

BOOSTING FEW-SHOT CLASSIFICATION WITH VIEW-LEARNABLE CONTRASTIVE LEARNING

Xu Luo^{1*}, Yuxuan Chen^{1*}, Liangjian Wen^{1†}, Lili Pan¹, Zenglin Xu^{2,3†}

¹University of Electronic Science and Technology of China, Chengdu, China

²Harbin Institute of Technology Shenzhen, Shenzhen, China

³ Pengcheng Lab, Shenzhen, China

Frank.Luox@outlook.com, {Edmondx.chen, wlj6816}@gmail.com,
lilipan@uestc.edu.cn, xuzenglin@hit.edu.cn

ABSTRACT

The goal of few-shot classification is to classify new categories with few labeled examples within each class. Nowadays, the excellent performance in handling few-shot classification problems is shown by metric-based meta-learning methods. However, it is very hard for previous methods to discriminate the fine-grained sub-categories in the embedding space without fine-grained labels. This may lead to unsatisfactory generalization to fine-grained sub-categories, and thus affects model interpretation. To tackle this problem, we introduce the contrastive loss into few-shot classification for learning latent fine-grained structure in the embedding space. Furthermore, to overcome the drawbacks of random image transformation used in current contrastive learning in producing noisy and inaccurate image pairs (i.e., views), we develop a learning-to-learn algorithm to automatically generate different views of the same image. Extensive experiments on standard few-shot learning benchmarks demonstrate the superiority of our method.

Index Terms— Few-shot learning, contrastive learning, meta-learning

1. INTRODUCTION

Few-shot learning has been widely studied to recognize unseen classes with limited samples for each novel class [1, 2, 3, 4]. Recently, metric-based meta-learning methods have attracted extensive attention in image classification due to their superior performance and simplicity [5, 6, 7]. For making inference, these methods compare the similarity between the feature embedding of query images and that of a few labeled images of each class. This therefore requires learning a flexible encoder, which can map the data points with similar semantics in the input space to locate closely in the embedding space. Meanwhile, those data with different semantic

meanings in the input space should disperse in the embedding space. Accordingly, a new sample from the novel class can be recognized directly through a simple distance metric within the learned embedding space. Indeed, the performance of recognizing novel classes in metric-based meta-learning extremely relies on the learned embedding space.

Despite the success of recognizing novel classes, existing metric-based few-shot approaches often fail to push the fine-grained sub-categories apart in the embedding space given no fine-grained labels in training. For illustration, we merge nine different sub-categories of dogs in the miniImageNet dataset into a coarse-grained class as a new label “dog” to train the Prototypical Network(PN) [6] without changing other classes. As shown in Figure 1(a), we visualize the features of the input data of three fine-grained sub-categories of “dog” using t-SNE. It is clearly revealed that features of fine-grained classes learned by the Prototypical Network cannot be separated without further label information. This means that these methods often do not generalize well to fine-grained sub-categories. Since labeling the fine-grained sub-categories requires strong expertise, the generalization ability to unseen fine-grained sub-categories is of critical importance.

In this paper, we try to alleviate this problem by seeking self-supervised learning to learn the fine-grained structure without given corresponding label information. As a powerful self-supervised representation learning paradigm, Contrastive Learning [8, 9] has outperformed over even supervised learning in many situations. The key insight of contrastive learning is to contrast semantically similar (positive) with dissimilar (negative) pairs of data points. Nice theoretical results for contrastive learning has been given in [10], by hypothesizing that semantically similar points are often sampled from the same latent class. Hence, contrastive learning has the potential to bring closer the representations from the same latent class and to separate those from different latent classes. Therefore, we introduce a contrastive loss into few-shot classification and learn latent fine-grained structures in the embedding space, which helps to cluster samples with similar

*These authors contributed equally to this work.

†Corresponding authors.

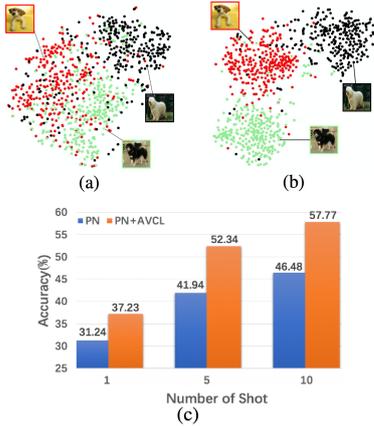


Fig. 1. (a) and (b): T-SNE visualization of feature vectors extracted from three out of nine fine-grained classes of dogs in miniImageNet using Prototypical Network and our proposed method, respectively. (c): Evaluation results on the fine-grained Stanford Dogs dataset [11].

representations to form similar sub-categories.

A critical issue in contrastive learning is generating a pair of semantically similar representations (views) of the same image for contrasting with those dissimilar ones in the embedding space. However, due to the limited number of training images, random image transformation may generate poor positive pairs with more substantial noises and less concept-relevant information when directly applied to few-shot learning. This may therefore make contrastive learning fail to learn fine-grained structure. To effectively improve fine-grained structure learning in the few shot learning setting, we propose view-learnable contrastive learning (VLCL) for metric-based meta-learning. Specifically we replace random image transformation of contrastive learning with spatial transformer network (STN) [12], a learned module that allows flexible spatial manipulation of images, and develop a learning-to-learn algorithm to adaptively generate different views of the same image. In detail, the parameters of STN are optimized through the contribution of the contrastive loss to few-shot image recognition.

To verify that our proposed approach can improve the generalization ability of unseen fine-grained sub-categories without corresponding fine-grained label information, we follow the experiment setting of the coarse-grained class “dog” for training and further test the learned model on Stanford Dogs, a fine-grain dataset containing 120 fine-grained dog classes [11]. As shown in Figure 1(b,c), our proposed method can learn a better embedding space and significantly improve the test accuracy on the Stanford Dogs dataset. For instance, our method obtains an improvement of 10.40% on the 5-way 5-shot task over the Prototypical Network.

2. RELATED WORKS

Few-shot image recognition: Few-shot image recognition was first proposed by [1], with the aim to solve the problem of classifying novel categories with few labeled images per class. Nowadays, two types of meta-learning methods are the mainstream methods to address this problem. One is gradient-based method that empower the model with ability to rapidly fine-tune to novel classes with limited labeled images [2, 3]. The other is metric-based method, which makes predictions based on a similar metric in a learned feature space between images with and without labels. Common similar metric used in this method includes cosine similarity [5], Euclidean similarity [6], relation module [7], and graph neural network [13].

In our work, we primarily consider improving the performance of metric-based meta-learning methods, especially Prototypical Network [6]. Recently, there are many methods proposed to improve the ability of metric-based meta-learning by constraining the structure of the feature space. Li et al. [14] introduce an extra margin loss that leverages external content information, e.g., pre-trained word embeddings, to generate adaptive margin between classes. This leads the feature space to have better semantic structure. Contrastively, our work does not import external information and utilizes semantic information extracted purely from images themselves. Works most relevant to ours is [15] which utilize self-supervised tasks to improve few-shot learning. However, their self-supervised pretext tasks are fixed at the training stage. Our proposed framework can progressively change the view of contrastive learning under a learning-to-learn paradigm.

Contrastive Learning: Contrastive learning is one of the most popular methods for unsupervised visual representation learning. View transformations and contrastive loss are two key parts of contrastive learning. This framework attains representations by optimizing contrastive loss that maximize agreement between transformed views of the same image and minimize agreement between transformed views of different images. Contrastive learning was first proposed by [16]. Recently, Wu et al. [17] consider instance discrimination and use noise contrastive loss to learn representations. Contrastive multiview Coding [18] utilizes contrastive learning to attain representations on multi-view setting. SimCLR [8] summarizes a standard framework of contrastive learning and shows the effect of different random view transformations. However, the performance of SimCLR relies on large batch size. Momentum Contrast (MoCo) [9] is proposed to alleviate this problem by constructing a queue to preserve immediate preceding samples. In our work, due to the relatively small batch size of metric-based meta-learning, we adopt MoCo in our model. Tian et al. [19] also point out the importance of view transformation for contrastive learning, and propose to learn view transformation via information bottleneck principle under unsupervised and semi-supervised settings. Our work employs a learning-to-learn framework to automatically learn

views for better performance for few-shot classification.

3. PROPOSED METHOD

3.1. Problem Setting

A few-shot classification task is characterized as N-way, K-shot only if the model adapt to classify new data after seeing K examples from each of the N classes. Meta-learning algorithms imitate this setting in each iteration by randomly sampling N classes and their corresponding images from a base dataset \mathbf{D}_B to construct a pseudo few-shot task or episode T . Specifically, the input of each episode T can be divided into two sets: (1) Observable support set $\mathbf{D}_S = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^{N \times K}$, formed by randomly selecting K samples from each of the N classes and (2) Unseen query set $\mathbf{D}_Q = \{(\mathbf{x}_j, \mathbf{y}_j)\}_{j=1}^M$, containing other M samples from the same N categories. Given \mathbf{D}_S , the pseudo task is employed by making predictions of \mathbf{D}_Q . Then the labels of query set are used to compute classification loss to guide the update of the model.

3.2. Model description

Figure 2 illustrates the framework of our proposed method. Two complementary classification tasks are employed simultaneously to learn the main encoder $F_\theta(\cdot)$, which is the key component that maps the input into a feature space. One path is metric-based meta-learning, which utilizes explicit label information to regularize the feature space. Another path is contrastive prediction task, which is a self-supervised instance-level classification task. This task is designed to identify latent fine-grained structure in the feature space by aggregating the representations of the same latent class and separating those of different latent classes at the same time.

In an episode T, sampled images in \mathbf{D}_S and \mathbf{D}_Q are used in the two paths. In the metric-based meta-learning path, the main encoder $F_\theta(\cdot)$ first maps all images into the feature space, and then all support features in the same class C_i aggregate into one vector h_i . Typically, this is accomplished by averaging all support features, i.e. $\mathbf{h}_i = \frac{1}{K} \sum_{(x,y) \in \mathbf{D}_S} \mathbb{1}_{[y=i]} F_\theta(x)$. It is followed by computing the similarities between query features and aggregated features in each class. The final classification loss \mathcal{L}^{meta} is defined as average cross entropy between true labels and predictions based on similarities. This can be formulated as:

$$\mathcal{L}^{meta}(\mathbf{D}_S, \mathbf{D}_Q, \theta) = -\frac{1}{M} \sum_{(x,y) \in \mathbf{D}_Q} \log \frac{e^{\text{sim}(F_\theta(x), h_y)}}{\sum_{i=1}^N e^{\text{sim}(F_\theta(x), h_i)}}, \quad (1)$$

where sim denotes a similarity metric.

At the core of the contrastive path is the process of producing two views from one image which lie in the same latent space, i.e. containing similar semantic contents. This process is accomplished by two differentiable auto-view modules $\mathcal{G}_{\gamma_1}(\cdot)$ and $\mathcal{G}_{\gamma_2}(\cdot)$ parameterized by γ_1 and γ_2 , respectively.

Each of them is a spatial transformer network(STN) [12] that allows flexible semantics-invariant spatial manipulation of images, including cropping, one dimensional scaling, As for translation, image deformation and proportional shrinkage, see appendix for more details of STN. They are applied to each image \mathbf{x}_i , producing two views $\mathbf{x}_i^{(1)}$ and $\mathbf{x}_i^{(2)}$. These views are then mapped into feature space by the main encoder $F_\theta(\cdot)$ and momentum encoder $F_\omega(\cdot)$, respectively. As mentioned in [9], the momentum encoder’s parameter ω is a moving average of θ , which makes two encoders behave similar, see the appendix for more details of Moco. Given $F_\theta(\mathbf{x}_i^{(1)})$, the contrastive loss aims to identify $F_\omega(\mathbf{x}_i^{(2)})$ in thousands of features $\{F_\omega(\mathbf{x}_k^{(2)})\}_{k \neq i}$, and can be formulated as:

$$\begin{aligned} \mathcal{L}^{con}(\mathbf{D}_S, \mathbf{D}_Q, \omega, \theta, \gamma) \\ = - \sum_{x \in \mathbf{D}_S \cup \mathbf{D}_Q} \log \frac{e^{\text{sim}(F_\theta(\mathbf{x}^{(1)}), F_\omega(\mathbf{x}^{(2)}))}}{\sum_{j=1}^r e^{\text{sim}(F_\theta(\mathbf{x}^{(1)}), F_\omega(\mathbf{x}_j^{(2)}))}}, \end{aligned} \quad (2)$$

where r denotes the number of negative samples, and $\gamma = [\gamma_1, \gamma_2]$ denotes the parameters of auto-view modules. By minimizing \mathcal{L}^{con} w.r.t θ , we force the main encoder $F_\theta(\cdot)$ to map views of one image which are semantically similar into closer points in the feature space, thus constructs a better fine-grained semantic structure.

3.3. Learning strategy

The optimization of our model during each iteration contains two stages. We denote $\theta^t, \omega^t, \gamma^t$ as parameters and $\mathbf{D}_S^t, \mathbf{D}_Q^t$ as support set and query set during iteration t, respectively. We first update two encoders based on the meta loss \mathcal{L}^{meta} and contrastive loss \mathcal{L}^{con} :

$$\mathcal{L}^{total} = \mathcal{L}^{meta}(\mathbf{D}_S^t, \mathbf{D}_Q^t, \theta^t) + \beta \mathcal{L}^{con}(\mathbf{D}_S^t, \mathbf{D}_Q^t, \omega^t, \theta^t, \gamma^t), \quad (3)$$

$$\theta^{t+1} = \theta^t - \alpha \nabla_{\theta^t} \mathcal{L}^{total}, \quad (4)$$

$$\omega^{t+1} = \epsilon \omega^t + (1 - \epsilon) \theta^{t+1}, \quad (5)$$

where β denotes the regularization hyperparameter weighting two losses, α denotes the learning rate of θ , and ϵ denotes momentum coefficient that controls the chasing speed of the momentum encoder $F_\omega(\cdot)$. Note that the value of θ^{t+1} relies on γ^t via the contrastive loss \mathcal{L}^{con} . This fact is crucial for the update of γ^t in the next stage.

The updating criterion for the auto-view modules $G_{\gamma_1}(\cdot)$ and $G_{\gamma_2}(\cdot)$ is to improve the positive effect of contrastive loss for meta-learning. Keeping this in mind, we use the updated encoder $F_{\theta^{t+1}}(\cdot)$ to compute meta loss again on the same task, and cast the loss as evaluation of update quality in the first stage. Then the loss is used to guide the update of the auto-view modules via gradient descent:

$$\gamma^{t+1} = \gamma^t - \eta \nabla_{\gamma^t} \mathcal{L}^{meta}(\mathbf{D}_S^t, \mathbf{D}_Q^t, \theta^{t+1}), \quad (6)$$

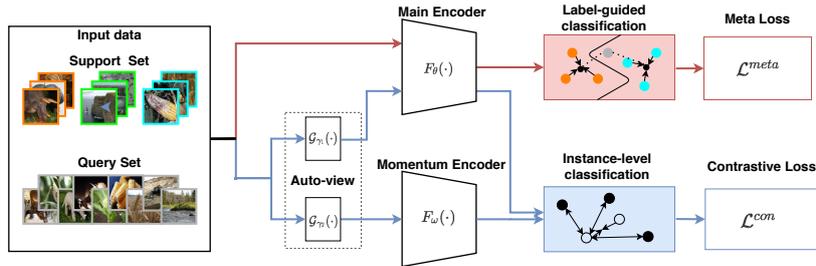


Fig. 2. A brief flow diagram of our model at training stage. Our framework consists of two main tasks. The red arrow is metric-based meta-learning, while blue arrow depicts instance-level contrastive learning. At meta-test stage, only the main encoder $F_{\theta}(\cdot)$ is held for evaluation under few-shot setting.

where η is the learning rate of the auto-view modules. We mention again that θ^{t+1} can be cast as a function of γ^t . Thus the loss depends on γ^t through the computation graph of θ^{t+1} in first stage. The update of γ^t adjusts produced views towards better update of encoder $F_{\theta}(\cdot)$ in the first stage. This indicates that learned views further pushes positive effect of contrastive loss for meta-learning.

4. EXPERIMENTAL RESULTS

In this section, we conduct experiments to demonstrate the effectiveness of our method. Our proposed method is evaluated on standard few-shot learning benchmarks. We also conduct ablative study about the auto-view module in our learning framework. To further verify the effectiveness of our method, we evaluate our method on three few-shot fine-grained datasets in appendix.

4.1. Experimental Setup

We conduct 5-way 5-shot and 5-way 1-shot classification for all datasets. The metric-based meta-learning method adopted in our model is Prototypical Network, one of the state-of-the-art metric-based meta-learning methods for few-shot learning.

Datasets: We follow the general few-shot image recognition settings and evaluate our method on two benchmarks: **miniImageNet** [5] and **CUB-200-2011** [20]. miniImageNet is selected from the well-known ImageNet[21] dataset. CUB-200-2011 was originally proposed for fine-grained bird classification.

Auto-view transformation architectures: We use 4-layer convolutional nets with 64 channels as localisation networks in STN modules that receive images of size 80×80 and output 4-dimensional vectors. The vectors are used as parameters of two diagonal affine transformations which apply to images and produce two views. When implementing without auto-view module, we replace the module with a random cropping function followed by resize to original size. Regardless of having this module or not, we randomly apply color

jitter, gaussian blur and horizontal flip before it, and the hyperparameters are the same as in [9].

Feature extractor architectures: In all our experiments, the main encoder and momentum encoder share the same architecture. Following prior work [22, 15], we use a 2-layer Wide Residual Network(WRN-28-10) that outputs 640-dimensional feature vectors after global pooling given images(or views) of size 80×80 . This feature space is directly used for metric-based meta-learning, but will be further mapped by a 2-layer mlp project head to a 128-dimensional hidden space for contrastive learning.

Training details: At training, each minibatch contains 4 tasks, and classes for each task are randomly selected from training set. The query set contains 4 samples during meta-training and 16 samples during meta-testing. All samples from query set and support set are used for computing contrastive loss. All learnable components of our model are trained for 60 epochs by SGD optimizer with Nesterov momentum 0.9 and weight decay 0.0005. The learning rate for STN was set to 0.00001, the same magnitude as in the original paper. The learning rate for other parts of model was initially set to 0.1, and then changed to 0.01 and 0.001 at epochs 20 and 40, respectively. Moreover, regularization hyperparameter β was set to 2.0. We use a queue containing 63000 negative samples for contrastive learning. Momentum coefficient ϵ for updating momentum encoder was set to 0.999, following [9].

4.2. Evaluation on Benchmarks

Below we report comparative results on two benchmarks for FSL: MiniImageNet and CUB in Table 1 and 2. In Table 1 we divide methods into two groups and compare them with our proposed method, respectively. The first group contains recent comparative few-shot learning methods. The second group contains baseline method(Prototypical Network) and other methods that aim at improving it. PN+CL denotes models that replace auto-view modules with random cropping functions. Analysis from the results, we can find that: (1)Our method consistently improves the baseline method (proto-

Table 1. Comparative results for 5-way classification on miniImageNet. Average accuracies on the meta-test set with 95 confidence interval are reported. † denotes methods using external text information. ‡ denotes result reported in [15].

Model	1-shot	5-shot
MAML [2]	48.70 ± 1.84	63.11 ± 0.92
Matching Network [5]	43.56 ± 0.84	55.31 ± 0.73
Relation Networks [7]	50.44 ± 0.82	65.32 ± 0.70
IDeMe-Net [23]	59.14 ± 0.86	74.63 ± 0.74
PPA [24]	59.60 ± 0.41	73.74 ± 0.19
PN‡ [6]	55.85 ± 0.48	68.72 ± 0.36
PN+TRAML† [14]	60.31 ± 0.48	77.94 ± 0.57
SEN PN [25]	-	72.3
PN+rotation [15]	58.28 ± 0.49	72.13 ± 0.38
PN + CL(ours)	59.54 ± 0.47	74.46 ± 0.52
PN + VLCL(ours)	61.75 ± 0.43	76.32 ± 0.49

Table 2. Comparative results for 5-way classification on CUB. Average accuracies on the meta-test set with 95 confidence interval are reported.

Model	1-shot	5-shot
AFHN [26]	70.53 ± 1.01	83.95 ± 0.63
FEAT [27]	68.87 ± 0.22	82.90 ± 0.15
MAML [28]	67.28 ± 1.08	83.47 ± 0.59
cosine classifier [28]	68.00 ± 0.83	84.50 ± 0.51
Relationnet [28]	66.20 ± 0.99	82.30 ± 0.58
DEML [27]	66.95 ± 1.06	77.11 ± 0.78
PN [6]	66.08 ± 0.54	78.79 ± 0.23
PN + CL(ours)	70.45 ± 0.41	82.67 ± 0.46
PN+VLCL(ours)	71.21 ± 0.43	85.08 ± 0.36

typical Network). For instance, our model boosts performance of Prototypical Network on miniImageNet by 5.90% and 7.60% under the 1-shot and 5-shot settings, respectively. This verifies that our method can indeed improve model performance by refining fine-grained semantic structure of the feature space. (2) Our model outperforms recent comparable few-shot learning methods and also outperforms other approaches that aim at improving Prototypical Network. Moreover, we achieve competitive performance with the method using external text information (PN+TRAML). This further gives evidence of the superiority of our learned feature space for FSL. (3) Our auto-view module can indeed improve the quality of views, thus reach a better performance. Compared to random views, our VLCL method obtains 2.21% and 1.86% performance gains under the 5-way 1-shot and 5-shot settings on miniImageNet, respectively.

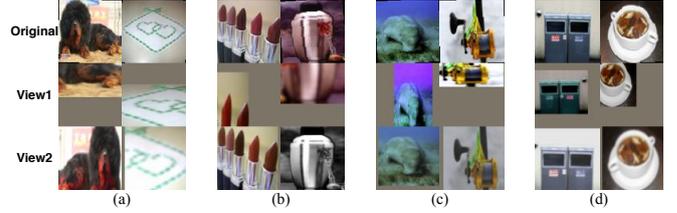


Fig. 3. Four types of views learned by our method: (a) Local-to-local, (b) Global-to-local, (c) One-dimension Scaling, and (d) Proportional zooming.

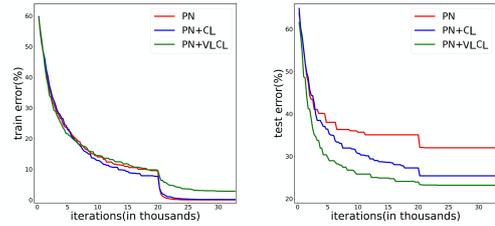


Fig. 4. Training errors (left) and test errors (right) on miniImageNet. The auto-view module significantly decreases test errors, while keeping not overfit to the training set.

4.3. Evaluation of view-learnable learning

In Fig. 3 we show four distinct types of views produced by our auto-view module. We first notice that every pair of learned views are different but both captures the main information of the image, which is expected [19]. We further notice that, aside from local-to-local and global-to-local views which can also be accomplished by traditional random cropping, our auto-view module additionally allows one dimension scaling, translation transformation, image deformation and proportional shrinkage. This flexibility can enrich semantics-invariant transformations applied to the images, forcing the encoder to extract essential content of the image. This allows the samples in the feature space to be distributed according to their semantics. Thus images from the same novel category can be mapped to close points in the feature space, which greatly improves generalization capability.

We additionally show training and test errors during training on miniImageNet in Figure 4. It can be observed that the curves of training error are similar, while the curves of test errors are different. While contrastive regularization helps the model generalize better, our auto-view module further improves it. This strongly supports our motivation that encoding fine-grained semantic contents can help metric-based meta-learning generalize better to novel classes in FSL.

5. CONCLUSION

In this paper, we propose view-learnable contrastive learning to improve few-shot image recognition. In particular, we design a learning-to-learn algorithm to adaptively learn the views. We carry out two paths of tasks, one is label-guided metric-based meta-learning, another is instance-level classification for exploring fine-grained semantic structure of feature space. Extensive experiments on benchmarks demonstrate that our method effectively boosts performance of metric-based meta-learning.

Acknowledgments

This work was partially supported by the National Key Research and Development Program of China (No. 2018AAA0100204) and a fundamental Program of Shenzhen Science and Technology Innovation Commission (No. ZX20210035).

6. REFERENCES

- [1] F. Li, R. Fergus, and P. Perona, “One-shot learning of object categories,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 28, no. 4, pp. 594–611, 2006.
- [2] C. Finn, P. Abbeel, and S. Levine, “Model-agnostic meta-learning for fast adaptation of deep networks,” in *ICML*, vol. 70, 2017, pp. 1126–1135.
- [3] K. Lee, S. Maji, A. Ravichandran, and S. Soatto, “Meta-learning with differentiable convex optimization,” in *CVPR*, 2019, pp. 10 657–10 665.
- [4] H. Tseng, H. Lee, J. Huang, and M. Yang, “Cross-domain few-shot classification via learned feature-wise transformation,” in *ICLR*, 2020.
- [5] O. Vinyals, C. Blundell, T. Lillicrap, K. Kavukcuoglu, and D. Wierstra, “Matching networks for one shot learning,” in *NeurIPS*, 2016, pp. 3630–3638.
- [6] J. Snell, K. Swersky, and R. S. Zemel, “Prototypical networks for few-shot learning,” in *NeurIPS*, 2017, pp. 4077–4087.
- [7] F. Sung, Y. Yang, L. Zhang, T. Xiang, P. H. S. Torr, and T. M. Hospedales, “Learning to compare: Relation network for few-shot learning,” in *CVPR*, 2018, pp. 1199–1208.
- [8] T. Chen, S. Kornblith, M. Norouzi, and G. E. Hinton, “A simple framework for contrastive learning of visual representations,” *CoRR*, vol. abs/2002.05709, 2020.
- [9] K. He, H. Fan, Y. Wu, S. Xie, and R. B. Girshick, “Momentum contrast for unsupervised visual representation learning,” in *CVPR*, 2020, pp. 9726–9735.
- [10] N. Saunshi, O. Plevrakis, S. Arora, M. Khodak, and H. Khandeparkar, “A theoretical analysis of contrastive unsupervised representation learning,” in *ICML*, vol. 97, 2019, pp. 5628–5637.
- [11] A. Khosla, N. Jayadevaprakash, B. Yao, and L. Fei-Fei, “Novel dataset for fine-grained image categorization,” in *CVPR, workshop*, Colorado Springs, CO, June 2011.
- [12] M. Jaderberg, K. Simonyan, A. Zisserman, and K. Kavukcuoglu, “Spatial transformer networks,” in *NeurIPS*, 2015, pp. 2017–2025.
- [13] V. G. Satorras and J. B. Estrach, “Few-shot learning with graph neural networks,” in *ICLR*, 2018.
- [14] A. Li, W. Huang, X. Lan, J. Feng, Z. Li, and L. Wang, “Boosting few-shot learning with adaptive margin loss,” in *CVPR*, 2020, pp. 12 573–12 581.
- [15] S. Gidaris, A. Bursuc, N. Komodakis, P. Pérez, and M. Cord, “Boosting few-shot visual learning with self-supervision,” in *ICCV*, 2019, pp. 8058–8067.
- [16] R. Hadsell, S. Chopra, and Y. LeCun, “Dimensionality reduction by learning an invariant mapping,” in *CVPR*, 2006, pp. 1735–1742.
- [17] Z. Wu, Y. Xiong, S. X. Yu, and D. Lin, “Unsupervised feature learning via non-parametric instance discrimination,” in *CVPR*, 2018, pp. 3733–3742.
- [18] Y. Tian, D. Krishnan, and P. Isola, “Contrastive multi-view coding,” *CoRR*, vol. abs/1906.05849, 2019.
- [19] Y. Tian, C. Sun, B. Poole, D. Krishnan, C. Schmid, and P. Isola, “What makes for good views for contrastive learning,” *CoRR*, vol. abs/2005.10243, 2020.
- [20] P. Welinder, S. Branson, T. Mita, C. Wah, F. Schroff, S. Belongie, and P. Perona, “Caltech-ucsd birds 200,” *Technical Report CNS-TR-2010-001*, California Institute of Technology, 2010.
- [21] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. S. Bernstein, A. C. Berg, and F. Li, “Imagenet large scale visual recognition challenge,” *Int. J. Comput. Vis.*, vol. 115, no. 3, pp. 211–252, 2015.
- [22] S. Gidaris and N. Komodakis, “Generating classification weights with GNN denoising autoencoders for few-shot learning,” in *CVPR*, 2019, pp. 21–30.
- [23] Z. Chen, Y. Fu, Y. Wang, L. Ma, W. Liu, and M. Hebert, “Image deformation meta-networks for one-shot learning,” in *CVPR*, 2019, pp. 8680–8689.
- [24] S. Qiao, C. Liu, W. Shen, and A. L. Yuille, “Few-shot image recognition by predicting parameters from activations,” in *CVPR*, 2018, pp. 7229–7238.
- [25] D. R. Kampffmeyer and R. Jenssen, “Sen: A novel feature normalization dissimilarity measure for prototypical few-shot learning networks,” in *ECCV*, 2020.
- [26] K. Li, Y. Zhang, K. Li, and Y. Fu, “Adversarial feature hallucination networks for few-shot learning,” in *CVPR*, 2020, pp. 13 467–13 476.
- [27] H. Ye, H. Hu, D. Zhan, and F. Sha, “Learning embedding adaptation for few-shot learning,” *CoRR*, vol. abs/1812.03664, 2018.
- [28] W. Chen, Y. Liu, Z. Kira, Y. F. Wang, and J. Huang, “A closer look at few-shot classification,” in *ICLR*, 2019.

- [29] J. Krause, M. Stark, J. Deng, and L. Fei-Fei, “3d object representations for fine-grained categorization,” in *ICCV*, 2013, pp. 554–561.
- [30] B. Zhou, À. Lapedriza, A. Khosla, A. Oliva, and A. Torralba, “Places: A 10 million image database for scene recognition,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 40, no. 6, pp. 1452–1464, 2018.
- [31] G. V. Horn, O. M. Aodha, Y. Song, Y. Cui, C. Sun, A. Shepard, H. Adam, P. Perona, and S. J. Belongie, “The inaturalist species classification and detection dataset,” in *CVPR*, 2018, pp. 8769–8778.

A. LEARNING PROCEDURE

The pseudo code of our learning procedure is shown in Algorithm 1,

Algorithm 1 View-Learnable Contrastive Learning (VLCL) for metric-based meta-learning

- 1: **Require:** Base dataset \mathbf{D}_B , learning rate α and η , weight hyperparameter β , momentum coefficient ϵ , and maximum iteration number t_{\max}
 - 2: Random initialization for $\theta^0, \omega^0, \gamma^0$
 - 3: **for** $t = 0$ to t_{\max} **do**
 - 4: /* Sample tasks */
 Randomly sample N classes from \mathbf{D}_B .
 Randomly sample K images from each class in \mathbf{D}_B to form \mathbf{D}_S^t
 Randomly sample other M images from the same N classes in \mathbf{D}_B to form \mathbf{D}_Q^t
 - 5: /* First forward pass */
 Using θ^t and ω^t to compute $\mathcal{L}^{\text{meta}}$ and \mathcal{L}^{con} through Eq. (1) and Eq. (2)
 - 6: /* Optimize main encoder $F_\theta(\cdot)$, project head $g_\theta(\cdot)$ and momentum encoder $F_\omega(\cdot)$ */
 Update (θ^t, ω^t) to $(\theta^{t+1}, \omega^{t+1})$ through Eq. (3) and Eq.(4), and retain computational graph.
 - 7: /* Second forward pass */
 Using θ^{t+1} to compute $\mathcal{L}^{\text{meta}}$ through Eq. (1)
 - 8: /* Optimize spatial transformation module $G_{\gamma_1}(\cdot)$ and $G_{\gamma_2}(\cdot)$ */
 Update (γ_1^t, γ_2^t) to $(\gamma_1^{t+1}, \gamma_2^{t+1})$ through Eq. (5)
 - 9: **end for**
-

B. SPATIAL TRANSFORMER NETWORKS

In spatial transformer networks [12], the input source image \mathbf{x}^s is first fed into a localisation net $G_\gamma(\cdot)$ and outputs six affine transformation parameters. This parameters form a 2×3 matrix which defines a affine transformation mapping each pixel coordinates (u_i^t, v_i^t) in the output \mathbf{x}^t to a source coordinates (u_i^s, v_i^s) in the input. In our setting, we constrain the matrix to be diagonal so as to avoid skewing which could

Table 3. Results for 5-way few-shot classification on three fine-grained datasets: Cars, Places and Plantae. Average accuracies on the meta-test set are reported.

Model	Cars		Places		Plantae	
	1-shot	5-shot	1-shot	5-shot	1-shot	5-shot
PN	62.11	75.83	62.19	76.67	53.59	67.98
PN+VLCL	76.93	87.04	64.50	78.96	59.80	75.37

possibly change the semantics of images: $\{\lambda_{i,j}\}_{i,j \in \{1,2,3\}}$.

$$\begin{pmatrix} u_i^s \\ v_i^s \end{pmatrix} = \tau_\lambda(u_t^s, v_t^s) = \begin{bmatrix} \lambda_{11} & 0 & \lambda_{13} \\ 0 & \lambda_{22} & \lambda_{23} \end{bmatrix} \begin{pmatrix} u_i^t \\ v_i^t \\ 1 \end{pmatrix}$$

Finally, the values of each pixels in \mathbf{x}^t is determined by bilinear interpolation at their corresponding coordinates in the source images, called differentiable image sampling.

C. MOMENTUM CONTRAST

Different from the standard framework in SimCLR [8], momentum contrast [9] framework introduces a queue q preserving negative samples and a momentum encoder $F_\omega(\cdot)$, to alleviate the problem of need for very large batch size for contrastive learning. In each iteration, immediate preceding features in the queue encoded by $F_\omega(\cdot)$ could be reused as negative samples to compute the contrastive loss. At the end of each iteration, features of current mini-batch is enqueued to the queue, and earliest features in the queue are removed. The update of the encoder $F_\omega(\cdot)$ is intractable by back-propagation. To maintain consistency, MoCo updates ω as a moving average of the main encoder’s parameter θ , as shown in eq. (4).

D. RESULTS FOR FEW-SHOT FINE-GRAINED CLASSIFICATION

Fine-grained categories are distinguished by subtle and local semantic differences, which makes few-shot fine-grained classification more difficult. We experimentally show that such difficulty can be largely addressed by our method. Table 1 presents 5-way mean accuracy on three datasets with fine-grained categories: Cars [29], Places [30] and Plantae [31]. It can be observed that our method improved performance of Prototypical Network by a large margin under both 5-shot and 1-shot settings. For instance, our method obtains 14.83% and 11.21% gains under 1-shot and 5-shot settings on Cars, respectively. This verifies that metric-based meta-learning benefits from better fine-grained semantic structure learnt by our method.

Table 4. Effects of the β value in VLCL on model performances.

	5-way Acc.	
	1-shot	5-shot
$\beta = 0.5$	$58.84 \pm 0.39\%$	$74.21 \pm 0.55\%$
$\beta = 1.0$	$59.45 \pm 0.36\%$	$75.64 \pm 0.71\%$
$\beta = 2.0$	$61.75 \pm 0.43\%$	$77.19 \pm 0.51\%$
$\beta = 5.0$	$58.73 \pm 0.50\%$	$74.37 \pm 0.77\%$

E. EFFECT OF REGULARIZATION HYPERPARAMETER

We perform ablation study w.r.t. the regularization hyperparameter β which controls the magnitude of contrastive loss. Table 2 shows the accuracies of 5-way few-shot learning on miniImageNet. We can observe that when β is small, the contrastive loss cannot thoroughly explore the semantics inside data, thus cannot boost the performance much. When the value of β is too large, the accuracy also decreases. This implies that supervised information is somewhat ignored, which is still important.