL methods follow a similar procedure which was first introduced in [1]. An example of this procedure is depicted in Figure 1. Let $\mathbf{X} = \{x_i\}_{i=1}^N$ be a minibatch with $N$ samples, $V$ the number of augmented views, $T(\cdot)$ a function that applies randomly selected augmentations to an image, $f$ a visual encoder network, and $g$ a projector network.

First, each image $x_i \in \mathbf{X}$ is paired with augmented versions of itself, making the augmented dataset $\mathbf{X}_{\text{aug}} = \{T_1(x_i), \cdots, T_V(x_i)\}_{i=1}^N = \{X_1, \cdots, X_V\}$ With ideal augmentations, $(X_1)_{i,:}$ and $(X_2)_{i,:}$ have identical *semantic* content and different *style* content. Note that typically $V = 2$, but we make no such assumptions. For each augmentation, the embedding set is given by $Y_v = \{f((X_v)_{i,:})\}_{i=1}^N$ and projection set $Z_v = \{g((Y_v)_{i,:})\}_{i=1}^N$. Finally, an SSL objective is computed on

the projections and backpropagated through both networks. The goal of the objective is to ensure that encoded augmentations for the same image are mapped close together by the projector, i.e. $(Z_a)_{i,:}$ and $(Z_b)_{i,:}$ are close in some sense of distance for all $a, b = 1, 2, \ldots, V$. At the same time, projections should capture the variability among images. Thus the goal of SSL is to train the networks $f$ and $g$ to extract *semantic* features that are invariant to any augmentations induced by $T(\cdot)$. In the following, we take a closer look at choices for the SSL objective.

## 2.2    The Three Families of Joint Embedding SSL Objectives

Objective functions for joint embedding self-supervised learning can be divided into three families. The first family consists of *sample-contrastive* methods [1, 2, 3, 6, 12] which use a contrastive loss to learn a representation that maps positive samples (augmentation of the same image) close together while pushing negative samples (different images) apart. These methods avoid collapse at the expense of making comparisons between positive and negative samples.

The second family consists of *asymmetric network* methods [4, 5, 13] which place restrictions on the architecture of the mapping network used, including asymmetrical encoders [4, 5], momentum encoders [6], and stop gradients [4, 14]. While these methods can achieve great results, they are rooted in implementation details and there is no clear theoretical understanding of how they avoid collapse [10].

The third, and most recent, family are the *dimension-contrastive* methods, which are sometimes called negative-free contrastive [15] or feature decorrelation methods [16]. These methods operate by reducing the redundancy in feature dimensions. Instead of examining *where* samples live in feature space, these methods examine *how* feature dimensions are being used. Methods in this family can avoid the use of negative samples while also not requiring restrictions in the network architecture to prevent collapse. Barlow Twins objective pushes the normalized cross-covariance between views towards the identity matrix [7]. VICReg consists of three terms: the invariance term enforces similarity in embeddings across views, while the variance/covariance terms regularize the covariance matrices of each view to prevent collapse [10]. W-MSE whitens and projects embeddings to the unit sphere before maximizing cosine similarity between positive samples [17]. I-VNE maximizes the von Neumann entropy of the embedding covariance matrices [18]. Finally, CorInfoMax maximizes the log det entropy of both views while minimizing mean-squared error [19].

## 2.3    A Framework for Dimension-Contrastive Methods

Many recent works in dimension-contrastive SSL, whether explicitly or implicitly, consist of a combination of two competing objectives: an augmentation **invariance** term that pulls different augmentations from the same image close together, and a **variance** term that avoids collapse of the mapping by regulating

**Table 1:** Taxonomy of dimension-contrastive SSL methods describing how they avoid informational collapse and achieve augmentation invariance in the $D_{\text{inv}}$ and $D_{\text{var}}$ framework of Section 2.3.

| Method | Variance $D_{\text{var}}(\Sigma_v \| \mathbf{I})$ | Invariance $D_{\text{inv}}(Z_v, Z_r)$ |
|---|---|---|
| VICReg | (Variance) Hinge loss on auto-covariance diagonal | |
| | (Covariance) covariance off-diagonals per view | |
| | $\sum\limits_{k}^{D} \max\left(0, 1 - \sqrt{(\Sigma_v)_{k,k} + \epsilon}\right) + \nu \left\|\Sigma_v - \Sigma_v \odot \mathbf{I}\right\|_F^2$ | MSE |
| W-MSE | Implicit through whitening that | $\frac{1}{N} \|Z_v - Z_r\|_2^2$ |
| | $D_{\text{var}}(\Sigma_v \| \mathbf{I}) = 0$ since $\Sigma_v = \mathbf{I}$ for all $v$ | |
| CorInfoMax | Log Det Divergence: $D_{\log \det}(A\|B) = \text{trace}(AB^{-1}) - D - \log \det(AB^{-1})$ | |
| | $D_{\log \det}(\Sigma_v + \epsilon \mathbf{I} \| \mathbf{I}) = \text{trace}(\Sigma_v + \epsilon \mathbf{I}) - D - \log \det(\Sigma_v + \epsilon \mathbf{I})$ | |
| I-VNE | von Neumann Relative Entropy: $S_1(A\|B) = \text{trace}(A(\ln A - \ln B))$ | Cosine Similarity |
| | $S_1(\Sigma_v \| \mathbf{I}) = \text{trace}(\Sigma_v \ln \Sigma_v)$ | |
| FroSSL (ours) | 2-Order Petz-Rényi Relative Entropy: $S_2(A\|B) = \log \text{trace}(A^2 B^{-1})$ | MSE |
| | $S_2(\Sigma_v \| \mathbf{I}) = \ln\left(\sum_{i=1}^{N} \sigma_i^2\right) = \ln \|\Sigma_v\|_F^2$ | $\frac{1}{N} \|Z_v - Z_r\|_2^2$ |

variance. Below, we unify dimension-contrastive methods into a general framework that is parameterized by choices of two distances. By carefully selecting these distances, specific dimension-contrastive methods can be recovered.

Let $Z_v \in \mathbb{R}^{N \times D}$ be a batch of projections and $\Sigma_v = \frac{1}{N} \hat{Z}_v^T \hat{Z}_v$ the corresponding covariance, where $\hat{Z}_v$ are the centered projections. A dimension-contrastive objective can be written as follows:

$$\min \sum_{v=1}^{V-1} \sum_{r=v+1}^{V} D_{\text{inv}}(Z_v, Z_r) + \gamma \sum_{v=1}^{V} D_{\text{var}}(\Sigma_v \| \mathbf{I}). \tag{2}$$

The first term of (2) is the **invariance** term which minimizes the distance $D_{\text{inv}} : \mathbb{R}^{N \times D} \times \mathbb{R}^{N \times D} \mapsto \mathbb{R}_{\geq 0}$ between all pairs of augmentations. The second term of (2) is a **variance** factor that forces the covariance of each augmentation to be close to identity according to a dissimilarity $D_{\text{var}} : \mathbb{R}^{D \times D} \times \mathbb{R}^{D \times D} \mapsto \mathbb{R}_{\geq 0}$. For instance, in VICReg [10], $D_{\text{inv}}(Z_v, Z_r) = \|Z_v - Z_r\|_F^2$ and $D_{\text{var}}(\Sigma_v \| \mathbf{I}) = \sum_{k}^{D} \max\left(0, 1 - \sqrt{(\Sigma_v)_{k,k} + \epsilon}\right) + \nu \|\Sigma_v - \Sigma_v \odot \mathbf{I}\|_F^2$. Similarly, in CorInfoMax [19] $D_{\text{inv}}(Z_v, Z_r)$ is the same as VICReg, but $D_{\text{var}}(\Sigma_v \| \mathbf{I})$ can be related to the log det divergence $D_{\log \det}(A\|B) = \text{trace}(AB^{-1}) - D - \log \det(AB^{-1})$ setting $A = \Sigma_v + \epsilon \mathbf{I}$ and $B$ to a scaling of identity due to the normalization step in their projector. In Table 1, we show the dimension-contrastive methods which fit into this framework. We provide derivations in Appendix B.

**Multiview Invariance Term** In (2) the invariance term requires $V(V-1)/2$ comparisons which scales quadratically with the number of views. However, if $D_{\text{inv}}(Z_v, Z_r) = \|Z_v - Z_r\|_F^2$, then the invariance term may be simplified to

$$\sum_{v=1}^{V-1} \sum_{r=v+1}^{V} D_{\text{inv}}(Z_v, Z_r) = V \sum_{v=1}^{V} D_{\text{inv}}\left(Z_v, \overline{Z}\right) \tag{3}$$

**Table 2:** Taxonomy of dimension-contrastive SSL methods showing which desirable criteria they fulfill.

|  | Invariant to Projection Rotations | Manipulates Eigenvalues Explicitly | Quadratic in Batch Size and Dimension | Linear in Views |
|---|---|---|---|---|
| Barlow Twins | ✗ | ✗ | ✓ | ✗ |
| VICReg | ✗ | ✗ | ✓ | ✓ |
| W-MSE | ✓ | ✗ | ✗ | ✓ |
| CorInfoMax | ✓ | ✓ | ✗ | ✓ |
| I-VNE | ✓ | ✓ | ✗ | ✓ |
| MMCR | ✓ | ✓ | ✗ | ✓ |
| FroSSL (ours) | ✓ | ✓ | ✓ | ✓ |

where $\overline{Z} = \frac{1}{V} \sum_{i=1}^{V} Z_i$ is the average projection across all views. If a method has a $D_{\text{inv}}$ that can be rewritten this way, we say the method scales linearly with views. All methods displayed in Table 1 have this property.

## 3   The FroSSL Objective

To motivate FroSSL, we begin by posing four desirable criteria of dimension-contrastive methods.

1. **Invariant to Projection Rotations** We argue that dimension-contrastive methods should be invariant to rotations in the projections because the orientation of the covariance does not affect the relationships between principal components. In other words, redundancy in the embedding dimensions is invariant to the rotation of the embeddings. Thus the choices of $D_{\text{var}}$ and $D_{\text{inv}}$ should be rotationally invariant as well.
2. **Manipulates Eigenvalues Explicitly** Several works have shown that regularizing projection covariance eigenvalues in SSL can lead to reduced training time and improved downstream performance [14, 18, 20]. We provide empirical support for this in Section 4.1.
3. **Scales Quadratically in Batch Size and Dimension** The time complexity of the objective function scale *at most* quadratically with respect to $N$ and $D$. This is often in opposition to the prior criteria which typically requires cubic eigendecomposition.
4. **Scales Linearly in Views** The time complexity of the objective function should be linear in the number of views $V$. This is advantageous because recent work has shown that using more views can reduce training time and improve downstream performance [12, 20]. We provide empirical support for this in Sections 4.2 and 5. Any dimension-contrastive method with $D_{\text{inv}}$ that satisfies Equation (3) fulfills this criterion.

As shown in Table 2, no prior method meets all four criteria. We provide explanations in Appendix C. To construct a method that fulfills all criteria, we

modify the I-VNE objective function:

$$\max \mathcal{L}_{\text{I-VNE}} = \sum_{v=1}^{V} \text{Tr} \, \Sigma_v \ln \Sigma_v + \sum_{v=1}^{V-1} \sum_{r=v+1}^{V} \frac{Z_v^T Z_r}{\|Z_v\|_2 \|Z_r\|_2} \tag{4}$$

The invariance term maximizes the cosine similarity between views. The variance term maximizes the von Neumann entropy of each view covariance matrix. The only criteria that I-VNE does not fulfill is being subcubic in batch size and dimension. This is due to the eigendecomposition needed to compute the matrix logarithm for the entropy. To begin addressing this, we first notice that the von Neumann entropy is a limit case of matrix-based $\alpha$-order entropy [21, 22, 23]:

$$S_\alpha \left( \Sigma_v \right) = \frac{1}{1-\alpha} \log \left[ \text{Tr} \left( \Sigma_v^\alpha \right) \right] = \frac{1}{1-\alpha} \log \left( \sum_i^{\min(N,D)} \lambda_i^\alpha(\Sigma_v) \right) \tag{5}$$

Here, we do not require trace$(\Sigma_v) = 1$ as is typically required by $\alpha$-order entropy. The von Neumann entropy is equivalent to $S_1(\Sigma_v)$ in the limit. Another special case is collision entropy, given by $S_2(\Sigma_v)$ below:

$$S_2 \left( \Sigma_v \right) = -\log \left( \sum_{i=1}^{\min(N,D)} \lambda_i^2(\Sigma_v) \right) = -\log(\|\Sigma_v\|_F^2) = -\log \sum_i \sum_j (\Sigma_v)_{ij}^2 \tag{6}$$

Notice how the left-hand side in the above equation explicitly uses the eigenvalues, while the right-hand side only uses matrix elements. This is made possible by the Frobenius norm, which offers an equivalency between a sum over eigenvalues and a sum over matrix elements. This has immediate impacts on the loss time complexity by relaxing the $O(\min(D, N)^3)$ eigendecomposition to the $O(\min(D, N)^2)$ Frobenius norm computation. The case of 2-order entropy is the only matrix-based $\alpha$-entropy which does not require eigendecomposition. One potential downside is $S_2(\Sigma_v)$ does not penalize outlying eigenvalues as heavily as in $S_1(\Sigma_v)$. This is akin to the difference between mean-absolute error and mean-squared error. However, this has no significant impact in our experiments.

The variance term $D_{\text{var}}$ for our objective will minimize the log Frobenius norm of normalized embeddings, causing the embeddings to spread out equally in all directions. Normalizing the embeddings is crucial because otherwise, minimizing the Frobenius norm will lead to trivial collapse. For the invariance term $D_{\text{inv}}$, we opt to use the mean-squared error between views. Our objective function FroSSL is given below:

$$\text{minimize } \mathcal{L}_{\text{FroSSL}} = \sum_{v=1}^{V} \log(\|\Sigma_v\|_F^2) + \gamma \left\| Z_v - \overline{Z} \right\|_F^2 \tag{7}$$

Note we simplify the pairwise mean-squared error via Equation (3). Because the Frobenius norm is invariant to transposition, we can choose to compute either $\left\| Z_v^T Z_v \right\|_F^2$ or $\left\| Z_v Z_v^T \right\|_F$ depending on if $D > N$. The former has time complexity $O(ND^2)$ while the latter has complexity $O(N^2D)$. For consistency, we always use the former in our experiments. We provide pseudocode in Appendix A.

## 3.1   The Role of the Logarithm

The log in Equation (7) ensures that the contributions of the variance term to the gradient of the objective function become self-regulated ($\frac{d \log f(x)}{dx} = \frac{1}{f(x)} \frac{df(x)}{dx}$) with respect to the invariance term. We later compare the experimental performance of Equation (7) with and without the logarithms, showing that using logarithms leads to a gain in performance. Prior work has shown that Equation (7) with no logarithms causes dead neurons in the final encoder layer [18].

## 3.2   FroSSL is both Sample-contrastive and Dimension-contrastive

It can be shown, up to an embedding normalization, that FroSSL is both dimension-contrastive and sample-contrastive. First, we provide formal definitions of dimension-contrastive and sample-contrastive SSL, following Garrido *et al.* [24].

**Definition 1 (Dimension-Contrastive Method).** *An SSL method is said to be dimension-contrastive if it minimizes the non-contrastive criterion* $\mathcal{L}_{nc}(Z) = \left\| Z^T Z - diag(Z^T Z) \right\|_F^2$, *where* $Z \in \mathbb{R}^{N \times D}$ *is a matrix of embeddings as defined above. This may be interpreted as penalizing the off-diagonal terms of the embedding covariance.*

**Definition 2 (Sample-Contrastive Method).** *An SSL method is said to be sample-contrastive if it minimizes the contrastive criterion* $\mathcal{L}_c(Z) = \left\| Z Z^T - diag(Z Z^T) \right\|_F^2$. *This may be interpreted as penalizing the similarity between pairs of different images.*

Next, we use the duality of the Frobenius norm, given by $\left\| Z^T Z \right\|_F = \left\| Z Z^T \right\|_F$, to show that FroSSL satisfies the qualifying criteria of both dimension-contrastive and sample-contrastive methods.

**Proposition 1.** *If every embedding dimension is normalized to have equal variance, then FroSSL is a dimension-contrastive method. Proof in Appendix E.1.*

**Proposition 2.** *If every embedding is normalized to have equal norm, then FroSSL is a sample-contrastive method. Proof in Appendix E.2.*

**Proposition 3.** *If the embedding matrices are doubly stochastic, then FroSSL is simultaneously dimension-contrastive and sample-contrastive.*

Proposition 3 allows for interpreting FroSSL as either a sample-contrastive or dimension-contrastive method, up to a normalization of the data embeddings. The choice of normalization strategy is not important to the performance of an SSL method [24]. Unless otherwise specified, we only normalize the variance and not the embeddings. These same proof techniques can be used to show that TiCo, MMCR, I-VNE, and CorInfoMax also belong to both families [18, 19, 20, 25]. Additionally, variants of the dimension-contrastive VICReg have been
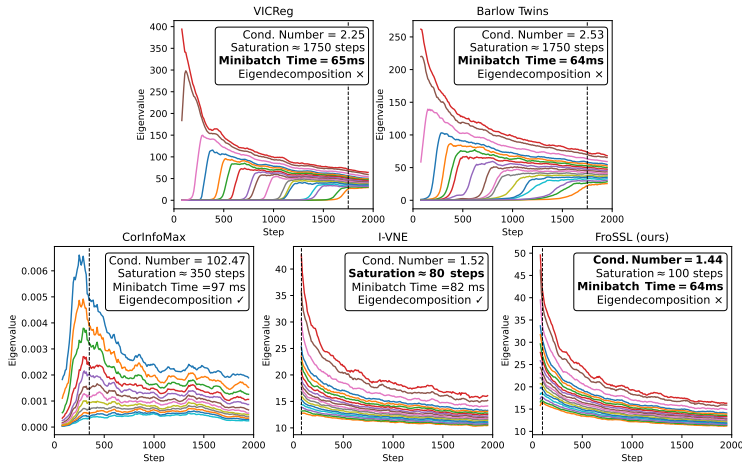
proposed [24] that allow it to be rewritten as the sample-contrastive SimCLR. However, VICReg cannot be rewritten in such a way due to the hinge loss.

While Proposition 3 is interesting theoretically, it also offers empirical benefits to FroSSL. We examine overall wall-training time to reach competitive accuracies (4.3), robustness to augmentations (5.1), and performance on little pretraining data (5.2). In all of these experiments, sample-contrastive methods outperform dimension-contrastive methods. FroSSL shares the advantages observed empirically in sample-contrastive methods.

# 4   On Efficiency in Self-Supervised Learning

It is well-known that traditional SSL algorithms need hundreds or thousands of epochs to reach competitive accuracies. To compare the efficiency of different SSL algorithms, we can borrow theoretical and practical tools from the broader field of algorithmic complexity. In the context of machine learning, there are two measurements of particular interest to practitioners: *wall-time* needed to reach a given accuracy and VRAM *space* used. The former can be decomposed into the atomic quantities of average *wall-time per minibatch* and the *number of epochs* needed to reach a given accuracy. To emphasize why this decomposition matters, consider a scenario where wall-time per minibatch differs between two methods but overall wall-time does not. In such a scenario, using the method with the slower minibatch wall-time is advantageous for using fewer disk reads and less distributed network traffic. This is not obvious without observing the atomic quantities. Note we are careful to specify "epochs to reach a given accuracy" rather than "epochs to convergence". One reason for this is that classical experiments in SSL train for a fixed number of epochs rather than until convergence. Another reason is algorithms that more quickly reach a target performance, such as FroSSL or I-VNE, do not necessarily converge in fewer epochs.

The design of an SSL algorithm is a balancing act between minibatch time, space, and number of epochs. While conversations involving minibatch time and space have been rarely discussed in the SSL literature, discussion about the number of epochs has seen renewed interest [9, 14, 26]. However, if SSL algorithm design is indeed a balancing act of the three quantities above, then space and minibatch time deserve discussion too. Methods that boast reductions to one quantity may come with significant penalties to a different quantity. For example, dimension-contrastive methods use less space in practice than sample-contrastive methods, which prefer large minibatch sizes, or asymmetric methods like BYOL, which need an additional prediction network. However, the improved space efficiency comes at the cost of requiring a higher number of epochs [9]. In Sections 4.1 and 4.2, we discuss the advantages and drawbacks of two approaches to reducing the number of epochs. In Section 4.3, we compare a variety of SSL algorithms and visualize their time, space, and epoch tradeoffs.

**Fig. 2: The choice of variance term, $D_{\mathbf{var}}(\Sigma_v \| \mathbf{I})$, has a significant impact on training dynamics.** Each subplot visualizes the trajectories of the top 20 eigenvalues of the embedding covariance matrix $\Sigma_1$ when trained with dimension-contrastive methods. These trajectories show how quickly $\Sigma_v$ converges to $\gamma I$, which has eigenvalues equal to $\frac{\gamma}{D}$. VICReg, Barlow Twins, and CorInfoMax converge slowly. FroSSL and I-VNE have similar training dynamics, but FroSSL has significantly lower computational complexity because it avoids explicitly computing the eigendecomposition.

### 4.1   Reducing the Number of Epochs with Eigenvalue Dynamics

Recent work has examined the training dynamics of SSL models [9]. In particular, they find that the eigenvalues of the covariance exhibit *stepwise* behavior, meaning that one eigendirection is learned at a time. This is readily seen in Figure 2 for VICReg and Barlow Twins. This phenomenon contributes to slowness in SSL optimization with the smallest eigendirections taking the longest to be learned. Other work shows that high-rank representations lead to better downstream performances [27]. It directly follows that if an SSL method requires a high number of epochs to learn high-rank representations, then it also needs a high number of epochs to learn useful representations.

We hypothesize that by carefully choosing the variance term $D_{\text{var}}(\Sigma_v \| \mathbf{I})$ to reduce stepwise eigenvalue dynamics, useful representations can be learned more quickly. Indeed, this behavior has already been observed in several SSL objectives already. CorInfoMax optimizes the log-determinant of each covariance $\Sigma_v$, which is defined as the log of the product of the $\Sigma_v$ eigenvalues [19]. IsoLoss uses $\Sigma_v$ eigenvalues as learning rate multipliers to equalize the convergence rate of different eigenmodes [14]. MMCR optimizes the nuclear norm of the average view embedding, which is defined as the sum of the singular value magnitudes [20]. I-VNE optimizes the von Neumann entropy of $\Sigma_v$, which is equal to the Shannon entropy of the $\Sigma_v$ eigenvalues [18].

It is straightforward to show that FroSSL also directly influences the covariance eigenvalue dynamics. However, FroSSL is unique from prior methods because it does so while avoiding explicit eigendecomposition. This can be seen from Equation 6. Additionally, using the Frobenius norm eliminates numerical instabilities typically associated with eigendecomposition [28].

To highlight the existence and remedy of stepwise phenomena in practical scenarios, we create an experiment similar to the one used by Simon *et al.* [9]. In Figure 2, we plot the trajectories of the top 20 eigenvalues of $\Sigma_1$ when trained with different dimension-contrastive objectives. For all SSL objectives, a ResNet-18 was trained for 5 epochs on STL-10 using SGD with $lr = 0.01$ and a batch size of 256. Further details are given in Appendix D.1.

The eigenvalues trajectories show how quickly $\Sigma_v$ is approaching $\gamma I$, which has eigenvalues equal to $\frac{\gamma}{D}$. We say that an objective is saturated once the stepwise learning phase is ended. This is marked as the step where $\lambda_{20}$ has increased from zero and started decreasing. It is clear to see that SSL objectives that directly influence eigenvalues, namely CorInfoMax, I-VNE, and FroSSL, saturate much quicker than the others. Interestingly, the condition number for CorInfo-Max, computed as $\frac{\lambda_1}{\lambda_{20}}$, is much larger than any other tested method. We hypothesize this is due to the choice of the $\epsilon$ hyperparameter for the regularization term when computing the determinant as $\det(\Sigma_1 + \epsilon I)$.

## 4.2    Reducing the Number of Epochs by Using More Views

**Multiview with 3 or More Views** In contrastive learning, using more views has been shown to have significant impacts on representation quality and downstream performance [29]. In SSL, using more augmentations for each image has the effect of averaging out noise from the mean embedding, which acts as a target for many SSL objectives as shown in Equation (3). This differs from increasing the batch size, which would instead average out noise across samples and not across targets. While using more views is promising, it has not seen widespread adoption in self-supervised learning. This is in part due to many sample-contrastive methods being quadratic in the number of views. However, this problem is circumvented for the dimension-contrastive methods shown in Table 1, which are instead linear in the number of views. One such method, W-MSE, has shown performance improvements when the number of views is increased from 2 to 4 [17]. Interestingly, MMCR is constant in the number of views because it operates only on the mean embedding [20].

**Multi-Patch and Multi-Crop Methods** An approach in a different vein is to extract and augment image patches to serve as views, rather than using full images. EMP-SSL has shown that this drastically reduces the number of epochs and overall wall-time needed to reach competitive accuracies by utilizing a bag-of-features model that embeds hundreds of small augmented patches per image [26, 30]. However, EMP-SSL comes at the cost of major penalties to time-per-minibatch and space in both training time and inference time.

**Table 3:** Comparison of the time/space/epoch tradeoffs for SSL algorithms trained on STL-10. FroSSL with 8 views achieves 80% top-1 accuracy in the least wall-time.

| | SimCLR | MoCo v2 | BYOL | VICReg | Barlow | CorInfoMax | MMCR | | | FroSSL (ours) | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Num. Views | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 4 | 8 | 2 | 4 | 8 |
| Loss Time Complexity | | $O(V^2N^2)$ | | $O(VD^2)$ | $O(V^2D^2)$ | $O(VD^3)$ | $O(\min(D,N)^3)$ | | | $O(V\min(D,N)^2)$ | | |
| VRAM Space (GB) | **1.6** | 2.8 | 2.0 | **1.6** | **1.6** | 1.7 | 1.7 | 2.9 | 5.3 | 1.7 | 2.9 | 5.3 |
| Minibatch Wall-time (ms) | **60** | 79 | 76 | 65 | 64 | 97 | 71 | 108 | 187 | 64 | 105 | 187 |
| Number of Epochs to 80% Acc | 347 | 180 | 187 | 360 | 370 | 405 | 380 | 211 | 63 | 290 | 144 | **55** |
| Wall-time to 80% Acc (hours) | 2.4 | 1.6 | 1.6 | 2.7 | 2.7 | 4.5 | 3.1 | 2.6 | 1.3 | 2.1 | 1.7 | **1.2** |

**Table 4:** Top-1 accuracies on STL-10 using an online linear classifier during training for specific numbers of epochs (left) and specific elapsed times (right).

| Method | Epochs | | | | | | Method | Training Wall-Time (min.) | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | 3 | 10 | 30 | 50 | 100 | | | 10 | 30 | 60 |
| SimCLR | 40.7 | 44.8 | 61.5 | 66.2 | 70.1 | | SimCLR | 61.5 | 68.8 | 73.9 |
| MoCo v2 | 24.6 | 45.0 | 63.8 | 69.4 | 75.2 | | MoCo v2 | 57.4 | 70.0 | 76.2 |
| BYOL | 28.8 | 32.7 | 59.6 | 64.7 | 70.6 | | BYOL | 50.2 | 65.3 | 75.0 |
| VICReg | 43.6 | 51.1 | 61.2 | 67.5 | 71.1 | | VICReg | 63.4 | 70.3 | 72.9 |
| Barlow Twins | 32.1 | 46.6 | 62.0 | 62.6 | 69.0 | | Barlow Twins | 55.5 | 66.1 | 68.0 |
| CorInfoMax | 39.0 | 49.1 | 58.0 | 62.5 | 66.2 | | CorInfoMax | 56.1 | 64.9 | 65.6 |
| MMCR (2 views) | 39.6 | 53.3 | 62.8 | 63.3 | 67.0 | | MMCR (2 views) | 54.7 | 68.4 | 70.9 |
| MMCR (4 views) | 46.0 | 61.5 | 70.2 | 71.5 | 75.7 | | MMCR (4 views) | **69.8** | **74.7** | 76.6 |
| MMCR (8 views) | **51.1** | 64.7 | 72.9 | 77.2 | 79.4 | | MMCR (8 views) | 64.7 | 73.5 | 78.0 |
| FroSSL (2 Views) | 44.8 | 56.9 | 64.8 | 67.1 | 72.0 | | FroSSL (2 Views) | 63.4 | 69.7 | 74.9 |
| FroSSL (4 Views) | 49.3 | 60.7 | 70.3 | 67.1 | 76.9 | | FroSSL (4 Views) | 68.5 | 73.7 | 76.3 |
| FroSSL (8 Views) | 47.6 | **65.5** | **74.5** | **78.4** | **81.8** | | FroSSL (8 Views) | 58.6 | 74.1 | **79.0** |

As an alternative to using full-sized images or tiny patches for each view, multi-crop methods strike a balance [12]. A certain number of views are full-sized images while the remaining views are smaller crops. A typical setup for ImageNet is using two $224 \times 224$ views and six $96 \times 96$ views. These approaches differ from our experiments which use all full-sized views with FroSSL. However, we expect that FroSSL should work well as an objective function for these paradigms too.

## 4.3   Exploring Time, Space, and Epoch Tradeoffs

We now compare the efficiency of different SSL algorithms. We train a ResNet-18 for 500 epochs on STL-10 and measure the number of epochs needed to reach a top-1 accuracy of 80%. This threshold of 80% was chosen because all methods achieve that accuracy within 500 epochs. We used $N = 256$ and $D = 1024$ for all methods. Because these models were trained on a distributed cluster, it is important to compensate for different compute when measuring minibatch wall-time. In particular, we measure minibatch time by averaging over 2000 iterations of training on one NVIDIA A5000 GPU. We measure VRAM space as the maximum space requested by the training script. We calculate wall-time to 80% accuracy by multiplying minibatch time, epochs, and iterations per epoch.

In Table 3 we show the resources needed for each SSL objective. There are several observations to glean from this table. First, increasing the number of views reduces epochs and overall wall-time, even though space and minibatch time become larger. FroSSL with 8 views reaches 80% top-1 accuracy faster than any other tested method. Second, asymmetric methods require the least overall

**Table 5:** Comparison of SSL methods on small datasets. All CorInfoMax and MMCR results are from our implementation. All Tiny ImageNet and STL-10 results are from our implementation. CIFAR-10 and CIFAR-100 results are reported from [17, 31]. IN-100 baseline results are from [31]. We observed negligible improvements from using more views for FroSSL on the CIFAR datasets; we used the 2-view CIFAR accuracies to compute 4/8 view averages. An asterisk (*) denotes Tiny ImageNet results where weak augmentations outperformed strong ones. Results within 0.5% of best are **bolded**.

| Method | CIFAR-10 | CIFAR-100 | STL-10 | Tiny-IN | IN-100 | Average |
|---|---|---|---|---|---|---|
| **Sample-Contrastive** | | | | | | |
| SimCLR | 91.8 | 65.8 | 85.9 | 41.9 | 77.6 | 72.6 |
| SwAV | 89.2 | 64.9 | 82.6 | 41.2 | 74.3 | 70.5 |
| MoCo v2 | **92.9** | 69.9 | 83.2 | 41.9 | 79.3 | 73.4 |
| **Asymmetric Network** | | | | | | |
| SimSiam | 90.5 | 66.0 | 88.5 | **45.6*** | 78.7 | 73.9 |
| BYOL | **92.6** | **70.5** | 88.7 | 40.1 | **80.3** | 74.4 |
| DINO | 89.5 | 66.8 | 78.9 | 34.9 | 78.9 | 69.8 |
| **Dimension-Contrastive** | | | | | | |
| VICReg | 92.1 | 68.5 | 85.9 | 37.5 | 79.4 | 65.8 |
| Barlow Twins | 92.1 | **70.9** | 85.0 | 45.3 | **80.2** | 74.7 |
| W-MSE 2 | 91.6 | 66.1 | 72.4 | 28.8* | 69.0 | 65.6 |
| CorInfoMax | **92.6** | 69.7 | 83.1 | 43.9 | 74.7 | 72.8 |
| I-VNE | 89.7 | 65.7 | 87.4 | **45.2** | 77.6 | 73.1 |
| MMCR (2 views) | 88.6 | 65.8 | 84.3 | 41.2 | 76.7 | 71.3 |
| MMCR (4 views) | 89.6 | 67.3 | 88.2 | 42.8 | 78.8 | 73.3 |
| MMCR (8 views) | 89.3 | 68.3 | 90.3 | 43.2 | **80.3** | 74.2 |
| FroSSL (2 views) [no log] | 88.9 | 62.3 | 82.4 | 36.4 | 78.3 | 69.7 |
| FroSSL (2 views) | **92.8** | **70.6** | 87.3 | 44.2 | 78.2 | 74.6 |
| FroSSL (4 views) | - | - | 90.0 | **45.3** | 79.4 | 75.6 |
| FroSSL (8 views) | - | - | **90.9** | **45.3** | 79.8 | **75.9** |

wall-time for any method using 2 views, at the cost of space. We hypothesize this is due to enhanced training stability from momentum encoders. Third, doubling the number of views does not necessarily double the minibatch wall-time. This is because some parts of the training script, such as data loading and logging, do not get slower as the number of views increases. In Table 4, we show top-1 accuracies over epochs and over time. In both scenarios, FroSSL with 8 views has the highest accuracy after training is finished.

The impact of loss time complexity on minibatch time is more apparent at extreme batch sizes and dimensionalities. We show the VRAM space and minibatch time for $N = 1024$ and $D = 1024$ in Table 8. In this setting, FroSSL has 33% and 14% faster minibatch times than MMCR when both have 2 or 8 views, respectively. Additionally, when compared to MoCo v2 and BYOL, FroSSL with 2 views has a 20% faster minibatch time and uses 20% less space.

## 5 Experimental Results

In this section, we train a ResNet-18 on CIFAR-10 [32], CIFAR-100, STL-10 [33], Tiny ImageNet [34], and ImageNet-100 [29]. Our implementation is based on the solo-learn SSL framework [31].

In Table 5, we show linear probe evaluation results on these datasets. It is readily seen that FroSSL learns competitive representations in comparison to other SSL methods. The implementation details can be summarized as:

**Table 6:** The top-1% accuracies after training on Tiny-ImageNet using weak or strong augmentations.

| Method | Weak | Strong | Δ% |
|---|---|---|---|
| SimCLR | 39.5 | 41.9 | 2.4 |
| SwAV | 39.9 | 41.2 | 1.5 |
| MoCo v2 | 40.9 | 41.9 | 1.0 |
| SimSiam | 45.6 | 39.7 | -5.9 |
| BYOL | 39.4 | 40.1 | 0.7 |
| DINO | 32.2 | 34.9 | 2.7 |
| VICReg | 18.1 | 37.5 | 19.6 |
| Barlow Twins | 36.8 | 45.3 | 8.5 |
| CorInfoMax | 33.1 | 43.9 | 10.8 |
| MMCR (2 views) | 24.2 | 41.2 | 17.0 |
| FroSSL (2 views) | 39.4 | 44.2 | 4.8 |

**Table 7:** Accuracies after pretraining on ImageNet-1k for 100 epochs with only 10% of data.

| | Top-1 | Top-5 |
|---|---|---|
| SimCLR | 31.1 | 56.6 |
| BYOL | 12.7 | 29.1 |
| SimSiam | 22.7 | 46.3 |
| Barlow Twins | 23.6 | 46.9 |
| FroSSL (2 Views) | 33.4 | 59.1 |
| FroSSL (8 Views) | 38.2 | 64.1 |

- **Optimizer** The backbone uses LARS optimizer [35] with an initial learning rate of 0.3, weight decay of 1e-6, and a warmup cosine learning rate scheduler. The linear probe uses the SGD optimizer [36] with an initial learning rate of 0.3, no weight decay, and a step learning rate scheduler with decreases at 60 and 80 epochs.
- **Epochs** For CIFAR-10 and CIFAR-100, we pretrain the backbone for 1000 epochs. For STL-10, we pretrain for 500 epochs. For Tiny ImageNet, we pretrain for 800 epochs. For ImageNet-100, we pretrain for 800 epochs. All linear probes were trained for 100 epochs.
- **Hyperparameters** A batch size of N=256 is used for all datasets, except for Tiny ImageNet which used N=512. For FroSSL, we used $\gamma = 1.4$ for 2 views and $\gamma = 2$ for 4 and 8 views. We used an MLP with output dimension $D = 1024$ for FroSSL. Details about augmentations and method-specific hyperparameters are given in Appendix D.2.

## 5.1  Robustness to Augmentations

We trained ResNet-18 models on Tiny ImageNet using both weak and strong augmentations. Weak augmentations had Gaussian blur probabilities (0.5, 0.5) and solarization probabilities (0, 0) for each view. Strong augmentations had Gaussian blur probabilities (1.0, 0.1) and solarization probabilities (0.2, 0). As shown in Table 6, FroSSL is more robust to changes in augmentations than any other dimension-contrastive method.

## 5.2  Performance In Low Data Regime

We trained ResNet-18 models on ImageNet-1K [37] for 100 epochs using only 10% of the data and evaluated them using the standard linear probe. Note that limited data was used in both pretraining and evaluation. As shown in Table 7, FroSSL achieves a better downstream performance on limited data than any other tested method.

# 6    Conclusion

We introduced FroSSL, a self-supervised learning method that can be seen as both sample- and dimension-contrastive. We showed that FroSSL enjoys the simplicity of dimension-contrastive methods while achieving the empirical advantages of sample-contrastive methods. In particular, we discovered that FroSSL can achieve substantially stronger performance than alternative SSL methods when trained with less overall wall-time. To better understand why this is happening, we presented empirical results based on eigenvalue trajectories. We also demonstrated the effectiveness of FroSSL through extensive experiments on standard datasets.

# Acknowledgements

# Bibliography

[1] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International Conference on Machine Learning*, pages 1597–1607. PMLR, 2020.

[2] Jeff Z HaoChen, Colin Wei, Adrien Gaidon, and Tengyu Ma. Provable guarantees for self-supervised deep learning with spectral contrastive loss. *Advances in Neural Information Processing Systems*, 34, 2021.

[3] Yao-Hung Hubert Tsai, Yue Wu, Ruslan Salakhutdinov, and Louis-Philippe Morency. Self-supervised learning from a multi-view perspective. In *International Conference on Learning Representations*, 2021.

[4] Xinlei Chen and Kaiming He. Exploring simple Siamese representation learning. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15750–15758, 2021.

[5] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent-a new approach to self-supervised learning. *Advances in Neural Information Processing Systems*, 33, 2020.

[6] 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*, pages 9729–9738, 2020.

[7] Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stéphane Deny. Barlow twins: Self-supervised learning via redundancy reduction. In *International Conference on Machine Learning*, pages 12310–12320, 2021.

[8] Yazhe Li, Roman Pogodin, Danica J Sutherland, and Arthur Gretton. Self-supervised learning with kernel dependence maximization. *Advances in Neural Information Processing Systems*, 34, 2021.

[9] James B. Simon, Maksis Knutins, Ziyin Liu, Daniel Geisz, Abraham J. Fetterman, and Joshua Albrecht. On the stepwise nature of self-supervised learning. In *International Conference on Machine Learning*, 2023.

[10] Adrien Bardes, Jean Ponce, and Yann LeCun. VICReg: Variance-invariance-covariance regularization for self-supervised learning. In *International Conference on Learning Representations*, 2022.

[11] Roger A. Horn and Charles R. Johnson. *Matrix Analysis: Second Edition.* Cambridge university press, 2013.

[12] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. *Advances in Neural Information Processing Systems*, 33, 2020.

[13] 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*, pages 9650–9660, 2021.

[14] Manu Srinath Halvagal, Axel Laborieux, and Friedemann Zenke. Implicit variance regularization in non-contrastive ssl. *Advances in Neural Information Processing Systems*, 36, 2023.

[15] Yao-Hung Hubert Tsai, Shaojie Bai, Louis-Philippe Morency, and Ruslan Salakhutdinov. A note on connecting barlow twins with negative-sample-free contrastive learning. *arXiv preprint arXiv:2104.13712*, 2021.

[16] Chenxin Tao, Honghui Wang, Xizhou Zhu, Jiahua Dong, Shiji Song, Gao Huang, and Jifeng Dai. Exploring the equivalence of siamese self-supervised learning via a unified gradient framework. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14431–14440, 2022.

[17] Aleksandr Ermolov, Aliaksandr Siarohin, Enver Sangineto, and Nicu Sebe. Whitening for self-supervised representation learning. In *International Conference on Machine Learning*, pages 3015–3024, 2021.

[18] Jaeill Kim, Suhyun Kang, Duhun Hwang, Jungwook Shin, and Wonjong Rhee. VNE: An effective method for improving deep representation by manipulating eigenvalue distribution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3799–3810, 2023.

[19] Serdar Ozsoy, Shadi Hamdan, Sercan Arik, Deniz Yuret, and Alper Erdogan. Self-supervised learning with an information maximization criterion. *Advances in Neural Information Processing Systems*, 35, 2022.

[20] Thomas Yerxa, Yilun Kuang, Eero Simoncelli, and SueYeon Chung. Learning efficient coding of natural images with maximum manifold capacity representations. *Advances in Neural Information Processing Systems*, 36, 2024.

[21] Luis Gonzalo Sanchez Giraldo, Murali Rao, and Jose C Principe. Measures of entropy from data using infinitely divisible kernels. *IEEE Transactions on Information Theory*, 61(1):535–548, 2015.

[22] Oscar Skean, Jhoan Keider Hoyos Osorio, Austin J Brockmeier, and Luis Gonzalo Sanchez Giraldo. DiME: Maximizing mutual information by a difference of matrix-based entropies. *arXiv preprint arXiv:2301.08164*, 2023.

[23] Jhoan K Hoyos-Osorio and Luis G Sanchez-Giraldo. The representation Jensen-Shannon divergence. *arXiv preprint arXiv:2305.16446*, 2023.

[24] Quentin Garrido, Yubei Chen, Adrien Bardes, Laurent Najman, and Yann LeCun. On the duality between contrastive and non-contrastive self-supervised learning. In *International Conference on Learning Representations*, 2023.

[25] Jiachen Zhu, Rafael M Moraes, Serkan Karakulak, Vlad Sobol, Alfredo Canziani, and Yann LeCun. TiCo: Transformation invariance and covariance contrast for self-supervised visual representation learning. *arXiv preprint arXiv:2206.10698*, 2022.

[26] Shengbang Tong, Yubei Chen, Yi Ma, and Yann Lecun. EMP-SSL: Towards self-supervised learning in one training epoch. *arXiv preprint arXiv:2304.03977*, 2023.

[27] Quentin Garrido, Randall Balestriero, Laurent Najman, and Yann Le-cun. RankMe: Assessing the downstream performance of pretrained self-supervised representations by their rank. In *International Conference on Machine Learning*, pages 10929–10974, 2023.

[28] Zheng Dang, Kwang Moo Yi, Yinlin Hu, Fei Wang, Pascal Fua, and Mathieu Salzmann. Eigendecomposition-free training of deep networks with zero eigenvalue-based losses. In *European Conference on Computer Vision*, pages 768–783, 2018.

[29] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive multiview coding. In *European Conference on Computer Vision*, pages 776–794, 2020.

[30] Yubei Chen, Adrien Bardes, ZENGYI LI, and Yann LeCun. Bag of image patch embedding behind the success of self-supervised learning. *Transactions on Machine Learning Research*, 2023.

[31] Victor Guilherme Turrisi da Costa, Enrico Fini, Moin Nabi, Nicu Sebe, and Elisa Ricci. solo-learn: A library of self-supervised methods for visual representation learning. *Journal of Machine Learning Research*, 23(56):1–6, 2022.

[32] Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. Technical report, University of Toronto, 2009.

[33] Adam Coates, Andrew Ng, and Honglak Lee. An analysis of single-layer networks in unsupervised feature learning. In *International Conference on Artificial Intelligence and Statistics*, pages 215–223, 2011.

[34] Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. *CS 231N*, 7(7):3, 2015.

[35] Yang You, Igor Gitman, and Boris Ginsburg. Large batch training of convolutional networks. *arXiv preprint arXiv:1708.03888*, 2017.

[36] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

[37] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale image database. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255. Ieee, 2009.

[38] Yifan Zhang, Zhiquan Tan, Jingqin Yang, Weiran Huang, and Yang Yuan. Matrix information theory for self-supervised learning. In *International Conference on Machine Learning*, 2024.

# A  Pseudocode

## A.1  Pseudocode for FroSSL - 2 Views

```
for x in loader:
    # augment the image
    x_a, x_b = augment(x)

    # pass through network f to get embeddings
    z_a = f(x_a)
    z_b = f(x_b)
    N, D = Z_a.shape

    # center embeddings
    Z_a = Z_a - Z_a.mean(0)
    Z_b = Z_b - Z_b.mean(0)

    # normalize dimensions to sqrt(D) std.
    Z_a = (D**0.5) * (Z_a / Z_a.norm())
    Z_b = (D**0.5) * (Z_b / Z_b.norm())

    # calculate invariance (MSE) term
    invariance_loss = MSELoss(Z_a, Z_b)

    # calculate variance (Frobenius norm) term
    frobenius_a = torch.log(torch.norm(Z_a.T @ Z_a, ord='fro'))
    frobenius_b = torch.log(torch.norm(Z_b.T @ Z_b, ord='fro'))
    variance_loss = frobenius_a + frobenius_b

    # FroSSL loss
    loss = gamma*invariance_loss + variance_loss
    loss.backward()
    optimizer.step()
```

## A.2    Pseudocode for FroSSL - V Views

```
for x in loader:
    # augment the image
    X_list = [augment(x) for v in range(V)]

    # pass through network f to get embeddings
    Z_list = [f(x) for x in X_list]
    N, D = Z_list[0].shape

    # normalize embeddings
    Z_list = [F.normalize(z) for z in Z_list]
    Z_mean = torch.mean(torch.stack(Z_list), dim=0)

    total_loss = 0
    for Z_v in Z_list:
        # calculate invariance (MSE) term
        invariance_loss = MSELoss(Z_v, Z_mean)

        # calculate variance (Frobenius norm) term quickly
        if N > D:
            cov = view_embeddings.T @ view_embeddings
        else:
            cov = view_embeddings @ view_embeddings.T
        cov = cov / torch.trace(cov)
        fro_norm = torch.linalg.norm(cov, ord='fro')

        # bring frobenius square outside log
        variance_loss = -2*torch.log(fro_norm)


        loss = gamma*invariance_loss + variance_loss
        total_loss += loss

    total_loss.backward()
    optimizer.step()
```

# B    Derivations for Table 1

Here we derive the $D_{var}$ and $D_{inv}$ terms for each row in Table 1. We would like to note that while it is tempting to include works like Barlow Twins and MMCR in the framework of Table 1, they do not fit neatly as they are calculated between views rather than on individual views. A concurrent work, Matrix-SSL [38], introduces a loss function that could also be represented in the framework of Table 1. However, we leave that analysis to future work.

## B.1    FroSSL

In FroSSL, the minimization of the Frobenius norm brings the embedding covariance eigenvalues towards $\frac{c}{D}$, where c is some scaling factor. These eigenvalues are the same as the eigenvalues of $cI$. Thus the goal of the variance term in FroSSL is to bring each view covariance towards a unitary transformation of $I$. In other words, we are minimizing some distance between $I$ and the covariance. This distance is, in fact, the Renyi relative entropy of order 2. Below we start with the 2-order Petz-Renyi relative entropy and derive the FroSSL variance term.

$$S_2(A||B) = \ln\left(\mathrm{trace}(A^2 B)\right)$$

$$S_2(\Sigma_v||I) = \ln\mathrm{trace}(\Sigma_v^2) = \ln\|\Sigma_v\|_F^2$$

## B.2    I-VNE

The intuition for I-VNE follows similarly to FroSSL. The only difference is the choice of relative entropy distance between $\Sigma_v$ and $I$. Here $S_1$, or von Neumann relative entropy, instead of $S_2$ as in FroSSL.

## B.3    VICReg

The distance between the covariance and identity is explicitly defined in the VICReg objective function. We simply use its variance terms and invariance terms for $D_{\mathrm{var}}$ and $D_{\mathrm{inv}}$, respectively.

## B.4    CorInfoMax

First, notice that if we use variance-normalized features, we have $\mathrm{trace}(\Sigma_v + \epsilon\mathbf{I}) = D$. If we instead use sample-normalized features, then we have $\mathrm{trace}(\Sigma_v + \epsilon\mathbf{I}) = N$. In either case, the trace term remains constant and thus is not relevant for the optimization. The only non-constant term that changes is $\log\det(\Sigma_v + \epsilon\mathbf{I})$ which directly corresponds to the CorInfoMax loss.

### B.5    W-MSE

Because the projections are explicitly whitened, we automatically get $D_{var} = 0$. The choice of $D_{inv}$ comes straightforwardly from the loss.

## C    Desirable Dimension-Contrastive Criteria

In this section, we go over each dimension-contrastive method and discuss why it meets or fails to meet the four desirable criteria laid out in Section 3.

### C.1    Barlow Twins

**Invariant to Projection Rotation** ✕ The objective function explicitly tries to make the embedding cross-covariance identity. Rotations of identity are not suitable, even though they offer the same low redundancies as identity. Thus this criteria is not met.

**Manipulates Eigenvalues Explicitly** ✕ The objective function is almost, but not quite, a Frobenius norm. It is not clear how the objective function could be written in terms of eigenvalues.

**Scales Quadratically in Batch Size and Dimension** ✓ The objective function is a sum of matrix elements, which scales in $O(D^2)$.

**Scales Linearly in Views** ✕ Because Barlow Twins is based on the cross-covariance, it must compute all pairwise cross-covariances as the number of views increases. Thus it scales quadratically in views.

### C.2    VICReg

**Invariant to Projection Rotation** ✕ Same reasoning as with Barlow Twins, but with covariance matrices instead of cross-correlations.

**Manipulates Eigenvalues Explicitly** ✕ Same reasoning as Barlow Twins.

**Scales Quadratically in Batch Size and Dimension** ✓ Same reasoning as Barlow Twins.

**Scales Linearly in Views** ✓ Uses the covariance for each view independently, rather than computing all pairwise cross-correlation matrices as in Barlow Twins. The mean-square error can be reduced to linear by Equation 3 in the main text.

## C.3    W-MSE

**Invariant to Projection Rotation** ✓  The whitening operation achieves this.

**Manipulates Eigenvalues Explicitly** ×  The objective function only tries to minimize the distance between the views.

**Scales Quadratically in Batch Size and Dimension** ×  The whitening operation is computed in $O(D^3)$.

**Scales Linearly in Views** ✓  The whitening operation is done for each view, and the mean-square error can be reduced to linear by Equation 3 in the main text.

## C.4    CorInfoMax

**Invariant to Projection Rotation** ✓  Determinant is invariant to rotation.

**Manipulates Eigenvalues Explicitly** ✓  The objective function uses the log of the covariance determinant, which is defined as the product of eigenvalues.

**Scales Quadratically in Batch Size and Dimension** ×  Determinant can only be computed in $O(D^3)$.

**Scales Linearly in Views** ✓  Computes a determinant for each view. The mean-square error can be reduced to linear by Equation 3 in the main text.

## C.5    I-VNE

**Invariant to Projection Rotation** ✓  The von Neumann entropy is invariant to rotations.

**Manipulates Eigenvalues Explicitly** ✓  The von Neumann entropy of a matrix is defined as the Shannon entropy of its eigenvalues.

**Scales Quadratically in Batch Size and Dimension** ×  The von Neumann entropy requires explicit eigendecomposition, which is $O(\min(N, D)^3)$.

**Scales Linearly in Views** ✓  Computes an entropy for each view. The cosine similarity can be reduced to linear by Equation 3 in the main text.

### C.6    MMCR

**Invariant to Projection Rotation** ✓ The nuclear norm is invariant to rotations.

**Manipulates Eigenvalues Explicitly** ✓ The nuclear norm is defined as the sum of the magnitudes of singular values.

**Scales Quadratically in Batch Size and Dimension** ✗ The nuclear norm requires explicit eigendecomposition, which is $O(\min(N, D)^3)$.

**Scales Linearly in Views** ✓ Interestingly, MMCR is constant in the number of views. The objective function is computed on the average view embedding, rather than each view independently.

### C.7    FroSSL

**Invariant to Projection Rotation** ✓ The Frobenius norm is invariant to rotations.

**Manipulates Eigenvalues Explicitly** ✓ The Frobenius norm can be defined as the sum of the squared eigenvalues.

**Scales Quadratically in Batch Size and Dimension** ✓ The Frobenius norm can be computed as a sum over matrix elements in $O(\min(D^2, N^2))$.

**Scales Linearly in Views** ✓ A Frobenius norm is computed for each view. The mean-square error can be reduced to linear by Equation 3 in the main text.

## D    Experimental Details

### D.1    Stepwise Convergence Experimental Details

We trained ResNet-18 on STL10 using SGD with $lr = 0.01$ and a batch size of 256. Training occurred for only 5 epochs ($\approx$2000 steps) because we were interested in stepwise behaviors early during training. For all methods except CorInfoMax, a projector dimensionality of $D = 1024$ was used.

- **Barlow Twins** We used $\lambda = 0.05$ as recommended by [7].
- **VICReg** We used $\lambda = 25$, $\mu = 25$, $\nu = 1$ as recommended by [10].
- **CorInfoMax** We used temperature $\alpha = 500$, $\epsilon = 1e-6$, $\mu = 1e-2$, $D = 256$ following the STL-10 recommendations of [19].
- **I-VNE** The original work does not mention choices of tradeoffs between the invariance and variance terms [18]. Therefore, we equally weight the invariance and variance terms. We find this works well in practice
- **FroSSL** We used an invariance weight $\gamma = 1.4$.

### D.2   Linear Probe Experimental Details

**Hyperparameters** We follow the guidance of [31], [19], and [20] for selecting baseline hyperparameters. Our code contains .yaml files with the exact hyperparameter configurations we used on each dataset.

**Augmentations** For datasets other than Tiny ImageNet, we use the default symmetric augmentations provided in [31] for FroSSL. These:

- Random Resized Crop - Crop scale ranges from 0.08 to 1.0, then resized to 32x32 for CIFAR, 96x96 for STL-10, 224x224 for ImageNet-100
- Color Jitter with probability 0.8 and (brightness, contrast, saturation, hue) values of (0.8, 0.8, 0.8, 0.2)
- Grayscale with probability 0.2
- Gaussian blur with probability 0.5
- Horizontal flip with probability 0.5

For Tiny ImageNet, we use the asymmetric augmentations provided in [19]. These are shown below. Because these were originally designed for methods using 2 views, we used equal numbers of each view type for multiview methods.

- Random Resized Crop - Crop scale ranges from 0.08 to 1.0, then resized to 64x64
- Color Jitter with probability 0.8 and (brightness, contrast, saturation, hue) values of (0.4, 0.4, 0.2, 0.1)
- Grayscale with probability 0.2
- Gaussian blur with probability 1.0 in view 1 and probability 0.1 in view 2
- Solarization with probability 0 in view 1 and probability 0.2 in view 2
- Horizontal flip with probability 0.5

## E   Proofs

### E.1   Proof of Proposition 1

WLOG, we prove the following for 2 views. We start with rewriting the arg min of Equation 7 as such:

$$\underset{Z_1,Z_2}{\arg\min}\,\mathcal{L}_{\mathrm{FroSSL}} = \underset{Z_1,Z_2}{\arg\min}\left[\log(\left\|Z_1^T Z_1\right\|_F^2) + \log(\left\|Z_2^T Z_2\right\|_F^2) + \mathcal{L}_{\mathrm{MSE}}(Z_1, Z_2)\right]$$

$$= \underset{Z_1,Z_2}{\arg\min}\left[\left\|Z_1^T Z_1\right\|_F^2 + \left\|Z_2^T Z_2\right\|_F^2 + \mathcal{L}_{\mathrm{MSE}}(Z_1, Z_2)\right]$$

Assume the normalization step normalizes each embedding dimension to have unit variance. Then both covariance matrices have 1 in each diagonal element.

$$= \underset{Z_1,Z_2}{\arg\min}\left[\left\|Z_1^T Z_1 - \mathrm{diag}(Z_1^T Z_1)\right\|_F^2 + \left\|Z_2^T Z_2 - \mathrm{diag}(Z_2^T Z_2)\right\|_F^2 + 2D + \mathcal{L}_{\mathrm{MSE}}(Z_1, Z_2)\right]$$

$$= \underset{Z_1,Z_2}{\arg\min}\left[\mathcal{L}_{nc}(Z_1) + \mathcal{L}_{nc}(Z_2) + 2D + \mathcal{L}_{\mathrm{MSE}}(Z_1, Z_2)\right]$$

Thus we have that the embeddings that minimize FroSSL also minimize the non-contrastive losses $\mathcal{L}_{nc}$ for both views.

## E.2    Proof of Proposition 2

WLOG, we prove the following for 2 views. With the duality of the Frobenius norm, we rewrite Equation 7 to use Gram matrices rather than covariance matrices:

$$\mathcal{L}_{\text{FroSSL}} = \log(\left\|Z_1^T Z_1\right\|_F^2) + \log(\left\|Z_2^T Z_2\right\|_F^2) + \mathcal{L}_{\text{MSE}}(Z_1, Z_2)$$

$$= \log(\left\|Z_1 Z_1^T\right\|_F^2) + \log(\left\|Z_2 Z_2^T\right\|_F^2) + \mathcal{L}_{\text{MSE}}(Z_1, Z_2)$$

Assuming that each embedding is normalized to have unit norm, then both Gram matrices have 1 in each diagonal element. Then the rest of the proof then follows similarly to Proposition 1.

# F    Minibatch Time for Larger Batch Size and Dimensionality

**Table 8:** Comparison of the time/space tradeoffs for SSL algorithms trained on STL-10. We used a batch size of N=1024 and dimensionality D=1024. The advantages of the reduced time complexity of FroSSL are more apparent here than in Table 3 of the main text due to the larger batch size and dimensionality. We encourage the reader to constrast this with Table 3.

| | SimCLR | MoCo v2 | BYOL | VICReg | Barlow | CorInfoMax | MMCR | | | FroSSL (ours) | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Loss Time Complexity | | $O(V^2 N^2)$ | | $O(VD^2)$ | $O(V^2 D^2)$ | $O(VD^3)$ | $O(\min(D,N)^3)$ | | | $O(V\min(D,N)^2)$ | | |
| Num. Views | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 4 | 8 | 2 | 4 | 8 |
| VRAM Space (GB) | 5.5 | 6.7 | 6.3 | 5.4 | 5.4 | 5.5 | 5.4 | 10.4 | 19.6 | 5.4 | 10.2 | 19.6 |
| Minibatch Wall-time (ms) | 165 | 220 | 210 | 168 | 165 | 231 | 252 | 400 | 700 | 168 | 318 | 601 |