

AN ADAPTIVE PLUG-AND-PLAY NETWORK FOR FEW-SHOT LEARNING

Hao Li¹, Li Li^{1*}, Yunmeng Huang¹, Ning Li¹, Yongtao Zhang²

¹School of Electronic and Information Engineering, Beihang University

²College of Mechanical & Electrical Engineering, Nanjing University of Aeronautics and Astronautics

ABSTRACT

Few-shot learning (FSL) requires a model to classify new samples after learning from only a few samples. While remarkable results are achieved in existing methods, the performance of embedding and metrics determines the upper limit of classification accuracy in FSL. The bottleneck is that deep networks and complex metrics tend to induce overfitting in FSL, making it difficult to further improve the performance. Towards this, we propose plug-and-play model-adaptive resizer (MAR) and adaptive similarity metric (ASM) without any other losses. MAR retains high-resolution details to alleviate the overfitting problem caused by data scarcity, and ASM decouples the relationship between different metrics and then fuses them into an advanced one. Extensive experiments show that the proposed method could boost existing methods on two standard dataset and a fine-grained datasets, and achieve state-of-the-art results on *mini*-ImageNet and *tiered*-ImageNet.

Index Terms— Few-shot learning, adaptive resizer, learnable metric

1. INTRODUCTION

The purpose of FSL is to learn a classifier that can recognize unseen categories from a small number of samples. A mainstream mechanism [1, 2] is to learn transferable knowledge from seen categories, and then use this transferable knowledge to build a classifier, and finally apply the classifier to unseen categories. Following this mechanism, recently great advances have been made in metric-based methods for FSL, achieving state-of-the-art results.

The pipeline of metric-based methods is to learn advanced deep embedding module and then measure the similarity between query sample and a few support samples [3]. The seminal work could date back to [4], which propose Siamese Network for obtaining image representations with L_1 distance as similarity metric. Then Matching Network [5] introduce attention mechanism and cosine metric. Prototypical Network [3] propose mean vector as the corresponding class prototype representation and adopt Euclidean distance. On the basis of the above methods, recent approaches can be broadly divided

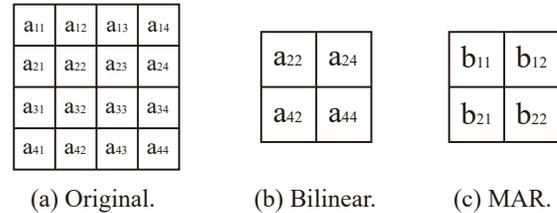


Fig. 1. Diagram of bilinear and MAR resizing. As shown in (b), bilinear resizing is equivalent to sampling the original image in (a). However, as shown in (c), MAR resizing yields new pixels.

into two categories. The first exploits Gaussian distribution [6, 7], discrete probability [8, 9], Ridge regression feature reconstruction [10], and attention mechanism [11, 12, 13] for advanced embedding representations. And the second category aims to obtain progressive similarity metric by learnable metric modules [14, 15], covariance matrices [16, 17], nearest neighbors [18], and multiple-metric fusion [19]. Remarkable results have been achieved in above methods. However, deep networks are prone to overfitting when data scarcity, making it a challenge to continue improving performance of FSL. To solve this challenge, the proposed approach presents a pioneering solution consisting of two ways.

The first is adding learnable preprocessing module in front of backbone network. Learnable preprocessing modules have been early shown could improve visual recognition [20]. Most approaches explore from super resolution [21, 22], image decompression [23], image denoising [24], image de-hazing [25] to improve the performance of specific task of visual recognition. [26] makes great progress, proposing a generic framework effective on multiple tasks. It cannot be used directly for FSL because it requires the introduction of specific loss and do not take into account FSL characteristics. For FSL, due to data scarcity, it is crucial for the framework to obtain enough discriminative details. However, as shown in Fig. 1, traditional resizing methods directly downsample high-resolution images to a uniform size, thus details with discriminative information in high-resolution images may be lost. To solve this problem, we propose MAR, a learnable model-adaptive resizer. MAR incorporates the spatial in-

* Corresponding author.

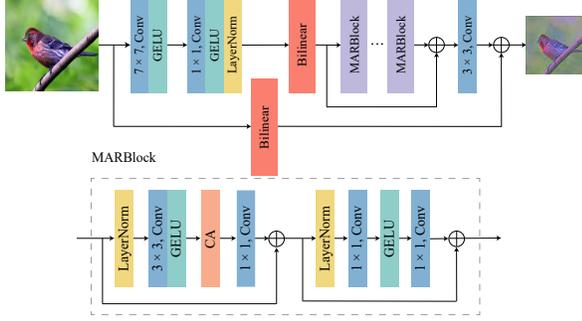


Fig. 2. The architecture of the proposed MAR. CA denotes channel attention module.

formation of surrounding pixels during resizing, so that the details of high-resolution images can be effectively preserved. In addition, MAR does not introduce any other losses thus can be flexibly instantiated in existing methods.

The second is to adaptively fuse existing metrics to obtain advanced discriminative one. Existing methods mainly focus on finding advanced discriminative higher-order metrics, and there has been little research on metric fusion. As the first method to fuse multiple measures into one, [19] proves that multiple simple metrics outperform a few single high-order metrics by manually presetting the weights. However, manual presets are based on sparse priors, which cannot deeply mine the relationship among various metrics. To solve the problem, we propose an adaptive metric ASM. ASM can adapt to different combinations of metrics by learning the weights. Extensive experiments show that the proposed ASM outperforms the method in [19].

2. PROPOSED METHOD

In this section, we first outline metric-based FSL, and then introduce the details of proposed framework.

2.1. FSL setting

In standard FSL, the dataset \mathcal{C} consists of training set \mathcal{C}_{train} , validation set \mathcal{C}_{val} and test set \mathcal{C}_{test} . Their classes do not have intersection. A specific task \mathcal{T} in FSL is denoted as N -way K -shot classification, which represents a random sampling of N classes from the support set $\mathcal{D}_{support}$ with K samples per class. Given a query set \mathcal{D}_{query} , the purpose is to correctly classify unlabeled samples in it by learning the knowledge in the support set $\mathcal{D}_{support}$. We denote the support set as $\mathcal{D}_{support} = (\mathbf{x}_i, \mathbf{y}_i)_{i=1}^{N,K}$, with the instance as $\mathbf{x}_i \in \mathbb{R}^D$ and one-hot vector label $\mathbf{y}_i = \{0, 1\}^N$. The value of N is generally 1 and 5.

The first step of metric-based methods is to obtain a classifier by regular classification training on \mathcal{C}_{train} , and then take specific tasks from \mathcal{C}_{train} to fine-tune. The final goal is to

obtain a classifier $f_{\theta}(\cdot)$ that performs well on unseen classes \mathcal{C}_{test} . Mathematically,

$$\hat{\mathbf{y}}_j = f_{\theta}(\mathbf{x}_j; \mathcal{C}_{train}), \mathbf{x}_j \in \mathcal{C}_{test} \quad (1)$$

where \mathbf{x}_j denotes the unseen sample and $\hat{\mathbf{y}}_j$ denotes the predicted label.

2.2. Model-adaptive resizer

The proposed MAR is shown in Fig. 2. For large receptive field and less parameters, the convolution kernels of the first two layers are 7×7 and 1×1 . And GELU is used for non-linear activation. Then features are resized by bilinear to a uniform size. Output features are fed into n MARBlocks, where MARBlock is a module for feature enhancement. After the features plus residual are convolved with a 3×3 kernel and then summed with original downsampled residual, the MAR output is obtained. As shown in Fig. 1, the conventional resizer is only able to extract partial information of the high-resolution image. However, the proposed MAR can extract partial information while preserving some key details in high-resolution images.

2.3. Adaptive similarity metric

In general, inputs are all transformed into vectors of the same length after embedding, and prediction is achieved by similarity measurement on the vectors from $\mathcal{D}_{support}$ and \mathcal{D}_{query} . And the widely adopted metrics are Euclidean distance and cosine similarity in existing methods. Almost all methods select one of them, and few have explored the fusion of them. As unique to our knowledge, BSNet [19] innovatively explores metric fusion roughly by manually presetting the weights. To further extend the work of BSNet from an adaptive perspective, we propose ASM. Specifically, given a FSL task \mathcal{T} , \mathcal{C}_k is used to denote the sample set of class k in the task. The prototype of class k is represented as

$$\omega_k = \frac{1}{\mathcal{C}_k} \sum_{(\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{C}_k} \Psi(\mathbf{x}_i) \quad (2)$$

where $\Psi(\cdot)$ represents embedding module. For the query sample \mathbf{x}_j in task, similarity with \mathcal{C}_k can be denoted as

$$\text{sim}(\Psi(\mathbf{x}_j), \omega_k) = \alpha \cdot \text{euc}(\Psi(\mathbf{x}_j), \omega_k) + \beta \cdot \cos(\Psi(\mathbf{x}_j), \omega_k) \quad (3)$$

where α and β denotes learned weights, $\text{sim}(\cdot)$ denotes ASM, $\text{euc}(\cdot)$ denotes Euclidean distance and $\cos(\cdot)$ denotes cosine similarity, respectively.

2.4. Implementation of proposed method

In our method, as shown in Fig. 3, training is divided into three stages: backbone training, joint training and fine-tuning.

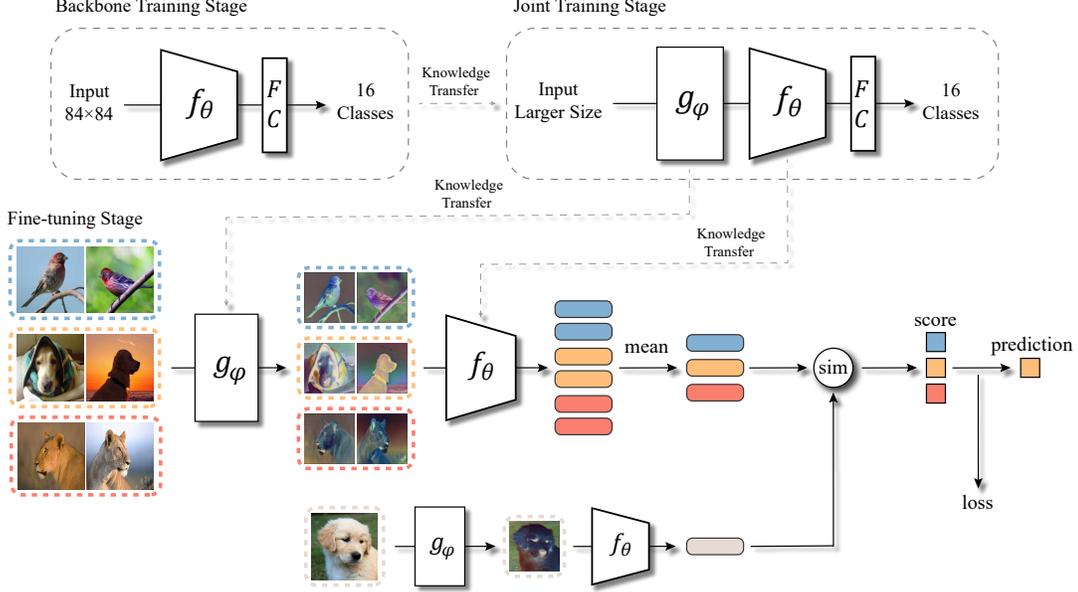


Fig. 3. Illustration of the proposed framework in 3-way 2-shot setting. $f_{\theta}(\cdot)$ represents backbone network, $g_{\phi}(\cdot)$ represents MAR and $\text{sim}(\cdot)$ is ASM. Training consists of three stages: backbone training, joint training and fine-tuning.

In backbone training stage, we only train the backbone network $f(\cdot)$ in traditional classification way to learn the parameters θ_1 , obtaining the preliminary model $f_{\theta_1}(\cdot)$. Then, in joint training stage, MAR, denoted as $g(\cdot)$ is instantiated in front of $f_{\theta_1}(\cdot)$. This stage learns parameters ϕ_1 and θ_2 in the same way as first stage. The obtained model are denoted as $g_{\phi_1}(\cdot)$ and $f_{\theta_2}(\cdot)$, respectively. Finally, in fine-tuning stage, we train the model obtained in the second stage in meta manners to obtain the final parameters ϕ and θ , i.e., the embedding module $\Psi(\cdot) = f_{\theta} \circ g_{\phi}(\cdot)$.

Specifically, in traditional classification way, training is done by adding a fully-connected (FC) layer after CNN. Then FC layer is removed and we randomly sample several tasks \mathcal{T} from \mathcal{D}_{query} in validation set \mathcal{C}_{val} to save the best model. In meta manner, i.e. fine-tuning stage, a large number of $(\mathbf{x}_i, \mathbf{y}_i)$ are taken from $\mathcal{D}_{support}$. Then \mathbf{x}_i are fed to resizer $g_{\phi_1}(\cdot)$ to get the resized \mathbf{x}_i (denoted as $\tilde{\mathbf{x}}_i$), and $\tilde{\mathbf{x}}_i$ passes backbone $f_{\theta_2}(\cdot)$ to get the embedding for similarity measure and loss calculation. Mathematically, backbone training stage can be written as

$$\arg \min_{\theta_1} \sum_{(\mathbf{x}_j, \mathbf{y}_j) \in \mathcal{D}_{query}} l(\mathbf{y}_j, f_{\theta_1}(\mathbf{x}_j; \mathcal{D}_{support})) \quad (4)$$

where $(\mathbf{x}_j, \mathbf{y}_j)$ denotes the sample pairs taken from validation set \mathcal{C}_{val} , and $l(\cdot)$ represents the loss between the prediction and the actual label. Similarly, joint training stage can be marked as

$$\arg \min_{\theta_2, \phi_1} \sum_{(\mathbf{x}_j, \mathbf{y}_j) \in \mathcal{D}_{query}} l(\mathbf{y}_j, f_{\theta_2} \circ g_{\phi_1}(\mathbf{x}_j; \mathcal{D}_{support})) \quad (5)$$

where θ_2, ϕ_1 are learned parameters, and other terms are consistent with Eq. 4. Fine-tuning stage is the same as the joint training, except that the learned objectives are θ and ϕ .

Overall, α, β, θ and ϕ are updated with cross-entropy loss. Finally, predictions are mapped from embedding space to probabilities via softmax. Concretely, the probability of query sample \mathbf{x}_j belonging to category k is

$$p(\hat{\mathbf{y}}_j = k | \mathbf{x}_j) = \frac{\exp(\text{sim}(\Psi(\mathbf{x}_j), \omega_k))}{\sum_{\mathbf{x}_j \in \mathcal{D}_{query}} \exp(\text{sim}(\Psi(\mathbf{x}_j), \omega_k))} \quad (6)$$

and the loss can be calculated by

$$\mathcal{L} = -\frac{1}{|\mathcal{D}_{query}|} \sum_{(\mathbf{x}_j, \mathbf{y}_j) \in \mathcal{D}_{query}} \mathbf{y}_j \log(p(\hat{\mathbf{y}}_j = k | \mathbf{x}_j)) \quad (7)$$

3. EXPERIMENTS

3.1. Datasets and experiment settings

Our experiments are conducted on datasets *mini-ImageNet* [5], *tiered-ImageNet* [27] and CUB [28]. For *mini-ImageNet* and *tiered-ImageNet* we adopt ResNet-12 as the backbone network, and for CUB we employ Conv-4 and ResNet-12 (Conv-4 is the same setting as in [19]). In training of ResNet-12, we utilize Adam optimizer with weight decay 0.0005, and learning rate starts from 0.0001 and drops by a factor of 0.1. For Conv-4, we use SGD optimizer with a momentum of 0.9 with learning rate starting at 0.0002 and decreasing by a factor of 0.5.

Model	<i>mini</i> -ImageNet		<i>tiered</i> -ImageNet	
	1-shot	5-shot	1-shot	5-shot
DeepEMD [8]	65.91	82.41	71.16	86.03
BML [29]	67.04	83.63	68.99	85.49
Baseline++ [30] [†]	63.15	81.77	67.97	84.43
ProtoNet [3] [†]	62.26	80.19	67.47	82.07
FEAT [12] [†]	66.67	82.01	70.97	84.87
Meta-Baseline [31] [†]	64.91	80.78	65.77	83.21
Meta DeepBDC [9] [†]	67.37	84.39	72.68	87.87
Ours	65.19	82.38	69.01	85.17
FEAT + Ours	67.41	82.98	71.74	85.81
Meta-Baseline + Ours	65.11	81.67	68.81	85.34
Meta DeepBDC + Ours	68.24	84.45	73.71	87.91

Table 1. Results on *mini*-ImageNet and *tiered*-ImageNet at 5-way setting. The best results are marked in **bold black**.

Model	Resizer’s Input Size	Resizer’s Output Size	5-way 1-shot	5-way 5-shot
Baseline	original	84 × 84	62.26	80.19
ASM	original	84 × 84	64.41	82.34
MAR	112 × 112	84 × 84	64.84	80.28
MAR	126 × 126	84 × 84	64.45	80.23
MAR	168 × 168	84 × 84	65.02	80.31
MAR+ASM	112 × 112	84 × 84	64.10	82.00
MAR+ASM	126 × 126	84 × 84	63.98	81.97
MAR+ASM	168 × 168	84 × 84	65.19	82.38

Table 2. Results of different combinations of MAR and ASM on *mini*-ImageNet.

Model	Number of MARBlocks	5-way 1-shot (%)	5-way 5-shot (%)
Ours	2	63.73 ± 0.49	81.73 ± 0.34
Ours	4	65.19 ± 0.52	82.38 ± 0.33

Table 3. Results of different number of MARBlocks on *mini*-ImageNet. The input size of MAR is 168 × 168.

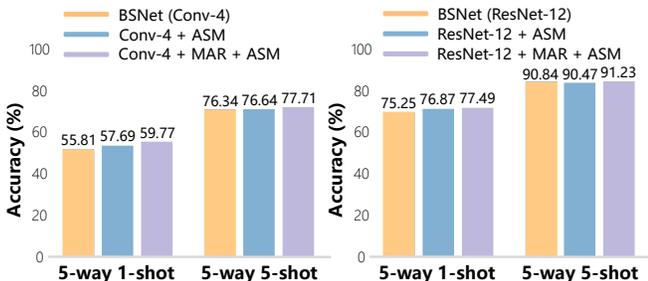


Fig. 4. Results of fine-grained FSL experiments compared with BSNet [19] on CUB.

3.2. Experimental results

For a fair comparison, a baseline of the proposed method has been built with ResNet-12 as backbone. Input size of MAR is set to 168 × 168 and weights of ASM, α and β , are initialized to 1.24 and 0.1. Results can be seen from Table 1. On *mini*-ImageNet and *tiered*-ImageNet, this baseline improves 2.93% and 2.19%, 1.54% and 3.1% compared with [3], and 0.28% and 1.6%, 3.24% and 1.96% compared with [31], at 5-way 1-shot and 5-shot settings, respectively. To further explore the plug-and-play property of the proposed model, we instantiate it in FEAT, Meta-Baseline and Meta DeepBDC. The results are shown in Table 1. It improves FEAT by 0.74% and 0.77%, Meta-Baseline by 0.2% and 3.04%, Meta DeepBDC by 0.87% and 1.03% at 5-way 1-shot, 0.97% and 0.94%, 0.89% and 2.13%, 0.06% and 0.04% at 5-way 5-shot, on *mini*-ImageNet and *tiered*-ImageNet, respectively.

3.3. Ablation analysis

Extensive experiments have been done to demonstrate the effectiveness of the proposed method. Specifically, we build a baseline with ResNet-12. Three different configurations, MAR, ASM and MAR + ASM, are obtained by instantiating MAR and ASM. 112 × 112, 126 × 126, 168 × 168 are explored as different input sizes for MAR. The details are shown in Table 2. Compared with baseline, the addition of MAR and ASM improved the model at 5-way 1-shot and 5-way 5-shot by up to 2.93% and 2.19% on *mini*-ImageNet. And MAR has a large improvement for 5-way 1-shot and ASM for 5-way 5-shot, reaching 2.76% and 2.15%, respectively. To further explore the proposed method, we investigate the impact of different number of MARBlocks on the performance. Results can be seen in the Table 3. It shows that 4 MARBlocks are superior to 2. Moreover, we compare the proposed framework with BSNet [19]. Similar to ResNet-12, we also build a baseline with Conv-4. Based on them, we explore the performance of different configurations on CUB. Results are shown in Fig. 4. At 5-way 1-shot and 5-way 5-shot, our method improves up to 3.96% and 1.37% for Conv-4 and 2.24% and 0.39% for ResNet-12. And learned α and β are 1.24 and 0.1, 1.91 and 0.1, respectively, which further explores BSNet. Since CUB is a fine-grained dataset, the results also demonstrate that a generalization of our approach to other domains.

4. CONCLUSION

In this paper, we propose a plug-and-play framework for FSL. A learnable resizer (MAR) is adopted to enhance the embeddings, and adaptive metric (ASM) is used to obtain advanced discrimination capabilities. MAR alleviates the loss of detailed information in traditional preprocessing pipeline, and ASM converts simple metrics into an advanced and efficient one. In addition, we open a novel view to improve FSL from

the perspective of input data and adaptive similarity measurement. Extensive experiments demonstrate that the proposed adaptive framework is a potential direction and could effectively boost existing methods.

5. REFERENCES

- [1] Li Fei-Fei, Robert Fergus, and Pietro Perona, “One-shot learning of object categories,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 28, no. 4, pp. 594–611, 2006.
- [2] Brenden Lake, Ruslan Salakhutdinov, Jason Gross, and Joshua Tenenbaum, “One shot learning of simple visual concepts,” in *Proceedings of the annual meeting of the cognitive science society*, 2011.
- [3] Jake Snell, Kevin Swersky, and Richard S. Zemel, “Prototypical networks for few-shot learning,” in *NIPS*, 2017.
- [4] Gregory Koch, Richard Zemel, Ruslan Salakhutdinov, et al., “Siamese neural networks for one-shot image recognition,” in *ICML deep learning workshop*, 2015.
- [5] Oriol Vinyals, Charles Blundell, Tim Lillicrap, Koray Kavukcuoglu, and Daan Wierstra, “Matching networks for one shot learning,” in *NIPS*, 2016.
- [6] Wenbin Li, Lei Wang, Jing Huo, Yinghuan Shi, Yang Gao, and Jiebo Luo, “Asymmetric distribution measure for few-shot learning,” in *IJCAI*, 2020.
- [7] Shuo Yang, Lu Liu, and Min Xu, “Free lunch for few-shot learning: Distribution calibration,” in *ICLR*, 2021.
- [8] Chi Zhang, Yujun Cai, Guosheng Lin, and Chunhua Shen, “Deepemd: Few-shot image classification with differentiable earth mover’s distance and structured classifiers,” in *CVPR*, 2020.
- [9] Jiangtao Xie, Fei Long, Jiaming Lv, Qilong Wang, and Peihua Li, “Joint distribution matters: Deep brownian distance covariance for few-shot classification,” in *CVPR*, 2022.
- [10] Davis Wertheimer, Luming Tang, and Bharath Hariharan, “Few-shot classification with feature map reconstruction networks,” in *CVPR*, 2021.
- [11] Carl Doersch, Ankush Gupta, and Andrew Zisserman, “Crosstransformers: spatially-aware few-shot transfer,” in *NIPS*, 2020.
- [12] Han-Jia Ye, Hexiang Hu, De-Chuan Zhan, and Fei Sha, “Few-shot learning via embedding adaptation with set-to-set functions,” in *CVPR*, 2020.
- [13] Yangji He, Weihan Liang, Dongyang Zhao, Hong-Yu Zhou, Weifeng Ge, Yizhou Yu, and Wenqiang Zhang, “Attribute surrogates learning and spectral tokens pooling in transformers for few-shot learning,” in *CVPR*, 2022.
- [14] Flood Sung, Yongxin Yang, Li Zhang, Tao Xiang, Philip HS Torr, and Timothy M Hospedales, “Learning to compare: Relation network for few-shot learning,” in *CVPR*, 2018.
- [15] Hongguang Zhang and Piotr Koniusz, “Power normalizing second-order similarity network for few-shot learning,” in *WACV*, 2019.
- [16] Wenbin Li, Jinglin Xu, Jing Huo, Lei Wang, Yang Gao, and Jiebo Luo, “Distribution consistency based covariance metric networks for few-shot learning,” in *AAAI*, 2019.
- [17] Davis Wertheimer and Bharath Hariharan, “Few-shot learning with localization in realistic settings,” in *CVPR*, 2019.
- [18] Wenbin Li, Lei Wang, Jinglin Xu, Jing Huo, Yang Gao, and Jiebo Luo, “Revisiting local descriptor based image-to-class measure for few-shot learning,” in *CVPR*, 2019.
- [19] Xiaoxu Li, Jijie Wu, Zhuo Sun, Zhanyu Ma, Jie Cao, and Jing-Hao Xue, “Bsnet: Bi-similarity network for few-shot fine-grained image classification,” *IEEE Transactions on Image Processing*, 2020.
- [20] Matthew D Zeiler, Dilip Krishnan, Graham W Taylor, and Rob Fergus, “Deconvolutional networks,” in *CVPR*, 2010.
- [21] Vinay P Nambodiri, Vincent De Smet, and Luc Van Gool, “Systematic evaluation of super-resolution using classification,” in *VCIP*. IEEE, 2011.
- [22] Muhammad Haris, Greg Shakhnarovich, and Norimichi Ukita, “Task-driven super resolution: Object detection in low-resolution images,” in *International Conference on Neural Information Processing*. Springer, 2021, pp. 387–395.
- [23] Xiyang Luo, Hossein Talebi, Feng Yang, Michael Elad, and Peyman Milanfar, “The rate-distortion-accuracy tradeoff: JPEG case study,” in *Data Compression Conference, DCC*. 2021, IEEE.
- [24] Steven Diamond, Vincent Sitzmann, Stephen P. Boyd, Gordon Wetzstein, and Felix Heide, “Dirty pixels: Optimizing image classification architectures for raw sensor data,” *CoRR*, vol. abs/1701.06487, 2017.
- [25] Boyi Li, Xiulian Peng, Zhangyang Wang, Jizheng Xu, and Dan Feng, “Aod-net: All-in-one dehazing network,” in *ICCV*, 2017.
- [26] Zhuang Liu, Tinghui Zhou, Zhiqiang Shen, Bingyi Kang, and Trevor Darrell, “Transferable recognition-aware image processing,” *CoRR*, vol. abs/1910.09185, 2019.
- [27] Mengye Ren, Eleni Triantafillou, Sachin Ravi, Jake Snell, Kevin Swersky, Joshua B. Tenenbaum, Hugo Larochelle, and Richard S. Zemel, “Meta-learning for semi-supervised few-shot classification,” in *ICLR*, 2018.
- [28] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie, “The caltech-ucsd birds-200-2011 dataset,” 2011.
- [29] Ziqi Zhou, Xi Qiu, Jiangtao Xie, Jianan Wu, and Chi Zhang, “Binocular mutual learning for improving few-shot classification,” in *ICCV*, 2021.
- [30] Wei-Yu Chen, Yen-Cheng Liu, Zsolt Kira, Yu-Chiang Frank Wang, and Jia-Bin Huang, “A closer look at few-shot classification,” in *ICLR*, 2019.
- [31] Yinbo Chen, Zhuang Liu, Huijuan Xu, Trevor Darrell, and Xiaolong Wang, “Meta-baseline: Exploring simple meta-learning for few-shot learning,” in *ICCV*, 2021.