

Adaptive Soft Contrastive Learning

Chen Feng

School of Electronic Engineering and Computer Science
Queen Mary University of London
London, UK

Ioannis Patras

School of Electronic Engineering and Computer Science
Queen Mary University of London
London, UK

Abstract—Self-supervised learning has recently achieved great success in representation learning without human annotations. The dominant method – that is contrastive learning, is generally based on instance discrimination tasks, i.e., individual samples are treated as independent categories. However, presuming all the samples are different contradicts the natural grouping of similar samples in common visual datasets, e.g., multiple views of the same dog. To bridge the gap, this paper proposes an adaptive method that introduces soft inter-sample relations, namely Adaptive Soft Contrastive Learning (*ASCL*). More specifically, *ASCL* transforms the original instance discrimination task into a multi-instance soft discrimination task, and adaptively introduces inter-sample relations. As an effective and concise plug-in module for existing self-supervised learning frameworks, *ASCL* achieves the best performance on several benchmarks in terms of both performance and efficiency. Code is available at https://github.com/MrChenFeng/ASCL_ICPR2022.

I. INTRODUCTION

Self-supervised learning learns meaningful representation information through label-independent tasks, achieving performance that approaches or even exceeds that of supervised learning models in many tasks. Early self-supervised learning methods are often based on heuristic tasks, such as the prediction of image rotation angles, while the current mainstream methods are generally based on instance discrimination tasks, i.e., treating each individual instance as a separate semantic class. Methods in this category usually share the same framework, named as *contrastive learning*. For a specific view of a specific instance, they define as positives other views of it and negatives views from other instances, and minimize its distance to positives while maximizing its distance to negatives. Meanwhile, a large number of works have been done to improve this framework, such as using a momentum encoder and memory bank to increase the number of negatives [1].

In this paper, we focus on an inherent deficiency of contrastive learning, namely “class collision” [2], [3]. The instance discrimination hypothesis violates the natural grouping in visual datasets and induces false negatives, e.g., the representation of two similar dogs should be close to each other rather than pushed away. To bridge the gap, we need to introduce meaningful inter-sample relations in contrastive learning.

Debiased contrastive learning [4] proposes a theoretical unbiased approximation of contrastive loss with the simplified hypothesis of the dataset distribution, however, does not address the issue of real false negatives. Some works [5], [6] apply a progressive mechanism to identify and remove false

negatives in the training. NNCLR [7] tries to define extra positives for each specific view by ranking and extracting the top- K neighbors in the learned feature space. Considering soft inter-sample relations, Co2 [8] introduces a consistency regularization enforcing relative distribution consistency of different positive views to all negatives. Clustering-based approaches [9], [10] also provide additional positives, but assuming the entire cluster is positive early in the training is problematic and clustering has an additional computational cost. In addition, all these methods rely on a manually set threshold or a predefined number of neighbors, which is often unknown or hard to determine in advance.

In this work, we propose *ASCL*, an efficient and effective module for current contrastive learning frameworks. We reformulate the contrastive learning problem and introduce inter-sample relations in an adaptive style. To make the training more stable and the inter-sample relationships more accurate, we use weakly augmented views to compute the relative similarity distribution and obtain the sharpened soft label information. Based on the uncertainty of the similarity distribution, we adaptively adjust the weights of the soft labels. In the early stages of training, due to the random initialization, the weights of the soft labels are low and the training of the model will be similar to the original contrastive learning. As the features mature and the soft labels become more concentrated, the model will learn stronger inter-sample relations.

The main contributions of this work are summarized as follows:

- We propose a novel adaptive soft contrastive learning (*ASCL*) method which smoothly alleviates the false negative and over-confidence in the instance discrimination task, and reduces the gap between instance-based learning with cluster-based learning;
- We show that weak augmentation strategies help to stabilize the contrastive training, both in our method and classical contrastive frameworks such as MoCo;
- We show that *ASCL* keeps a high learning speed in the initial epochs compared to two other variants in our work that both try to introduce inter-sample relations in a hard style;
- Our method achieved the state of the art results in various benchmarks with very limited additional cost.

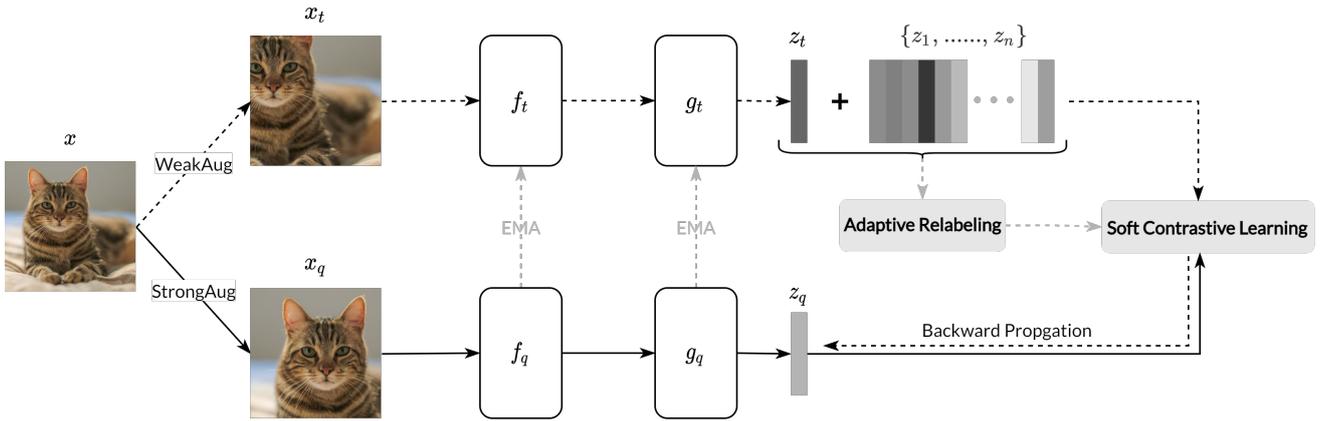


Fig. 1. Structure of ASCL. When we remove the adaptive relabelling step (indicated in light grey), ASCL can be considered as a general contrastive learning framework such as MoCo.

II. RELATED WORKS

A. Early methods on self-supervised learning

Most of the early self-supervised learning methods rely on carefully designed heuristic tasks, such as jigsaw puzzles [11], patch localization [12], image inpainting [13] and rotation prediction [14]. However, these pre-defined tasks lack enough relations to subsequent tasks such as image classification, and are outperformed by contrastive learning methods by a large margin currently.

B. Contrastive learning with instance discrimination

The idea of contrastive learning is first proposed in Hadsell et al. [15] and has become popular again as the dominant self-supervised learning method in recent years, achieving great performance close to or even exceeding that of fully supervised methods on many datasets. Exemplar network [16] proposes to construct a K -way classifier for a dataset with K images by treating every single image and its transformations as a unique surrogate class. Instance discrimination [17] replaces the K -way classifier in [16] with a non-parametric one consisting of all samples' representations and saved in a memory bank, to solve the problem of excessive memory requirement. MoCo [1] is an important baseline for current contrastive learning methods, which reuses the memory bank since samples in a single mini-batch may lead to insufficient negative pairs, and proposes a momentum encoder to update the memory bank in real-time to avoid outdated data representation. SimCLR [18] is another important baseline that finds that setting the mini-batch size to be large enough can eliminate the need for a memory bank. They also used stronger augmentation strategies and replaced the original single linear layer with an MLP.

Meanwhile, several works explore contrastive learning without negative samples. BYOL [19] proposes an asymmetric network structure with a predictor and batch normalization to avoid mode collapse without explicit negative samples, while SimSiam [20] shows that even a momentum encoder is not necessary.

C. Introducing inter-sample relations

Most related to our work are recent works that explore how to introduce inter-sample relations into the original instance discrimination task. NNCLR [7] builds on SimCLR by introducing a memory bank and searches for nearest neighbors to replace the original positive samples, while MeanShift [21] relies on the same idea but builds on BYOL. Co2 [8] proposes an extra regularization term to ensure the relative consistency of both positive views with negative samples, while ReSSL [22] validates that the consistency regularization term itself is enough to learn meaningful representations.

III. ADAPTIVE SOFT CONTRASTIVE LEARNING

Current self-supervised learning methods focus on the instance discrimination task, more specifically, learning by considering each image instance as a separate semantic class. In this work, we follow the representative structure in MoCo [1]. More specifically, given a specific sample x , and two different transformed views of it, as query x_q and target x_t , we want to minimize the distance of the corresponding representation projection z_q and z_t while maximizing the distance of z_q and representations of other samples cached in a memory bank $\{z_1, \dots, z_n\}$. Here $z_- = g(f(x_-))$. The learned representation $f(x_-)$ will be fixed and utilized in subsequent tasks such as image classification with an extra linear classifier [Fig. 1]. With the encoders f_q, f_t and projectors g_q, g_t , we optimize the infoNCE loss:

$$L = -\log \frac{\exp(z_q^T z_t / \tau)}{\exp(z_q^T z_t / \tau) + \sum_{i=1}^n \exp(z_q^T z_i / \tau)} \quad (1)$$

Where τ is a temperature hyperparameter that controls the feature density.

A. Soft contrastive learning

Combining z_t and memory bank $\{z_1, \dots, z_n\}$ together as $\{z'_1, z'_2, \dots, z'_{n+1}\} \triangleq \{z_t, z_1, \dots, z_n\}$ ¹, we can easily rewrite

¹For the convenience, we may use these two notations interchangeably in the following.

eq. 1 below:

$$L = - \sum_{j=1}^{n+1} y_j \log p_j \quad (2)$$

where

$$p_j = \frac{\exp(z_q^T z'_j / \tau)}{\sum_{i=1}^{n+1} \exp(z_q^T z'_i / \tau)} \quad (3)$$

$$y_j = \begin{cases} 1, & j = 1 \\ 0, & \text{otherwise} \end{cases} \quad (4)$$

Here $\mathbf{y} = [y_1, \dots, y_{n+1}]$ is the one-hot **pseudo label** while $\mathbf{p} = [p_1, \dots, p_{n+1}]$ is the corresponding prediction probability vector. Recalling normal supervised learning, prediction over-confidence has inspired research on label smoothing and knowledge distillation. Similarly in self-supervised learning, this problem is more pronounced due to the fact that the distance between individual samples is smaller compared to that between categories, especially when there are duplicate samples or extremely similar samples in the dataset, i.e., the false negatives described earlier. By modifying **pseudo label**, especially the part regarding with other samples, we can convert original contrastive learning problem as a soft contrastive learning problem, with the optimization goal in eq. 2.

B. Adaptive Relabelling

As mentioned above, the **pseudo label** in infoNCE loss ignores the inter-sample relations which will result in false negatives. To address this problem, we propose to modify the **pseudo label** based on the neighboring relations in the feature space. We first calculate the cosine similarity d between self positive view z'_1 and other representations in memory bank $\{z'_2, z'_3, \dots, z'_{n+1}\}$:

$$d_j = \frac{z'_1{}^T z'_j}{\|z'_1\|_2 \|z'_j\|_2}, \quad j = 2, \dots, n+1 \quad (5)$$

1) *Hard relabelling*: According to $d_j, i = 2, \dots, n+1$, we define the top- K nearest neighbors set \mathcal{N}_K in the memory bank of z'_1 as extra positives for z_q . The new **pseudo label** \mathbf{y}_{hard} will be defined as below:

$$y_j = \begin{cases} 1, & j = 1 \text{ or } z_j \in \mathcal{N}_K \\ 0, & \text{otherwise} \end{cases} \quad (6)$$

Intuitively speaking, we consider not only z'_1 as positive for z_q but also the top- K nearest neighbors of z'_1 .

2) *Adaptive hard relabelling*: However, it is risky to recklessly assume that the top- K nearest neighbors are positive, and, especially early in the training, some hard samples may have fewer close neighbors compared to others. To alleviate these problems of \mathbf{y}_{hard} , we propose an adaptive mechanism that automatically modifies the confidence of the **pseudo label**. More specifically, with cosine similarity d we build the

relative distribution \mathbf{q} between self positive view z'_1 and other representations in memory bank $\{z'_2, z'_3, \dots, z'_{n+1}\}$:

$$q_j = \frac{\exp(d_j / \tau')}{\sum_{l=2}^{n+1} \exp(d_l / \tau')}, \quad j = 2, \dots, n+1 \quad (7)$$

To quantify the uncertainty of relative distribution, i.e., how confident when we extract the neighbors, we define a confidence measure as the normalized entropy of the distribution \mathbf{q} :

$$c = 1 - \frac{H(\mathbf{q})}{\log(n)} \quad (8)$$

Here $H(\mathbf{q})$ is the Shannon entropy of \mathbf{q} . We further use $\log(n)$ to normalize c into $[0, 1]$. We then get the adaptive hard label \mathbf{y}_{ahcl} by augmenting \mathbf{y}_{hard} with c :

$$y_j = \begin{cases} 1, & j = 1 \\ c, & j \neq 1 \text{ and } z_j \in \mathcal{N}_K \\ 0, & j \neq 1 \text{ and } z_j \notin \mathcal{N}_K \end{cases} \quad (9)$$

3) *Adaptive soft relabelling*: Moreover, instead of using top- K neighbors for the extra positives, we also propose using the distribution \mathbf{q} itself as soft labels. Intuitively speaking, a more concentrated distribution yields a higher degree of confidence, implying a more reliable neighboring relationship for the sample. We then define the adaptive soft label \mathbf{y}_{ascl} as:

$$y_j = \begin{cases} 1, & j = 1 \\ \min(1, c \cdot K \cdot q_j), & j \neq 1 \end{cases} \quad (10)$$

Here, c is defined in eq. 8 to weight the soft labels, and K is the number of neighbors in \mathcal{N}_K . Please note, that we put an upper bound of one – that means that the most confident positive neighbor is not more confident than a view of the sample itself, i.e., than z'_1 . Here, we

Finally, \mathbf{y}_{ascl} , \mathbf{y}_{ahcl} and \mathbf{y}_{hard} are then normalized, that is:

$$y_j = \frac{y_j}{\sum_{-} y_{-}} \quad (11)$$

For simplicity, we use the same notation for the normalized **pseudo label** as the unnormalized ones. By default we use \mathbf{y}_{ascl} for training — this is the **ASCL** method. We call the training method that uses \mathbf{y}_{ahcl} as **AHCL**, and the one with \mathbf{y}_{hard} as **Hard**. When we set K as zero, the method degenerates to the original MoCo framework.

C. Distribution sharpening

The temperature τ in infoNCE loss (eq. (1)) controls the density of the learned representations. Motivated by current semi-supervised learning works, to filter out possible noisy relations in the feature space we set a smaller temperature τ' (eq. (7)) for relative distribution \mathbf{q} than τ in soft contrastive learning. We $\tau = 0.1, \tau' = 0.05$ by default.

D. Augmentation strategies

Motivated by [21], [22], we also explore different augmentation strategies in *ASCL*. Intuitively, strongly augmented samples have greater randomness and have larger errors in describing inter-sample relationships, while the use of weak augmentation leads to purer nearest-neighbor relationships, which in turn makes training more stable. In our approach, we use weak augmentation for the momentum encoder and memory bank, and strong augmentation for the online encoder.

E. ASCL without negative samples and memory bank

In this section, we propose using ASCL in a contrastive learning framework that does not require a memory bank or use explicit negative samples. More specifically, so as to show that *ASCL* is a flexible framework, we also apply *ASCL* with BYOL [19], which is a representative self-supervised learning work. BYOL uses an extra predictor h and learns by enforcing the consistency between $h(z_q)$ and z_t . We extend this by considering all samples in the batch, labeling them softly according to the distance from the sample in question, and optimizing the consistency according to the soft labels. For more details, please refer to Appendix A.

IV. EXPERIMENTS AND RESULTS

A. Experiment settings

1) Datasets:

a) *CIFAR10 and CIFAR100*: Both CIFAR10 and CIFAR100 consist of 50K training images and 10K test images with 32×32 pixel image resolution. CIFAR10 has 10 classes while CIFAR100 has 100 classes.

b) *STL10*: STL10 consists of 100K unlabeled images, 5K labeled training images, and 8K test images with 96×96 image resolution and 10 classes.

c) *Tiny ImageNet*: Tiny ImageNet consists of 100K training images and 10K validation images with 200 classes. Tiny ImageNet is a lite version of ImageNet with image resizing to 64×64 pixels.

d) *ImageNet-1k*: ImageNet-1K is a large dataset with almost 1.3M images in the training set and 50K images in the validation set, also known as the ILSVRC-2012 dataset.

2) *Implementation details*: For small-scale datasets: CIFAR10, CIFAR100, STL10, and Tiny ImageNet, we apply ResNet-18 as the backbone. To adapt to the low image resolution, we modify the ResNet-18 structure by modifying the first convolutional layer and removing the max-pooling layer. For ImageNet-1k, we apply ResNet-50 as backbone.

For small-scale datasets we set the memory bank size as 4096 while for ImageNet-1k as 65536. The updating momentum for memory encoder is 0.99. For data augmentations, we apply strong augmentations consisting of a random resized crop (range from 0.2 to 1.0, for CIFAR10 and CIFAR100, size 32×32 , for STL10 and Tiny ImageNet, size 64×64 , for ImageNet-1K, size 224×224), horizontal flip (with probability=0.5), color distortion (with strength=0.8), Gaussian blur (with probability=0.5) and grayscale (with probability=0.2). For weak augmentations we only keep random resized crop

and horizontal flip. For model hyperparameters, we set $K = 1$, $\tau = 0.1$ and $\tau' = 0.05$ by default.

We train the network with SGD optimizer for 200 epochs with a momentum of 0.9 and weight decay of $1e-4$. The initial learning rate is 0.06 and is controlled by a cosine annealing scheduler. The batchsize is fixed as 256. Optional batchsize should be equipped with learning rate = $0.06 \times \text{batchsize}/256$.

3) *Evaluation protocol*: We freeze the parameters of encoder f_- and train a linear classifier without the projector g_- . We train the classifier with SGD optimizer for 100 epochs with a momentum of 0.9 and no weight decay. The initial learning rate is 10 and reduced to 1 and 0.1 at the 60th and 80th epochs, respectively. Especially for STL10, we pretrain with both 100K unlabeled data and 5K labeled training images, while for evaluation, we train with only the 5K labeled training images and test on the 8K test images. For online KNN evaluation, we extract the representation for all train samples and apply a distance-weighted KNN classification on the test samples.

B. Results on small-scale and medium-scale datasets

We evaluate our method on small-scale datasets and compare with baselines, in Table I. The results of other methods are copied from the recent work [22] with their best results. For fair comparison, we also report the reproduced results of [22], noted with star marker. It is clear that we perform the best compared to the baselines.

C. Ablations study

1) *Effect of augmentation strategies*: Strong augmentation strategies help model to learn transformation-invariant representations. Some works explore even stronger augmentations for further improvement [23]. By contrast, in *ASCL* we apply a weak augmented view for momentum encoder and memory bank to stabilize the inter-sample relations. We find a stable memory bank with weak augmentations benefits learning. Since the lower performance may come from the slower convergence induced by strong augmentation, we train the model for more epochs to further validate the effect of the weak augmentations.

In Table II, we can see that for both *ASCL* and original MoCo, a more stable memory bank is always beneficial to bring consistent improvement, while *ASCL* always performs better than MoCo in all settings.

2) *Learning speed analysis*: In Fig. 2 we show the online KNN evaluation accuracy of the original MoCo and three variants of our method. As expected, all three variants surpass the original MoCo on both CIFAR10 and CIFAR100 datasets considering the KNN accuracy. However, previous works [7], [21] usually experience a slower learning speed as they described that introducing neighbors makes the task harder. This is validated by the KNN accuracy trend of the *Hard* variant of our method. However, both *ASCL* and *AHCL* keep a high learning speed similar to the original MoCo which implies the adaptive mechanism makes the framework more reliable and efficient, bringing risk-free improvement even when trained with few epochs.

TABLE I
RESULTS ON SMALL-SCALE AND MEDIUM-SCALE DATASETS.

Method	BackProp	EMA	CIFAR-10	CIFAR-100	STL-10	Tiny ImageNet
Supervised	-	No	94.22	74.66	82.55	59.26
SimCLR [18]	2x	No	84.92	59.28	85.48	44.38
BYOL [19]	2x	Yes	85.82	57.75	87.45	42.70
SimSiam [20]	2x	No	88.51	60.00	87.47	37.04
MoCo [1]	1x	Yes	86.18	59.51	85.88	43.36
ReSSL [22]	1x	Yes	90.20	63.79	88.25	46.60
ReSSL(*)	1x	Yes	90.23	64.31	87.69	45.94
ASCL(Ours)	1x	Yes	90.55	65.27	89.54	48.36

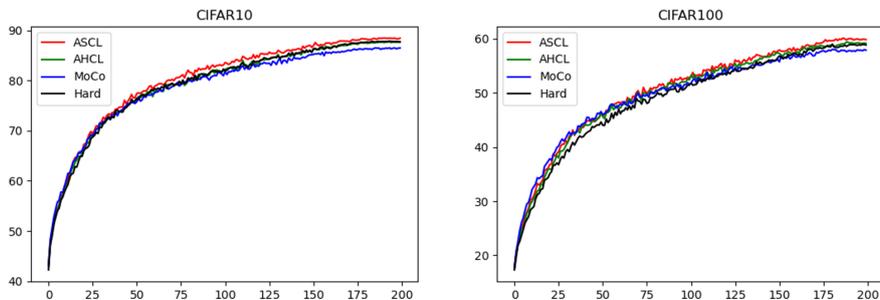


Fig. 2. KNN online evaluation accuracy.

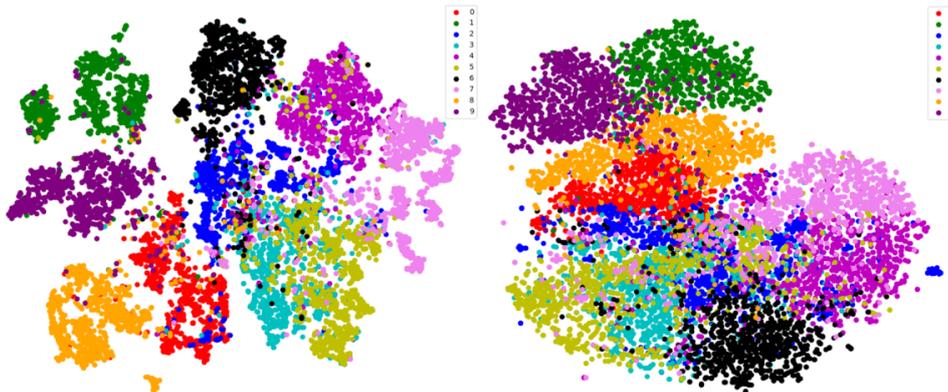


Fig. 3. t-SNE visualization of learned features on CIFAR10, classes indicated by different colors. Left: ASCL, Right: MoCo.

TABLE II
EFFECT OF WEAK AUGMENTATIONS.

Method		CIFAR10		CIFAR100	
		200 epochs	800 epochs	200 epochs	800 epochs
MoCo	Strong	83.48	89.61	57.98	64.00
	Weak	88.55	91.98	64.02	68.16
ASCL	Strong	85.97	91.45	59.34	66.85
	Weak	90.00	92.84	65.46	69.19

3) *Representation visualizations*: In Fig. 3, we compare the representations learned by our method with those learned by MoCo. Different classes of samples are mixed with each other in the learned presentation space of MoCo while we can find clear boundaries in the ASCL representations. It is clear

that the relative relationship of different classes is significantly improved with ASCL.

4) *Appropriate sharpening of inter-sample relations is beneficial*: In most contrastive learning algorithms, the temperature parameter τ is very critical. In ASCL, we use an extra τ' to sharpen the relative distribution (eq. 7), thus to remove the possible noisy inter-sample relations but focus on the most important one. In Table V, we fix $\tau = 0.1$ and perform extensive experiments on CIFAR10 and CIFAR100 with $\tau' = \{0.01, 0.02, 0.05, 0.08, 0.1\}$. It is clear that a too low or too high τ' is not optimal, but still better than the original MoCo (88.55% for CIFAR10, 64.02% for CIFAR100), which again shows the robustness of ASCL. Intuitively speaking, when $\tau' \rightarrow 0$, ASCL is equivalent to the *Hard* with $K = 1$, and when τ' increases, our method is equivalent to introducing uniform distribution, similar to label smoothing in supervised

TABLE III
EFFECT OF DIFFERENT NUMBER OF NEIGHBORS K .

Method	CIFAR10				CIFAR100			
	$K = 1$	$K = 2$	$K = 5$	$K = 10$	$K = 1$	$K = 2$	$K = 5$	$K = 10$
MoCo	88.55				64.02			
Hard	89.52	89.99	90.06	89.99	64.19	64.08	63.82	63.71
AHCL	89.60	89.85	89.92	89.89	64.98	64.72	64.43	63.94
ASCL	90.00	90.09	90.55	90.27	65.46	65.35	64.83	64.39

TABLE IV

Method	Architecture	BackProp	EMA	Batch Size	Epochs	Top-1 Acc
Supervised	ResNet50	1x	No	256	120	76.5
InstDisc [17]	ResNet50	1x	No	256	200	58.5
LocalAgg [24]	ResNet50	1x	NO	128	200	58.8
MoCo [1]	ResNet50	1x	Yes	256	200	67.5
CO2 [8]	ResNet50	1x	No	256	200	68.0
PCL [3]	ResNet50	1x	Yes	256	200	67.6
ReSSL [22]	ResNet50	1x	Yes	256	200	69.9
ASCL(Ours)	ResNet50	1x	Yes	256	200	71.5
SimCLR [18]	ResNet50	2x	No	4096	200	66.8
NNCLR [7]	ResNet50	2x	No	4096	200	70.7
CLSA [23]	ResNet50	2x	Yes	256	200	69.4
SwAV [25]	ResNet50	2x	No	4096	200	69.1
SimSiam [20]	ResNet50	2x	No	256	200	70.0
BYOL [19]	ResNet50	2x	Yes	4096	200	70.6

TABLE V
EFFECT OF DIFFERENT TEMPERATURE τ' .

Dataset	$\tau' = 0.01$	$\tau' = 0.02$	$\tau' = 0.05$	$\tau' = 0.08$	$\tau' = 0.1$
CIFAR10	89.67	89.82	90.00	88.91	88.58
CIFAR100	64.94	64.69	65.46	64.26	64.32

learning.

5) *Robustness to number of neighbors*: In Table III we evaluate our methods with different K — number of neighbors. Please note, for a fair comparison, we report MoCo with weak augmentations for the momentum encoder here. We can find that for *Hard* and *AHCL*, too many neighbors ($K = 10$ for example) result in reduced performance compared to MoCo (which equals introducing no neighbors). The adaptive mechanism helps as *AHCL* is always better than *Hard* while *ASCL* is more robust achieving the best results in all conditions, with further soft pseudo labels.

D. Results on ImageNet-1k

We also evaluate *ASCL* on ImageNet-1k in Table IV. With all methods pretrained for 200 epochs, *ASCL* outperforms the current state-of-the-art methods. Also, please note that *ASCL* requires only one backpropagation pass, which reduces a significant amount of computational cost compared to methods such as BYOL, SimCLR, etc.

E. ASCL without negative pairs and memory bank

In Table VI, we compare the performance of BYOL and BYOL with *ASCL* (appendix A). *ASCL* gets better and worse performance on CIFAR10 and CIFAR100, respectively.

Intuitively, we conjecture that the risk of introducing false positives outweighs the benefits of extra positives. Specifically, for CIFAR100 dataset, we set the batchsize to 256 which is relatively small considering there are 100 different semantic classes. To verify this hypothesis, we increase the batch size to 512 and find that the gap between the two was significantly reduced.

TABLE VI
PERFORMANCE OF ASCL WITH BYOL.

Method	CIFAR10 (batchsize = 256)	CIFAR100 (batchsize = 256)	CIFAR100 (batchsize = 512)
BYOL [19]	86.95	62.31	61.03
BYOL + ASCL	90.49	58.51	60.37

V. CONCLUSIONS

In this work, we propose *ASCL*, a reliable and efficient framework based on the current contrastive learning framework. We utilize a sharpened inter-sample distribution to introduce extra positives and adaptively adjust its confidence based on the entropy of the distribution. Our method achieves the state of the art in various benchmarks, with a negligible extra computational cost. We also show the potential of our method with self-supervised learning methods requiring no memory bank and explicit negative pairs.

Acknowledgments: This work was supported by the EU H2020 AI4Media No. 951911 project.

REFERENCES

- [1] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, “Momentum contrast for unsupervised visual representation learning,” in *Proceedings of the*

- IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 9729–9738.
- [2] S. Arora, H. Khandeparkar, M. Khodak, O. Plevrakis, and N. Saunshi, “A theoretical analysis of contrastive unsupervised representation learning,” *arXiv preprint arXiv:1902.09229*, 2019.
 - [3] J. Li, P. Zhou, C. Xiong, and S. C. Hoi, “Prototypical contrastive learning of unsupervised representations,” *arXiv preprint arXiv:2005.04966*, 2020.
 - [4] C.-Y. Chuang, J. Robinson, L. Yen-Chen, A. Torralba, and S. Jegelka, “Debiased contrastive learning,” *arXiv preprint arXiv:2007.00224*, 2020.
 - [5] T. Huynh, S. Kornblith, M. R. Walter, M. Maire, and M. Khademi, “Boosting contrastive self-supervised learning with false negative cancellation,” in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2022, pp. 2785–2795.
 - [6] T.-S. Chen, W.-C. Hung, H.-Y. Tseng, S.-Y. Chien, and M.-H. Yang, “Incremental false negative detection for contrastive learning,” *arXiv preprint arXiv:2106.03719*, 2021.
 - [7] D. Dwivedi, Y. Aytar, J. Tompson, P. Sermanet, and A. Zisserman, “With a little help from my friends: Nearest-neighbor contrastive learning of visual representations,” *arXiv preprint arXiv:2104.14548*, 2021.
 - [8] C. Wei, H. Wang, W. Shen, and A. Yuille, “Co2: Consistent contrast for unsupervised visual representation learning,” 2020.
 - [9] M. Caron, P. Bojanowski, A. Joulin, and M. Douze, “Deep clustering for unsupervised learning of visual features,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 132–149.
 - [10] Y. M. Asano, C. Rupprecht, and A. Vedaldi, “Self-labelling via simultaneous clustering and representation learning,” *arXiv preprint arXiv:1911.05371*, 2019.
 - [11] M. Noroozi and P. Favaro, “Unsupervised learning of visual representations by solving jigsaw puzzles,” in *European conference on computer vision*. Springer, 2016, pp. 69–84.
 - [12] C. Doersch, A. Gupta, and A. A. Efros, “Unsupervised visual representation learning by context prediction,” in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1422–1430.
 - [13] R. Zhang, P. Isola, and A. A. Efros, “Colorful image colorization,” in *European conference on computer vision*. Springer, 2016, pp. 649–666.
 - [14] S. Gidaris, P. Singh, and N. Komodakis, “Unsupervised representation learning by predicting image rotations,” *arXiv preprint arXiv:1803.07728*, 2018.
 - [15] R. Hadsell, S. Chopra, and Y. LeCun, “Dimensionality reduction by learning an invariant mapping,” in *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR’06)*, vol. 2. IEEE, 2006, pp. 1735–1742.
 - [16] A. Dosovitskiy, J. T. Springenberg, M. Riedmiller, and T. Brox, “Discriminative unsupervised feature learning with convolutional neural networks,” *Advances in neural information processing systems*, vol. 27, pp. 766–774, 2014.
 - [17] Z. Wu, Y. Xiong, S. X. Yu, and D. Lin, “Unsupervised feature learning via non-parametric instance discrimination,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 3733–3742.
 - [18] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, “A simple framework for contrastive learning of visual representations,” in *International conference on machine learning*. PMLR, 2020, pp. 1597–1607.
 - [19] J.-B. Grill, F. Strub, F. Altché, C. Tallec, P. H. Richemond, E. Buchatskaya, C. Doersch, B. A. Pires, Z. D. Guo, M. G. Azar *et al.*, “Bootstrap your own latent: A new approach to self-supervised learning,” *arXiv preprint arXiv:2006.07733*, 2020.
 - [20] X. Chen and K. He, “Exploring simple siamese representation learning,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 15 750–15 758.
 - [21] S. A. Koohpayegani, A. Tejankar, and H. Pirsiavash, “Mean shift for self-supervised learning,” *arXiv preprint arXiv:2105.07269*, 2021.
 - [22] M. Zheng, S. You, F. Wang, C. Qian, C. Zhang, X. Wang, and C. Xu, “Resl: Relational self-supervised learning with weak augmentation,” *Advances in Neural Information Processing Systems*, vol. 34, 2021.
 - [23] X. Wang and G.-J. Qi, “Contrastive learning with stronger augmentations,” *arXiv preprint arXiv:2104.07713*, 2021.
 - [24] C. Zhuang, A. L. Zhai, and D. Yamins, “Local aggregation for unsupervised learning of visual embeddings,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 6002–6012.
 - [25] M. Caron, I. Misra, J. Mairal, P. Goyal, P. Bojanowski, and A. Joulin, “Unsupervised learning of visual features by contrasting cluster assignments,” *arXiv preprint arXiv:2006.09882*, 2020.

A. Details of ASCL with BYOL

We here explain how to apply *ASCL* with BYOL in detail. For implementation details of the BYOL, please refer to the original paper [19]. With a mini-batch of samples as x_1, x_2, \dots, x_b and respective feature projections $z_t^1, z_t^2, \dots, z_t^b$ (from momentum encoder) and $z_q^1, z_q^2, \dots, z_q^b$ (from online encoder), BYOL optimizes the following loss for each sample x_i :

$$L = 2 - 2 \cdot \frac{z_t^{iT} h(z_q^i)}{\|z_t^i\|_2 \|h(z_q^i)\|_2} \quad (12)$$

We can easily reformulate it as below:

$$L = 2 - 2 \cdot \sum_{j=1}^b y_j \frac{z_t^{iT} h(z_q^j)}{\|z_t^i\|_2 \|h(z_q^j)\|_2} \quad (13)$$

With the *pseudo label* y denoting inter-sample relations as:

$$y_j = \begin{cases} 1, & j = i \\ 0, & j = \{1, \dots, b\} \setminus i \end{cases} \quad (14)$$

We then apply *ASCL* to build inter-sample relations directly based on the samples of each mini-batch. For sample x_i , we calculate the cosine similarity between z_t^i and all other projections $\{z_t^1, z_t^2, \dots, z_t^b\} \setminus z_t^i$ in the mini-batch.

$$d_j = \frac{z_t^{iT} z_t^j}{\|z_t^i\|_2 \|z_t^j\|_2}, \quad j = \{1, \dots, b\} \setminus i \quad (15)$$

Similarly, with cosine similarity $d_j, j = \{1, \dots, b\} \setminus i$ we build the relative distribution q :

$$q_j = \frac{\exp(d_j/\tau)}{\sum_{l \in \{1, \dots, b\} \setminus i} \exp(d_l/\tau)}, \quad j = \{1, \dots, b\} \setminus i \quad (16)$$

To quantify the uncertainty of relative distribution, i.e., how confident when we extract the neighbors, we define a similar confidence measure c for in-batch samples:

$$c = 1 - \frac{H(q)}{\log(b-1)} \quad (17)$$

and the corresponding adaptive soft label as:

$$y_j = \begin{cases} 1, & j = i \\ \max(1, c \cdot K \cdot p_k), & j \neq i \end{cases} \quad (18)$$