

Multi-level Similarity Learning for Low-Shot Recognition

Hongwei Xu, Xin Sun, Junyu Dong, Shu Zhang, Qiong Li

Department of Computer Science and Technology

Ocean University of China

Qingdao, China

{xhw, liqiong}@stu.ouc.edu.cn, {sunxin, dongjunyu, zhangshu}@ouc.edu.cn

Abstract—Low-shot learning indicates the ability to recognize unseen objects based on very limited labeled training samples, which simulates human visual intelligence. According to this concept, we propose a multi-level similarity model (MLSM) to capture the deep encoded distance metric between the support and query samples. Our approach is achieved based on the fact that the image similarity learning can be decomposed into image-level, global-level, and object-level. Once the similarity function is established, MLSM will be able to classify images for unseen classes by computing the similarity scores between a limited number of labeled samples and the target images. Furthermore, we conduct 5-way experiments with both 1-shot and 5-shot setting on Caltech-UCSD datasets. It is demonstrated that the proposed model can achieve promising results compared with the existing methods in practical applications.

Index Terms—Multi-level, Low-shot learning, Similarity learning

I. INTRODUCTION

Deep convolutional neural networks (ConvNets) have achieved great success on visual recognition tasks, such as face recognition [1], object detection [2] and many others. Among those tasks, the success of visual recognition systems greatly relies on the large quantities of annotated data, e.g., ImageNet [3]. It usually needs thousands of labeled examples for each class to saturate ConvNets' performance [4]. However, in practice, it might be extremely expensive or infeasible to obtain sufficient labeled images. On the other hand, the human perception system can easily understand unseen concepts with little knowledge training, especially for the domain-specific low-shot tasks (e.g., Ornithologists will recognize new birds more quickly than ordinary people). This challenge of learning new concepts from a very limited number of labeled samples often is referred to as *low-shot learning* or *few-shot learning*. And this is what this paper focuses on.

One way to address this problem is to utilize the strategy of transfer learning—a fine-tuning network using a few labeled samples for a new classification category. However, the whole network can be broken down by very few data due to the overfitting. With an urgent need on low-shot learning in practice, there are two research approaches fall under the umbrella of *meta-learning* [5], [6] and *metric learning* [7],

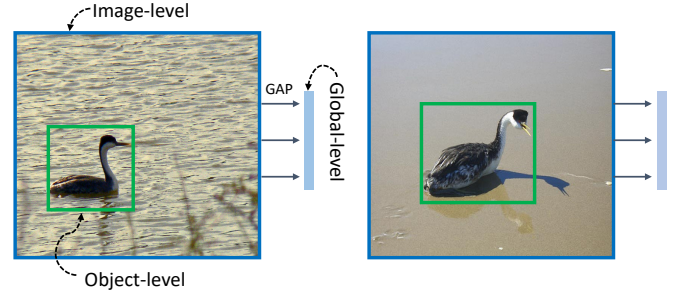


Fig. 1. Three different levels of similarity between images.

[8]. Approaches of meta-learning often train a meta-learner, which learns to directly map between the training sets and the testing examples for classification. Meanwhile, there are other meta-learning methods such as Meta-SGD [9] generate an adaptive learning rate by adjusting differentiable learner in one or few steps for training the classifier. Metric-learning tries to solve the low-shot problems by learning an embedding function to map a few numbers of labeled images into an embedding space. Then the learned space can classify test images via a nearest neighbor algorithm based on the distance measurements such as Euclidean or cosine distance.

In this paper, we adopt the metric learning method. The proposed strategy is based on the following facts. When a person observes the correlation between two images, they first directly get the concept of the overall similarity between the two images (image-level information as shown in Fig. 1). Then his/her attention will move to the target object existed in the images (object-level information as shown in Fig. 1) to measure the target object similarity. In addition, the global-level information is a compressed version of overall image-information, it is a complement to image-level information. Inspiringly, we propose a Multi-level Similarity Model (MLSM) that performs a low-shot recognition by learning how to compare the test images against low-shot labeled samples. The element-sum operation will be performed on these three levels of similarities. The similar function is achieved by a fully connected network.

The main contribution of this paper is to propose a multi-level relation model obtained from training classes. It com-

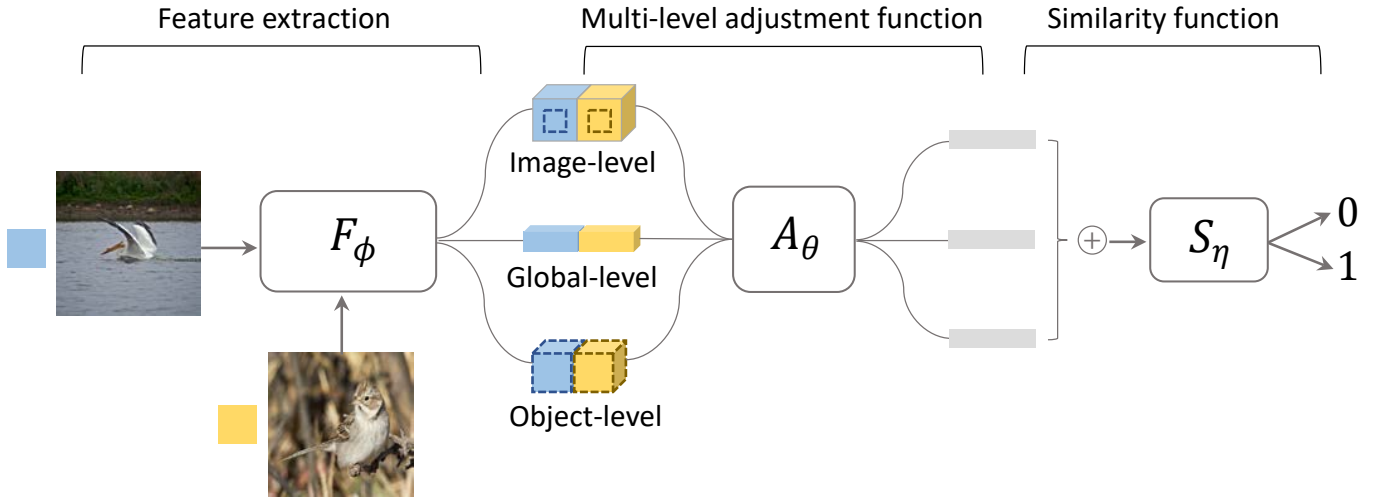


Fig. 2. The overall architecture of MLSM. The blue and yellow small squares represent different images, respectively, they can be the same category or not.

putes relation scores between the test images and the limited number of examples. The experiment shows that the proposed model could achieve promising result compare with the state-of-the-art methods on Caltech-UCSD Birds datasets.

II. RELATED WORK

A. Few-shot Learning

Amongst representation learning approaches, metric-learning is widely adopted, such as siamese networks [10] or triplet loss [11]. They are used to minimize the distance between the inputs and the samples with the same label while keeping the distances of samples with different labels far apart. Such loss function has shown benefits in face identification and low-shot learning [10], [12]. More generally, metric-learning based approaches, such as Matching Networks employs a differentiable nearest neighbor classifier over the learned representations of the training examples in order to classify an unlabeled example. Prototypical Networks [13] handles the problems by calculating the Euclidean-distance between the embedding points of test samples and prototype representation of labeled samples. Relation Networks [14] uses concatenated feature maps from the query and labeled images to distinguish the similarity and dissimilarity. Those methods including the proposed one follow the training strategy of meta-learning, which is to sample small training sets and query sets from base-train classes and feed the sampled training set to the learner for a classifier. It then computes the loss of the classifier on the sampled query set.

B. Visualizing CNNs

In addition, the object-level information is essentially the focused-regions in the image. The focused-region is usually captured by some methods like visualizing ConvNets(CNN). A number of previous works [15], [16] have visualized CNN predictions by highlighting 'important' pixels (i.e., usually some areas with discriminative features, such as the area

containing the object). Class Activation Mapping (CAM) [15] proposed to localize object by modifying CNN architectures. It replaced the fully-connected layers with convolutional layers and global average pooling [17] to obtain the focused-regions. Gradient-weighted Class Activation Mapping (Grad-CAM) [16] was applicable to a significantly broader range of CNNs without modifying models anymore. Grad-CAM utilized the gradients of any target concept (for example the logits for a caption) and fed it into the final convolutional layer to produce a coarse localization map, which highlighted the important regions in the image for prediction tasks. The object-level module in our work is inspired by the localization concept from Grad-CAM.

C. Similarity Learning

Amongst similarity learning, the ability to compare between images is one of the most fundamental operations among all of computing. Classic per-pixel measures, such as L_2 Euclidean-distance are insufficient for assessing structured images. How to measure two similar images in a way that coincides with human judgment is a longstanding goal [18]. Richard Zhang et al. [19] systematically evaluated the deep features that are obtained by averaging the information across spatial dimension and across all layers. For the proposed model, different levels of deep features are utilized to measure the similarities rather than averaging all layers' outputs.

III. MULTI-LEVEL SIMILARITY MODEL

A. Problem Definition

Assuming that there is a base train dataset \mathcal{D}_{base} , a novel test dataset \mathcal{D}_{novel} , where $\mathcal{D}_{base} \cap \mathcal{D}_{novel} = \emptyset$. Following the meta-learning normal training strategy, in each training iteration, we randomly sample C classes with K labeled samples (called C -way K -shot) from \mathcal{D}_{base} . These samples serve as support set $\mathcal{S} = \{(x_i, y_i)\}_{i=1}^{K \times C}$. Then a fraction of the remainder in those C classes is sampled to serve as a query set

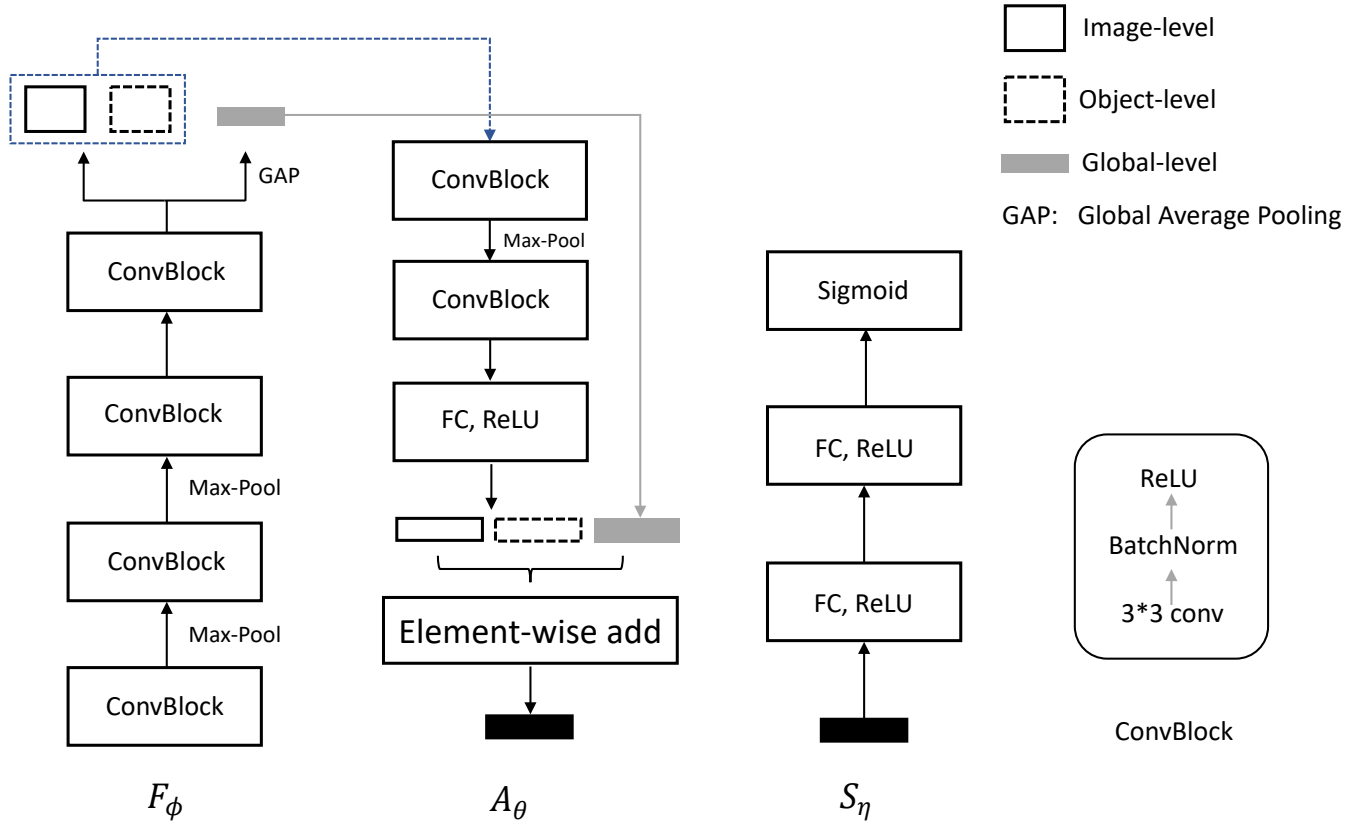


Fig. 3. Details of the architecture of MLSM.

$\mathcal{Q} = \{(x_i, y_i)\}_{i=1}^N$. The loss of the classifier is computed on the \mathcal{Q} to optimize the proposed model by back-propagation. During the testing of the proposed model, the \mathcal{S} and \mathcal{Q} are sampled from \mathcal{D}_{novel} similarly to the above training procedure. The \mathcal{Q} of \mathcal{D}_{novel} can be classified without further updating the model.

As shown in Fig. 2, the proposed model consists of three parts: a feature extraction module F_ϕ , an adjustment function module A_θ and a similarity function module S_η . Given two images $\mathbf{a} \in \mathcal{S}$ and $\mathbf{b} \in \mathcal{Q}$ (or $\mathbf{b} \in \mathcal{S}$), the feature extractor F_ϕ will produce a pair of image-level feature maps $\mathcal{I}_\mathbf{a}$ and $\mathcal{I}_\mathbf{b}$. Object-level of \mathbf{a} and \mathbf{b} are also fed into the feature extractor F_ϕ . Then $\mathcal{I}_\mathbf{a}$ and $\mathcal{I}_\mathbf{b}$ will generate a pair of global-level feature vectors $\mathcal{G}_\mathbf{a}$ and $\mathcal{G}_\mathbf{b}$ through the global average pooling (GAP). After processed by the adjustment function module A_θ , three-level feature vectors will be obtained through the unified channels with the element-sum operation to integrate the multi-level feature representations. Then $\mathcal{C}(\mathcal{I}_\mathbf{a} \oplus \mathcal{O}_\mathbf{a} \oplus \mathcal{G}_\mathbf{a})$, $(\mathcal{I}_\mathbf{b} \oplus \mathcal{O}_\mathbf{b} \oplus \mathcal{G}_\mathbf{b})$ will be sent to S_η to learn similar scores (1 represents the same class while 0 represents the different), where $\mathcal{C}(\cdot, \cdot)$ means the concatenation of feature vectors in depth. For K -shot, where $K > 1$, the feature vectors of K labeled samples are averaged to serve as the representation of this class. The detailed structure of the three parts is shown in Fig. 3.

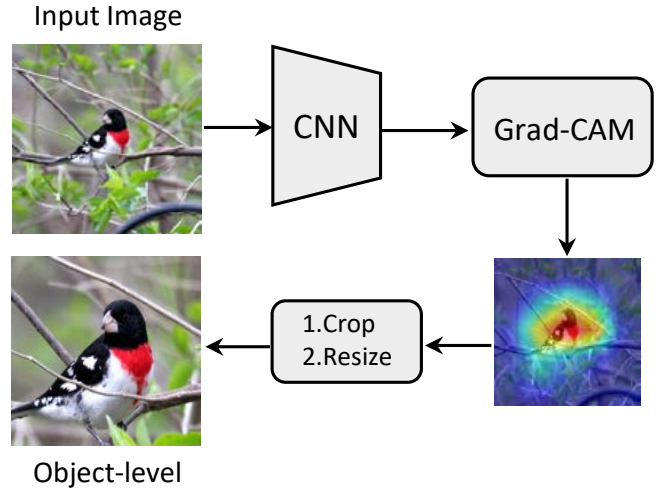


Fig. 4. A flowchart for generating object-level areas.

B. Image-level Similarity

Image-level information refers to the feature map of the whole image obtained by feature extractor F_ϕ . For the proposed model, F_ϕ consists of four convolutional blocks, as shown in Fig. 3. More specifically, each convolutional block contains a 3×3 convolution with 64 filters, batch normalization

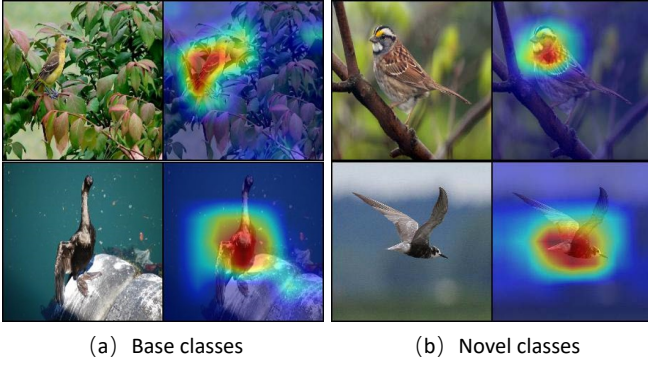


Fig. 5. The object-level areas generated by Grad-CAM. We can observe that the novel classes in (b) also can locate the coarse focus regions even it has never seen.

and a non-linear activation function (ReLU). Both training and testing samples are uniformly resized into 84×84 . Feature maps at this level contain the highest-order semantic information of the image. Therefore, image-level is the best choice for similarity learning between images. A simple rule is to concatenate feature maps from a and b pair-wisely, i.e.

$$\mathcal{C}(\mathcal{I}_a, \mathcal{I}_b) = \{\text{Concat}(\mathcal{I}_a^i, \mathcal{I}_b^i)\}_{i=1}^{M*N} \quad (1)$$

where $M * N$ represents the size of the last feature map. Once two feature maps are concatenated together, the channel information contained in each pixel of the last feature map will be compared, like Relation Networks did [14].

C. Object-level Similarity

The final goal of low-shot learning is to classify different types of objects. The existing studies have shown that learning from object-areas could help for recognition tasks at image-level [20], [21]. Inspiredly, we assume that there may exist some discriminative regions in the images that are beneficial to low-shot tasks. The discriminative region is the focused-area in the image. It usually contains all the information of the objects rather than the background. We call the focus-area as object-level area. To acquire the object-level area, Li et al. [21] proposed a zoom network which utilized the candidate region to crop the original images. Wei et al. [22] adopted the unsupervised object discovery and co-localization mechanism by deep descriptor transformation to discover the object-level area. Differently, to be more efficient, we utilize Grad-CAM to get the object-level area. A coarse localization map highlighting the important regions in the image for predicting the concept could be calculated by:

$$L_{Grad-CAM}^c = ReLU\left(\sum_k \alpha_k^c A^k\right) \quad (2)$$

where $\alpha_k^c = \frac{1}{Z} \sum_i \sum_j \frac{\partial y^c}{\partial A_{ij}^k}$ denotes the weight of the k -th feature map for category c and Z is the number of feature map's pixels. A_{ij}^k represents the pixel value at the location of (i, j) of the k -th feature map, and y^c is the classification score corresponding to the c -th class.

For this step, as shown in Fig. 4, we first train a ConvNets with the existing base data \mathcal{D}_{base} . The ConvNets is only used to generate object-level areas rather than extracting any features. Then we use the Grad-CAM to generate a coarse object-level area. We further scale and zoom the target areas into a size of 84×84 , which is the same as the original images. Nevertheless, there are still some problems when applying the proposed method, i.e. we do not know the label (corresponds to c in equation (2)) of new categories to generate the object-level areas. In this case, we utilize Grad-CAM directly and look for the closest class to \mathcal{D}_{base} to locate the coarse object-level area, as shown in Fig. 5.

D. Global-level Similarity

For the image-level feature, each of the learned filters operates with a local receptive field. It is thus unable to exploit contextual information outside of this region. In order to capture the global-level similarity, global spatial information is applied to the channel descriptor. This is achieved by the feature extractor F_ϕ , which employs global average pooling after the last convolutional layer. The output features of a (or b) are as follows:

$$\mathcal{G}_a(\mathcal{G}_b) = \frac{1}{H * W} \sum_{x,y} \alpha_{(x,y)}. \quad (3)$$

where $\alpha_{(x,y)}$ denotes the C -dimensional slice of α at spatial location (x, y) . The global-level similarity of images a and b can be directly calculated by:

$$\mathcal{C}(\mathcal{G}_a, \mathcal{G}_b) = \{\text{Concat}(\mathcal{G}_a^i, \mathcal{G}_b^i)\}_{i=1}^K \quad (4)$$

where K denotes K -dimensional feature vector after GAP. In addition, due to the three different levels, the feature channel sizes are not unified. As shown in Fig. 3, we then use an adjustment function A_θ to align the feature vectors of three levels for the element-sum operation. A_θ consists of two convolutional blocks and a fully-connected layer.

Overall, the multi-level similarity is calculated as follows:

$$\text{Similarity} = S_\eta(\text{Concat}((\mathcal{I}_a \oplus \mathcal{O}_a \oplus \mathcal{G}_a), (\mathcal{I}_b \oplus \mathcal{O}_b \oplus \mathcal{G}_b))) \quad (5)$$

where S_η denotes the similarity function. It uses Sigmoid activation function to force the network to learn the concept of similarity and dissimilarity.

IV. EXPERIMENTS

A. Datasets and Experimental Procedure

Datasets: we evaluate the proposed MLSM on the dataset of Caltech-UCSD 200-2011 [24]. It consists of 11788 bird images in 200 categories. The samples of this dataset are shown in Fig. 6. Each category has some semantic descriptions. We further divide the dataset into three parts: 100 classes for training (i.e. \mathcal{D}_{base}), 50 for validation, and 50 for testing (i.e. \mathcal{D}_{novel}). For the proposed model, the input images (including object-level areas) are resized to 84×84 pixels without any data augmentations (e.g. randomized rotation, and horizontal

TABLE I
AVERAGE TEST SET CLASSIFICATION ACCURACY ON CALTECH-UCSD BIRDS.

Methods	Caltech-UCSD Birds (%)			
	<i>Fine Tune</i>	<i>5-way 1-shot</i>	<i>5-way 5-shot</i>	<i>Distance</i>
MAML [5]	Y	38.43	59.15	-
META-LEARN LSTM [9]	N	40.43	49.65	-
Matching Nets [7]	N	49.34	59.31	Cosine
PROTO-Nets [13]	N	45.27	56.35	Euclid.
RELATION-Nets [14]	N	46.69	55.86	Image-level deep metric
MACO [23]	N	60.76	74.96	-
MLSM	N	64.50	70.50	Multi-level deep metric

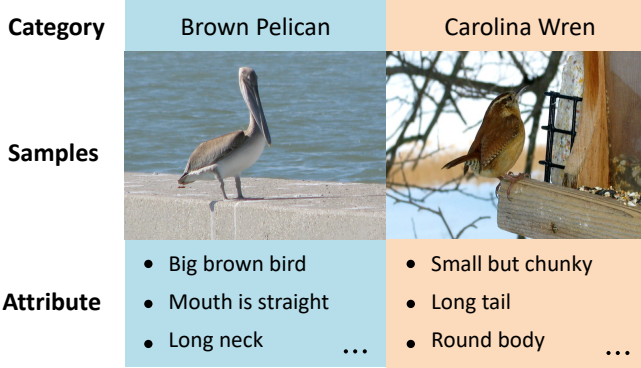


Fig. 6. Two categories randomly sampled from the Caltech-UCSD Birds dataset.

TABLE II
PERFORMANCE COMPARISON ON DIFFERENT LEVELS.

Different level	1-shot (%)	5-shot (%)
\mathcal{I}	50.0	59.0
$\mathcal{I} + \mathcal{G}$	51.0	61.0
$\mathcal{I} + \mathcal{G} + \mathcal{O}$	64.5	70.5

flipping). We use Adam solver [25] with an initial learning rate of 0.001, and annealed by half for every 100,000 episodes.

B. Experiments on CUB-200-2011

As shown in Table I, it can be seen that the proposed model achieves promising results compared with other methods. The accuracy is computed by averaging over 100 randomly generated episodes from the testing set. Each episode contains 5 classes with 200 query images. Our model could achieve state-of-the-art performance on the 5-way 1-shot setting. At the same time, there are also some gaps between our method and MACO [23] on 5-way 5-shot setting. Compared with the image-level deep metric method RelationNets [26], it can be noticed that the multi-level metric could be beneficial to learn a better similarity function with no additional training set. We also compared the impact of different levels metric on the accuracy, from table II, we can see that object-level

representations bring great improvement both on 1-shot and 5-shot setting.

V. CONCLUSION

This paper introduces a novel multi-level similarity learning model (MLSM) for low-shot tasks. MLSM learns a similarity function from image-level, global-level, and object-level which coincides with human judgment to some extent. The performance is evaluated on Caltech-UCSD 200-2011 dataset by comparing query and support samples. It achieves promising results compared with the existing methods. We believe that our method can act as a valuable complement to low-shot learning.

ACKNOWLEDGMENT

This work was supported in part by the National Natural Science Foundation of China (No. 41741007, 41576011, U1706218) the Key Research and Development Program of Shandong Province (no.GG201703140154) and Applied Basic Research Programs of Qingdao (no. 18-2-2-38-jch). [27]

REFERENCES

- [1] A. Calefati, M. K. Janjua, S. Nawaz, and I. Gallo, "Git Loss for Deep Face Recognition," *BMVC*, pp. 1–12, 2018.
- [2] T. Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature pyramid networks for object detection," in *CVPR*, 2017.
- [3] Jia Deng, Wei Dong, R. Socher, Li-Jia Li, Kai Li, and Li Fei-Fei, "ImageNet: A large-scale hierarchical image database," in *CVPR*, 2009.
- [4] Y. X. Wang, R. Girshick, M. Hebert, and B. Hariharan, "Low-Shot Learning from Imaginary Data," *CVPR*, pp. 7278–7286, 2018.
- [5] C. Finn, P. Abbeel, and S. Levine, "Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks," *International Conference on Machine Learning (ICML)*, 2017.
- [6] S. Ravi and H. Larochelle, "Optimization as a model for few-shot learning," *ICLR*, 2017.
- [7] O. Vinyals, C. Blundell, T. Lillicrap, K. Kavukcuoglu, and D. Wierstra, "Matching Networks for One Shot Learning," *NIPS*, 2016.
- [8] H. Qi, M. Brown, and D. G. Lowe, "Low-Shot Learning with Imprinted Weights," *CVPR*, 2018.
- [9] Z. Li, F. Zhou, F. Chen, and H. Li, "Meta-SGD: Learning to Learn Quickly for Few-Shot Learning," *arXiv preprint arXiv:1707.09835*, 2017.
- [10] G. Koch, R. Zemel, and R. Salakhutdinov, "Siamese Neural Networks for One-shot Image Recognition," *International Conference on Machine Learning*, 2015.

- [11] F. Schroff, D. Kalenichenko, and J. Philbin, "FaceNet: A unified embedding for face recognition and clustering," in *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2015, pp. 815–823.
- [12] X. Zhang, F. Zhou, Y. Lin, and S. Zhang, "Embedding Label Structures for Fine-Grained Feature Representation," *CVPR*, 2016.
- [13] J. Snell, K. Swersky, and R. S. Zemel, "Prototypical Networks for Few-shot Learning," *NIPS*, 2017.
- [14] 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," *CVPR*, 2017.
- [15] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, and A. Torralba, "Learning Deep Features for Discriminative Localization," *CVPR*, 2016.
- [16] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, "Grad-CAM: Visual Explanations from Deep Networks via Gradient-Based Localization," in *CVPR*, 2017, pp. 618–626.
- [17] M. Lin, Q. Chen, and S. Yan, "Network In Network," *International Conference of Learning Representation*, 2014.
- [18] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: From error visibility to structural similarity," *IEEE Transactions on Image Processing*, 2004.
- [19] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang, "The Unreasonable Effectiveness of Deep Features as a Perceptual Metric," *CVPR*, pp. 586–595, 2018.
- [20] J. Fu, H. Zheng, and T. Mei, "Look closer to see better: Recurrent attention convolutional neural network for fine-grained image recognition," *CVPR*, pp. 4476–4484, 2017.
- [21] Y. Li, J. Zhang, J. Zhang, and K. Huang, "Discriminative Learning of Latent Features for Zero-Shot Recognition," *CVPR*, pp. 7463–7471, 2018.
- [22] X. S. Wei, C. L. Zhang, J. Wu, C. Shen, and Z. H. Zhou, "Unsupervised object discovery and co-localization by deep descriptor transformation," *CVPR*, 2018.
- [23] N. Hilliard, L. Phillips, S. Howland, A. Yankov, C. D. Corley, and N. O. Hodas, "Few-Shot Learning with Metric-Agnostic Conditional Embeddings," *arXiv preprint arXiv:1802.04376*, 2018.
- [24] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie, "The Caltech-UCSD Birds-200-2011 Dataset," Tech. Rep., 2011.
- [25] D. P. Kingma and J. Ba, "Adam: A Method for Stochastic Optimization," *ICLR*, 2015.
- [26] C. Sun, F. Li, H. Lu, and G. Hua, "Visual Tracking via Joint Discriminative Appearance Learning," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 27, no. 12, pp. 2567–2577, 2017.
- [27] H. Xiao, K. Rasul, and R. Vollgraf, "Fashion-MNIST: a Novel Image Dataset for Benchmarking Machine Learning Algorithms," pp. 1–6, 2017. [Online]. Available: <http://arxiv.org/abs/1708.07747>