

# SSOLE: RETHINKING ORTHOGONAL LOW-RANK EMBEDDING FOR SELF-SUPERVISED LEARNING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Self-supervised learning (SSL) aims to learn meaningful representations from unlabeled data. Orthogonal Low-rank Embedding (OLE) shows promise for SSL by enhancing intra-class similarity in a low-rank subspace and promoting inter-class dissimilarity in a high-rank subspace, making it particularly suitable for multi-view learning tasks. However, directly applying OLE to SSL poses significant challenges: (1) the virtually infinite number of "classes" in SSL makes achieving the OLE objective impractical, leading to representational collapse; and (2) low-rank constraints may fail to distinguish between positively and negatively correlated features, further undermining learning. To address these issues, we propose **SSOLE** (Self-Supervised Orthogonal Low-rank Embedding), a novel framework that integrates OLE principles into SSL by (1) decoupling the low-rank and high-rank enforcement to align with SSL objectives; and (2) applying low-rank constraints to feature deviations from their mean, ensuring better alignment of positive pairs by accounting for the signs of cosine similarities. Our theoretical analysis and empirical results demonstrate that these adaptations are crucial to SSOLE's effectiveness. Moreover, SSOLE achieves competitive performance across SSL benchmarks without relying on large batch sizes, memory banks, or dual-encoder architectures, making it an efficient and scalable solution for self-supervised tasks.

## 1 INTRODUCTION

Self-supervised learning (SSL) (Chen et al., 2020a; He et al., 2020; Bardes et al., 2022), learns meaningful representations from unlabeled data by exploiting the intrinsic structure within the data, reducing dependence on costly labeled datasets. SSL has become crucial in fields like computer vision, natural language processing, and speech recognition, enabling models to harness vast amounts of unannotated data.

Orthogonal Low-rank Embedding (OLE) (Qiu & Sapiro, 2015; Lezama et al., 2018), originally developed for supervised image classification, constructs a feature space where same-class samples reside in low-rank subspaces, enhancing intra-class similarity, while different-class samples are orthogonal, promoting inter-class dissimilarity. OLE operates at the matrix level, optimizing the nuclear norm as a surrogate for rank, making it well-suited for leveraging multiple views or augmentations of data in SSL. Integrating OLE into SSL offers several benefits: (1) Its matrix-level operation is well-suited for leveraging multiple views or augmentations inherent in SSL; (2) Enforcing low-rank constraints on positive pairs not only brings their representations closer but also minimizes the number of factors controlling the image representation, leading to more compact features; (3) Enforcing high-rank constraints on negative pairs not only pushes their representations apart but also prevents dimensional collapse, maximizing the representational capacity of the embedding space.

However, OLE has not been successfully applied as a standalone metric in SSL. Previous works have only partially utilized OLE's components. For example, LORAC (Wang et al., 2022) incorporates low-rank embedding as a regularization term within the MoCo framework (He et al., 2020) but relies heavily on contrastive loss. [Methods like VICReg \(Bardes et al., 2022\) and Total Coding Rate \(TCR\) \(Yu et al., 2020; Tong et al., 2023\) enforce orthogonality among negative pairs via soft regularization, while W-MSE \(Ermolov et al., 2021\) employs feature whitening.](#) They align with OLE's high-rank constraints, but neglecting the low-rank embedding of positive pairs. This gap prompts the question: Why has OLE not been fully integrated into SSL?

Our investigation reveals two key challenges. First, in SSL, each instance acts as its own class, creating an infinite number of "classes." The OLE objective requires enforcing orthogonality between classes and minimizing intra-class ranks, but orthogonality is mathematically impossible with infinite classes in a limited feature space, leading to representational collapse. Second, enforcing low-rank constraints on feature vectors can fail to distinguish between positively and negatively correlated features. The nuclear norm ignores the signs of cosine similarities, which is critical for aligning positive pairs in SSL, where supervision is absent.

To address these issues, we propose **SSOLE** (Self-Supervised Orthogonal Low-rank Embedding), a framework that effectively integrates OLE into SSL. Our approach involves two key strategies: (1) **Decoupling low-rank and high-rank enforcement**: We separately manage positive pair attraction and negative pair repulsion, eliminating the dependence of low-rank enforcement on inter-class orthogonality. We optimize both the lower and upper bounds of the contrastive loss using nuclear norms to adapt OLE for SSL. (2) **Applying low-rank constraints to feature deviations**: Instead of directly applying low-rank constraints to the feature vectors, we apply them to the deviations from their mean, ensuring that the nuclear norm captures the signs of cosine similarities and aligns positive pairs correctly.

Through theoretical analysis and empirical evaluations, we demonstrate that these adaptations are essential for SSOLE's success in SSL. Our method addresses the identified challenges, achieving state-of-the-art performance on benchmark datasets. SSOLE leverages OLE's strengths while avoiding its limitations in SSL, offering an efficient solution for representation learning.

In summary, our contributions are: (1) identifying key challenges in applying OLE to SSL and offering insights into its limitations; (2) proposing SSOLE, which decouples low-rank and high-rank enforcement and applies low-rank constraints to feature deviations; (3) validating SSOLE's superior performance through theoretical and empirical analysis; and (4) opening new avenues for matrix-level operations and rank constraints in SSL, offering potential for future advancements.

## 2 BACKGROUND

### 2.1 ORTHOGONAL LOW-RANK EMBEDDING (OLE)

Orthogonal Low-rank Embedding (OLE) (Lezama et al., 2018) was introduced as a geometric loss function to improve deep network representations by simultaneously minimizing intra-class variance and maximizing inter-class separability. Unlike traditional softmax-based classification losses, OLE enforces that samples from the same class lie in a low-rank subspace while ensuring that different class subspaces are orthogonal. This approach encourages compact and well-separated representations in the embedding space, leading to improved generalization and robustness.

The OLE loss is defined as:

$$\mathcal{L}_{\text{OLE}} = \sum_{c=1}^K \|\Phi(\mathbf{X}_c)\|_* - \|\Phi(\mathbf{X})\|_*, \quad (1)$$

where  $\mathbf{X}_c$  denotes the set of samples from class  $c$ ,  $\mathbf{X}$  is the entire data set from  $K$  classes,  $\Phi(\mathbf{X})$  represents the transformation applied to the entire dataset, and  $\|\cdot\|_*$  is the nuclear norm, which approximates the rank of a matrix. The first term minimizes the intra-class rank, while the second term maximizes the inter-class separability by encouraging orthogonality between different classes.

OLE achieves a non-negative value and reaches zero when the representations of all classes are orthogonal, implying maximal inter-class separation. At the same time, intra-class samples are compressed into a low-dimensional subspace, preserving intra-class similarity. The nuclear norm serves as a convex relaxation of the rank function, facilitating efficient optimization through gradient-based methods.

## 3 LIMITATIONS OF OLE FOR SSL

Although OLE has been successful in supervised settings, its application to SSL poses significant challenges. [These challenges stem directly from two inherent limitations of OLE, which are manageable in the presence of labels but become problematic in SSL due to the lack of supervision.](#) Firstly, in

OLE, low-rank (intra-class similarity) and high-rank (inter-class dissimilarity) constraints are deeply intertwined, making it difficult to optimize them independently. Secondly, the nuclear norm used in OLE cannot differentiate between positively and negatively correlated vectors. In the following subsections, we elaborate on the resulted challenges from them.

### 3.1 VIRTUALLY INFINITE NUMBER OF "CLASSES" IN SSL

A key challenge of applying OLE to self-supervised learning (SSL) arises from the fact that, in SSL, each instance is treated as its own class. This results in a virtually infinite number of "classes," making it impossible to achieve orthogonality among all class representations. The entanglement of low-rank and high-rank enforcement exacerbates this issue, as the original OLE objective depends on achieving inter-class orthogonality before minimizing intra-class ranks.

In supervised settings with finite classes, OLE works by first ensuring orthogonality between different class subspaces and then minimizing the intra-class rank, leading to compact and discriminative features. However, in SSL, where orthogonality is not feasible, OLE "cheats" by reducing vector length  $l$  and the angle  $\theta$  between vectors to minimize the nuclear norm, resulting in representational collapse.

To explain this, consider the model generating random representations from a  $d$ -dimensional Gaussian distribution. For  $m$  samples drawn from this distribution, the expected nuclear norm of the matrix  $\mathbf{V}$  of these samples can be approximated by the following lemma:

**Lemma 3.1.** *Let  $\mathbf{V}$  be an  $m \times d$  matrix, with  $d \gg 1$ , whose rows are sampled from a  $d$ -dimensional Gaussian distribution  $\mathcal{N}\left(\frac{s\mathbf{1}_d}{\sqrt{d}}, \frac{s^2\sigma^2\mathbf{I}_d}{d}\right)$ . The expected nuclear norm of  $\mathbf{V}$  satisfies:*

$$\mathbb{E}[\|\mathbf{V}\|_*] = \begin{cases} s \cdot (\sigma \cdot \mathcal{O}(m) + \mathcal{O}(\sqrt{m})), & \text{if } m \ll d, \\ s \cdot (\sigma \cdot \mathcal{O}(\sqrt{md}) + \mathcal{O}(\sqrt{m})), & \text{if } m \gg d. \end{cases}$$

Lemma 3.1 follows the random matrix theory, and the proof is provided in Appendix A.1. Extending it to OLE loss over  $B$  images, each with  $N$  views, we can approximate the OLE loss.

**Theorem 3.2.** *If the model's output representation conforms to  $\mathcal{N}\left(\frac{s\mathbf{1}_d}{\sqrt{d}}, \frac{s^2\sigma^2\mathbf{I}_d}{d}\right)$ , then for  $B$  images, each with  $N$  views, where  $1 < N \ll d \ll BN$ , the OLE loss satisfies:*

$$\mathcal{L}_{OLE} = s\sqrt{BN} \left( \sigma \cdot \mathcal{O}(\sqrt{BN}) + \mathcal{O}(\sqrt{B}) \right).$$

The proof is provided in Appendix A.2. Theorem 3.2 shows that  $\mathcal{L}_{OLE}$  is negatively correlated with the scale  $s$  and  $\sigma$ . The analysis in Appendix A.3 also shows how  $s$  and  $\sigma$  relate the average vector length  $l$  and angle  $\theta$  between vectors. So the OLE loss decreases as  $l$  or  $\theta$  becomes smaller. In SSL, where orthogonality is unachievable, OLE reduces  $l$  and  $\theta$  to minimize the loss, leading to collapsed, trivial representations. Once this collapse occurs, further training becomes difficult, as the nuclear norm is dominated by the reduction in  $l$  and  $\theta$ , preventing the model from learning separable features.

Even normalization techniques in SSL, which prevent  $l$  from shrinking, do not solve the problem. The OLE loss can still minimize  $\theta$ , leading to collapse. Additionally, imposing a bound on the intra-class nuclear norm, as done in the original OLE, does not resolve this. If the bound is not reached, the model continues to shrink both  $l$  and  $\theta$ . Once the bound is reached, the model oscillates between increasing and shrinking  $l$  and  $\theta$ , causing instability.

This challenge highlights a fundamental limitation for OLE. While OLE performs well in supervised settings by enforcing orthogonality, it struggles in SSL, where orthogonality cannot be achieved. Despite efforts like normalization and bounded nuclear norms, there is currently no effective solution to this issue. To address this, Section 4.1 proposes to decouple low-rank and high-rank enforcements, allowing independent optimization of them to avoid representational collapse. It also explores the connection between nuclear norm optimization and contrastive objectives to refine these enforcements.

### 3.2 LOW-RANK REPRESENTATIONS MAY NOT BE GOOD REPRESENTATIONS

Another significant challenge arises from the nuclear norm's inability to distinguish between positively and negatively correlated features. In SSL, this limitation leads to misalignment between positive pairs

because the nuclear norm treats aligned and anti-aligned vectors similarly, degrading the quality of learned representations. Without distinguishing between aligned and anti-aligned vectors, the nuclear norm can treat opposite directions similarly, causing misaligned positive pairs and undermining the learning process. This issue is particularly problematic in SSL, where the lack of labels prevents the model from correcting these misalignment.

Consider the nuclear norm for a 2-row matrix  $\mathbf{V}$  where the rows are two unit vectors separated by an angle  $\theta$ . Its nuclear norm is given by  $\|\mathbf{V}\|_* = \sqrt{1 + \cos \theta} + \sqrt{1 - \cos \theta}$ . This shows that the nuclear norm only depends on the magnitude of  $\cos \theta$ , regardless of whether the cosine is positive or negative. Whether the vectors are aligned ( $\cos \theta > 0$ ) or anti-aligned ( $\cos \theta < 0$ ), the nuclear norm remains unchanged. This inability to distinguish between the directions of the vectors is problematic for SSL, where positive pairs should be aligned and negative pairs should remain distinct.

More generally, this limitation can be formalized. Let  $\mathbf{V}$  be an  $m$ -row matrix where each row is a  $d$ -dimensional unit vector, and let  $\mathbf{P}$  be an  $m \times m$  diagonal matrix whose diagonal elements are either 1 or  $-1$ , which invert or preserve the signs of the rows of  $\mathbf{V}$ . The nuclear norm of  $\mathbf{P}\mathbf{V}$  remains unchanged, demonstrating the unitarily invariant property of the nuclear norm. Further discussion is provided in Appendix A.4.

Currently, there are no effective methods to resolve this limitation in SSL. As a result, the inability of the nuclear norm to differentiate vector directions remains a significant obstacle for OLE-based methods in SSL. To address this challenge, Section 4.2 proposes using deviation matrices for low-rank enforcement. This ensures alignment of positive pairs by focusing on deviations from the mean and circumventing the nuclear norm’s insensitivity to cosine similarity signs.

## 4 METHOD

We present how to address the key challenges of applying OLE to SSL. To handle the issue of infinite classes in SSL, we decouple low-rank and high-rank enforcement, allowing independent control of positive pair alignment and negative pair separation. We also modify the nuclear norm enforcement to account for the signs of cosine similarities, preventing representational collapse and ensuring proper alignment of feature spaces. Our framework, SSOLE (Self-Supervised Orthogonal Low-rank Embedding), integrates these solutions to effectively adapt OLE for SSL tasks.

### 4.1 ADDRESSING CHALLENGE 1: DECOUPLING LOW-RANK AND HIGH-RANK ENFORCEMENT

To address the challenge of infinite number of "classes" in SSL, we propose to decouple low-rank and high-rank enforcement to ensure that each term operates independently, allowing the model to enforce both alignment within positive pairs and uniformity across negative pairs more effectively.

**Normalizing feature vectors:** First, we alleviate some of the collapse problems by normalizing the feature vectors. By ensuring that all feature vectors have a fixed unit length, we prevent the model from shrinking vector lengths to minimize the loss artificially.

**Decoupling low-rank and high-rank enforcement:** To further address the challenge, we modify the original OLE objective to separate the low-rank enforcement for positive pairs from the high-rank enforcement for negative pairs. We propose the following loss function:

$$\mathcal{L} = \frac{1}{B} \sum_{b=1}^B h_1(\|\mathbf{Z}_{b,:}\|_*, N) + \lambda \frac{1}{N} \sum_{n=1}^N h_2(\|\mathbf{Z}_{:,n}\|_*, B), \quad (2)$$

where  $\mathbf{Z} \in \mathbb{R}^{B \times N \times d}$  is the representation tensor for  $B$  images in a batch, each with  $N$  augmented views. The matrix  $\mathbf{Z}_{b,:}$  represents the views of the  $b^{th}$  image, and  $\mathbf{Z}_{:,n}$  represents the  $n^{th}$  view across all images.  $h_1(\cdot)$  and  $h_2(\cdot)$  are strictly increasing and decreasing functions respectively and used to control the nuclear norm-based enforcement for low-rank and high-rank. The term  $\lambda$  balances the two objectives.

This formulation decouples the two terms and allows each objective to be enforced more appropriately, reducing the entanglement between intra-class and inter-class representations that causes collapse.

However, directly applying the nuclear norm as the function  $h_1(\cdot)$  or  $h_2(\cdot)$  still carries the risk of collapse. To address this, we explore relationships between nuclear norms and cosine similarities,

drawing inspiration from prior work (Wang & Isola, 2020), which decomposes SSL objectives into alignment (for positive pairs) and uniformity (for negative pairs).

**Studying the nuclear norms.** To guide the choice of  $h_1(\cdot)$  and  $h_2(\cdot)$ , we develop the following theorem, which relates nuclear norms to the cosine similarities between vectors.

**Theorem 4.1.** *For a set of  $N$   $d$ -dimensional unit vectors  $\{\mathbf{v}_i\}$ , where  $d > N$ , the nuclear norm of the matrix  $\mathbf{V}$ , containing the vectors as columns, is bounded by:*

$$\sqrt{\frac{N}{|\overline{\cos(\theta)}|}} \leq \|\mathbf{V}\|_* \leq \sqrt{N} \cdot \sqrt{\overline{\cos(\theta)}} + \sqrt{N(N-1)} \cdot \sqrt{1 - \overline{\cos(\theta)}},$$

where  $\overline{\cos(\theta)} = \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N \cos(\theta_{i,j}) = \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N \mathbf{v}_i^T \mathbf{v}_j$  is the average cosine similarity between every two vectors, and  $|\overline{\cos(\theta)}|$  is the average absolute cosine similarity.

The lower bound is attained when the vectors can be equally grouped into  $\frac{1}{|\overline{\cos(\theta)}|}$  groups such that within each group, the vectors are either perfectly aligned or anti-aligned, while the vectors from different groups are orthogonal. The upper bound holds when the angles are identical for all vector pairs  $(\mathbf{v}_i, \mathbf{v}_j)$  where  $i \neq j$ .

The proof is provided in Appendix A.5. Theorem 4.1 shows that the nuclear norm is both lower and upper bounded by a function of the average (absolute) cosine similarity between vectors. When  $\overline{\cos(\theta)} = |\overline{\cos(\theta)}| = \frac{1}{N}$ , both the lower and upper bound attain the maximum of  $N$ . Conversely, when  $\overline{\cos(\theta)} = |\overline{\cos(\theta)}| = 1$ , both the lower and upper bound attain the minimum of  $\sqrt{N}$ . These bounds allow us to design functions  $h_1(\cdot)$  and  $h_2(\cdot)$ .

Inverting the bounds in Theorem 4.1, we can derive lower and upper bounds for the average absolute cosine similarity  $|\overline{\cos(\theta)}|$  in terms of the nuclear norm:

$$\frac{N}{\|\mathbf{V}\|_*^2} \leq |\overline{\cos(\theta)}|; \quad \overline{\cos(\theta)} \leq \cos^2 \left( \arcsin\left(\frac{\|\mathbf{V}\|_*}{N}\right) - \arcsin\left(\frac{1}{\sqrt{N}}\right) \right). \quad (3)$$

This insight enables us to control the cosine similarities between vectors through nuclear norm minimization or maximization.

**Caution in Low-rank Enforcement:** However, care must be taken when using these bounds for low-rank enforcement. The lower bound on  $|\overline{\cos(\theta)}|$  ignore the signs of the cosine similarities, and the upper bound has a local minimum at  $\overline{\cos(\theta)} = 0$  where all off-diagonal cosine similarities are  $\frac{-1}{N-1}$  which is negative. This is not an issue for high-rank enforcement, as the goal is to achieve orthogonality. However, for low-rank enforcement, there may be cases where the absolute cosine similarity is maximized, but the actual cosine similarity averages to zero, resulting in misaligned vectors. Thus, we need to be cautious in designing  $h_1(\cdot)$  to ensure it aligns positive pairs effectively.

## 4.2 ADDRESSING CHALLENGE 2: LOW-RANK ENFORCEMENT VIA DEVIATION MATRICES

To address [the challenge of misalignment between positive pairs](#), we propose enforcing low-rank constraints on a **deviation matrix**. The core issue lies in the nuclear norm’s insensitivity to cosine similarity signs. Applying the nuclear norm directly to positive pairs’ representations risks misaligning them. By subtracting the mean vector from each individual vector, we create a deviation matrix that captures how each vector deviates from the average. The rank of the deviation matrix approximates the rank of the original matrix (with at most a difference of 1), so low-rank enforcement on the deviation matrix remains valid. Additionally, deviations in opposite directions are acceptable in SSL, meaning this approach circumvents the cosine similarity sign issue.

We formalize this approach with the following theorem:

**Theorem 4.2.** *For a set of  $N$   $d$ -dimensional unit vectors  $\{\mathbf{v}_i\}$ , where  $d > N$ , the nuclear norm of  $\tilde{\mathbf{V}} = \mathbf{V} - \bar{\mathbf{v}}\mathbf{1}^T$ , where  $\bar{\mathbf{v}} := \frac{1}{N} \sum_{i=1}^N \mathbf{v}_i$  and  $\mathbf{1} \in \mathbb{R}^N$  is a vector of ones, is bounded by:*

$$\sqrt{N} \cdot \sqrt{1 - \overline{\cos(\theta)}} \leq \|\tilde{\mathbf{V}}\|_* \leq \sqrt{N(N-1)} \cdot \sqrt{1 - \overline{\cos(\theta)}}.$$

The lower bound is reached when the vectors  $\{\mathbf{v}_i\}$  can be grouped into two perfectly aligned sets, and the upper bound holds when the angles are identical for all vector pairs  $(\mathbf{v}_i, \mathbf{v}_j)$  where  $i \neq j$ .

The proof is provided in Appendix A.6. The bounds provided by Theorem 4.2 depend on  $\overline{\cos(\theta)}$ , which incorporates the signs of the cosine similarities. This ensures that low-rank enforcement will not confuse with anti-aligned vectors, unlike in the original matrix where such signs may be ignored.

Further, when  $\overline{\cos(\theta)}$  reaches its maximum value of 1 (perfect alignment), both the lower and upper bounds become 0, indicating a perfectly low-rank representation. Conversely, when  $\overline{\cos(\theta)}$  reaches its minimum value of 0 (orthogonal or anti-aligned vectors), the lower bound becomes  $\sqrt{N}$ , corresponding to the lowest nuclear norm of the original matrix, while the upper bound reaches  $\sqrt{N(N-1)}$ , the highest possible nuclear norm for the deviation matrix.

By inverting the bounds in Theorem 4.2, we derive lower and upper bounds for  $1 - \overline{\cos(\theta)}$  in terms of the nuclear norm of the deviation matrix:

$$\frac{\|\tilde{\mathbf{V}}\|_*^2}{N(N-1)} \leq 1 - \overline{\cos(\theta)} \leq \frac{\|\tilde{\mathbf{V}}\|_*^2}{N}. \quad (4)$$

These bounds offer insights into the cosine similarities based on the nuclear norm of the deviation matrix. Enforcing low-rank representations on the deviation matrix resolves the issues caused by applying the nuclear norm directly to the feature matrix, ensuring that low-rank enforcement does not collapse into suboptimal solutions. This approach allows SSL to maintain the benefits of low-rank enforcement while preserving meaningful feature representations.

### 4.3 TRAINING OBJECTIVE

Given the improvements to low-rank enforcement, we need to modify Equation (2) to derive the training objective for SSOLE as follows:

$$\mathcal{L}_{SSOLE} = \frac{1}{B} \sum_{b=1}^B h_1(\|\tilde{\mathbf{z}}_{b,:}\|_*, N) + \lambda \frac{1}{N} \sum_{n=1}^N h_2(\|\mathbf{z}_{:,n}\|_*, B), \quad (5)$$

where  $\tilde{\mathbf{z}}_{b,:}$  denotes the deviation matrix of  $\mathbf{z}_{b,:}$ .

The loss  $\mathcal{L}_{SSOLE}$  leverages both lower and upper bounds for intra-class and inter-class cosine similarity derived from Theorem 4.1 and Theorem 4.2. These bounds provide a framework for optimizing a contrastive-like objective based on nuclear norms, which aligns with the goals of maximizing intra-class similarity (alignment) and inter-class dissimilarity (uniformity).

For low-rank enforcement, we derive  $h_1(\cdot)$  from Equation (4), basing on optimizing the lower bound of the nuclear norm. Since the lower and upper bounds are equivalent up to a factor of  $\frac{1}{N-1}$ , optimizing the lower bound is sufficient to optimize both bounds. Then  $h_1(\cdot)$  is given by<sup>1</sup>:

$$h_1(\|\tilde{\mathbf{z}}_{b,:}\|_*, N) = \frac{\|\tilde{\mathbf{z}}_{b,:}\|_*^2}{(N-1)^2}. \quad (6)$$

For high-rank enforcement, we derive  $h_2(\cdot)$  from Equation (3), optimizing the average of the lower and upper bounds of average (absolute) cosine similarity from the nuclear norm.  $h_2(\cdot)$  is given by:

$$h_2(\|\mathbf{z}_{:,n}\|_*, B) = \frac{B}{2(B-1)} \left( \frac{B}{\|\mathbf{z}_{:,n}\|_*^2} + \cos^2(\arcsin(\frac{\|\mathbf{z}_{:,n}\|_*}{B}) - \arcsin(\frac{1}{\sqrt{B}})) \right) - \frac{1}{B-1}. \quad (7)$$

The SSOLE training objective combines the advantages of OLE for SSL while optimizing both a lower and upper bound for a contrastive-like loss using nuclear norms. By utilizing these bounds, the model can enforce alignment and uniformity in a more efficient manner, avoiding the representational collapse seen in vanilla OLE approaches.

<sup>1</sup>Note that in  $h_1(\cdot)$  and  $h_2(\cdot)$ , we need to multiply a factor of  $\frac{N-1}{N}$  or  $\frac{B-1}{B}$  to obtain the estimated average (absolute) cosine similarities for all vector pairs  $(\mathbf{v}_i, \mathbf{v}_j)$  where  $i \neq j$ .



## 5 RELATED WORK

### 5.1 SELF-SUPERVISED LEARNING

Self-supervised learning has advanced significantly in the realm of image recognition by leveraging various innovative techniques. Methods such as BYOL (Grill et al., 2020) and MoCo (He et al., 2020) use bootstrapping and dynamic dictionaries to enhance representation learning. SimSiam (Chen & He, 2021) explores learning representations without negative pairs, while CPCv2 (Henaff, 2020) emphasizes data-efficient image recognition. The principles of alignment and uniformity on a hypersphere have been analyzed by Wang & Isola (2020). Meanwhile, the impact of view (augmentation) selection in contrastive learning is investigated by Tian et al. (2020). DINO (Caron et al., 2021) stands out as a robust self-supervised method, having evolved its training protocol to achieve competitive results. MoCo v3 (Chen et al., 2021) builds upon momentum contrast for training Vision Transformers. VicReg (Bardes et al., 2022) introduces an approach based on variance, invariance, and covariance, while Wang et al. (2021) address inefficiencies in representation learning. iBOT (Zhou et al., 2022) focuses on Image BERT pre-training and RELIC v2 (Tomasev et al., 2022) ambitiously aims to outperform supervised learning on ImageNet without labels. Furthermore, the work of He et al. (2021) introduces Masked Autoencoders (MAE), a scalable vision learner that benefits from the reconstruction of masked image patches. This method implicitly utilizes multiple views by treating visible and masked patches differently during the learning process. Wu et al. (2018) propose a non-parametric approach to instance discrimination.

Multi-View Self-Supervised Learning (MV-SSL) has recently emerged as a potent paradigm to harness the information from various augmentations or views of the same data. This approach has led to significant advancements in SSL by promoting more generalized feature representations. SwAV (Caron et al., 2020) introduces a unique "swapped prediction" task to SSL, utilizing cluster assignments as pseudo-labels to encourage consistency across different augmentations or views. It employs multiple views of an image to compute these assignments, promoting invariance across views and improving the learned representations. LORAC (Wang et al., 2022) extends the principles of MoCo by incorporating low-rank embedding as a prior, which is particularly beneficial for SSL. It leverages multiple views to enforce consistency. EMP-SSL (Tong et al., 2023) takes a different approach by generating an extremely large number of patches or views from the input images, significantly reducing the training epochs to converge.

This work is related to MV-SSL and LORAC. While LORAC incorporates low-rank constraints as a regularization term within contrastive learning frameworks, it does not fully exploit the potential of OLE as a central metric. In contrast, our proposed SSOLE redefines the role of OLE in SSL by using it as an intrinsic metric for both positive alignment and negative separation. By strategically adapting OLE to address SSL-specific challenges, SSOLE fully leverages the unique properties of multi-view data. This represents a significant departure from prior approaches and unlocks the full potential of OLE in SSL.

### 5.2 ORTHOGONAL LOW-RANK EMBEDDING AND RELATED TECHNIQUES

The concept of learning low-dimensional, structured representations through low-rank constraints has been extensively studied across multiple domains, including Principal Component Analysis (PCA), Linear Discriminant Analysis (LDA) (Hastie et al., 2009), face recognition (Yang et al., 2016; Xue et al., 2017; Lezama et al., 2017; Xue et al., 2019), and image classification (Zhang et al., 2013; Jiang et al., 2014; Zhang et al., 2016). High-rank regularization has also been explored for applications such as learning orthogonal projections in deep networks (Vorontsov et al., 2017), improving recurrent network performance (Bansal et al., 2018), and capsule subspace projection (Zhang et al., 2018).

In deep learning, Orthogonal Low-Rank Embedding (OLE) (Lezama et al., 2018) extended these ideas by introducing low-rank and orthogonal constraints into supervised classification tasks. By enforcing intra-class low-rank and inter-class orthogonality, OLE achieves compact and discriminative representations. Inspired by these principles, subsequent works like LORAC (Wang et al., 2022) incorporated low-rank priors into SSL, albeit as regularizations rather than direct metrics.

Further extensions of OLE include data-dependent regularizations (Zhu et al., 2019), which aim to enhance pattern discovery and prevent overfitting, and Meta-OLE (Wang et al., 2023), which adapts OLE for meta-learning.

While these works establish the utility of low-rank and high-rank constraints, they primarily focus on supervised settings. Their direct application to SSL faces unique challenges. Our proposed method, SSOLE, builds upon OLE by addressing these challenges through two key innovations. These adaptations enable SSOLE to effectively extend OLE principles to SSL.

## 6 EXPERIMENTS

### 6.1 ABLATION STUDIES

We focus on exploring effective adaptations of OLE for SSL. The ResNet-18 (He et al., 2016) architecture is employed on the ImageNet100 (Deng et al., 2009) dataset. Detailed information on data augmentation and training procedures is provided in Appendix D.

#### 6.1.1 ADAPTING OLE TO SSL

We introduce a baseline method of InfoNCE-M, an extension of the InfoNCE loss adapted for MV-SSL. InfoNCE-M uses the mean of all views as the anchor for each image, computing the InfoNCE loss for each view and averaging these values. The formula for InfoNCE-M is

$$\mathcal{L}_{\text{InfoNCE-M}} = -\frac{1}{BN} \sum_{i=1}^B \sum_{j=1}^N \log \frac{e^{\text{sim}(\mathcal{Z}_{i,j}, \mathbf{m}_i)/\tau}}{\sum_{k=1}^B e^{\text{sim}(\mathcal{Z}_{i,j}, \mathbf{m}_k)/\tau}}, \quad (8)$$

where  $\mathbf{m}_i = \frac{1}{N} \sum_{j=1}^N \mathcal{Z}_{i,j}$  is the mean embedding (anchor) of all views for the  $i$ -th image, and  $\tau$  is the temperature scaling parameter. The function  $\text{sim}(\cdot, \cdot)$  computes the cosine similarity.

We then explore adaptations of OLE for SSL. The results of these adaptations are summarized in Table 1. The standard  $\mathcal{L}_{\text{InfoNCE-M}}$  with a hyperparameter ( $\tau = 0.2$ ) sets a baseline with a respectable Top-1 accuracy of 76.4% and Top-5 accuracy of 93.0%. Then we observe that the direct application of  $\mathcal{L}_{\text{OLE}}$  faced convergence issues. When we normalize the representation vectors, the training collapses. This indicates inherent challenges in applying OLE to SSL without suitable modifications. Decoupling the low-rank and high-rank enforcement, but without using the deviation matrix for intra-class alignment, helps to stabilize the training but has bad performance, and it showed high sensitivity to the hyperparameter  $\lambda$ . Specifically, training tended to collapse to constants for  $\lambda \leq 2.10$  and to random values for  $\lambda \geq 2.20$ . A temporary stabilization occurred at  $\lambda = 2.15$ , but the training loss eventually diverged after about 30 epochs, leading to SVD errors. This instability indicates the difficulties in directly applying low-rank constraints to the original embedding matrix.

Table 1: Studying adaptations of OLE for SSL using various strategies. We use 5 views per image. all models are trained for 100 epochs.

Objective	h.param.	Top-1	Top-5
$\mathcal{L}_{\text{InfoNCE-M}}$	$\tau = 0.2$	76.4	93.0
$\mathcal{L}_{\text{OLE}}$	-	failed to converge	
+ normalization	-	collapse	
+ loss decoupling	$\lambda = 2.15$	43.2	71.2
+ enhanced low-rank ( $\mathcal{L}_{\text{SSOLE}}$ )	$\lambda = 0.7$	<b>78.5</b>	<b>94.4</b>

Table 2: Impact of  $\lambda$  on linear probing top-1 Accuracy (%).

$\lambda$	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1.0	1.1	1.2	1.3	1.4	1.5	1.6	1.7	1.8	1.9	2.0	2.5	3.0
Acc.	72.1	73.9	76.5	77.5	78.4	78.5	78.5	78.4	78.5	78.0	78.3	78.6	78.3	78.3	78.1	78.3	78.4	78.3	78.4	78.0	77.3	77.3

In contrast, the  $\mathcal{L}_{\text{SSOLE}}$  loss with further enhanced low-rank enforcement demonstrated remarkable insensitivity to the value of  $\lambda$ , which was varied between 0.5 and 1.9, shown in Table 2. For smaller values of  $\lambda \in [0.1, 0.4]$ , the model prioritizes intra-class compactness, which can lead to under-penalized inter-class overlap, slightly degrading performance. For larger values of  $\lambda \in [2.0, 3.0]$ , the model emphasizes inter-class separability, sometimes at the expense of intra-class consistency, resulting in over-dispersed features. Optimal performance was achieved around  $\lambda = 0.7$ , with top-1 accuracy of 78.5% and top-5 accuracy of 94.4%. This relative robustness of  $\mathcal{L}_{\text{SSOLE}}$  highlights its suitability for SSL tasks and verifies the effectiveness of our approach in adapting OLE to the SSL framework.

More analysis of nuclear norm of matrices is available in Appendix B.1.



### 6.1.2 STUDIES ON THE NUMBER OF VIEWS

Understanding the impact of the number of views on model performance is crucial in SSL. As shown in Figure 1, both models exhibit an increase in Top-1 accuracy with more views, but SSOLE consistently outperforms InfoNCE-M. SSOLE’s performance notably improves up to 8 views, after which it plateaus, suggesting that the intrinsic rank of views from the same instance is likely less than 8. This plateau indicates diminishing returns beyond that point. SSOLE’s consistently superior performance highlights its effectiveness in utilizing additional views, while InfoNCE-M shows more modest gains, underscoring its relative inefficiency in leveraging extra views.

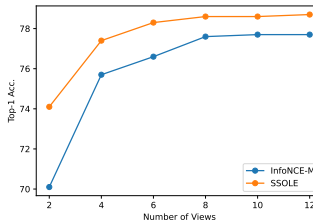


Figure 1: Comparison between InfoNCE-M and SSOLE with various numbers of views.

### 6.2 COMPARISON ON IMAGENET100

Table 3 provides a detailed comparison of various methods on the ImageNet100 dataset, using ResNet-18 as the backbone.

Our proposed method, SSOLE, achieves the highest Top-1 accuracy of 82.5%, outperforming established methods such as MoCo-M, SwAV, VICReg, BYOL, LORAC, and EMP-SSL. This result highlights SSOLE’s superior performance, especially considering its modest computational demands.

In contrast to methods like BYOL and VICReg, which require large batch sizes, SSOLE achieves higher performance with a batch size of only 128. Additionally, SSOLE employs a single-branch encoder, while methods like MoCo, BYOL, and LORAC use an EMA teacher model. **Lastly, unlike MoCo and LORAC, which rely on large memory banks, SSOLE avoids additional memory requirements.** We utilize 8 crops per image, consisting of 4 large and 4

Table 3: Comparative performance on ImageNet100. The table includes batch size (bs), number of epochs, number of crops, number of forward passes, and Top-1 accuracy (%).

Method	bs	#epochs	#crops	#forwards	Top-1
MoCo-M (Wang et al., 2022)	128	100	8	10	77.0
SwAV (Caron et al., 2020)	256	400	8	8	74.0
VICReg (Bardes et al., 2022)	2048	400	2	2	79.2
BYOL (Grill et al., 2020)	4096	400	2	4	80.2
LORAC (Wang et al., 2022)	128	100	8	10	78.7
EMP-SSL (Tong et al., 2023)	100	10	200	200	78.9
INTL (Weng et al., 2024)	128	400	2	2	81.7
SSOLE (Ours)	128	100	8	8	<b>82.5</b>

small crops. Details on our multi-crop strategy are provided in Appendix D.1. For a fair comparison, we train SSOLE for just 100 epochs, ensuring that the number of forward passes per iteration ( $\#forwards$ )  $\times$  the number of epochs ( $\#epochs$ ) is equal to or less than those of other methods.

### 6.3 EVALUATION ON FULL IMAGENET-1K

In this subsection, we assess the performance of our proposed SSOLE method when evaluated on the full ImageNet-1k dataset. We follow a similar evaluation protocol as used for ImageNet100, with the addition of semi-supervised learning settings where only a fraction of the labels are available.

#### 6.3.1 SELF-SUPERVISED AND SEMI-SUPERVISED PERFORMANCE

Table 4 showcases the performance of various methods on the full ImageNet dataset, with a special emphasis on the SSOLE method here proposed. The table reports both linear probing and semi-supervised learning accuracies, highlighting the efficacy of SSOLE in different learning regimes.

SSOLE demonstrates superior overall performance, particularly when compared to LORAC, which also leverages multiple views and employs low-rank embedding as a regularization prior. While LORAC achieves a respectable Top-1 accuracy of 73.2% and Top-5 accuracy of 91.6%, SSOLE surpasses it with a Top-1 accuracy of 73.9% and a Top-5 accuracy of 91.7%. Although SSOLE slightly underperforms INTL, which uses an EMA teacher, it is important to note that SSOLE does not rely on dual-branch encoders or EMA. Moreover, SSOLE significantly outperforms INTL when EMA is not applied. In semi-supervised learning, the advantage of SSOLE becomes even more pronounced. SSOLE outperforms LORAC by over 5 percentage points in the 1% labeled data setting and approximately 2 percentage points in the 10% labeled data setting. These improvements underscore the effectiveness of integrating OLE directly as a metric, rather than using it merely as a regularization term, as is the case with LORAC.

Table 4: Performance on full ImageNet of different methods. The table reports Top-1 and Top-5 accuracies (%) for linear probing and semi-supervised settings with 1% and 10% labeled data.

Method	BS	#Epochs	#Crops	#Forwards	Linear Probing		Semi-supervised (1%)		Semi-supervised (10%)	
					Top-1	Top-5	Top-1	Top-5	Top-1	Top-5
Supervised	-	-	-	-	76.5	-	25.4	56.4	48.4	80.4
MoCo v2 (Chen et al., 2020b)	256	200	2	2	71.1	-	-	-	-	-
SimCLR (Chen et al., 2020a)	4096	1000	2	2	69.3	89.0	48.3	75.5	65.6	87.8
SimSiam (Chen & He, 2021)	256	800	2	2	71.3	-	-	-	-	-
Barlow Twins (Zbontar et al., 2021)	2048	1000	2	2	73.2	91.0	55.0	79.2	69.7	89.3
VICReg (Bardes et al., 2022)	2048	1000	2	2	73.2	91.1	54.8	79.4	69.5	89.5
BYOL (Grill et al., 2020)	4096	400	2	4	73.2	-	-	-	-	-
SwAV (Caron et al., 2020)	256	200	8	10	72.7	91.5	49.6	76.1	67.7	88.7
LORAC (Wang et al., 2022)	256	200	8	10	73.2	91.6	50.0	76.3	68.0	88.9
MEC (Liu et al., 2022)	256	400	2	4	73.5	-	-	-	-	-
MatrixL-SSL (Zhang et al., 2024)	256	400	2	4	73.6	-	-	-	-	-
INTL (Weng et al., 2024)	512	800	2	2	73.1	-	55.0	<b>80.8</b>	69.4	89.8
INTL (EMA) (Weng et al., 2024)	256	800	2	2	<b>74.3</b>	-	-	-	-	-
SSOLE (Ours)	256	200	8	8	73.9	<b>91.7</b>	<b>55.4</b>	79.6	<b>70.3</b>	<b>90.3</b>

### 6.3.2 TRANSFERRING TO OTHER DATASETS

In this section, we evaluate the adaptability and robustness of the SSOLE framework through transfer learning. We apply the feature extractor trained on ImageNet and fine-tune on [MS-COCO \(Lin et al., 2015\)](#) for [object detection and instance segmentation tasks](#). We also evaluate transfer learning to linear classification tasks on datasets including CIFAR10 (Krizhevsky & Hinton, 2009), CIFAR100 (Krizhevsky & Hinton, 2009), Aircraft (Maji et al., 2013), DTD (Cimpoi et al., 2014), and Flowers (Nilsback & Zisserman, 2008), each offering unique image content and complexity challenges. This diverse set enables a comprehensive assessment of how SSOLE’s learned representations generalize across visual domains.

[Table 5 shows that SSOLE outperforms all other methods on COCO object detection and instance segmentation, highlighting its robustness and adaptability](#); [Table 6 illustrates SSOLE’s superior performance in transfer settings compared to state-of-the-art methods like MoCov2, SwAV, and LORAC](#). SSOLE achieves the highest accuracy on CIFAR10, Aircraft, DTD, and Flowers, indicating its capability to capture generalizable and robust features. SSOLE’s performance on CIFAR 10/100 highlights its ability to handle complex small-image classifications. Its success in the fine-grained classification on the Aircraft dataset and in diverse recognition tasks on DTD and Flowers further showcases its adaptability. Overall, SSOLE’s consistent effectiveness across these varied datasets attests to its versatility and efficiency as a feature extractor.

## 7 CONCLUSION

In this paper, we presented SSOLE, which integrates OLE into the SSL paradigm. Our method addresses two critical challenges in applying OLE to SSL: the difficulty of enforcing orthogonality in the presence of an infinite number of classes, and the nuclear norm’s inability to distinguish between positive and negative correlations. By decoupling low-rank and high-rank enforcement and applying low-rank constraints to feature deviations, SSOLE effectively adapts OLE for self-supervised tasks. Through extensive experiments, we demonstrated that SSOLE achieves state-of-the-art performance across linear probing, semi-supervised, and transfer learning tasks, all while maintaining computational efficiency. Notably, SSOLE achieves these results without relying on large batch sizes, memory banks, or complex architectures. SSOLE sets a new benchmark for integrating orthogonal low-rank representations into SSL, opening up promising directions for future research in SSL.

Table 5: Transfer Learning on object detection and instance segmentation on MS-COCO.

Model	Object Detection			Instance Segmentation		
	AP <sub>50</sub>	AP	AP <sub>75</sub>	AP <sub>50</sub> <sup>mk</sup>	AP <sup>mk</sup>	AP <sub>75</sub> <sup>mk</sup>
SimCLR	57.7	37.9	40.9	54.6	33.3	35.3
MoCo v2	58.9	39.3	42.5	55.8	34.4	36.5
BYOL	57.8	37.9	40.9	54.3	33.2	35.0
SwAV	58.6	38.4	41.3	55.2	33.8	35.9
SimSiam	59.3	39.2	42.1	56.0	34.4	36.7
Barlow Twins	59.0	39.2	42.5	56.0	34.3	36.5
VICReg	-	40.0	-	-	-	36.7
MEC	59.8	39.8	43.2	56.3	34.7	36.8
Matrix-SSL	60.8	41.0	44.2	57.5	35.6	38.0
INTL	61.0	41.0	44.5	57.7	35.6	37.8
SSOLE (Ours)	<b>61.5</b>	<b>41.3</b>	<b>44.8</b>	<b>58.0</b>	<b>35.9</b>	<b>38.4</b>

Table 6: Transfer Learning on linear classification on various datasets.

Method	CIFAR10	CIFAR100	Aircraft	DTD	Flowers
Supervised	90.0	73.4	42.6	68.8	89.7
MoCov2	89.8	71.0	39.3	69.2	87.4
SwAV	90.8	73.4	45.5	72.2	88.9
LORAC	91.8	<b>75.3</b>	47.8	72.7	89.5
SSOLE (Ours)	<b>92.2</b>	74.4	<b>48.0</b>	<b>73.3</b>	<b>90.0</b>

## REFERENCES

- 540  
541  
542 Nitin Bansal, Xiaohan Chen, and Zhangyang Wang. Can we gain more from orthogonality regu-  
543 larizations in training deep networks? *Advances in Neural Information Processing Systems*, 31,  
544 2018.
- 545 Adrien Bardes, Jean Ponce, and Yann LeCun. Vicreg: Variance-invariance-covariance regularization  
546 for self-supervised learning. In *International Conference on Learning Representations*, 2022.  
547
- 548 Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin.  
549 Unsupervised learning of visual features by contrasting cluster assignments. *Advances in Neural*  
550 *Information Processing Systems*, 33:9912–9924, 2020.
- 551 Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and  
552 Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the*  
553 *IEEE/CVF International Conference on Computer Vision*, pp. 9650–9660, 2021.  
554
- 555 Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for  
556 contrastive learning of visual representations. In *International conference on machine learning*, pp.  
557 1597–1607. PMLR, 2020a.
- 558 Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *Proceedings of*  
559 *the IEEE/CVF conference on computer vision and pattern recognition*, pp. 15750–15758, 2021.  
560
- 561 Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum  
562 contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020b.  
563
- 564 Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised vision  
565 transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp.  
566 9640–9649, 2021.
- 567 Mircea Cimpoi, Subhansu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. Describ-  
568 ing textures in the wild. In *Proceedings of the IEEE conference on computer vision and pattern*  
569 *recognition*, pp. 3606–3613, 2014.  
570
- 571 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale  
572 hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*,  
573 pp. 248–255. Ieee, 2009.
- 574 Aleksandr Ermolov, Aliaksandr Siarohin, Enver Sangineto, and Nicu Sebe. Whitening for self-  
575 supervised representation learning. In *International conference on machine learning*, pp. 3015–  
576 3024. PMLR, 2021. ISBN 2640-3498.  
577
- 578 Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena  
579 Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, and Mohammad Gheshlaghi Azar.  
580 Bootstrap your own latent—a new approach to self-supervised learning. *Advances in Neural*  
581 *Information Processing Systems*, 33:21271–21284, 2020.
- 582 Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The elements of statistical learning: data*  
583 *mining, inference, and prediction*, 2009.  
584
- 585 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
586 recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*,  
587 pp. 770–778, 2016.
- 588 Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for  
589 unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on*  
590 *computer vision and pattern recognition*, pp. 9729–9738, 2020.  
591
- 592 Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked  
593 autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer*  
*vision and pattern recognition*, pp. 16000–16009, November 2021.

- 594 Olivier Henaff. Data-efficient image recognition with contrastive predictive coding. In *Proceedings*  
595 *of the 37th International Conference on Machine Learning*, pp. 4182–4192. PMLR, November  
596 2020.
- 597 Ziheng Jiang, Ping Guo, and Lihong Peng. Locality-constrained low-rank coding for image classifi-  
598 cation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 28, 2014.
- 600 Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. 2009.
- 601 José Lezama, Qiang Qiu, and Guillermo Sapiro. Not afraid of the dark: Nir-vis face recognition via  
602 cross-spectral hallucination and low-rank embedding. In *Proceedings of the IEEE conference on*  
603 *computer vision and pattern recognition*, pp. 6628–6637, 2017.
- 605 José Lezama, Qiang Qiu, Pablo Musé, and Guillermo Sapiro. Ole: Orthogonal low-rank embedding-  
606 a plug and play geometric loss for deep learning. In *Proceedings of the IEEE Conference on*  
607 *Computer Vision and Pattern Recognition*, pp. 8109–8118, 2018.
- 608 Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays, Pietro  
609 Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollár. Microsoft coco: Common objects  
610 in context, February 2015.
- 612 Xin Liu, Zhongdao Wang, Ya-Li Li, and Shengjin Wang. Self-supervised learning via maximum  
613 entropy coding. *Advances in Neural Information Processing Systems*, 35:34091–34105, 2022.
- 614 Subhransu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. Fine-grained  
615 visual classification of aircraft. *arXiv preprint arXiv:1306.5151*, 2013.
- 616 Kanti V. Mardia and Peter E. Jupp. *Directional statistics*. John Wiley & Sons, 2009.
- 618 George Marsaglia. Choosing a point from the surface of a sphere. *The Annals of Mathematical*  
619 *Statistics*, 43(2):645–646, 1972.
- 620 Mervin E. Muller. A note on a method for generating points uniformly on  $n$ -dimensional spheres.  
621 *Communications of the ACM*, 2(4):19–20, April 1959. ISSN 0001-0782, 1557-7317. doi: 10.1145/  
622 377939.377946.
- 624 Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number  
625 of classes. In *2008 Sixth Indian conference on computer vision, graphics & image processing*, pp.  
626 722–729. IEEE, 2008.
- 627 Qiang Qiu and Guillermo Sapiro. Learning transformations for clustering and classification. *J. Mach.*  
628 *Learn. Res.*, 16(1):187–225, 2015.
- 630 Mark Rudelson and Roman Vershynin. Smallest singular value of a random rectangular matrix.  
631 *Communications on Pure and Applied Mathematics*, 62(12):1707–1739, December 2009. ISSN  
632 0010-3640, 1097-0312. doi: 10.1002/cpa.20294.
- 633 Yonglong Tian, Chen Sun, Ben Poole, Dilip Krishnan, Cordelia Schmid, and Phillip Isola. What  
634 makes for good views for contrastive learning? *Advances in neural information processing systems*,  
635 33:6827–6839, 2020.
- 636 Nenad Tomasev, Ioana Bica, Brian McWilliams, Lars Buesing, Razvan Pascanu, Charles Blundell,  
637 and Jovana Mitrovic. Pushing the limits of self-supervised resnets: Can we outperform supervised  
638 learning without labels on imagenet? *arXiv preprint arXiv:2201.05119*, 2022.
- 640 Shengbang Tong, Yubei Chen, Yi Ma, and Yann Lecun. Emp-ssl: Towards self-supervised learning  
641 in one training epoch. *arXiv preprint arXiv:2304.03977*, 2023.
- 642 Eugene Vorontsov, Chiheb Trabelsi, Samuel Kadoury, and Chris Pal. On orthogonality and learning  
643 recurrent networks with long term dependencies. In *International Conference on Machine Learning*,  
644 pp. 3570–3578. PMLR, 2017.
- 646 Guangrun Wang, Keze Wang, Guangcong Wang, Philip HS Torr, and Liang Lin. Solving ineffi-  
647 ciency of self-supervised representation learning. In *Proceedings of the IEEE/CVF International*  
*Conference on Computer Vision*, pp. 9505–9515, 2021.

- 648 Tongzhou Wang and Phillip Isola. Understanding contrastive representation learning through align-  
649 ment and uniformity on the hypersphere. In *International Conference on Machine Learning*, pp.  
650 9929–9939. PMLR, 2020. ISBN 2640-3498.
- 651 Yu Wang, Jingyang Lin, Qi Cai, Yingwei Pan, Ting Yao, Hongyang Chao, and Tao Mei. A low rank  
652 promoting prior for unsupervised contrastive learning. *IEEE Transactions on Pattern Analysis and*  
653 *Machine Intelligence*, 45(3):2667–2681, 2022. ISSN 0162-8828.
- 654 Ze Wang, Yue Lu, and Qiang Qiu. Meta-ole: Meta-learned orthogonal low-rank embedding. In  
655 *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 5305–  
656 5314, 2023.
- 657 Xi Weng, Yunhao Ni, Tengwei Song, Jie Luo, Rao Muhammad Anwer, Salman Khan, Fahad Shahbaz  
658 Khan, and Lei Huang. Modulate your spectrum in self-supervised learning. In *International*  
659 *Conference on Learning Representations*, 2024.
- 660 Zhirong Wu, Yuanjun Xiong, Stella X. Yu, and Dahua Lin. Unsupervised feature learning via  
661 non-parametric instance discrimination. In *Proceedings of the IEEE conference on computer vision*  
662 *and pattern recognition*, pp. 3733–3742, 2018.
- 663 Niannan Xue, Yannis Panagakis, and Stefanos Zafeiriou. Side information in robust principal  
664 component analysis: Algorithms and applications. In *Proceedings of the IEEE international*  
665 *conference on computer vision*, pp. 4317–4325, 2017.
- 666 Niannan Xue, Jiankang Deng, Shiyang Cheng, Yannis Panagakis, and Stefanos Zafeiriou. Side  
667 information for face completion: a robust pca approach. *IEEE transactions on pattern analysis*  
668 *and machine intelligence*, 41(10):2349–2364, 2019.
- 669 Jian Yang, Lei Luo, Jianjun Qian, Ying Tai, Fanlong Zhang, and Yong Xu. Nuclear norm based  
670 matrix regression with applications to face recognition with occlusion and illumination changes.  
671 *IEEE transactions on pattern analysis and machine intelligence*, 39(1):156–171, 2016.
- 672 Yaodong Yu, Kwan Ho Ryan Chan, Chong You, Chaobing Song, and Yi Ma. Learning diverse and  
673 discriminative representations via the principle of maximal coding rate reduction. In *Advances in*  
674 *Neural Information Processing Systems*, volume 33, pp. 9422–9434, 2020.
- 675 Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stéphane Deny. Barlow twins: Self-supervised  
676 learning via redundancy reduction. In *International Conference on Machine Learning*, pp. 12310–  
677 12320. PMLR, 2021. ISBN 2640-3498.
- 678 Liheng Zhang, Marzieh Edraki, and Guo-Jun Qi. Cappronet: Deep feature learning via orthogonal  
679 projections onto capsule subspaces. *Advances in neural information processing systems*, 31, 2018.
- 680 Yangmuzi Zhang, Zhuolin Jiang, and Larry S. Davis. Learning structured low-rank representations  
681 for image classification. In *Proceedings of the IEEE conference on computer vision and pattern*  
682 *recognition*, pp. 676–683, 2013.
- 683 Yifan Zhang, Zhiqian Tan, Jingqin Yang, Weiran Huang, and Yang Yuan. Matrix information theory  
684 for self-supervised learning. In *International Conference on Machine Learning*, 2024.
- 685 Zhao Zhang, Fanzhang Li, Mingbo Zhao, Li Zhang, and Shuicheng Yan. Joint low-rank and  
686 sparse principal feature coding for enhanced robust representation and visual classification. *IEEE*  
687 *Transactions on Image Processing*, 25(6):2429–2443, 2016.
- 688 Jinghao Zhou, Chen Wei, Huiyu Wang, Wei Shen, Cihang Xie, Alan Yuille, and Tao Kong. ibot: Image  
689 bert pre-training with online tokenizer. In *International Conference on Learning Representations*,  
690 2022.
- 691 Wei Zhu, Qiang Qiu, Bao Wang, Jianfeng Lu, Guillermo Sapiro, and Ingrid Daubechies. Stop  
692 memorizing: A data-dependent regularization framework for intrinsic pattern learning. *SIAM*  
693 *Journal on Mathematics of Data Science*, 1(3):476–496, January 2019. ISSN 2577-0187. doi:  
694 10.1137/19M1236886.



702 A PROOFS AND ANALYSIS

703  
704 A.1 PROOF FOR LEMMA 3.1

705  
706 *Proof.* We begin by centering the matrix  $\mathbf{V}$  to create a zero-mean matrix  $\tilde{\mathbf{V}}$  as follows:

707  
708 
$$\tilde{\mathbf{V}} = \mathbf{V} - \frac{s\mathbf{1}_d^T}{\sqrt{d}}.$$

709  
710 This gives rows of  $\tilde{\mathbf{V}}$  sampled from  $\mathcal{N}\left(0, \frac{s^2\sigma^2\mathbf{I}_d}{d}\right)$ . We first compute the nuclear norm for  $\tilde{\mathbf{V}}$  and  
711 then correct for the non-centered matrix.

712 **Cosine similarity of vectors.**

713  
714 Let  $\tilde{\mathbf{v}}_i$  and  $\tilde{\mathbf{v}}_j$  be two distinct rows of  $\tilde{\mathbf{V}}$ , sampled from  $\mathcal{N}\left(0, \frac{s^2\sigma^2\mathbf{I}_d}{d}\right)$ . The cosine similarity between  
715 these vectors is given by:

716  
717 
$$\cos(\theta_{i,j}) = \frac{\langle \tilde{\mathbf{v}}_i, \tilde{\mathbf{v}}_j \rangle}{\|\tilde{\mathbf{v}}_i\| \|\tilde{\mathbf{v}}_j\|} = \left\langle \frac{\tilde{\mathbf{v}}_i}{\|\tilde{\mathbf{v}}_i\|}, \frac{\tilde{\mathbf{v}}_j}{\|\tilde{\mathbf{v}}_j\|} \right\rangle.$$

718  
719 Then both  $\frac{\tilde{\mathbf{v}}_i}{\|\tilde{\mathbf{v}}_i\|}$  and  $\frac{\tilde{\mathbf{v}}_j}{\|\tilde{\mathbf{v}}_j\|}$  follow a uniform probability distribution over the unit sphere  $S^{d-1}$  (Muller,  
720 1959; Marsaglia, 1972).

721  
722 Mardia & Jupp (2009) shows that  $\frac{1+\cos(\theta_{i,j})}{2}$  follows a Beta distribution,  $\text{Beta}\left(\frac{d-1}{2}, \frac{d+1}{2}\right)$ , and

723  
724 
$$\mathbb{E}[\cos(\theta_{i,j})] = 0,$$
  
725  
726 
$$\text{Var}[\cos(\theta_{i,j})] = \frac{1}{d}.$$

727  
728 Using Chebyshev's inequality, we bound the tail probability of  $|\cos(\theta_{i,j})|$  for a given threshold  $\frac{1}{d}$ :

729  
730 
$$P\left(|\cos(\theta_{i,j})| > \frac{C}{\sqrt{d}}\right) \leq \frac{1}{C^2}, \quad (9)$$

731  
732 where  $C$  is a constant and satisfies  $1 \ll C \ll \sqrt{d}$ .

733 **Expected nuclear norm for  $m \ll d$ .**

734  
735 The nuclear norm of  $\tilde{\mathbf{V}}$  depends on the norm of each row and cosine similarity between them.

736  
737 When the rows of  $\tilde{\mathbf{V}}$  are perfectly orthogonal,  $\tilde{\mathbf{V}}$  achieves the maximum nuclear norm for given row  
738 norms, then  $\tilde{\mathbf{V}}\tilde{\mathbf{V}}^T$  is an  $m \times m$  diagonal matrix where the  $i^{\text{th}}$  diagonal element is  $\tilde{\mathbf{v}}_i\tilde{\mathbf{v}}_i^T = \|\tilde{\mathbf{v}}_i\|^2$ .  
739 Then for  $i = 1, \dots, m$ ,  $\tilde{\mathbf{V}}\tilde{\mathbf{V}}^T$  has an eigenvalue of  $\|\tilde{\mathbf{v}}_i\|^2$ ;  $\tilde{\mathbf{V}}$  has a corresponding singular value of  
740  $\|\tilde{\mathbf{v}}_i\|$ . This determines the upper bound for the expected nuclear norm:

741  
742 
$$\mathbb{E}[\|\tilde{\mathbf{V}}\|_*] \leq m \cdot \mathbb{E}[\|\tilde{\mathbf{v}}\|].$$

743  
744 Since  $\tilde{\mathbf{v}}$  follows  $\mathcal{N}\left(0, \frac{s^2\sigma^2\mathbf{I}_d}{d}\right)$ ,  $\|\frac{\sqrt{d}}{s\sigma}\tilde{\mathbf{v}}\|$  follows the chi distribution with  $d$  degrees of freedom, *i.e.*  
745  $\mathcal{X}_d$ . Therefore,

746  
747 
$$\mathbb{E}\left[\left\|\frac{\sqrt{d}}{s\sigma}\tilde{\mathbf{v}}\right\|\right] = \sqrt{2}\frac{\Gamma\left(\frac{d+1}{2}\right)}{\Gamma\left(\frac{d}{2}\right)},$$
  
748  
749 
$$\text{Var}\left[\left\|\frac{\sqrt{d}}{s\sigma}\tilde{\mathbf{v}}\right\|\right] = d - \mathbb{E}^2[\|\tilde{\mathbf{v}}\|].$$

750  
751  
752 Using the recurrence relation for the Gamma function and Stirling's approximation for large  $d$ , we  
753 can bound  $\mathbb{E}\left[\left\|\frac{\sqrt{d}}{s\sigma}\tilde{\mathbf{v}}\right\|\right]$  as:

754  
755 
$$\mathbb{E}\left[\left\|\frac{\sqrt{d}}{s\sigma}\tilde{\mathbf{v}}\right\|\right] = \sqrt{d}\left(1 - \frac{1}{4d} + \mathcal{O}\left(\frac{1}{d^2}\right)\right) = \sqrt{d} + \mathcal{O}\left(\frac{1}{\sqrt{d}}\right).$$

Further,

$$\text{Var}[\|\frac{\sqrt{d}}{s\sigma}\tilde{\mathbf{v}}\|] = d - \left(d - \frac{1}{2} + \mathcal{O}\left(\frac{1}{d}\right)\right) = \frac{1}{2} + \mathcal{O}\left(\frac{1}{d}\right).$$

Thus,

$$\mathbb{E}[\|\tilde{\mathbf{v}}\|] = s\sigma + \mathcal{O}\left(\frac{1}{d}\right), \quad (10)$$

$$\text{Var}[\|\tilde{\mathbf{v}}\|] = \frac{s^2\sigma^2}{2d} + \mathcal{O}\left(\frac{1}{d^2}\right). \quad (11)$$

Therefore,

$$\mathbb{E}[\|\tilde{\mathbf{V}}\|_*] \leq s \cdot \sigma \cdot m + \mathcal{O}\left(\frac{1}{d}\right). \quad (12)$$

Now, we derive the lower bound.  $\|\tilde{\mathbf{V}}\|_*$  negatively relates to the absolute cosine similarity between its rows. Applying Equation (9), we derive the below probability:

$$P\left(\max_{i,j,i \neq j} |\cos(\theta_{i,j})| \leq \frac{C}{\sqrt{d}}\right) \geq \left(1 - \frac{1}{C^2}\right)^{\frac{m(m-1)}{2}} = 1 - \mathcal{O}\left(\frac{1}{C^2}\right). \quad (13)$$

On the other hand,  $\|\tilde{\mathbf{V}}\|_*$  positively relates to the norm of its rows.

Using Chebyshev's inequality with Equation (10) and Equation (11), we have:

$$P\left(\|\tilde{\mathbf{v}}\| - s \cdot \sigma > \frac{C \cdot s \cdot \sigma}{\sqrt{d}}\right) \leq \frac{1}{2C^2} + \mathcal{O}\left(\frac{1}{d^2}\right). \quad (14)$$

Further,

$$P\left(\min_i \|\tilde{\mathbf{v}}\| \geq s \cdot \sigma \cdot \left(1 - \frac{C}{\sqrt{d}}\right)\right) \geq \left(1 - \frac{1}{2C^2} + \mathcal{O}\left(\frac{1}{d^2}\right)\right)^m = 1 - \mathcal{O}\left(\frac{1}{C^2}\right). \quad (15)$$

Combining Equation (13) and Equation (15), we have: with probability at least  $(1 - \mathcal{O}(\frac{1}{C^2})) \cdot (1 - \mathcal{O}(\frac{1}{C^2})) = 1 - \mathcal{O}(\frac{1}{C^2})$ , every two distinct rows has absolute cosine similarity at most  $\frac{C}{\sqrt{d}}$ , and each row has norm at least  $s \cdot \sigma \cdot \left(1 - \frac{C}{\sqrt{d}}\right)$ .

When every two distinct rows of  $\tilde{\mathbf{V}}$  has cosine similarity of  $\frac{C}{\sqrt{d}}$ , and each row has norm of  $s \cdot \sigma \cdot \left(1 - \frac{C}{\sqrt{d}}\right)$ , then its nuclear norm is:

$$\begin{aligned} & s\sigma \left(1 - \frac{C}{\sqrt{d}}\right) \left(\sqrt{1 + (m-1)\frac{C}{\sqrt{d}}} + (m-1)\sqrt{1 - \frac{C}{\sqrt{d}}}\right) \\ &= s\sigma \left(1 - \frac{C}{\sqrt{d}}\right) \left(1 + \frac{1}{2}(m-1)\frac{C}{\sqrt{d}} + \mathcal{O}\left((m-1)^2\frac{C^2}{d}\right) + (m-1)\left(1 - \frac{C}{2\sqrt{d}} + \mathcal{O}\left(\frac{C^2}{d}\right)\right)\right) \\ &= s\sigma \left(1 - \frac{C}{\sqrt{d}}\right) \left(m + \mathcal{O}\left(\frac{C^2}{d}\right)\right) \\ &= s \cdot \sigma \cdot m - \mathcal{O}\left(\frac{C}{\sqrt{d}}\right). \end{aligned}$$

Therefore,

$$\mathbb{E}[\|\tilde{\mathbf{V}}\|_*] \geq \left(1 - \mathcal{O}\left(\frac{1}{C^2}\right)\right) \cdot \left(s \cdot \sigma \cdot m - \mathcal{O}\left(\frac{C}{\sqrt{d}}\right)\right) = s \cdot \sigma \cdot m - \mathcal{O}\left(\frac{C}{\sqrt{d}}\right) - \mathcal{O}\left(\frac{1}{C^2}\right). \quad (16)$$

Combining Equation (12) and Equation (16), we obtain:

$$\mathbb{E}[\|\tilde{\mathbf{V}}\|_*] = \mathcal{O}(s \cdot \sigma \cdot m). \quad (17)$$

810 **Expected nuclear norm for  $m \gg d$ .**

811 When  $m \gg d$ ,  $\tilde{\mathbf{V}}$  is full rank, with  $d$  singular values approximately uniformly distributed.

812 The singular values of  $\tilde{\mathbf{V}}$ , denoted as  $\lambda_i$ , for  $i = 1, \dots, d$ , are the square roots of the eigenvalues of  
813  $\tilde{\mathbf{V}}^T \tilde{\mathbf{V}}$ .

814 Since the trace of  $\tilde{\mathbf{V}}^T \tilde{\mathbf{V}}$  is the sum of all its eigenvalues, we have:

$$815 \sum_{i=1}^d \lambda_i^2 = \text{Tr}(\tilde{\mathbf{V}}^T \tilde{\mathbf{V}}) = \sum_{j=1}^m \|\tilde{\mathbf{v}}_j\|^2. \quad (18)$$

816 Therefore,

$$817 \sum_{i=1}^d \lambda_i \leq d \cdot \sqrt{\frac{\sum_{i=1}^d \lambda_i^2}{d}} = \sqrt{d \cdot \sum_{j=1}^m \|\tilde{\mathbf{v}}_j\|^2}. \quad (19)$$

818 Since  $\frac{\sqrt{d \cdot \sum_{j=1}^m \|\tilde{\mathbf{v}}_j\|^2}}{s\sigma}$  follows a chi distribution with  $md$  degrees of freedom, i.e.  $\mathcal{X}_{md}$ ,

$$819 \mathbb{E}\left[\sqrt{d \cdot \sum_{j=1}^m \|\tilde{\mathbf{v}}_j\|^2}\right] = s \cdot \sigma \cdot \sqrt{md} - \mathcal{O}\left(\frac{1}{\sqrt{md}}\right). \quad (20)$$

820 Thus,

$$821 \mathbb{E}[\|\tilde{\mathbf{V}}\|_*] = \mathbb{E}\left[\sum_{i=1}^d \lambda_i\right] \leq s \cdot \sigma \cdot \sqrt{md} - \mathcal{O}\left(\frac{1}{\sqrt{md}}\right). \quad (21)$$

822 On the other hand, Rudelson & Vershynin (2009) show that:

$$823 P\left(\min_i \frac{\sqrt{d}\lambda_i}{s\sigma} \leq \sqrt{m} - \sqrt{d} - t\right) \leq e^{-\frac{t^2}{2}}, \quad t > 0. \quad (22)$$

824 Let  $t = \sqrt{2 \log C}$ , where  $1 \ll C \ll e^m$ . Then, with probability at least  $1 - \frac{1}{C}$ :

$$825 \|\tilde{\mathbf{V}}\|_* \geq d \cdot \frac{s\sigma}{\sqrt{d}} \left(\sqrt{m} - \sqrt{d} - \sqrt{2 \log C}\right) = s \cdot \sigma \left(\sqrt{md} - \mathcal{O}(d)\right). \quad (23)$$

826 Therefore,

$$827 \mathbb{E}[\|\tilde{\mathbf{V}}\|_*] \geq \left(1 - \frac{1}{C}\right) \cdot s \cdot \sigma \left(\sqrt{md} - \mathcal{O}(d)\right) = s \cdot \sigma \cdot \sqrt{md} - \mathcal{O}(d) - \mathcal{O}\left(\frac{1}{C}\right). \quad (24)$$

828 Combining Equation (21) and Equation (24), we obtain:

$$829 \mathbb{E}[\|\tilde{\mathbf{V}}\|_*] = \mathcal{O}\left(s \cdot \sigma \cdot \sqrt{md}\right). \quad (25)$$

830 **Correction for non-centered  $\mathbf{V}$ .**

831 Finally, for the non-centered matrix  $\mathbf{V}$ , we have  $\|\mathbf{V}\|_* \leq \|\tilde{\mathbf{V}}\|_* + \sqrt{m} \|\frac{s\mathbf{1}_d^T}{\sqrt{d}}\|_*$ ; and  $\|\mathbf{V}\|_* \geq$   
832  $\|\tilde{\mathbf{V}}\|_*$ ,  $\|\mathbf{V}\|_* \geq \sqrt{m} \|\frac{s\mathbf{1}_d^T}{\sqrt{d}}\|_*$ . When  $\sigma$  is large,  $\|\mathbf{V}\|_*$  is dominated by  $\|\tilde{\mathbf{V}}\|_*$ ; whereas when  $\sigma$  is  
833 small,  $\|\mathbf{V}\|_*$  is dominated by  $\sqrt{m} \|\frac{s\mathbf{1}_d^T}{\sqrt{d}}\|_*$ .

834 We add a correction term of  $\mathcal{O}(s \cdot \sqrt{m})$  to  $\|\tilde{\mathbf{V}}\|_*$ , yielding:

$$835 \mathbb{E}[\|\mathbf{V}\|_*] = \begin{cases} s \cdot (\sigma \cdot \mathcal{O}(m) + \mathcal{O}(\sqrt{m})), & \text{if } m \ll d, \\ s \cdot (\sigma \cdot \mathcal{O}(\sqrt{md}) + \mathcal{O}(\sqrt{m})), & \text{if } m \gg d. \end{cases}$$

836  $\square$

## 864 A.2 PROOF FOR THEOREM 3.2

865 *Proof.* The OLE loss is based on the difference between intra-class compactness and inter-class  
866 separation. For each class with  $N$  augmentations, the nuclear norm of the intra-class representations  
867 follows the behavior for  $m \ll d$ :

$$868 \mathbb{E}[\|\mathbf{V}_{\text{intra}}\|_*] = s \cdot \left( \sigma \mathcal{O}(N) + \mathcal{O}(\sqrt{N}) \right).$$

871 For  $B$  classes, we sum over the nuclear norms:

$$872 \mathcal{L}_{\text{intra}} = s \cdot B \left( \sigma \mathcal{O}(N) + \mathcal{O}(\sqrt{N}) \right).$$

874 For inter-class separation across all  $BN$  samples, the nuclear norm follows the behavior for  $m \gg d$ :

$$875 \mathbb{E}[\|\mathbf{V}_{\text{all}}\|_*] = s \cdot \left( \sigma \mathcal{O}(\sqrt{BNd}) + \mathcal{O}(\sqrt{BN}) \right).$$

878 Thus, the total OLE loss is:

$$\begin{aligned} 879 \mathcal{L}_{OLE} &= s \cdot B \left( \sigma \mathcal{O}(N) + \mathcal{O}(\sqrt{N}) \right) - s \cdot \left( \sigma \mathcal{O}(\sqrt{BNd}) + \mathcal{O}(\sqrt{BN}) \right) \\ 880 &= s\sqrt{BN} \left( \sigma \cdot \left( \mathcal{O}(\sqrt{BN}) - \mathcal{O}(\sqrt{d}) \right) + \mathcal{O}(\sqrt{B}) - 1 \right) \\ 881 &= s\sqrt{BN} \left( \sigma \cdot \mathcal{O}(\sqrt{BN}) + \mathcal{O}(\sqrt{B}) \right). \end{aligned}$$

885  $\square$

## 886 A.3 ANALYSIS OF RANDOM GAUSSIAN MATRIX

888 To analyze the structure of the matrix  $\mathbf{V}$  in Lemma 3.1, we start by examining the length of the rows  
889  $\mathbf{v}_i$ . Each row of  $\mathbf{V}$  is sampled from  $\mathcal{N}\left(\frac{s\mathbf{1}_d}{\sqrt{d}}, \frac{s^2\sigma^2\mathbf{I}_d}{d}\right)$ , so we first compute the expected norm of  $\mathbf{v}_i$ .

### 891 **Expected length of each vector:**

892 The squared norm of a row  $\mathbf{v}_i$  is:

$$893 \mathbb{E}[\|\mathbf{v}_i\|^2] = \mathbb{E}\left[\left\|\frac{s}{\sqrt{d}}\mathbf{1}_d + \tilde{\mathbf{v}}_i\right\|^2\right].$$

897 Expanding this, we get:

$$898 \mathbb{E}[\|\mathbf{v}_i\|^2] = s^2 + s^2\sigma^2,$$

899 where  $\mathbb{E}[\|\tilde{\mathbf{v}}_i\|^2] = s^2\sigma^2$  due to the variance of the Gaussian distribution. Thus, the expected norm is:

$$900 \mathbb{E}[\|\mathbf{v}_i\|] = s\sqrt{1 + \sigma^2}.$$

### 902 **Cosine similarity between two vectors:**

903 The cosine similarity between two rows  $\mathbf{v}_i$  and  $\mathbf{v}_j$  is:

$$904 \cos(\theta_{i,j}) = \frac{\langle \mathbf{v}_i, \mathbf{v}_j \rangle}{\|\mathbf{v}_i\| \|\mathbf{v}_j\|}.$$

908 Since  $\mathbf{v}_i$  and  $\mathbf{v}_j$  are independent, we have:

$$909 \mathbb{E}[\langle \mathbf{v}_i, \mathbf{v}_j \rangle] = s^2,$$

911 and the expected cosine similarity is:

$$912 \mathbb{E}[\cos(\theta_{i,j})] = \frac{s^2}{s^2(1 + \sigma^2)} = \frac{1}{1 + \sigma^2}.$$

913 For large  $\sigma$ , the cosine similarity approaches zero and  $\theta$  approaches  $\frac{\pi}{2}$ , meaning the vectors are  
914 approximately orthogonal; for small  $\sigma$ , the cosine similarity approaches one and  $\theta$  approaches 0,  
915 meaning the vectors are nearly identical. Both  $\sigma$  and  $s$  control the vector length, and are positively  
916 correlated.

#### 918 A.4 DISCUSSION OF NUCLEAR NORM'S UNITARY INVARIANCE

919  
920 The nuclear norm of a matrix  $\mathbf{A}$ , denoted as  $\|\mathbf{A}\|_*$ , is defined as the sum of its singular values. The  
921 singular values of a matrix  $\mathbf{A}$  are the square roots of the eigenvalues of  $\mathbf{A}^\top \mathbf{A}$ . Therefore, for any  
922 matrix  $\mathbf{V} \in \mathbb{R}^{m \times d}$ , the nuclear norm is given by:

$$923 \quad \|\mathbf{V}\|_* = \sum_{i=1}^r \sigma_i(\mathbf{V}),$$

924 where  $\sigma_i(\mathbf{V})$  are the singular values of  $\mathbf{V}$ , and  $r$  is the rank of  $\mathbf{V}$ .

925  
926 Now, consider the matrix  $\mathbf{P}\mathbf{V}$ , where  $\mathbf{P}$  is a diagonal matrix with diagonal elements  $\pm 1$ , which flips  
927 the signs of the rows of  $\mathbf{V}$ . To compute the nuclear norm  $\|\mathbf{P}\mathbf{V}\|_*$ , we need to determine its singular  
928 values. The singular values are the square roots of the eigenvalues of  $(\mathbf{P}\mathbf{V})^\top (\mathbf{P}\mathbf{V})$ :

$$929 \quad (\mathbf{P}\mathbf{V})^\top (\mathbf{P}\mathbf{V}) = \mathbf{V}^\top \mathbf{P}^\top \mathbf{P} \mathbf{V}.$$

930 Since  $\mathbf{P}$  is a diagonal matrix with  $\pm 1$  entries, we have  $\mathbf{P}^\top \mathbf{P} = \mathbf{I}_m$ , where  $\mathbf{I}_m$  is the identity matrix  
931 of size  $m$ . Thus, the expression simplifies to:

$$932 \quad (\mathbf{P}\mathbf{V})^\top (\mathbf{P}\mathbf{V}) = \mathbf{V}^\top \mathbf{V}.$$

933 This shows that the matrix  $\mathbf{P}\mathbf{V}$  has the same singular values as  $\mathbf{V}$ , because the eigenvalues of  $\mathbf{V}^\top \mathbf{V}$   
934 are unchanged by the multiplication with  $\mathbf{P}$ .

935 Therefore, the nuclear norm of  $\mathbf{P}\mathbf{V}$  is equal to the nuclear norm of  $\mathbf{V}$ :

$$936 \quad \|\mathbf{P}\mathbf{V}\|_* = \|\mathbf{V}\|_*.$$

#### 937 A.5 PROOF FOR THEOREM 4.1

938  
939 *Proof.* The nuclear norm of  $\mathbf{V}$  is defined as the sum of the singular values of  $\mathbf{V}$ , denoted  
940  $\sigma_1, \sigma_2, \dots, \sigma_N$ , so that:

$$941 \quad \|\mathbf{V}\|_* = \sum_{i=1}^N \sigma_i.$$

942 We will use the properties of the singular values to establish bounds on  $\|\mathbf{V}\|_*$ .

943 First, we know that the sum of the squared singular values is equal to the Frobenius norm of  $\mathbf{V}$ ,  
944 which is:

$$945 \quad \sum_{i=1}^N \sigma_i^2 = \|\mathbf{V}\|_F^2 = N.$$

946 This is because each column of  $\mathbf{V}$  is a unit vector, so the total squared length of the matrix is  $N$ .

947 We know that the square of singular values of  $\mathbf{V}$  are eigenvalues of  $\mathbf{V}^\top \mathbf{V}$ , then we have

$$948 \quad \sum_{i=1}^N \sigma_i^4 = \|\mathbf{V}^\top \mathbf{V}\|_F^2 = N^2 \cdot \overline{\cos^2(\theta)}.$$

949 Combining these, we have

$$950 \quad \begin{cases} \|\mathbf{V}\|_* &= \sum_{i=1}^N \sigma_i \\ N &= \sum_{i=1}^N \sigma_i^2 \\ N^2 \cdot \overline{\cos^2(\theta)} &= \sum_{i=1}^N \sigma_i^4 \end{cases}$$

951 Since  $f_1(x) = \sqrt{x}$  is a concave function, and  $f_2(x) = x^2$  is a convex function;  $\|\mathbf{V}\|_* = \sum_{i=1}^N f_1(\sigma_i^2)$   
952 is negatively related to  $N^2 \cdot \overline{\cos^2(\theta)} = \sum_{i=1}^N f_2(\sigma_i^2)$  with fixed  $\sum_{i=1}^N \sigma_i^2$ .

953 To find the lower and upper bounds of  $\|\mathbf{V}\|_*$ . We need to get the max and min value of  $\overline{\cos^2(\theta)}$ .

954 **Minimization of  $\overline{\cos^2(\theta)}$ :** Using the Cauchy-Schwarz inequality, we have the following chain of  
955 inequalities for  $\cos(\theta_{i,j})$  where  $i \neq j$ :

$$956 \quad \overline{\cos^2(\theta_{i,j})} \geq |\overline{\cos(\theta_{i,j})}|^2 \geq \overline{\cos(\theta_{i,j})}^2.$$



Since  $\cos(\theta_{i,i}) = 1$  for all  $i$ , we focus on the off-diagonal elements  $\cos(\theta_{i,j})$  for  $i \neq j$ . The value of  $\overline{\cos^2(\theta)}$  is minimized for a fixed  $\overline{\cos(\theta)}$  when all the off-diagonal cosine similarities  $\cos(\theta_{i,j})$  are identical, i.e., when the vectors form an equally spaced configuration in high-dimensional space.

For such a configuration, the matrix  $\mathbf{V}^T \mathbf{V}$  has one eigenvalue corresponding to the eigenvector  $\mathbf{1} = (1, 1, \dots, 1)^\top$ , with eigenvalue  $\sqrt{N \overline{\cos(\theta)}}$ . The remaining  $N - 1$  eigenvalues have the same value  $\frac{\sqrt{N - N \overline{\cos(\theta)}}}{N-1}$  and correspond to eigenvectors orthogonal to  $\mathbf{1}$ .

Therefore,

$$\|\mathbf{V}\|_* \leq \sqrt{N} \cdot \sqrt{\overline{\cos(\theta)}} + \sqrt{N(N-1)} \cdot \sqrt{1 - \overline{\cos(\theta)}}.$$

**Maximization of  $\overline{\cos^2(\theta)}$ :** We note that  $\cos^2(\theta_{i,j}) \leq |\cos(\theta_{i,j})|$ . Equality holds when  $|\cos(\theta_{i,j})|$  is either 1 or 0, meaning that the vectors are either perfectly aligned, completely anti-aligned, or orthogonal to each other. This scenario occurs when the vectors can be grouped into  $G$  distinct groups, such that within each group, the vectors are either aligned or anti-aligned, while vectors from different groups are orthogonal.

Using this setup, the matrix  $\mathbf{V}^T \mathbf{V}$  will have  $G$  non-zero singular values, each corresponding to a group of vectors. To maximize the nuclear norm, we apply the Cauchy-Schwarz inequality to obtain:

$$\|\mathbf{V}\|_* \geq \sqrt{GN},$$

where  $G$  represents the number of groups.

**Minimizing  $G$ :** We now need to find the minimum possible value of  $G$ . The minimum value of  $G$  occurs when the vectors are partitioned into the fewest number of groups, while ensuring that the vectors within each group are either aligned or anti-aligned, and vectors between different groups are orthogonal. The smallest  $G$  is given by:

$$\frac{N^2 \cdot \overline{\cos^2(\theta)}}{G} \geq \left(\frac{N}{G}\right)^2 \implies G \geq \frac{1}{|\overline{\cos(\theta)}|}.$$

Substituting this into the earlier inequality for the nuclear norm, we obtain:

$$\|\mathbf{V}\|_* \geq \sqrt{\frac{N}{|\overline{\cos(\theta)}|}}.$$

Combining the results from the maximum and minimum cases, we obtain the desired bounds for the nuclear norm:

$$\sqrt{\frac{N}{|\overline{\cos(\theta)}|}} \leq \|\mathbf{V}\|_* \leq \sqrt{N} \cdot \sqrt{\overline{\cos(\theta)}} + \sqrt{N(N-1)} \cdot \sqrt{1 - \overline{\cos(\theta)}}.$$

□

## A.6 PROOF FOR THEOREM 4.2

*Proof.* Let  $\mathbf{V} \in \mathbb{R}^{d \times N}$  be a matrix where the  $i$ -th column is the unit vector  $\mathbf{v}_i \in \mathbb{R}^d$ , and define  $\bar{\mathbf{v}} := \frac{1}{N} \sum_{i=1}^N \mathbf{v}_i$  as the mean vector. The deviation matrix is defined as  $\tilde{\mathbf{V}} = \mathbf{V} - \bar{\mathbf{v}} \mathbf{1}^\top$ , where  $\mathbf{1} \in \mathbb{R}^N$  is a vector of ones.

The Frobenius norm of the deviation matrix is:

$$\|\tilde{\mathbf{V}}\|_F^2 = \sum_{i=1}^N \|\mathbf{v}_i - \bar{\mathbf{v}}\|^2.$$

For each vector  $\mathbf{v}_i$ , we can expand  $\|\mathbf{v}_i - \bar{\mathbf{v}}\|^2$  as:

$$\|\mathbf{v}_i - \bar{\mathbf{v}}\|^2 = 1 - 2\mathbf{v}_i^\top \bar{\mathbf{v}} + \|\bar{\mathbf{v}}\|^2.$$

We calculate  $\|\bar{v}\|^2$  as:

$$\|\bar{v}\|^2 = \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N \mathbf{v}_i^\top \mathbf{v}_j = \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N \cos(\theta_{i,j}),$$

where  $\theta_{i,j}$  is the angle between vectors  $\mathbf{v}_i$  and  $\mathbf{v}_j$ . Thus, the Frobenius norm of the deviation matrix is:

$$\|\tilde{\mathbf{V}}\|_F^2 = N \left(1 - \overline{\cos(\theta)}\right),$$

where  $\overline{\cos(\theta)}$  is the average cosine similarity.

The nuclear norm is the sum of the singular values of the matrix. The sum of the squared singular values equals the Frobenius norm,  $\|\tilde{\mathbf{V}}\|_F^2 = N(1 - \overline{\cos(\theta)})$ .

The nuclear norm is bounded between two extremes. The lower bound occurs when the vectors can be grouped into two perfectly aligned sets. In this case, the deviation vectors are either identical or the opposite. The deviation from the mean will be minimized, and the nuclear norm satisfies:

$$\|\tilde{\mathbf{V}}\|_* \geq \sqrt{N} \cdot \sqrt{1 - \overline{\cos(\theta)}}.$$

The upper bound is attained when the angles between all vectors are identical. In this configuration, the deviation vectors are orthogonal to each other, yielding:

$$\|\tilde{\mathbf{V}}\|_* \leq \sqrt{N(N-1)} \cdot \sqrt{1 - \overline{\cos(\theta)}}.$$

Therefore, the nuclear norm of the deviation matrix is bounded by:

$$\sqrt{N} \cdot \sqrt{1 - \overline{\cos(\theta)}} \leq \|\tilde{\mathbf{V}}\|_* \leq \sqrt{N(N-1)} \cdot \sqrt{1 - \overline{\cos(\theta)}}.$$

□

## B ADDITIONAL ANALYSIS

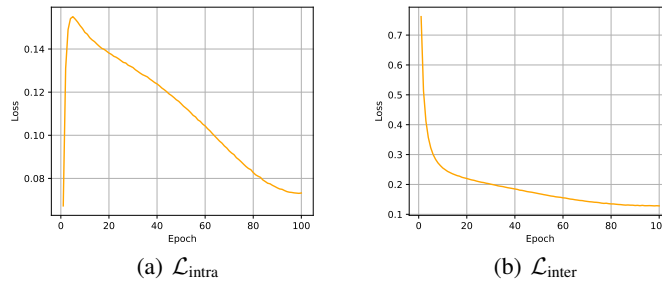


Figure 2:  $\mathcal{L}_{\text{intra}}$  and  $\mathcal{L}_{\text{inter}}$  throughout the training process of SSOLE trained on the ImageNet100 dataset. Batch size is 128, 5 views for each instances are used, and  $\lambda = 0.7$ . The model is trained for 100 epochs.

### B.1 NUCLEAR NORMS OF MATRIX

Figure 2 and Figure 3 provide a visual representation of the evolution of different losses and nuclear norms during the training of SSOLE on the ImageNet100 dataset. Analyzing the training dynamics depicted in Figure 2, we observe a pronounced initial phase where  $\mathcal{L}_{\text{intra}}$  sharply rises and  $\mathcal{L}_{\text{inter}}$  drastically falls. This abrupt behavior at the onset of training likely indicates a quick shift away from a suboptimal starting point where the model’s representations were converging towards a

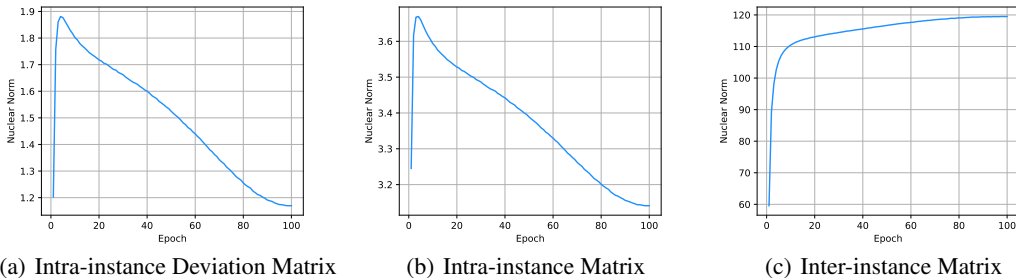


Figure 3: Nuclear norms of different matrix throughout the training process of SSOLE trained on the ImageNet100 dataset. Batch size is 128, 5 views for each instances are used, and  $\lambda = 0.7$ . The model is trained for 100 epochs.

trivial solution — one that maps disparate instances to indistinguishable points in the embedding space. As the training progresses past this phase, both  $L_{intra}$  and  $L_{inter}$  settle into a steady decline, suggesting that the model begins to develop more differentiated and stable representations, effectively separating instances in the embedding space, which is a desirable behavior in SSL to avoid collapse to uninformative features.

Corresponding to the changes in  $L_{intra}$  and  $L_{inter}$ , in Figure 3, the nuclear norms of the Intra-instance Deviation Matrix and the Intra-instance Matrix initially exhibit a sharp increase, followed by a steady descent. This pattern indicates an alignment with the loss behavior, reinforcing that the model is refining its internal representation to reduce intra-class variance while increasing inter-class separation. The bounded nature of these norms, with the Intra-instance Deviation Matrix approaching zero and the Intra-instance Matrix lower-bounded by  $\sqrt{5}$ , reflects a constrained yet effective optimization towards a low-rank within-class structure and a high-rank between-class structure. The Inter-instance Matrix’s nuclear norm, which aligns with  $L_{inter}$ , shows a sustained increase post-initial adjustment, capped by the batch size of 128. This increasing trend is a positive indication that the model effectively diversifies the feature space across instances, achieving the high-rank separation necessary for robust SSL.

In summary, these observations throughout the training process suggest that SSOLE not only quickly corrects initial misalignments but also steadily matures towards achieving its goal of a structured embedding space. The consistent optimization patterns, combined with the nuclear norm bounds, demonstrate SSOLE’s capability to navigate the challenges of MV-SSL with stability and efficacy.

### C LIMITATIONS

While this study introduces SSOLE as a robust and efficient training objective for Multi-view Self-Supervised Learning (MV-SSL), integrating Orthogonal Low-rank Embedding (OLE) concepts into SSL, it does have certain limitations. Our experiments validate SSOLE’s performance using ResNet-18 and ResNet-50 on ImageNet100 and the full ImageNet(1K) dataset. However, the potential of SSOLE with more complex architectures, such as wide ResNet-50, Vision Transformers (ViT), and Swin Transformers, remains unexplored. These architectures have shown promise in SSL, achieving state-of-the-art results. Future work could explore the applicability of SSOLE with these advanced architectures. Additionally, investigating SSOLE as a complementary regularization technique alongside other SSL methods presents another promising avenue for future research.

### D IMPLEMENTATION DETAILS

This section outlines the configurations for our experiments in ablation studies, focusing on InfoNCE-M and SSOLE methods.

## D.1 INCORPORATION OF MULTIPLE CROPS

The multiple crops technique, introduced in SwAV Caron et al. (2020), draws inspiration from discussions in Wang et al. (2022) on handling outlier views and balancing low-rank and high-rank constraints. In SSL, the diversity of views during training has a significant impact on the effectiveness of models. As highlighted in Wang et al. (2022), incorporating high-variance views such as small crops requires careful consideration because they can deviate from the principal components of larger views. These deviations challenge the low-rank assumptions central to SSL, necessitating a balanced approach in multi-view learning.

To address this, we propose a two-step approach to low-rank and high-rank enforcement for large and small crops. Our strategy ensures that gradients from small crops do not interfere with those from large crops, which represent the primary views used for low-rank enforcement. The steps are detailed as follows:

Let  $\mathbf{Z}_L \in \mathbb{R}^{B \times N^L \times d}$  represent the feature matrix of the large crops, where  $N^L$  is the number of large crops, and  $d$  is the feature dimension. First, we compute the deviation matrix  $\tilde{\mathbf{Z}}_{L,b,:}$  by subtracting the mean of the large crops’ feature vectors:

$$\tilde{\mathbf{Z}}_{L,b,:} = \mathbf{Z}_{L,b,:} - \frac{1}{N^L} \sum_{n=1}^{N^L} \mathbf{Z}_{L,b,n}.$$

The low-rank loss for the large crops is then applied to this deviation matrix as:

$$\mathcal{L}_{intra}^{(L)} = h_1 \left( \|\tilde{\mathbf{Z}}_{L,b,:}\|_*, N^L \right).$$

Let  $\mathbf{Z}_{all} \in \mathbb{R}^{B \times N^{all} \times d}$  represent the feature matrix for all crops, where  $N^{all}$  includes both large (L) and small (S) crops. To account for the variance introduced by small crops, we calculate the deviation matrix for all crops, but the mean is replaced by the mean of the large crops:

$$\tilde{\mathbf{Z}}_{all,b,:} = \mathbf{Z}_{all,b,:} - \text{sg} \left( \frac{1}{N^L} \sum_{n=1}^{N^L} \mathbf{Z}_{L,b,n} \right).$$

$\text{sg}()$  denotes the stop-gradient operation. To further prevent the gradients of small crops from affecting the large crops, we stop the gradient flow on the large crops in this step. The low-rank loss for all crops is then computed as:

$$\mathcal{L}_{intra}^{(All)} = h_1 \left( \|\tilde{\mathbf{Z}}_{all,b,:}\|_*, N^{all} \right).$$

Then the low-rank loss for small crops is then computed as:

$$\mathcal{L}_{intra}^{(S)} = 2 * \mathcal{L}_{intra}^{(All)} - \text{sg}(\mathcal{L}_{intra}^{(L)}).$$

For high-rank enforcement, we perform the same operation on all crops, including both large and small views. Let  $\mathbf{Z}_{:,n}$  represent the feature matrix of the  $n$ -th view. The high-rank loss is computed as:

$$\begin{aligned} \mathcal{L}_{inter}^L &= h_2 (\|\mathbf{Z}_{L:,n}\|_*, B), \\ \mathcal{L}_{inter}^S &= h_2 (\|\mathbf{Z}_{S:,n}\|_*, B). \end{aligned}$$

To adjust the influence of the small crops, we introduce a weighting factor  $\beta$ . The low-rank and high-rank losses for the small crops are scaled by  $\beta$ , ensuring that the small crops’ contribution is balanced against that of the large crops:

$$\begin{aligned} \mathcal{L}_{intra} &= \mathcal{L}_{intra}^{(L)} + \beta * \mathcal{L}_{intra}^{(S)}, \\ \mathcal{L}_{inter} &= \mathcal{L}_{inter}^{(L)} + \beta * \mathcal{L}_{inter}^{(S)}. \end{aligned}$$

This approach ensures that diverse crops are incorporated effectively while preserving the benefits of low-rank and high-rank enforcement. The factor  $\beta$  allows for fine-tuning of the influence of small crops, optimizing the balance between view diversity and stability in SSL.

## 1188 D.2 PRETRAINING

1189 **ImageNet100 Experiments:** We employ ResNet-18 as the backbone ( $f_\theta$ ) with a three-layer MLP  
 1190 (4096- $d$  hidden layer with ReLU, followed by normalization) as a projector, yielding a final embedding  
 1191 dimension of  $d = 4096$ . Training utilizes a batch size of  $B = 128$  across 4 GPUs, SGD optimizer  
 1192 with a base learning rate of  $lr = 2.0$ , and a cosine decay to 0.002.  $\lambda = 0.7$ . The experiment uses  
 1193  $N_L = 4$  full views and  $N_S = 4$  small views with  $\beta = 0.6$ .

1194 **Full ImageNet (1K) Experiments:** For the full ImageNet dataset, ResNet-50 is used as the backbone  
 1195 with an enhanced three-layer MLP (8192- $d$  with ReLU and normalization) in the projector, leading  
 1196 to an embedding size of  $d = 8192$ . The batch size is set at  $B = 256$ , evenly distributed over 8 GPUs.  
 1197 The SGD optimizer is used with a base learning rate of  $lr = 1.0$ , decaying to 0.001 following a  
 1198 cosine rule.  $\lambda = 0.7$ . The experiment uses  $N_L = 4$  full views and  $N_S = 4$  small views with  $\beta = 0.6$ .

1200 We adopt the “multi-crop” data augmentation strategy from SwAV, an exemplary multi-view training  
 1201 algorithm. This method enriches the diversity of input data, crucial for effective multi-view self-  
 1202 supervised learning. During each iteration of training, we generate a combination of views:

- 1203 • *Full Views:* We create  $N_L$  full views for each image, each of size  $224 \times 224$  pixels. The scale  
 1204 factor for these views varies within the range of  $[0.14, 1.0]$ , ensuring a wide representation  
 1205 of the original images.
- 1206 • *Small Views:* Alongside full views,  $N_S$  small views of size  $96 \times 96$  pixels are generated, with  
 1207 a scale factor ranging from  $[0.05, 0.33]$ . These smaller views focus on different segments of  
 1208 the images, introducing further variance.

1209 Each generated view undergoes a series of augmentation techniques to enhance model robustness:

- 1210 • *Random Horizontal Flip:* Applied with a probability of 0.5 to introduce horizontal variability.
- 1211 • *Color Distortion:* This includes color jittering (brightness, contrast, saturation, and hue  
 1212 adjustments with respective strengths of 0.8, 0.8, 0.8, and 0.2) applied with a probability of  
 1213 0.8, and color dropping (conversion to grayscale) with a probability of 0.2.
- 1214 • *Gaussian Blur:* Each view is subjected to Gaussian blur, having a standard deviation in the  
 1215 range of  $[0.1, 2.0]$ , to simulate focus variability.
- 1216 • *Random Solarization:* Applied with a probability of 0.2 to further diversify the visual input.

## 1221 D.3 LINEAR PROBING

1222 Linear probing evaluates the representational quality of our SSOLE model. We outline distinct  
 1223 training protocols for ImageNet100 and the full ImageNet dataset.

### 1225 Training Protocols:

- 1226 • *ImageNet100:*  
 1227 – Batch Size:  $B = 256$ .  
 1228 – Optimizer: SGD with a momentum of 0.9, no weight decay.  
 1229 – Learning Rate: Base rate  $lr = 0.1$ , cosine schedule over 100 epochs.
- 1230 • *Full ImageNet:*  
 1231 – Batch Size:  $B = 2048$ .  
 1232 – Optimizer: SGD, momentum 0.9, no weight decay.  
 1233 – Learning Rate: Starting at  $lr = 0.6$ , reduced by 0.3 every 20 epochs, across 100  
 1234 epochs.

### 1237 Data Augmentation:

- 1238 • *Training Images:* Random cropping and resizing to  $224 \times 224$ , plus random horizontal flips  
 1239 (probability 0.5).
- 1240 • *Test Images:* Resize to  $256 \times 256$ , then center crop to  $224 \times 224$ .

1241



1242 **Implementation Note:** For the full ImageNet, we implemented SwAV based on the available code  
1243 and compared our results with those reported in Wang et al. (2022) (referenced in Table 4).  
1244

#### 1245 D.4 SEMI-SUPERVISED LEARNING ON IMAGENET 1246

1247 We conduct fine-tuning experiments with limited labeled data, specifically 1% and 10% subsets of  
1248 the ImageNet dataset. These subsets are the same as those used in Chen et al. (2020a).

1249 For the fine-tuning process, we employ an SGD optimizer with a batch size of  $B = 256$  and a  
1250 momentum of 0.9, without any weight decay. The fine-tuning is carried out for 20 epochs for both  
1251 the 1% and 10% labeled datasets.

1252 Following the learning rate scaling strategy from SwAV, we set different learning rates for the linear  
1253 layers and the backbone network weights. Specifically, the linear layers' learning rates are scaled up  
1254 by 250 times and 20 times for the 1% and 10% tasks, respectively. We determined the optimal base  
1255 learning rates for the linear layers to be 5.0 for the 1% task and 0.2 for the 10% task after conducting  
1256 a search in the range of 0.01 to 10. These learning rates are then reduced by a factor of 0.2 at the 12th  
1257 and 16th epochs during the training period.  
1258

#### 1259 D.5 TRANSFER LEARNING FROM IMAGENET TO OTHER DATASETS 1260

1261 For transfer learning tasks, we maintain consistent image pre-processing protocols during linear  
1262 classifier training and testing, as used in the linear evaluation on ImageNet. However, for CIFAR10  
1263 and CIFAR100, we adjust the image size to  $224 \times 224$  and apply random horizontal flipping to  
1264 training images with a probability of 0.5.

1265 In these experiments, our focus is on training linear classifiers on top of the frozen feature represen-  
1266 tations extracted from the pre-trained network. We employ an SGD optimizer with a batch size of  
1267  $B = 256$ , momentum of 0.9, and no weight decay. To determine the optimal base learning rate for  
1268 each algorithm, we conduct a search across seven logarithmically-spaced values ranging from 0.1 to  
1269 100. Once the optimal learning rate is identified for each algorithm, we apply a cosine decay rule for  
1270 its annealing.

1271 For the CIFAR10 and CIFAR100 datasets, the linear classifiers are trained for a total of 30,000  
1272 iterations. In contrast, for the Aircraft, DTD, and Flowers datasets, we limit the training to 5,000  
1273 iterations.

1274 We report top-1 accuracies (%) achieved by each method on CIFAR10, CIFAR100, and DTD datasets;  
1275 and *mean per class* on Aircraft and Flowers, following practices in Wang et al. (2022); Grill et al.  
1276 (2020).  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295