

MSVQ: Self-Supervised Learning with Multiple Sample Views and Queues

Chen Peng, Xianzhong Long*, Yun Li

School of Computer Science, Nanjing University of Posts and Telecommunications, Nanjing, 210023, China

Abstract

Self-supervised methods based on contrastive learning have achieved great success in unsupervised visual representation learning. However, most methods under this framework suffer from the problem of false negative samples. Inspired by the mean shift for self-supervised learning, we propose a new simple framework, namely Multiple Sample Views and Queues (MSVQ). We jointly construct three soft labels on-the-fly by utilizing two complementary and symmetric approaches: multiple augmented positive views and two momentum encoders that generate various semantic features for negative samples. Two teacher networks perform similarity relationship calculations with negative samples and then transfer this knowledge to the student network. Let the student network mimic the similarity relationships between the samples, thus giving the student network a more flexible ability to identify false negative samples in the dataset. The classification results on four benchmark image datasets demonstrate the high effectiveness and efficiency of our approach compared to some classical methods. Source code and pretrained models are available here.

Keywords:

Self-supervised learning, Contrastive learning, Knowledge distillation, Data augmentation, Momentum encoder.

1. Introduction

Self-supervised learning (SSL) has received sufficient attention and rapid progress in the computer vision community. This is due to its ability to learn rich semantic features using unlabeled data [1–8]. Early self-supervised learning methods were typically based on geometric transformations or heuristics to design the corresponding pretext task, such as image rotation [9]. The current mainstream SSL approach is based on the instance discrimination task [6] under the contrastive learning framework. Briefly, each image in the dataset is treated as a separate semantic class. In feature space, augmented views of the same image are pulled closer and views between other images are pushed away by the noise contrastive estimation (NCE) loss [10]. Meanwhile, there are some milestone works to continuously improve the methods under contrastive learning. For example, MoCo [8] introduced momentum encoders and queues to address the dilemma of memory consumption and untimely updates of data features. SimCLR [4] increased the difficulty of the self-supervised pre-training task by applying complex data augmentation and additional non-linear projectors.

The problem of false negative samples has severely impeded the ongoing development of contrastive self-supervised learning. False negative samples are defined as samples within the negative sample set that exhibit similar semantics or categories to the positive sample. Some works have attempted to address the false negative sample problem by introducing nearest neighbors (NN). For example, NNCLR [11] searches for the nearest neighbors of the query sample in its imported support set

and performs NCE loss with the positive sample. MSF [12] enriches the semantic information of the positive sample by searching for the top-K neighbors in the Memory Bank [6] and performing the Mean Squared Error loss with the query sample. CMSF [13] improves the semantic diversity of neighbor samples with an additional Memory Bank. SNCLR [14] leverages cross-attention scores to distinguish the contribution of different neighbor samples to the model. However, these methods often require a predefined number of neighboring samples, and determining this number in advance can be challenging.

In this study, we are interested in improving the reliability and coverage of models to identify false negative samples. Inspired by MSF, we propose a new simple Self-Supervised Learning with Multiple Sample Views and Queues (MSVQ). Our approach employs two complementary and symmetrical methods within the teacher networks to generate three distinct soft labels for the student network. Firstly, we create multiple augmented views of the positive sample within the teacher networks and perform consistent similarity distillation with negative samples from the same queue. Secondly, we introduce two separate queues into the model to generate diverse semantic features for the negative samples within these queues. This is accomplished by utilizing two momentum encoders with different update coefficients.

Our main contributions are summarized as follows:

- More augmented views of the positive sample. In the teacher networks, we apply weak data augmentation to the positive sample multiple times to enhance feature diversity. It is then distilled for consistency similarity with the negative samples in the queues to improve the reliability of the model in identifying false negative samples.

*Corresponding author

Email address: 1xz@njupt.edu.cn (Xianzhong Long)

- Feature diversity of negative samples. We form the variability of features on the same negative samples by leveraging encoders with different momentum coefficients. Our intuition is that the embedding diversity of negative samples can reduce the risk of omission for the model to identify false negative samples in feature space.

2. Related Work

2.1. Self-supervised learning

Self-supervised learning is an approach for learning generic semantic features from data by solving pretext tasks using large amounts of unlabeled data. In the early SSL, common pretext tasks included rotations [9, 15], grayscale coloring [16], or cropping [17, 18]. However, this approach may result in semantic features containing noisy information related to specific pretext tasks, which can hinder generalization [9, 19].

In recent years, contrastive learning methods based on instance discrimination tasks [1–4, 8, 20–22] have achieved rapid development in SSL. The core idea is that the semantic features of augmented views generated from the same image should be invariant. There are also some classical works under this framework. For instance, MoCo is designed to mitigate rapid memory consumption caused by too large batchsize, and performance decreases due to data features in the memory bank not being updated in time. This is achieved by employing the momentum encoder and the queue. SimCLR enhances the generalization of data features by introducing complex data augmentation and additional nonlinear projectors. Un-Mix [23] uses the idea of mix-up [24] to make the model predict less confidently. Meanwhile, several recent studies have demonstrated the efficiency of contrastive learning by leveraging feature similarities among positive samples. For example, BYOL [7] prevents model collapse without using negative samples based on the introduction of predictor and Mean Squared Error (MSE) loss. SimSiam [5] demonstrates that a simple siamese framework with stop-grad can achieve favorable properties on its own without requiring additional components.

2.2. Knowledge distillation

Knowledge distillation [25, 26] is a method where a model (teacher network) that has learned rich semantic features is used to transfer its knowledge to another model (student network). While in most cases, the teacher network has a more complex model structure compared to the student network. In some work, the teacher network and the student network can also be the same model structure [27, 28].

Recent contrastive self-supervised learning can be seen as a structure that contains a student network and a teacher network. For example, in MoCo, the student network corresponds to an encoder with stochastic gradient descent (SGD) update and the teacher network corresponds to an encoder with momentum update. Compared to previous work [29, 30], the student network in the proposed MSVQ can learn rich semantic features from multiple teacher networks with different knowledge. The student network has a more flexible ability to identify

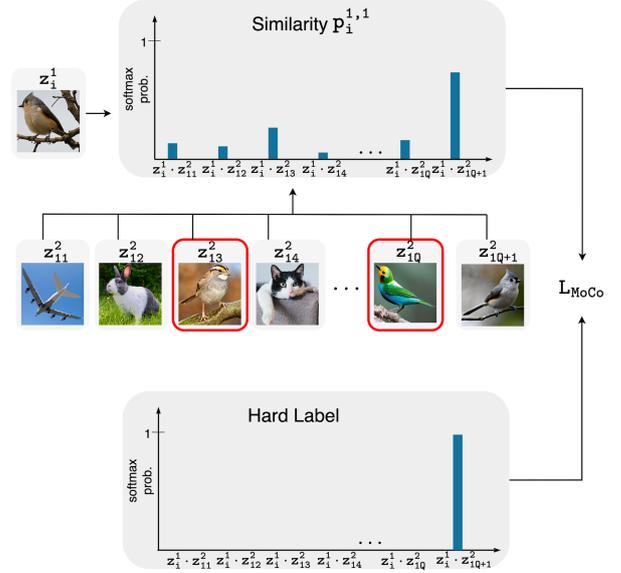


Fig. 1. Illustration of MoCo. While the similarity distribution $p_i^{1,1}$ can search for negative samples (the samples marked by red blocks, i.e., false negative samples) in the queue that are semantically similar to the query sample z_i^1 . But an artificial one-hot label ignores the semantic relationship between them.

false negative samples in the dataset. MOKD [31] is the closest related work to our approach that utilizes multiple teacher networks to teach a student network. However, our approach differs in that we use different momentum update coefficients to construct teacher networks, rather than introducing multiple heterogeneous models (e.g., a combination of ResNet [32] and ViT [33]).

2.3. Consistency regularization

Consistency regularization is an approach to make the output of a model consistent or similar under small perturbations in the input [34] or model parameters [35]. It allows the model to learn the most possible diversity of semantic features of the input data. Inspired by this idea, some works have been applied to contrastive self-supervised learning to solve the false negative sample problem. For example, MSF introduces top-K neighbors of the positive sample and minimizes the distance between the query sample and its nearest neighbors. This helps the network learn the semantic diversity of the query sample. On the SimCLR-based framework, NNCLR attempts to search for nearest neighbors in the support set and performs NCE loss with the positive sample. However, the performance of these methods is sensitive to the number of nearest neighbors K. Meanwhile, due to the random initialization, it is meaningless for the model to find the nearest neighbors in the early training stage.

CO2 [36] is based on MoCo by adding a consistency regularization term to distill various augmented views of the same image with negative samples. However, ReSSL [37] shows that the regularization term itself can learn meaningful information when the appropriate temperature parameter and data augmentation are utilized in the teacher network. In addition, SCE

[38] incorporates hard labels into ReSSL to enhance the discriminative ability of the model. Our approach utilizes multiple augmented views of the positive sample along with consistent similarity distillation of negative samples within the queues. Inspired by the concept of symmetry, we incorporate two momentum encoders to generate distinct semantic features for the same negative samples. Both of these complementary approaches are employed to jointly improve the ability of the model to identify false negative samples.

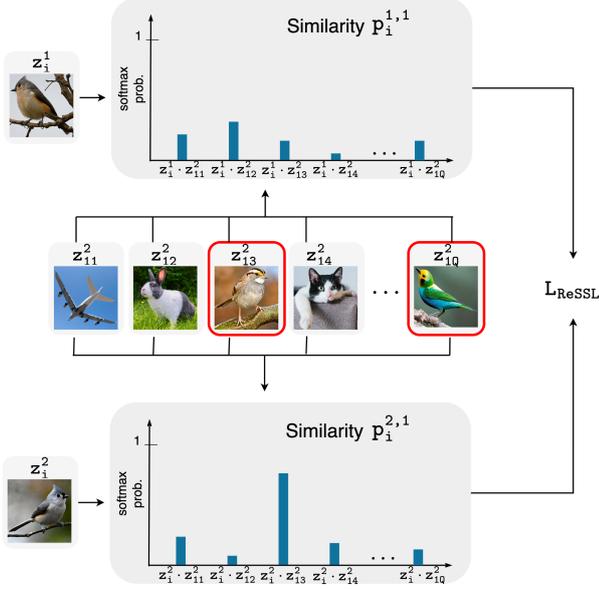


Fig. 2. Illustration of ReSSL. We can correspond the upper and lower branches to the student network and the teacher network, respectively. Since the teacher network uses an appropriate temperature parameter and weak data augmentation, $p_i^{2,1}$ highlights the important instance relationships and filters out some trivial connections. Nevertheless, depending on a single teacher network to accurately and consistently represent the similarity between instance samples is challenging.

3. Methodology

In this section, first, we briefly review our baselines: MoCo [8, 39] and ReSSL [37]. Then we introduce our proposed MSVQ framework. Meanwhile, the algorithm and implementation details of MSVQ will also be explained.

3.1. Previous contrastive self-supervised learning and relational learning

Let $X \in \mathbb{R}^{N \times H \times W \times C}$ be N training samples with height H , width W , and number of channels C . The queues $Queue_1 = \{z_{1j}^1\}_1^Q$ and $Queue_2 = \{z_{2j}^2\}_1^Q$ each contain a set of Q random embeddings of other samples. Meanwhile, $f_{\cdot}(\cdot)$ is a backbone network (e.g., ResNet [32]) and $g_{\cdot}(\cdot)$ is a nonlinear projector. First, we generate two different views by data augmentation ($T(\cdot)$) on the training samples: $X^1 = T_1(X)$ and $X^2 = T_2(X)$. Next, we input these views into the student network and the teacher network to obtain the corresponding embedding features $Z^1 = g_s(f_s(X^1))$ and $Z^2 = g_t(f_t(X^2))$, respectively. In

the case of MoCo (as depicted in Fig. 1), it employs the NCE loss, which is defined as follows:

$$L_{MoCo} = -\frac{1}{N} \sum_{i=1}^N \log \frac{\exp(\text{sim}(z_i^1, z_i^2)/\tau)}{\sum_{j=1}^{Q+1} \exp(\text{sim}(z_i^1, z_{1j}^2)/\tau)} \quad (1)$$

where $z_{1Q+1}^2 \triangleq z_i^2$ and τ is the temperature parameter. Meanwhile, $\text{sim}(u, v)$ is the similarity measure of the two feature embeddings, such as cosine similarity:

$$\text{sim}(u, v) = \frac{u^T v}{\|u\|_2 \|v\|_2} \quad (2)$$

In the context of ReSSL (as illustrated in Fig. 2), the method utilizes distributions $p_{ij}^{1,1}$ and $p_{ij}^{2,1}$ to represent the similarity relationship among the instance samples:

$$p_{ij}^{1,1} = \frac{\exp(\text{sim}(z_i^1, z_{1j}^2)/\tau_s)}{\sum_{k=1}^Q \exp(\text{sim}(z_i^1, z_{1k}^2)/\tau_s)}, j = 1, 2, \dots, Q \quad (3)$$

$$p_{ij}^{2,1} = \frac{\exp(\text{sim}(z_i^2, z_{1j}^1)/\tau_t)}{\sum_{k=1}^Q \exp(\text{sim}(z_i^2, z_{1k}^1)/\tau_t)}, j = 1, 2, \dots, Q \quad (4)$$

where τ_s and τ_t represent the temperature hyperparameters used in calculating the relationship distributions within the student network and the teacher network, respectively. The loss function of ReSSL is the KL divergence of $P^{2,1}$ and $P^{1,1}$:

$$L_{ReSSL} = KL(P^{2,1} \| P^{1,1}) \quad (5)$$

Both MoCo and ReSSL employ a similar approach to update their teacher networks, as described by the following formula:

$$f_{t1} = m_1 f_{t1} + (1 - m_1) f_s, g_{t1} = m_1 g_{t1} + (1 - m_1) g_s \quad (6)$$

where m_1 indicates the momentum update coefficient. Since the teacher network in ReSSL is not updated directly by the loss function, its loss function can simply use cross entropy instead of KL divergence. We also employ the momentum update mechanism, queue storage for negative samples, and KL divergence to capture inter-sample relationships.

3.2. Self-supervised learning with multiple sample views and queues

In this work, we propose two symmetric and complementary ways to improve the reliability and coverage of the model to identify false negative samples.

3.2.1. Multiple sample views

Inspired by some work [11, 12, 37], we incorporate multiple augmented views of the positive sample into the teacher network to comprehensively represent its semantics. This approach aids in the precise identification of false negative samples within the queue $Queue_1$. Specifically, the training samples denoted as X undergo data augmentation via $T(\cdot)$, resulting in $X^3 = T_3(X)$. Subsequently, it is projected using $g_{t1}(f_{t1}(\cdot))$ to obtain Z^3 . Finally, we calculate the similarity between Z^3 and the negative samples in $Queue_1$, applying the softmax function

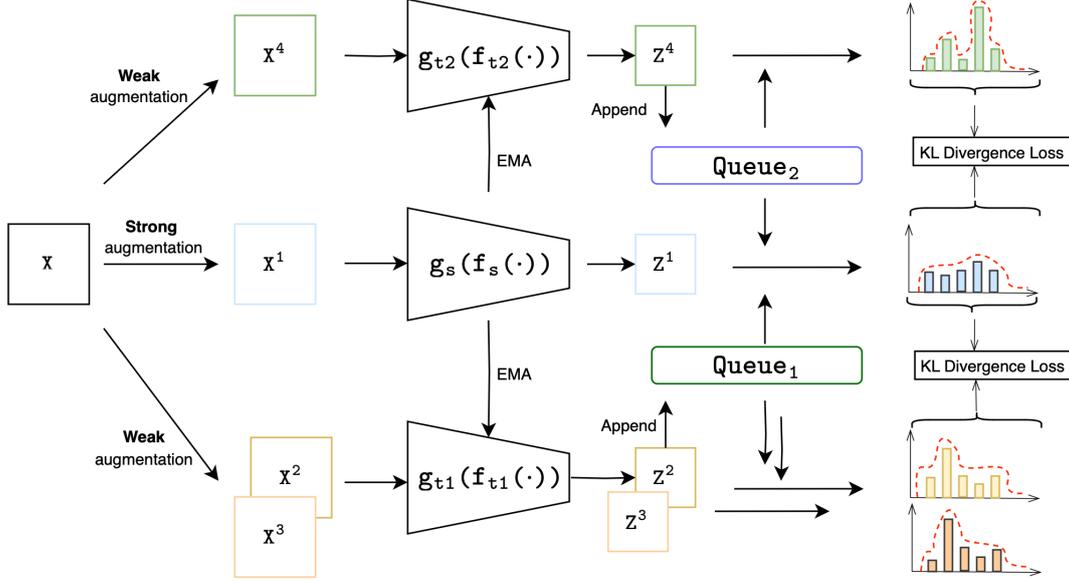


Fig. 3. The overall framework of MSVQ. We can consider the relationship distributions generated by the two teacher networks as three distinct soft labels, which serve as guidance for the student network in classifying the negative samples within the queues.

to derive a distribution that represents the relationships among these samples:

$$P_{ij}^{3,1} = \frac{\exp(\text{sim}(z_i^3, z_{1j}^2)/\tau_t)}{\sum_{k=1}^Q \exp(\text{sim}(z_i^3, z_{1k}^2)/\tau_t)}, j = 1, 2, \dots, Q \quad (7)$$

$P^{3,1}$ serves as an extra soft label to guide the student network. This guidance involves calculating the KL divergence between $P^{3,1}$ and $P^{1,1}$. To simplify, we average its loss function with ReSSL:

$$L_{MSV} = \frac{1}{2}(KL(P^{2,1}||P^{1,1}) + KL(P^{3,1}||P^{1,1})) \quad (8)$$

The update rule of the teacher network in Multiple Sample Views (MSV) is consistent with that of ReSSL.

3.2.2. Multiple queues

MSV can be viewed as a variant of previous work [11, 12], with the notable distinction that MSV enhances the semantics of the positive sample through multiple data augmentations rather than identifying its top-K neighbors. From a symmetric perspective, we also focus on the aspect of negative samples. Our intuition is that the diversity of both is complementary and can help achieve a more comprehensive and robust representation of the underlying data distribution.

To be specific, a teacher network $g_{t2}(f_{t2}(\cdot))$ with a momentum coefficient of m_2 is employed. The training images X are independently subjected to data augmentation using $T(\cdot)$ to obtain $X^4 = T_4(X)$, and then they are projected using $g_{t2}(f_{t2}(\cdot))$ to yield $Z^4 = \{z_i^4\}_i^N$. The following relationship distribution indicates the similarity between instances:

$$P_{ij}^{1,2} = \frac{\exp(\text{sim}(z_i^1, z_{2j}^4)/\tau_s)}{\sum_{k=1}^Q \exp(\text{sim}(z_i^1, z_{2k}^4)/\tau_s)}, j = 1, 2, \dots, Q \quad (9)$$

$$P_{ij}^{4,2} = \frac{\exp(\text{sim}(z_i^4, z_{2j}^4)/\tau_t)}{\sum_{k=1}^Q \exp(\text{sim}(z_i^4, z_{2k}^4)/\tau_t)}, j = 1, 2, \dots, Q \quad (10)$$

This part is the KL divergence of $P^{4,2}$ and $P^{1,2}$. For simplicity, we also just take the average value of its loss function with ReSSL:

$$L_{MQ} = \frac{1}{2}(KL(P^{2,1}||P^{1,1}) + KL(P^{4,2}||P^{1,2})) \quad (11)$$

The Multiple Queues (MQ) contains two teacher networks with different update coefficients. The first teacher network is updated in the identical way as ReSSL, while the second teacher network uses the following update mechanism:

$$f_{t2} = m_2 f_{t2} + (1 - m_2) f_s, g_{t2} = m_2 g_{t2} + (1 - m_2) g_s \quad (12)$$

3.2.3. Multiple sample views and queues

As illustrated in Fig. 3, a simple organic merging of these two symmetric ways yields our proposed method. We optimize the student network only by minimizing the following loss:

$$L_{MSVQ} = \frac{1}{3}(KL(P^{2,1}||P^{1,1}) + KL(P^{3,1}||P^{1,1}) + KL(P^{4,2}||P^{1,2})) \quad (13)$$

This method combines the strengths of both approaches to generate three distinct soft labels, $P^{2,1}$, $P^{3,1}$, and $P^{4,2}$, aiming to mitigate the under-detection of false negative samples by the model. These labels effectively transfer knowledge from multiple teacher networks to the student network. Further analysis and discussion of these three soft labels are presented in Sec. 4.6.

Notably, z_{ij}^- represents the semantic embedding of negative sample j in $Queue_t$, p_{ij}^- signifies the inter-instance similarity between the positive sample z_i^- and the negative sample z_{ij}^- . We have included a comprehensive list of important notations in the MSVQ framework along with their specific meanings in Table 1. The procedure of MSVQ is summarized in Algorithm 1.

3.2.4. Sharper distribution and friendly data augmentation

To effectively emphasize significant false negative samples within the teacher networks and simultaneously filter out potential noisy relationships in the queues, our model ensures that $\tau_t < \tau_s$. We conduct an analysis of the impact of different τ_t values on MSVQ in Sec. 4.5.1.

To reduce the impact of aggressive data augmentation [4] in the teacher networks, we employ a weaker data augmentation scheme [12, 37, 38] to generate more suitable soft labels. We also conduct an analysis of the impact of various data augmentation ways on MSVQ in Sec. 4.5.2.

Table 1: Notations in the MSVQ Framework.

Notation	Meaning
$X = \{x_i\}_1^N$	Comprising N training samples from the current batch, where each element is regarded as a positive sample.
$Queue_1 = \{z_{i1}^2\}_1^Q$ $(Queue_2 = \{z_{i2}^4\}_1^Q)$	Comprising Q negative sample features, with these features derived from the most recent previous batches of Z^2 (Z^4). The update method follows a first-in-first-out (FIFO) approach.
$T(\cdot)$	The data augmentation distribution is employed to generate various augmented views. Specifically, $T_1(\cdot)$ is utilized to apply strong data augmentation to the student network, while $T_i(\cdot), i \in \{2, 3, 4\}$ is employed for applying weak data augmentation to the teacher networks.
Z	$Z^i = \{z_j^i\}_{j=1}^N, i \in \{1, 2, 3, 4\}$ denote the 128-dimensional image features obtained by projecting $X^i, i \in \{1, 2, 3, 4\}$ through the corresponding network.
$p^{1,1} (p^{1,2})$	The similarity relationship exists between the features Z^1 of the positive samples and the features of the negative samples within $Queue_1$ ($Queue_2$).
$p^{2,1} (p^{3,1})$	The similarity relationship exists between the features Z^2 (Z^3) of the positive samples and the features of the negative samples within $Queue_1$.
$p^{4,2}$	The similarity relationship exists between the features Z^4 of the positive samples and the features of the negative samples within $Queue_2$.
$m_1 (m_2)$	The momentum update coefficient is used to update the teacher network $g_{r1}(f_{r1}(\cdot))$ ($g_{r2}(f_{r2}(\cdot))$). To introduce variability in the semantic features of different teacher networks, we set $m_1 \neq m_2$.
$\tau_s (\tau_t)$	Temperature hyperparameter for similarity distributions in the student network (teacher networks).

4. Experiments and Results

In this section, we will compare and analyze the proposed MSVQ with previous classical algorithms on four benchmark image datasets.

4.1. Datasets and device performance

Most SSL methods are typically evaluated using an ImageNet-1K dataset [40] containing almost 1.3M images. However, due to hardware limitations, implementing this evaluation can be challenging for most research labs. We evaluate the proposed method MSVQ with some classical contrastive self-supervised methods on four datasets.

- CIFAR-10 and CIFAR-100 [41]: Both datasets comprise 60,000 color images, consisting of 50,000 training images and 10,000 test images, all with a resolution of 32x32 pixels. However, they differ in the number of classes they contain: the CIFAR-10 dataset comprises 10 classes, while the CIFAR-100 dataset comprises 100 semantic classes.

Algorithm 1: PyTorch-style pseudocode for MSVQ

```

# Fs, Ft1, Ft2: encoder for student, teacher1 and
teacher2, F·  $\hat{=}$  g·(f·(·))
# queue1, queue2: two queues(CxQ)
# m1, m2: momentum for teacher1 and teacher2
#  $\tau_s, \tau_t$ : temperature for student and teacher
# CE: CrossEntropyLoss

Ft1.params = Fs.params # initialize teacher1
Ft2.params = Fs.params # initialize teacher2
# load a minibatch X with N samples
for X in loader:
    # random augmentation
    X1, X2 = strong_aug(X), weak_aug(X)
    X3, X4 = weak_aug(X), weak_aug(X)

    Z1, Z2 = Fs.forward(X1), Ft1.forward(X2) # NxN
    Z3, Z4 = Ft1.forward(X3), Ft2.forward(X4)
    # l2-normalize
    Z1, Z2, Z3, Z4 = normalize(Z1, Z2, Z3, Z4, dim=1)

    Z2, Z3, Z4 = Z2.detach(), Z3.detach(), Z4.detach()
    # mm: matrix multiplication
    logits11 = mm(Z1, queue1) # [NxN, CxQ] -> NxQ
    logits12 = mm(Z1, queue2)
    logits21, logits31 = mm(Z2, queue1), mm(Z3, queue1)
    logits42 = mm(Z4, queue2)

    loss1 = CE(logits11/ $\tau_s$ , softmax(logits21/ $\tau_t$ ))
    loss2 = CE(logits11/ $\tau_s$ , softmax(logits31/ $\tau_t$ ))
    loss3 = CE(logits12/ $\tau_s$ , softmax(logits42/ $\tau_t$ ))
    loss = (loss1+loss2+loss3)/3

    loss.backward()
    # SGD update: student
    update(Fs.params)
    # momentum update: teacher1 and teacher2
    Ft1.params = m1*Ft1.params+(1-m1)*Fs.params
    Ft2.params = m2*Ft2.params+(1-m2)*Fs.params
    # update two queues
    enqueue(queue1, Z2)
    enqueue(queue2, Z4)
    dequeue(queue1)
    dequeue(queue2)

```

Table 2: Parameters of the experiment.

	CIFAR-10	CIFAR-100	STL-10	Tiny ImageNet
Pre-training				
Epoch	200	200	200	200
Batch size	256	256	256	256
Warm up epoch	5	5	5	5
Base learning rate	0.06	0.06	0.06	0.06
(m_1, m_2)	(0.99, 0.95)	(0.99, 0.93)	(0.996, 0.99)	(0.996, 0.99)
(τ_s, τ_t)	(0.1, 0.04)	(0.1, 0.03)	(0.1, 0.04)	(0.1, 0.04)
Queue size	4096	4096	16384	16384
Weight decay	5e-4	5e-4	5e-4	5e-4
Cropped and Resized	32 × 32	32 × 32	64 × 64	64 × 64
Fine-tuning				
Epoch	100	100	100	100
Batch size	256	256	256	256
Base learning rate	1	1	1	1
Weight decay	0	0	0	0

Table 3: Data augmentation of the experiment. 'Resized Crops' specifies the lower and upper bounds of scale for cropping a random region based on the area of the original image. For the other data augmentations, they indicate the probability (denoted as ' ρ ') of being randomly applied. For instance, in the Strong category, the probability that "Horizontal Flip" is applied is $\rho = 0.5$.

	Resized Crops	Horizontal Flip	Color Jitter	GrayScale	Gaussian Blur
Strong	(0.2, 1.0)	0.5	0.8	0.2	0.5
Weak	(0.2, 1.0)	0.9			

- STL-10 [42] and Tiny ImageNet [43]: For the STL-10 dataset, the training set comprises 100,000 unlabeled color

images and 5,000 labeled color images. Additionally, the test set includes 8,000 labeled images. All of these images share a common resolution of 96x96 pixels, and the dataset encompasses a total of 10 categories. Tiny ImageNet comprises 120,000 images distributed across 200 classes. These images are resized to a dimension of 64x64 pixels. Specifically, each class includes 500 training images, along with 50 validation images and 50 test images.

A consistent hardware setup (1 Nvidia GTX 3090 GPU) was used for all algorithm experiments in this study. All algorithms are initially pre-trained using the training set. During the evaluation phase, we evaluate them using the test set, except in the case of Tiny ImageNet, where the validation set is used.

4.2. Pre-training

In all datasets, we use ResNet18 [32] as the backbone network $f(\cdot)$. Meanwhile, a nonlinear projector $g(\cdot)$ is added following the backbone network. All projectors within both the student network and teacher networks are composed of two fully-connected (FC) layers together with a linear rectification function (ReLU) layer between them, where the first FC layer is of size [512, 2048], and the second is of size [2048, 128].

To ensure a fair comparison, certain hyperparameters and data augmentations [44] in MSVQ were aligned with those in ReSSL. Regarding the model parameters, we employed the SGD optimizer with a momentum value of 0.9 and a weight decay of $5e-4$ for pre-training the model over 200 epochs. In terms of data augmentation, strong data augmentation was applied to the student network, while weak data augmentation was employed for the teacher networks. Additional details can be found in Table 2 and Table 3, respectively.

4.3. Fine-tuning

After the pre-training phase, we employed the widely adopted linear evaluation protocol to assess our model. In this protocol, we initially discarded the pre-trained nonlinear projector and fixed all parameters of the backbone network. Afterward, we added a linear classifier to the backbone of the student network with dimensions [512, cla], where cla represents the number of categories in the dataset. This linear classifier was used for the linear evaluation, enabling the model to perform classification based on the learned features. Finally, we fine-tuned the classifier for 100 epochs using the SGD optimizer with a momentum of 0.9 and weight decay of 0. Further details of these parameters can be found in Table 2.

4.4. Main results

4.4.1. Linear evaluation protocol

In Table 4, the results of other methods are copied from [37] for best results. For a fair comparison, we have also reproduced some recent work and marked it with *. It is clear that MSVQ outperforms the other classical methods on most benchmarks. Noticeably, MSVQ also has a significantly better performance compared to the MSV and the MQ trained alone. This indicates that the MSV and MQ methods learn complementary semantic features of images. Meanwhile, our method requires only a slight additional overhead without multiple backpropagations.

Table 4: Linear evaluation results. The optimal results are shown in bold, and the suboptimal results are underlined. Results marked with * were reproduced from the official code. †: Similar to the SCE [38] network framework, batch normalization [45] is incorporated after the first FC layer of the projector within MSVQ.

Method	BackProp	CIFAR-10	CIFAR-100	STL-10	Tiny ImageNet
Supervised	-	94.22	74.66	82.55	59.26
SimCLR [4]	2x	84.92	59.28	85.48	44.38
BYOL [7]	2x	85.82	57.75	87.45	42.70
SimSiam [5]	2x	88.51	60.00	87.47	37.04
MoCoV2 [39]	1x	86.18	59.51	85.88	43.36
ReSSL [37]	1x	90.20	63.79	88.25	46.60
ReSSL* [37]	1x	90.22	64.22	87.64	45.61
CMSF* [13]	2x	91.00	62.37	88.21	44.50
Un-Mix* [23]	2x	90.20	64.42	89.76	45.20
SNCLR* [14]	2x	88.50	62.40	88.24	46.14
SCE [38]	2x	90.34	65.45	89.94	51.90
MSV(Ours)	1x	90.92	65.02	89.35	46.68
MQ(Ours)	1x	90.91	65.15	89.58	46.32
MSVQ(Ours)	1x	91.46	66.44	90.41	48.09
SCE* [38]	2x	90.03	65.41	90.06	48.11
MSVQ†(Ours)	1x	91.28	65.82	89.71	49.51

4.4.2. Learning efficiency analysis

To reduce the influence of downstream task hyperparameters on the model and enhance evaluation efficiency, we also employ K Nearest Neighbors (KNN) classification to assess the pre-trained off-the-shelf features, with K set to 200 [6]. The online KNN classification results in Table 5 reflect that MSVQ can learn rich semantic features in the pre-training stage. In Fig. 4, the KNN (K=200) classification accuracy curves show that our method has reliable learning efficiency compared to ReSSL and MoCoV2.

Table 5: The online evaluation accuracy using KNN with K set to 200. The best results are shown in bold and the suboptimal results are underlined. * denotes results that were reproduced from the official code. †: Similar to the SCE [38] network framework, batch normalization is added after the first FC layer of the projector within MSVQ.

Method	CIFAR-10	CIFAR-100	STL-10	Tiny ImageNet
ReSSL* [37]	89.11	57.93	84.18	37.19
CMSF* [13]	89.30	55.57	84.11	36.79
Un-Mix* [23]	88.38	59.52	85.18	38.81
SNCLR* [14]	86.15	55.07	82.18	36.86
MSV(Ours)	89.64	59.33	85.58	38.78
MQ(Ours)	89.51	60.16	85.69	38.84
MSVQ(Ours)	90.16	60.66	86.51	40.26
SCE* [38]	88.41	59.30	85.45	40.17
MSVQ†(Ours)	90.23	61.63	86.56	41.81

4.5. Ablation studies

4.5.1. Sharper distribution

An appropriate temperature parameter τ_t can eliminate the noisy relationship between the positive sample and the negative samples in the queues, thus providing accurate soft labels for the student network. Table 6 provides valuable insights into the influence of temperature parameters on our method. Notably, it becomes evident that extremely low or high values of these temperature parameters are suboptimal. Intuitively, temperature τ_t controls the degree of smoothing of the three labels in the MSVQ. As τ_t approaches 0, MSVQ becomes analogous to using three artificial one-hot codings. This means that each label solely focuses on the most similar false negative sample

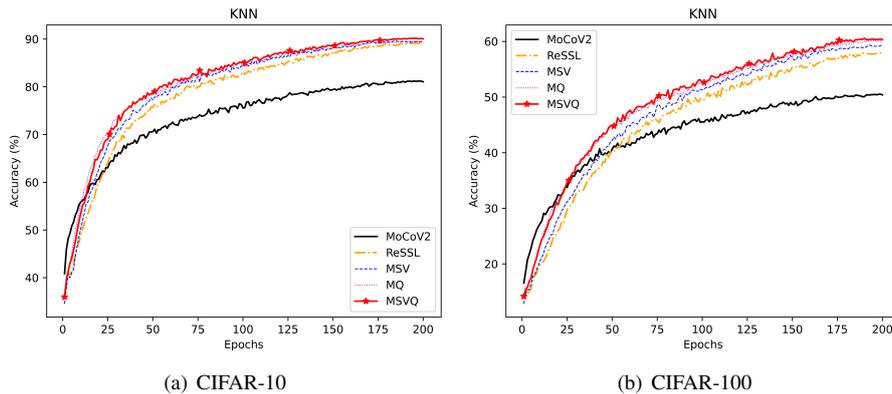


Fig. 4. KNN online evaluation accuracy curve.

within the queue. Conversely, as τ_t approaches 0.1, the student network is constrained to acquiring only trivial knowledge from the teacher networks.

Table 6: Effect of different τ_t for MSVQ.

Dataset	$\tau_t=0.1$	$\tau_t=0.01$	$\tau_t=0.02$	$\tau_t=0.03$	$\tau_t=0.04$	$\tau_t=0.05$	$\tau_t=0.06$	$\tau_t=0.07$
CIFAR-10	-	90.87	90.68	90.73	91.46	91.20	90.62	90.32
CIFAR-100	-	65.68	65.33	66.44	66.27	65.78	63.78	58.83
STL-10	-	88.68	89.21	89.64	90.41	89.56	88.96	87.75
Tiny ImageNet	-	47.26	47.81	47.96	48.09	46.99	45.72	44.47

4.5.2. Data augmentation

Unlike previous methods [4] that use aggressive data augmentation to encourage the learning of semantics invariant to geometric transformations, MSVQ adopts a different approach. Within the MSVQ framework, the teacher networks employ milder, weak data augmentation techniques with the specific goal of preserving image semantics to provide appropriate soft labels for the student network. In this section, we conduct a comprehensive study to assess the performance impact of various data augmentation approaches in teacher networks. Table 7 demonstrates that employing traditional data augmentation techniques for teacher networks can lead to semantic loss or distortion of the positive sample, resulting in incorrect soft labels that may mislead the student network.

Table 7: Effect of weak augmentation in teacher networks for MSVQ.

Teacher Aug	CIFAR-10	CIFAR-100	STL-10	Tiny ImageNet
Strong	88.33	59.55	86.14	39.44
Weak	91.46	66.44	90.41	48.09

4.5.3. Different momentum coefficients

Some may question the necessity of employing distinct momentum update coefficients for the two teacher networks within the MSVQ framework. As demonstrated in Table 8, when both teacher networks employ identical momentum update coefficients, a marginal reduction in performance is observed. We contend that $m_1 \neq m_2$ serves as a means to distinguish between $g_{r1}(f_{i1}(\cdot))$ and $g_{r2}(f_{i2}(\cdot))$. This differentiation enables a clear distinction between Z^2 (or Z^3) and Z^4 , as well as among the negative samples within $Queue_1$ and $Queue_2$. Consequently,

this approach maximizes the expression of semantic information within both the positive and negative samples.

Table 8: Effect of different m_2 for MSVQ.

Dataset	$m_1 = 0.99$	$m_2 = 0.93$	$m_2 = 0.95$	$m_2 = 0.97$	$m_2 = 0.99$
CIFAR-10	-	90.93	91.46	91.08	90.88
CIFAR-100	-	66.44	65.99	65.58	65.63
Dataset	$m_1 = 0.996$	$m_2 = 0.95$	$m_2 = 0.99$	$m_2 = 0.993$	$m_2 = 0.996$
STL-10	-	90.00	90.41	90.10	90.18
Tiny ImageNet	-	47.59	48.09	47.55	48.05

4.5.4. The size of queues

In the context of the MSVQ framework, the size of queues directly corresponds to the number of negative samples. Table 9 illustrates the linear evaluation accuracy corresponding to various queue sizes. It is evident that larger queue sizes lead to a substantial improvement in performance. We hypothesize that a larger queue size augments the probability of the positive sample discovering false negative samples within the queues that align with its semantic context. Consequently, the features acquired by the model become more generalizable. However, as we further increase the queue size, we observe performance instability. We suspect that this phenomenon arises from the presence of stale features within the queues, which are not promptly replaced when the queue size is enlarged. Therefore, it is necessary for the MSVQ to strike a balance between enhancing the chances of discovering false negative samples and ensuring the timely update of features within the queues.

Table 9: Analysis of the size of queues (Q).

Dataset	Q = 512	Q = 1024	Q = 2048	Q = 4096	Q = 8192	Q = 16384	Q = 32768
CIFAR-10	90.81	90.65	91.40	91.46	91.07	91.07	91.06
CIFAR-100	64.25	64.82	65.77	66.44	66.08	66.43	66.82
STL-10	88.10	89.30	89.60	89.69	89.76	90.41	90.11
Tiny ImageNet	44.82	45.93	47.04	47.31	47.74	48.09	48.34

4.6. Analysis and discussion

4.6.1. The position of augmented views

In addition to the default configuration described in Sec. 3.2.3, we conducted experiments exploring two variations of MSVQ: (i) relocating the augmented view X^3 from $g_{r1}(f_{i1}(\cdot))$ to $g_{r2}(f_{i2}(\cdot))$, and (ii) introducing a new view X^5 generated

Table 10: An analysis of the position of augmented views in MSVQ. The best results are indicated in bold, and the suboptimal results are underlined.

Method	X^2	X^3 in $g_{i1}(f_{i1}(\cdot))$	X^3 in $g_{i2}(f_{i2}(\cdot))$	X^4	X^5 in $g_{i2}(f_{i2}(\cdot))$	CIFAR-10	CIFAR-100	STL-10	Tiny ImageNet
MSVQ(Ours)	✓	✓		✓		91.46	<u>66.44</u>	<u>90.41</u>	<u>48.09</u>
	✓		✓	✓		91.05	66.16	90.11	47.45
	✓	✓		✓	✓	<u>91.16</u>	67.14	90.49	48.55
MSV(Ours)	✓	✓				90.92	65.02	89.35	46.68
MQ(Ours)	✓			✓		90.91	65.15	89.58	46.32
ReSSL [37]	✓					90.20	63.79	88.25	46.60

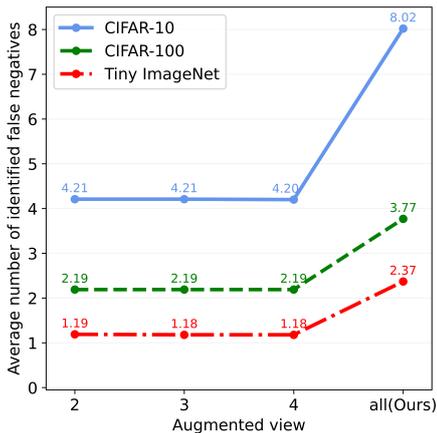
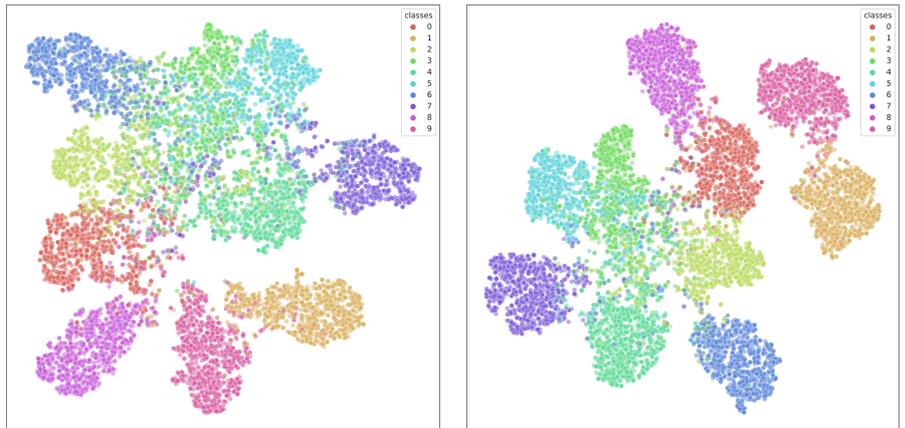


Fig. 5. The average number of false negative samples identified by different augmented views (i.e., X^2 , X^3 , and X^4) in teacher networks. Here, 'all' represents the cumulative effect of all three soft labels when utilized simultaneously.



(a) MoCoV2

(b) MSVQ

Fig. 6. t-SNE visualization of learned features on CIFAR-10, classes indicated by different colors. Best viewed in color.

through weak augmentation in $g_{i2}(f_{i2}(\cdot))$. The experimental results, as presented in Table 10, demonstrate that (ii) yielded slightly better performance improvements compared to (i). This suggests that the number of augmented views (i.e., distinct soft labels) may have a more significant impact on performance than the choice of the teacher network in which they are employed.

Furthermore, the default MSVQ settings yielded slight performance gains compared to variation (i) across all datasets. This was due to our experimental setup where we set the momentum update coefficient m_1 to be greater than m_2 . Larger momentum update coefficients reduce disturbances caused by inconsistencies among different batches of negative samples in the queue [8]. Given this simplicity and the optimal performance observed, we have chosen to retain the default settings of MSVQ.

4.6.2. Analyzing model reliability and coverage in identifying false negative samples

We treat the distribution of relationships between $\{X^i\}_{i=2}^4$ and the negative samples in the queues as three distinct soft labels. These labels provide guidance for the student network in classifying the negative samples within the queues. Ideally, false negative samples should receive higher prediction values in the student network, and vice versa.

To ensure the accuracy of these three soft labels, we employ

both weak data augmentation and lower temperature parameters in teacher networks. In this section, we investigate the reliability and scope of these three soft labels in the identification of false negative samples. These labels are expected to assign higher similarity values to the false negative samples. Fig. 5 illustrates the average number of false negative samples identified by the three soft labels. In detail, we begin by arranging each of these three distributions (i.e., $P^{2,1}$, $P^{3,1}$, and $P^{4,2}$) in descending order. Then, we calculate the average count of false negative samples that share the same labels as the positive sample within the *top 5* samples of each distribution. In this context, 'all' refers to the total number of distinct false negative samples identified by aggregating their respective sets of false negative samples when utilizing all three soft labels simultaneously. This observation implies that the three soft labels in our model can effectively identify distinct sets of false negative samples within the queues, resulting in the recognition of nearly twice as many false negatives compared to when each soft label is applied individually.

4.6.3. Visualization of features

As demonstrated by t-SNE visualization [46] in Fig. 6, our method exhibits more distinct class boundaries and a more compact internal arrangement of classes compared to MoCoV2. This suggests that MSVQ offers a greater ability to alleviate the

issue of false negative samples in the instance discrimination task.

5. Conclusion

In this work, we bring in the framework of MSVQ. We improve the reliability and coverage of false negative sample identification by introducing two complementary and symmetrical methods to generate three distinct soft labels within the teacher networks. The first method entails utilizing multiple weakly augmented views of the positive sample, while the second method involves employing two momentum encoders to generate distinct semantic features for negative samples. Our extensive experimental results on four benchmarks demonstrate the remarkable performance of MSVQ. In future research, our goal is to explore even more effective strategies for leveraging semantic diversity within the realm of SSL.

Acknowledgement

This work was supported by the National Natural Science Foundation of China under Grant No. 61906098.

References

- [1] Y. Tian, D. Krishnan, P. Isola, Contrastive multiview coding, in: European Conference on Computer Vision, 2020, pp. 776–794.
- [2] R. D. Hjelm, A. Fedorov, S. Lavoie-Marchildon, K. Grewal, P. Bachman, A. Trischler, Y. Bengio, Learning deep representations by mutual information estimation and maximization, in: International Conference on Learning Representations, 2019, pp. 1–24.
- [3] C.-I. Lai, Contrastive predictive coding based feature for automatic speaker verification, in: arXiv:1904.01575, 2019.
- [4] T. Chen, S. Kornblith, M. Norouzi, G. Hinton, A simple framework for contrastive learning of visual representations, in: International Conference on Machine Learning, 2020, pp. 1597–1607.
- [5] X. Chen, K. He, Exploring simple siamese representation learning, in: IEEE Conference on Computer Vision and Pattern Recognition, 2021, pp. 15750–15758.
- [6] Z. Wu, Y. Xiong, S. X. Yu, D. Lin, Unsupervised feature learning via non-parametric instance discrimination, in: IEEE Conference on Computer Vision and Pattern Recognition, 2018, pp. 3733–3742.
- [7] J.-B. Grill, F. Strub, F. Althé, C. Tallec, P. H. Richemond, E. Buchatskaya, C. Doersch, B. A. Pires, Z. D. Guo, M. G. Azar, B. Piot, K. Kavukcuoglu, R. Munos, M. Valko, Bootstrap your own latent: a new approach to self-supervised learning, in: Advances in Neural Information Processing Systems, 2020, pp. 21271–21284.
- [8] K. He, H. Fan, Y. Wu, S. Xie, R. Girshick, Momentum contrast for unsupervised visual representation learning, in: IEEE Conference on Computer Vision and Pattern Recognition, 2020, pp. 9729–9738.
- [9] N. Komodakis, S. Gidaris, Unsupervised representation learning by predicting image rotations, in: International Conference on Learning Representations, 2018, pp. 1–16.
- [10] M. Gutmann, A. Hyvärinen, Noise-contrastive estimation: A new estimation principle for unnormalized statistical models, in: Artificial Intelligence and Statistics, 2010, pp. 297–304.
- [11] D. Dwibedi, Y. Aytar, J. Tompson, P. Sermanet, A. Zisserman, With a little help from my friends: Nearest-neighbor contrastive learning of visual representations, in: International Conference on Computer Vision, 2021, pp. 9588–9597.
- [12] S. A. Koohpayegani, A. Tejankar, H. Pirsiavash, Mean shift for self-supervised learning, in: International Conference on Computer Vision, 2021, pp. 10326–10335.
- [13] K. Navaneet, S. Abbasi Koohpayegani, A. Tejankar, K. Pourahmadi, A. Subramanya, H. Pirsiavash, Constrained mean shift using distant yet related neighbors for representation learning, in: European Conference on Computer Vision, 2022, pp. 23–41.
- [14] C. GE, J. Wang, Z. Tong, S. Chen, Y. Song, P. Luo, Soft neighbors are positive supporters in contrastive visual representation learning, in: International Conference on Learning Representations, 2023, pp. 1–16.
- [15] H. Lee, S. J. Hwang, J. Shin, Self-supervised label augmentation via input transformations, in: International Conference on Machine Learning, 2020, pp. 5714–5724.
- [16] R. Zhang, P. Isola, A. A. Efros, Colorful image colorization, in: European Conference on Computer Vision, 2016, pp. 649–666.
- [17] M. Noroozi, P. Favaro, Unsupervised learning of visual representations by solving jigsaw puzzles, in: European Conference on Computer Vision, 2016, pp. 69–84.
- [18] C. Doersch, A. Gupta, A. A. Efros, Unsupervised visual representation learning by context prediction, in: International Conference on Computer Vision, 2015, pp. 1422–1430.
- [19] A. Kolesnikov, X. Zhai, L. Beyer, Revisiting self-supervised visual representation learning, in: IEEE Conference on Computer Vision and Pattern Recognition, 2019, pp. 1920–1929.
- [20] Y. Tian, C. Sun, B. Poole, D. Krishnan, C. Schmid, P. Isola, What makes for good views for contrastive learning?, in: Advances in Neural Information Processing Systems, 2020, pp. 6827–6839.
- [21] L. Huang, C. Zhang, H. Zhang, Self-adaptive training: Bridging supervised and self-supervised learning, in: IEEE Transactions on Pattern Analysis and Machine Intelligence, 2022, pp. 1–17.
- [22] Y. Kalantidis, M. B. Sariyildiz, N. Pion, P. Weinzaepfel, D. Larlus, Hard negative mixing for contrastive learning, in: Advances in Neural Information Processing Systems, 2020, pp. 21798–21809.
- [23] Z. Shen, Z. Liu, Z. Liu, M. Savvides, T. Darrell, E. Xing, Un-mix: Rethinking image mixtures for unsupervised visual representation learning, in: Proceedings of the AAAI Conference on Artificial Intelligence, 2022, pp. 2216–2224.
- [24] H. Zhang, M. Cisse, Y. N. Dauphin, D. Lopez-Paz, mixup: Beyond empirical risk minimization, in: International Conference on Learning Representations, 2018, pp. 1–13.
- [25] G. Hinton, O. Vinyals, J. Dean, Distilling the knowledge in a neural network, in: arXiv:1503.02531, 2015.
- [26] J. Ba, R. Caruana, Do deep nets really need to be deep?, in: Advances in Neural Information Processing Systems, 2014, pp. 2654–2662.
- [27] H. Bagherinezhad, M. Horton, M. Rastegari, A. Farhadi, Label refinery: Improving imagenet classification through label progression, in: arXiv:1805.02641, 2018.
- [28] T. Furlanello, Z. Lipton, M. Tschannen, L. Itti, A. Anandkumar, Born again neural networks, in: International Conference on Machine Learning, 2018, pp. 1607–1616.
- [29] A. Tejankar, S. A. Koohpayegani, V. Pillai, P. Favaro, H. Pirsiavash, Isd: Self-supervised learning by iterative similarity distillation, in: International Conference on Computer Vision, 2021, pp. 9609–9618.
- [30] Z. Fang, J. Wang, L. Wang, L. Zhang, Y. Yang, Z. Liu, (SEED): Self-supervised distillation for visual representation, in: International Conference on Learning Representations, 2021, pp. 1–21.
- [31] K. Song, J. Xie, S. Zhang, Z. Luo, Multi-mode online knowledge distillation for self-supervised visual representation learning, in: IEEE Conference on Computer Vision and Pattern Recognition, 2023, pp. 11848–11857.
- [32] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: IEEE Conference on Computer Vision and Pattern Recognition, 2016, pp. 770–778.
- [33] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, N. Houlsby, An image is worth 16x16 words: Transformers for image recognition at scale, in: International Conference on Learning Representations, 2021, pp. 1–21.
- [34] T. Miyato, S.-i. Maeda, M. Koyama, S. Ishii, Virtual adversarial training: a regularization method for supervised and semi-supervised learning, in: IEEE Transactions on Pattern Analysis and Machine Intelligence, 2018, pp. 1979–1993.
- [35] A. Tarvainen, H. Valpola, Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning re-

- sults, in: *Advances in Neural Information Processing Systems*, 2017, pp. 1195–1204.
- [36] C. Wei, H. Wang, W. Shen, A. Yuille, {CO}2: Consistent contrast for unsupervised visual representation learning, in: *International Conference on Learning Representations*, 2021, pp. 1–13.
- [37] M. Zheng, S. You, F. Wang, C. Qian, C. Zhang, X. Wang, C. Xu, Rssl: Relational self-supervised learning with weak augmentation, in: *Advances in Neural Information Processing Systems*, 2021, pp. 2543–2555.
- [38] J. Denize, J. Rabarisoa, A. Orcesi, R. Hérault, S. Canu, Similarity contrastive estimation for self-supervised soft contrastive learning, in: *Winter Conference on Applications of Computer Vision*, 2023, pp. 2706–2716.
- [39] X. Chen, H. Fan, R. Girshick, K. He, Improved baselines with momentum contrastive learning, in: *arXiv:2003.04297*, 2020.
- [40] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, L. Fei-Fei, Imagenet: A large-scale hierarchical image database, in: *IEEE Conference on Computer Vision and Pattern Recognition*, 2009, pp. 248–255.
- [41] A. Krizhevsky, G. Hinton, et al., Learning multiple layers of features from tiny images, 2009, pp. 1–60.
- [42] A. Coates, A. Ng, H. Lee, An analysis of single-layer networks in unsupervised feature learning, in: *Artificial Intelligence and Statistics*, 2011, pp. 215–223.
- [43] Y. Le, X. Yang, Tiny imagenet visual recognition challenge, in: *CS 231N*, 2015.
- [44] C. Feng, I. Patras, Adaptive soft contrastive learning, in: *International Conference on Pattern Recognition*, 2022, pp. 2721–2727.
- [45] S. Ioffe, C. Szegedy, Batch normalization: Accelerating deep network training by reducing internal covariate shift, in: *International Conference on Machine Learning*, 2015, pp. 448–456.
- [46] L. Van der Maaten, G. Hinton, Visualizing data using t-sne., in: *Journal of Machine Learning Research*, 2008, pp. 2579–2605.