# Few-Shot Learning with Geometric Constraints

Hong-Gyu Jung and Seong-Whan Lee, *Fellow, IEEE*

*Abstract*—In this paper, we consider the problem of few-shot learning for classification. We assume a network trained for base categories with a large number of training examples, and we aim to add novel categories to it that have only a few, e.g., one or five, training examples. This is a challenging scenario because (1) high performance is required in both the base and novel categories, and (2) training the network for the new categories with few training examples can contaminate the feature space trained well for the base categories. To address these challenges, we propose two geometric constraints to fine-tune the network with a few training examples. The first constraint enables features of the novel categories to cluster near the category weights, and the second maintains the weights of the novel categories far from the weights of the base categories. By applying the proposed constraints, we extract discriminative features for the novel categories while preserving the feature space learned for the base categories. Using public datasets for few-shot learning that are subsets of ImageNet, we demonstrate that the proposed method outperforms prevalent methods by a large margin.

*Index Terms*—Image Recognition, Few-Shot Learning, Deep Learning, Neural Network, Geometric Constraint.

## I. INTRODUCTION

DEEP networks have achieved state-of-the art performance in a variety of fields, such as visual recognition [1]–[5], object detection [6] and semantic segmentation [7]. Compared to previous studies [8]–[12], these impressive achievements are primarily owing to the availability of large numbers of training examples. For instance, current deep networks for image classification match human performance but require a large number of training examples [1]. Humans can quickly and accurately recognize novel categories using only one example based on previous knowledge [13]. Mimicking the cognitive abilities of humans is difficult, as the layers of deep networks are heavily optimized for the trained categories. In particular, the top layers tend to learn category-specific features [14], which are not flexible enough to learn novel categories without forgetting the parameters learned before. Thus, to continually add novel categories, the network must be re-trained by using a large number of training examples for the base and novel categories [15]. Apart from the complexity of training, building large numbers of training examples requires a significant amount of human effort and incurs high costs.

To solve this problem, few-shot learning, which can train neural networks using a few training examples, has attracted considerable research interest. In this paper, we consider a network problem where its base categories have been trained using a large number of training examples and novel categories with a few, e.g., one or five, training examples are to be added. This problem is challenging as the network needs to exhibit high performance in both the base and novel categories. In other words, an effective algorithm needs to be able to generate discriminative features and weights[1] for the novel categories without affecting those of the base categories.

One simple way to exploit training examples for a new category is to fine-tune a network by adding nodes to a classifier. However, this approach can destroy a feature space well constructed for base categories owing to the small number of training examples [17]. Instead of fine-tuning a network, recent studies have attempted to train it to become generalizable to unseen categories. For example, training examples for base categories can be used to train a weight predictor where the weights of both base and novel categories are calculated using activations from a feature extractor [18]. Similarly, the weights of novel categories can be generated by exploiting the weights of base categories with an attention mechanism [19]. To predict the weights of novel categories, these studies used activations or weights well learned for base categories with a large number of training examples. Nevertheless, the feature space for novel categories has not been explicitly considered, and whether the feature space trained for base categories is suitable for novel categories unseen during the training procedure is also uncertain.

In this paper, we propose a fine-tuning strategy that trains novel categories with a few training examples while not contaminating the feature space well learned for the base categories. The goal is to allow the base and novel categories to co-exist in a common space without interference. To this end, we first restrict the feature space to a high-dimensional sphere by normalizing the features and weights. This eliminates information related to magnitudes, and thus makes it easy to control the feature space using only angular information. Based on the high-dimensional sphere, we propose two geometric constraints. The first is called weight-centric feature clustering. For a given category, when used with cross-entropy loss, this reduces the angular distance between the features and the weights. The second is angular weight separation; this separates category weights based on the angular distance.

Fig. 1 shows our motivation for the proposed geometric constraints. From the CIFAR 100 dataset [16], we chose 10 base categories consisting of five pairs, each from a parent category. For example, the images of "Motorcycle" and "Bicycle" were included in the parent category, Vehicles

H.-G. Jung is with the Department of Brain and Cognitive Engineering and S.-W. Lee is with the Department of Artificial Intelligence, Korea University, Anam-dong, Seongbuk-ku, Seoul 02841, Korea. e-mail: {hkjung00, sw.lee}@korea.ac.kr (Corresponding author: Seong-Whan Lee).

[1]In this paper, we alternately denote the output of a feature extractor as features and activations. Weights are referred to as category weights used in a classifier.

(a) Trained for base categories  (b) Fine-tuned by Softmax  (c) Fine-tuned by geometric constraints
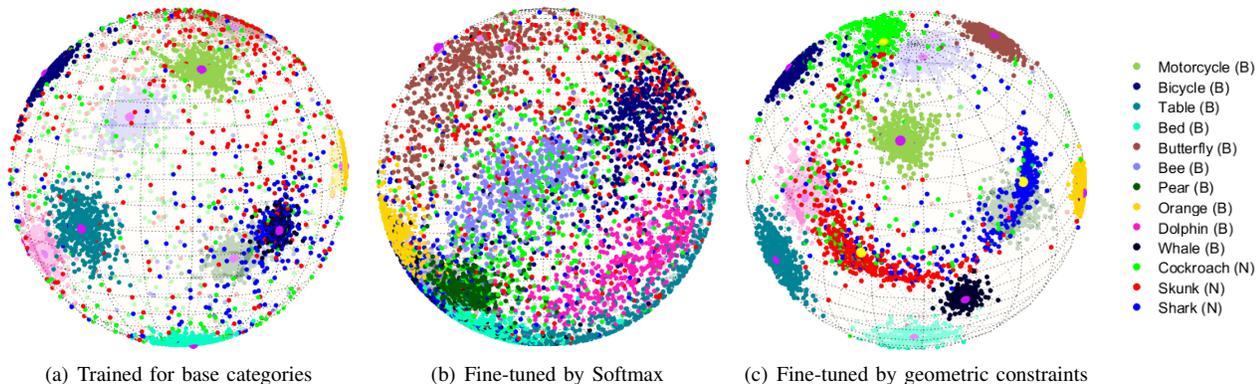
Fig. 1. (Best viewed in color) Examples of base and novel categories in the feature space. We trained 10 base categories chosen from the CIFAR 100 dataset [16] using 500 training examples for each base category. (a) shows features extracted from both base and novel categories. The base categories are clustered near the weights (hot pink). "Shark," one of the novel categories, is clustered well but its location overlaps with that of "Whale". The other two novel categories are spread out. (b) shows the results after the network was fine-tuned with five training examples for each novel category. For the network, we added classification weights for novel categories to the classifier and froze the classification weights for base categories during fine-tuning. We did not use any training examples for base categories. (c) demonstrates that our proposed geometric constraints can locate the features of the novel categories discriminatively and preserve the feature space of the base categories. The yellow circle means the weights of novel categories. The labels of the base (B) and novel (N) categories are shown on the far right. For all experiments, we used SGD with 300 iterations and the learning rates were (a) 1e-2, (b) 1e-4 and (c) 1e-3.

1. For novel categories, we chose "Cockroach" and "Skunk" which were not similar to any of the base categories, and "Shark" whose images looked similar to those in "Whale". Using 500 training examples for each base category, we first trained a shallow convolutional neural network (CNN) that had three nodes for the feature extractor and the softmax classifier followed by cross-entropy loss. We then extracted features of the examples for the base and novel categories. The extracted features are shown in Fig. 1. Features of "Shark" are located close to those of "Whale," which were well trained, and features of the other two novel categories are spread out. We can also see that merely fine-tuning the novel categories with five training examples contaminates the feature space that was well learned for the base categories. Furthermore, it is clear that our proposed method not only preserves features of the base categories but also constructs discriminative features of novel categories using only five training examples for each category. In Section V, we show that the proposed method can be applied to more complex datasets such as a subset of ImageNet built for few-shot learning.

To sum up, our contributions in this paper are threefold:

- Assuming only a few training examples for novel categories, a fine-tuning strategy is proposed to classify both base and novel categories. That is, we do not use any training examples of the base categories to fine-tune a network or for few-shot learning.
- We develop two geometric constraints that help classify the novel categories while maintaining the feature space previously learned for the base categories.
- We demonstrate that the proposed method achieves state-of-the art performance and we analyse it using several CNNs.

We discuss the overall model in Section III. The details of the geometric constraints are described in Section IV and the experimental results are detailed in Section V. We discuss the proposed method and the result in Section VI. Finally, we present the conclusions of this work in Section VII.

## II. RELATED WORK

A typical few-shot learning problem does not assume the existence of base categories with a large number of training examples, and its objective is to only classify novel categories with a few training examples. Meta-learning has recently attracted interest in the context of solving this problem [20]–[24]. In meta-learning, we use meta-training, meta-validation and meta-test datasets. Given a meta-training dataset, we randomly select a few training examples to obtain the parameters of the network and randomly choose test examples to generate a loss and train it. This process is repeated several times to make the network generalizable to unseen categories. Finally, the algorithms are evaluated in the meta-test set by the way they used in the meta-training set. Note that the three meta-datasets have disjoint categories.

Based on the concept of meta-learning, a matching network [23] was proposed to jointly learn two embedding functions for training examples and a test example. This method uses a deterministic distance metric to compare them. Unlike in a fixed manner, the proposed method in [24] learns a relation module to compare the results of the two embedding functions. A prototypical network [22] learns a prototype for each category, so that the examples discriminatively cluster around the prototypes corresponding to each category. Some studies [20], [25] developed memory modules that store useful information from training examples and exploited their memories when testing them. On the contrary, the relationship between activations and weights was established by [21]. The authors used $l2$ normalization layers to benefit from cosine similarity that renders activations and weights in the final layer symmetric. Thus, the activations of training examples for novel categories can be regarded as weights in the final layer.

However, the purpose of the above studies was to only classify novel categories with a few training examples, and they did not assume the existence of previously-trained base categories. In this respect, the most relevant works to ours are [18], [19], [26] and [27]. A weight predictor trained on

a large number of training examples for base categories was proposed by [18]. Following training, the weight predictor was used to calculate the weights of novel categories based on the activation of a few training examples. An attention-based weight generator [19] was developed to predict the weights of novel categories by exploiting those of the base categories. These two studies take advantage of activations or weights that are well learned for the base categories with a large number of training examples, and do not need to fine-tune a network for novel categories. However, it is unclear whether the feature space trained for the base categories is also suitable for novel categories hitherto unseen.

In terms of a generative model, Wang *et al.* [26] generated data that has similar characteristics to the training examples for novel categories. Similarly, Hariharan and Girshick [27] presented an example generation function where transformations learned from base categories are applied to the examples of novel categories. Then, they re-trained a classifier using a large dataset for base categories and the generated examples for novel categories. However, the complexity of this training procedure can be burdensome for few-shot learning. Instead, we propose a fine-tuning strategy that uses only a few training examples for novel categories. Moreover, fine-tuning the network does not affect the feature space learned for the base categories and generates discriminative features for novel categories.

Finally, a common concept exploiting marginal distance was proposed by [28]. However, the paper aims to classifying only novel categories without any distractors (e.g., base categories). This means their idea cannot be directly applied to our scenario. Specifically, they show the triplet (or contrastive) loss can be applied to few-shots. But, the losses should form triplets of anchor, positive and negative examples for training. Note that, in our scenario, generating triplets for training incurs very high complexity since we deal with both base and novel categories. Meanwhile, our method only uses a few training examples (not pairs of examples) and we will show our method works for prevalent datasets for few-shot image recognition.

## III. OVERALL FRAMEWORK

In this section, we define the few-shot learning problem, introduce a two-stage training procedure to classify base and novel categories, and explain how to test an example.

### A. Problem Definition

In our classification problem, the dataset $D_{Base} = \{D_{Base}^{Tr}, D_{Base}^{Val}, D_{Base}^{Test}\}$ is available to train, validate, and test the base categories. Similar to ImageNet [29], $D_{Base}$ is composed of a large number of examples for $C_{Base}$ categories. Then, (1) we train a network with $D_{Base}^{Tr}$ to classify $C_{Base}$ categories.

Based on the network, we aim to add novel categories to it. The dataset $D_{Novel}$ is available for novel categories, containing $C'_{Novel}$ categories that are disjoint to $C_{Base}$. (2) We randomly choose $C_{Novel} (\in C'_{Novel})$ categories, and $k$ training and $T_{Novel}$ test examples from each category. Then, the network is further trained using $k \times C_{Novel}$ examples.

This setting is called $C_{Novel}$-way $k$-shot learning. The typical number of $k$ is one or five.

Finally, (3) the network is evaluated using $T_{Base} (\in D_{Base}^{Test})$ and $T_{Novel}$ examples. Thus, given $C_{Both} = C_{Base} \cup C_{Novel}$, the goal is to correctly classify an example as belonging to one of $C_{Both}$ categories. Steps (2) and (3) are repeated several times to obtain a 95% confidence interval.

### B. Two-Stage Training Procedure

**Notation** We consider a network composed of a feature extractor and a classifier. We use $f_i$ to denote the feature extracted from the $i$-th example through the feature extractor. The weights of the base categories for the classifier are $W_B = [w_B^1 w_B^2 \cdots w_B^{n_B}]$, where $w_B^j$ is the column vector for the $j$-th category and $n_B$ the number of the base categories. Likewise, the weights of the novel categories are denoted by $W_N = [w_N^1 w_N^2 \cdots w_N^{n_N}]$. If $l2$ normalization is applied, we denote it by $\widetilde{v} = v/\|v\|$, where $v$ is a column vector and $\widetilde{V} = [\widetilde{v}^1 \widetilde{v}^2 \cdots \widetilde{v}^n]$, where $\widetilde{v}^j$ is a normalized column vector. We define scores in a classification layer as $S_{Base} = \{\widetilde{f}_i^T \widetilde{w}_B^1, \widetilde{f}_i^T \widetilde{w}_B^2, \cdots, \widetilde{f}_i^T \widetilde{w}_B^{n_B}\}$ and $S_{Novel} = \{\widetilde{f}_i^T \widetilde{w}_N^1, \widetilde{f}_i^T \widetilde{w}_N^2, \cdots, \widetilde{f}_i^T \widetilde{w}_N^{n_N}\}$.

**Training Stage 1** As shown in Fig. 2, a network is trained on $D_{Base}^{Tr}$ to classify $C_{Base}$ categories. The features and weights are normalized by $l2$-norm. For classification, we use the softmax layer followed by cross-entropy loss. This can be expressed as

$$L_{cls} = -\frac{1}{M} \sum_{i=1}^{M} log \frac{e^{s \widetilde{f}_i^T \widetilde{w}_B^{y_i}}}{\sum_{j=1}^{n_B} e^{s \widetilde{f}_i^T \widetilde{w}_B^j}}, \qquad (1)$$

where $M$ is a batch size, $y_i$ refers to the category of the $i$-th example, and $\widetilde{f}_i^T \widetilde{w}_B^j$ is the cosine similarity between the weight and the feature. $s$ is a learnable parameter.

Eq. 1 has two notable aspects. First, we normalize the features and weights by $l2$-norm. This locates the feature space in a high-dimensional sphere. In other words, we consider only angular information with respect to features and weights. This helps us in training stage 2 to geometrically control the locations of the features and weights in terms of angular distances. Second, we apply a scale parameter $s$. As addressed in [30], the range $[-1, 1]$ of $\widetilde{f}_i^T \widetilde{w}_B^j$ is too small to obtain a sufficient gradient for training. Thus, it is possible for the network to fail to converge. To solve this problem, we scale up the cosine similarity by using the learnable parameter $s$ [30].

In addition to Eq. 1, we also apply a loss called weight-centric feature clustering that will be introduced in Section IV. This is to place the weights for base categories at predictable locations, such as the center of features for each category.

**Training Stage 2** We now consider $D_{Novel}$ and the network trained in training stage 1. The parameters of bottom blocks are frozen. However, we duplicate the top convolution blocks in order to fine-tune them, and build a classifier consisting of $C_{Novel}$ nodes. Given a training example of
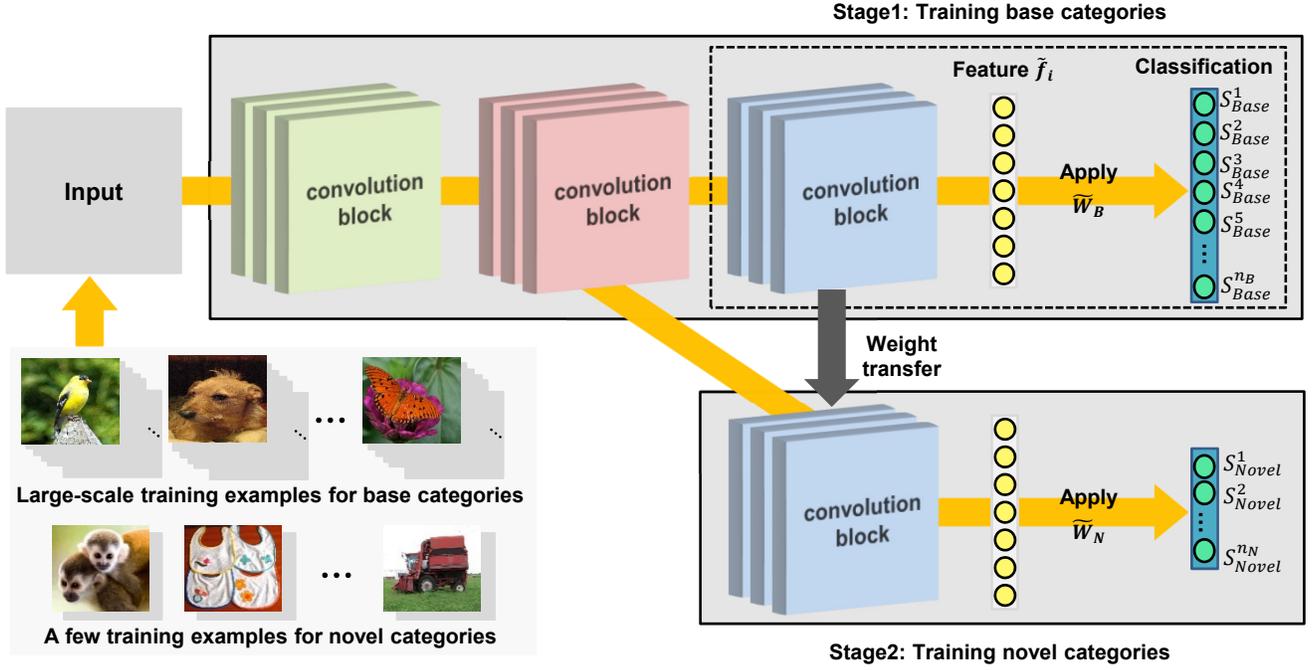
Fig. 2. Overall framework. In training stage 1, a network is trained for the base categories with a large number of training examples. In training stage 2, we duplicate parts of the convolution blocks to be fine-tuned and build a classification layer for the novel categories. The parameters of the bottom blocks are frozen. We forward the training examples for the novel categories to the bottom blocks, the duplicated top blocks and the novel classifier. In the test stage, an example is forwarded to all layers, and we calculate the largest value from the classification layers of the base and novel categories.

$D_{Novel}$, it is forwarded to the bottom blocks, the duplicated top layers, and the novel classifier. In training stage 2, we do not use the top convolution blocks and the classifier used for training the base categories, which are shown using broken lines in Fig. 2. Based on this network architecture, we use the loss functions proposed in Section IV to train examples of novel categories.

### C. Two-Stream Test Procedure

Once training stages 1 and 2 are complete, we use test examples from $D_{Base}^{Te}$ and $D_{Novel}$. We forward a test example $I_i$ to all layers used in training stages 1 and 2. The example is then classified into a category by

$$\underset{k}{\arg\max}(S_{Both}^k), \qquad (2)$$

where $S_{Both}^k$ is the $k$-th element of $S_{Base} \cup S_{Novel}$. To evaluate base categories or novel categories only, we substitute $S_{Both}^k$ into $S_{Base}^k$ or $S_{Novel}^k$, respectively.

### IV. LOSS FUNCTIONS

In this section, we elaborate on the loss functions used in training stage 2. As we use separate blocks and weights learned *independently* for the base and novel categories, it is unnatural for them to co-exist in a common feature space. To address this issue, we introduce three functions, cross-entropy loss and two geometric constraints, to extract discriminative features for novel categories while preserving the feature space trained for the base categories in training stage 1.

### A. Cross-Entropy Loss

We use a similar loss function to Eq. 1 but replace the weights for the base categories with those for the base and novel categories and use the fixed scale parameter $s$ trained in training stage 1.

### B. Weight-Centric Feature Clustering

Because novel categories are not trained in stage 1, their features may not be well clustered, and thus intra-class variation may be large. To solve this problem, we propose the weight-centric feature clustering (WCFC) defined as

$$L_{WCFC} = \sum_{i=1}^{C_{Novel}} -log(cos\ \theta_{g(f^i),\widetilde{w}_N^i}), \qquad (3)$$

where $cos\ \theta_{g(f^i),\widetilde{w}_N^i} = g(f^i)^T \cdot \dfrac{w_N^i}{\|w_N^i\|}$. We define a function $g(\cdot)$ as two types. For type 1,

$$g(f^i) = \frac{\bar{f}^i}{\|\bar{f}^i\|}, \qquad (4)$$

where $\bar{f}^i$ is the arithmetic mean of features of the training examples for the $i$-th category. As the initial weight $w_N^i$, we use $\bar{f}^i$. For type 2,

$$g(f^i) = \frac{\sum_i \widetilde{f}^i}{\|\sum_i \widetilde{f}^i\|}. \qquad (5)$$

When using this type, we set the initial weight $w_N^i$ to $\sum_i \widetilde{f}^i$.

(a) Feature space after training stage 1          (b) Feature space after training stage 2
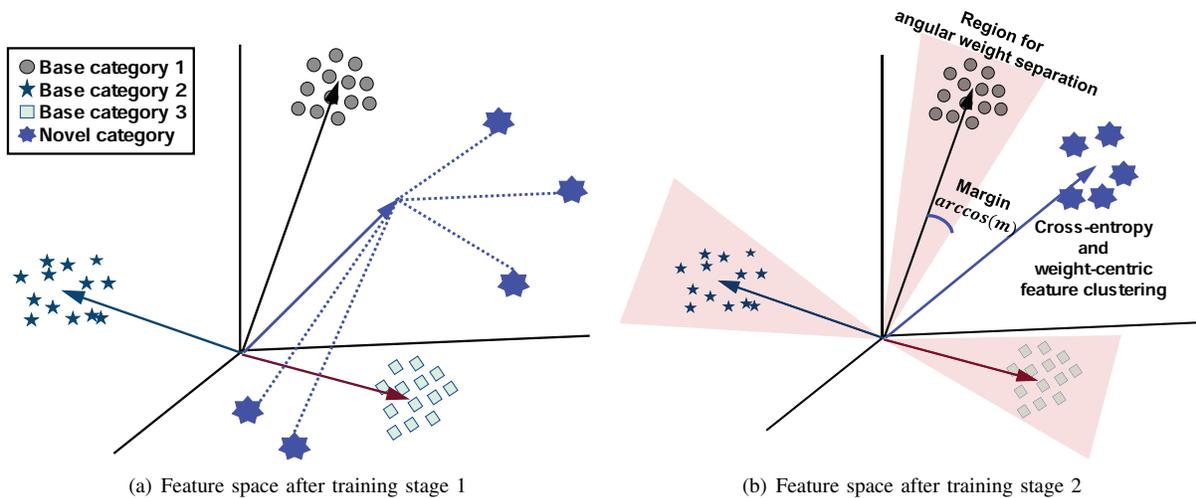
Fig. 3. Visual interpretation of training stages 1 and 2 in 3-dimensional feature space. (a) After the base categories have been trained, their features are clustered close to the weights. (b) By applying the fine-tuning strategy using the proposed loss functions, features of the novel categories are clustered near the weights, and are sufficiently distant from weights of other categories. Note that all features and weights are located on a high-dimensional sphere in practice.

Note that Eq. 4 is the normalization of averaged features and Eq. 5 is the normalization of the sum of normalized features. Thus, Eq. 5 has more degrees of freedom for magnitude and we have found that this is beneficial to novel categories for few-shot learning.

When used with cross-entropy loss, this weight-centric feature clustering states that features corresponding to a category should be located near the weight. In Eq. 3, maximizing the cosine similarity $cos\theta_{g(f^i),\widetilde{w}_N^i}$ is the core concept.

### C. Angular Weight Separation

Thus far, we have clustered features for novel categories near the weights to reduce intra-class variation. Our goal now is to ensure that the weights and features for the novel categories are sufficiently far from those for the base categories. To achieve this, we first define the angular distance $u_{ij}$ between weights as follows:

$$u_{ij} = \begin{cases} cos\ \theta_{\widetilde{w}^i,\widetilde{w}_N^j}, & \text{if } \widetilde{w}^i \not\equiv \widetilde{w}_N^j \\ 0, & \text{otherwise} \end{cases}, \quad (6)$$

where $\widetilde{w}^i$ is the $i$-th column vector of $\left[\widetilde{W}_B \widetilde{W}_N\right]$.

Then, we maximize $\theta_{\widetilde{w}^i,\widetilde{w}_N^j}$, so that $u_{ij}$ is smaller than margin $m$. This is called the angular weight separation (AWS) constraint and is defined as

$$L_{AWS} = \frac{\sum_{i,j} -log\left(-u_{ij} \cdot \mathbb{1}_M\left(u_{ij}\right) + 1\right)}{\sum_{i,j} \mathbb{1}_M\left(u_{ij}\right)}, \quad (7)$$

where $M = \{u_{ij} \mid u_{ij} > m,\ \forall i,j\}$ and $\mathbb{1}_M\left(u_{ij}\right) = \begin{cases} 1, & \text{if } u_{ij} \in M \\ 0, & \text{otherwise} \end{cases}$ is the indicator function. The logarithm and a factor of 1 are used to align the loss function with the scale of the other loss functions. During the training process, when all novel weights are separated from other weights in terms of the angular distance, we turn off the AWS loss.

### D. Regularization

We sum up all the losses for the final loss function.

$$L_{total} = \gamma L_{cls} + \alpha L_{WCFC} + \beta L_{AWS} \quad (8)$$

Using Eq. 8, we train the convolution blocks and weights for the novel categories in training stage 2. Fig. 3 shows the role of each loss function. For all experiments, we use $\gamma\alpha\beta = 111$. In Section V, we discuss how different settings of the hyper-parameters impact on the performance of the proposed method.

### V. EXPERIMENTS

To evaluate the proposed method, we used two common datasets for few-shot learning: *mini*ImageNet [23] and Bharath & Girshick's dataset [27]. The two datasets were built based on ImageNet [29].

### A. miniImageNet

*Mini*ImageNet [23] is the most commonly used dataset for few-shot learning. It is composed of 64 training, 16 validation, and 20 test categories. Each category contains 600 examples of size $84 \times 84$. We used the split provided by [31]. Most studies use *mini*ImageNet as follows: Using the training categories, $k$ training examples for each of $C_{Novel}$ categories are chosen to obtain the parameters of a network and $T$ test examples are randomly chosen to train it. Then, this process is repeated several times to warm-up or generalize the parameters of the network. Finally, the network is evaluated on the test categories in the same way as it was on the training categories. Therefore, this dataset has been primarily used to only classify novel categories in meta-learning.

Recently, Gidaris and Komodakis [19] extended the dataset by collecting 300 extra validation examples and 300 test examples for the 64 base categories. Thus, we used the original 600 examples to train the base categories and the additional examples for validation and testing. We trained and tested the base and novel categories as follows: In training stage 1,

TABLE I
CLASSIFICATION ACCURACY ON *mini*IMAGENET. THE RESULTS ARE AVERAGED TO OBTAIN A 95% CONFIDENCE INTERVAL. '-' DENOTES THAT THE PERFORMANCE WAS NOT REPORTED FOR THE MODEL. WE REFERRED TO [19] FOR THE ACCURACIES OF OTHER METHODS. † INDICATES THAT THE PROPOSED DATA AUGMENTATION TECHNIQUE WAS USED. THE HIGHEST NUMBERS ARE BOLDED.

| Models | Feature Extractors | 5-way 5-shot | | | 5-way 1-shot | | |
| --- | --- | --- | --- | --- | --- | --- | --- |
| | | Novel | Both | Base | Novel | Both | Base |
| Matching Nets [23] | C64F | 55.30 | - | - | 43.60 | - | - |
| Finn *et al.* [32] | C64F | 63.10$\pm$0.92 | - | - | 48.70$\pm$1.84 | - | - |
| Prototypical Nets [22] | C64F | 68.20$\pm$0.66 | - | - | 49.42$\pm$0.78 | - | - |
| Relation Nets [24] | C64F | 65.32$\pm$0.66 | - | - | 50.44$\pm$0.82 | - | - |
| Mishra *et al.* [33] | ResNetS | 68.88$\pm$0.92 | - | - | 55.71$\pm$0.99 | - | - |
| Gidaris and Komodakis [19]$^\dagger$ | C64F | 73.27$\pm$0.59 | 57.72 | 66.96 | 57.78$\pm$0.78 | 48.13 | 66.58 |
| | ResNetS | 70.32$\pm$0.66 | 55.93 | 79.58 | 56.76$\pm$0.80 | 49.68 | 79.43 |
| Proposed | C64F | 72.92$\pm$0.64 | 58.00 | 67.12 | 53.39$\pm$0.81 | 45.32 | 67.12 |
| | C64F-Dropout | 73.21$\pm$0.65 | 58.77 | 66.98 | 55.74$\pm$0.83 | 48.11 | 66.98 |
| | **ResNetS** | **78.00$\pm$0.61** | **68.16** | **79.78** | **58.52$\pm$0.82** | **56.05** | **79.78** |
| Ablation (w/o Fine-Tuning) | C64F | 72.85$\pm$0.63 | 54.43 | 67.12 | 53.43$\pm$0.81 | 43.17 | 67.12 |
| | C64F-Dropout | 71.32$\pm$0.63 | 51.71 | 66.98 | 54.86$\pm$0.80 | 38.52 | 66.98 |
| | ResNetS | 68.68$\pm$0.67 | 47.11 | 79.78 | 54.12$\pm$0.83 | 43.11 | 79.78 |

we trained 64 base categories with 600 examples. In training stage 2, we randomly chose $C_{Novel}$ novel categories and $k$ examples for each category to fine-tune the network using the proposed loss functions. In the test stage, we randomly chose $T$ examples from each of the 64 base categories and $C_{Novel}$ novel categories to test the network. We iterated training stage 2 and the test stage to obtain a 95% confidence interval.

For the feature extractor, we used a shallow CNN called C64F with four convolution blocks: each block had $3 \times 3$ convolutions, batch normalization [34], ReLU [35] and $2 \times 2$ max-pooling. The sizes of feature maps for the four convolution blocks were 64. For a deeper model, the small version of ResNet [1] proposed in [33] was used as in [19]. We call the network ResNetS. In all cases, we did not use ReLU for the last layer in the last convolutional block.

**Data Augmentation** As our method uses a fine-tuning strategy, it is important for training examples to generate a loss large enough to back-propagate the error. In this respect, we found that augmenting the training examples of each novel category helped improve performance. Given a training example, we performed zero-padding with 8 pixels on each border and randomly cropped it to make $84 \times 84$. It was then horizontally flipped with a probability of 0.5. This transformation is identical to that used when training the base categories in training stage 1. By iterating this process, we generated $nAug$ examples for each training example. In our experiments, we augmented examples to have 20 training examples for each novel category in total. Since this simple augmentation helps improve performance, more advanced example generation methods [27], [36] may further boost the performance of our algorithm.

**Effect of Dropout** Since our algorithm uses a fine-tuning strategy, it is crucial that a network has an ability to be generalizable to unseen categories. Typically, we rely on a large number of training examples to train a network

with a highly complicated distribution. For few-shot learning, however, we only have a few training examples, and thus other techniques for generalization should be considered. In our experiment, we have found that ResNetS significantly outperforms C64F. We assume that this is because not only the depth of the network but also the usage of Dropout [37] on the last fully connected layers. Thus, we add Dropout following a linear layer with $1,024$ neurons to C64F and set the dropout rate to 0.5. We call the network C64F-Dropout.

**Performance** In Table I, we compare our proposed method with prevalent methods [22]–[24], [32], [33] for few-shot learning and a model [19] considering both base and novel categories. It is clear that, as the feature extractor has a greater capacity to handle the training dataset, the performance of the proposed method improves significantly. For example, in training stage 1, ResNetS yields the accuracy of around 79% in the base categories, higher than the $65 \sim 68\%$ of C64F. In this case, our method outperforms others by a large margin. Therefore, the proposed fine-tuning strategy with geometric constraints is clearly advantageous when a deeper network is available. It is worth noting that the proposed strategy dose not affect the accuracy on base categories even after fine-tuning the network. For ResNetS, the performance improvements in both categories indicate that the geometric constraints not only preserve the feature space of base categories, but also extract discriminative features for novel categories, when compared to [19] that aims to obtain the weights of novel categories based on those of base categories. Note that Table I results from the choice of the weight-centric feature clustering as follows: we always use Eq. 5 for training stage 2. For training stage 1, we used Eq. 5 for C64F; and Eq. 4 for ResNetS. We used 0.6 for the margin value for the angular weight separation.

**Ablation Study** To assess the benefit of our fine-tuning strategy, we show the accuracy when not fine-tuning the network. In this case, the weights of novel categories are

TABLE II
TOP-5 CLASSIFICATION ACCURACY ON THE IMAGENET. THE RESULTS ARE AVERAGED TO OBTAIN A 95% CONFIDENCE INTERVAL. '-' DENOTES THAT THE PERFORMANCE WAS NOT REPORTED FOR THE MODEL. 'W/ H' IS SHORT FOR W/ HALLUCINATION WHERE THE ALGORITHMS USE EXAMPLE GENERATION TECHNIQUES. $^{\dagger}$ INDICATES THAT THE PROPOSED DATA AUGMENTATION TECHNIQUE WAS USED. THE HIGHEST NUMBERS ARE BOLDED. *RESULTS WERE REPORTED BY [26] AND WE REFERRED TO [19].

| Models | Novel | | | | | Both | | | | | Both with prior | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | k=1 | 2 | 5 | 10 | 20 | k=1 | 2 | 5 | 10 | 20 | k=1 | 2 | 5 | 10 | 20 |
| Prototypical Nets [22]* | 39.3 | 54.4 | 66.3 | 71.2 | 73.9 | 49.5 | 61 | 69.7 | 72.9 | 74.6 | 53.6 | 61.4 | 68.8 | 72 | 73.8 |
| Matching Nets [23]* | 43.6 | 54 | 66 | 72.5 | 76.9 | 54.4 | 61 | 69 | 73.7 | 76.5 | 54.5 | 60.7 | 68.2 | 72.6 | 75.6 |
| Logistic Regression [26]* | 38.4 | 51.1 | 64.8 | 71.6 | 76.6 | 40.8 | 49.9 | 64.2 | 71.9 | 76.9 | 52.9 | 60.4 | 68.6 | 72.9 | 76.3 |
| Logistic Regression w/ H [27]* | 40.7 | 50.8 | 62 | 69.3 | 76.5 | 52.2 | 59.4 | 67.6 | 72.8 | 76.9 | 53.2 | 59.1 | 66.8 | 71.7 | 76.3 |
| SGM w/ H [27] | - | - | - | - | - | 54.3 | 62.1 | 71.3 | 75.8 | 78.1 | - | - | - | - | - |
| Batch SGM [27] | - | - | - | - | - | 49.3 | 60.5 | 71.4 | 75.8 | 78.5 | - | - | - | - | - |
| Prototype Matching Nets w/ H [26]* | 45.8 | 57.8 | 69 | 74.3 | 77.4 | 57.6 | 64.7 | 71.9 | 75.2 | 77.5 | 56.4 | 63.3 | 70.6 | 74 | 76.2 |
| Prototype Matching Nets [26]* | 43.3 | 55.7 | 68.4 | 74 | 77 | 55.8 | 63.1 | 71.1 | 75 | 77.1 | 54.7 | 62 | 70.2 | 73.9 | 75.9 |
| Gidaris and Komodakis [19]$^{\dagger}$ | 49.08 ±.23 | 59.97 ±.15 | 70.77 ±.10 | 75.51 ±.06 | 78.26 ±.05 | 60.03 ±.14 | 66.71 ±.10 | 73.77 ±.06 | 76.94 ±.04 | 78.77 ±.03 | 58.65 ±.14 | 65.3 ±.09 | 72.37 ±.06 | 75.54 ±.04 | 77.34 ±.03 |
| Proposed | 45.98 ±.22 | 58.37 ±.16 | 70.17 ±.09 | 75.49 ±.07 | 78.62 ±.05 | 57.11 ±.14 | 64.53 ±.10 | 71.77 ±.07 | 75.20 ±.06 | 78.10 ±.03 | 56.24 ±.14 | 63.77 ±.10 | 71.43 ±.06 | 74.91 ±.05 | 77.08 ±.03 |
| Proposed-Dropout | **49.57** ±.24 | **60.89** ±.16 | **71.07** ±.09 | **76.11** ±.07 | **78.84** ±.05 | **60.39** ±.14 | **67.44** ±.10 | **74.22** ±.06 | **77.32** ±.05 | **79.05** ±.03 | **58.91** ±.14 | **65.84** ±.10 | **72.71** ±.06 | **75.96** ±.05 | **77.88** ±.04 |

set to the sum of normalized features of training examples for each category. The weights and the feature extractor for novel categories are not trained further. We observe that our fine-tuning strategy significantly improves the performance when evaluating both the base and novel categories.

### B. Bharath & Girshick's Dataset

Bharath & Girshick proposed a larger dataset for few-shot learning [27]. The dataset is composed of 193 base categories and 300 novel categories for cross-validation, and there are 196 base categories and 311 novel categories for evaluation. The categories have training, validation and test examples as in ImageNet [29]. We used the categories provided by the most recent work [19] and we show the performance using the evaluation categories. As a feature extractor, we exploit ResNet10 [1]. We further show the performance can be boosted with ResNet10-Dropout that adds Dropout [37] of the dropout rate 0.2 following a linear layer with 1,024 neurons.

**Performance** For this dataset, we provide the performance called Both with prior [26]. When computing a probability that an image $x$ belongs to a class $k$, $p_k(x) = p(y = k|x)$, this considers a prior probability of whether an example belongs to base categories or novel categories. By following [19], we set $p(y \in C_{base}|x) = 0.2$ and $p(y \in C_{novel}|x) = 0.8$ for Both with prior. Table II shows top-5 classification accuracy on the ImageNet. Our algorithm outperforms the state-of-the-art methods for all shots. It is worth noting that without any hallucination and the training examples of base categories [26], [27], our algorithm improves the performance only using a few training examples for novel categories. The result also verifies that a fine-tuning strategy for few-shot learning

TABLE III
CLASSIFICATION ACCURACY ON THE IMAGENET USING INCREMENTAL LEARNING. FOR THE PROPOSED-DROPOUT, WE USED 20-SHOTS AND FOR INCREMENTAL LEARNING, WE FINE-TUNED THE FEATURE EXTRACTOR GRADUALLY FROM 1-SHOT TO 20-SHOTS.

| Models | Novel | Both | Both with prior |
|---|---|---|---|
| Proposed-Dropout | 78.84±.05 | 79.05±.03 | 77.88±.04 |
| Incremental learning | 79.87±.05 | 79.44±.04 | 78.59±.04 |

can be boosted using a network with Dropout. We used 0.4 for the margin value for the angular weight separation.

### C. Incremental Learning

Although the literature of few-shot learning mainly considers the performance when $k$-shot training examples are available, we would like to provide a possible application of the proposed method. Since our method fine-tunes a network, one might be interested in the performance when training examples are given gradually. For this experiment, we fine-tuned the feature extractor for novel categories continually from 1-shot to 20-shots with ResNet10-Dropout. As shown in Table III, the performance of incremental learning is superior to 20-shots. This suggests that a recognition system can be continually upgraded by the proposed method as novel examples are continually provided.

### D. Graphical Analysis

We aim to extract discriminative features for the novel categories while preserving the feature space learned for the base categories. To show that our algorithm accords with the purpose, we visualize the feature space before and after applying our fine-tuning method on *mini*ImageNet with
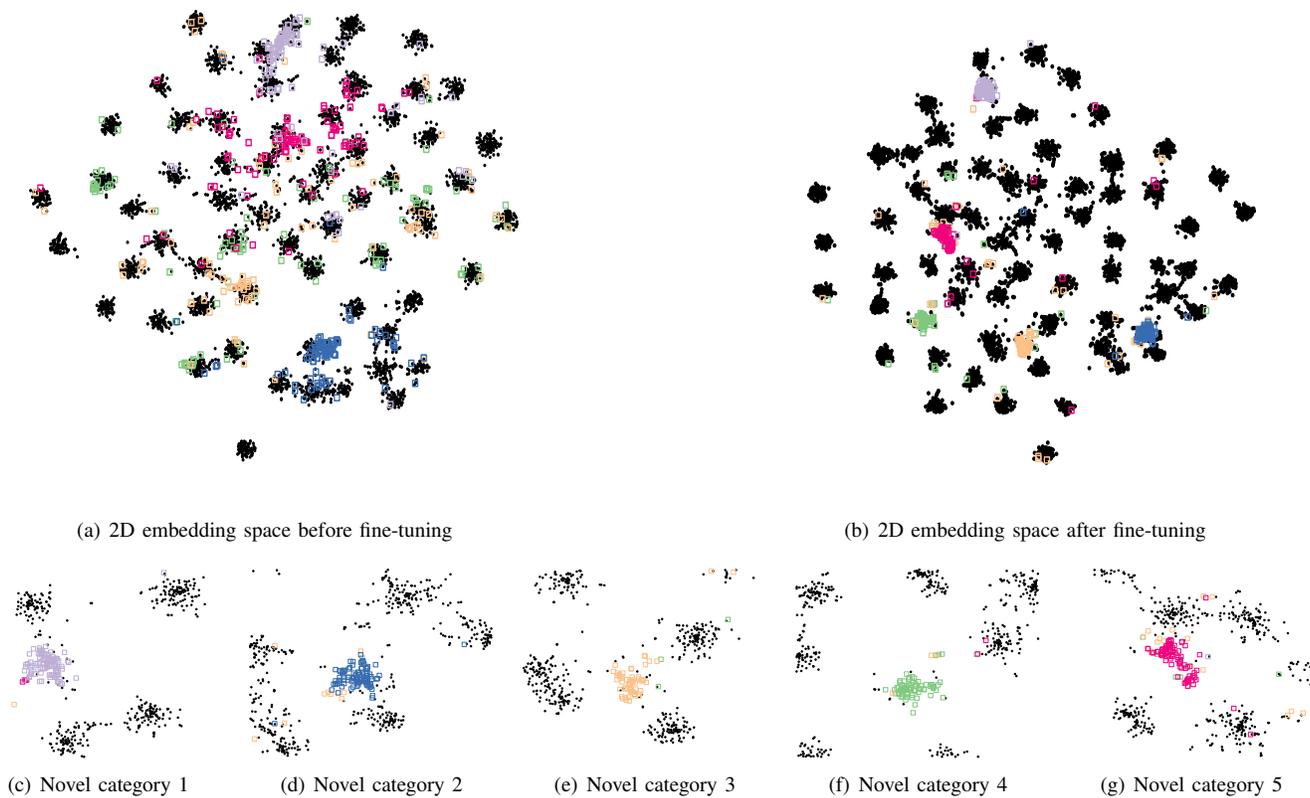
(a) 2D embedding space before fine-tuning

(b) 2D embedding space after fine-tuning

(c) Novel category 1    (d) Novel category 2    (e) Novel category 3    (f) Novel category 4    (g) Novel category 5

Fig. 4. (Best viewed in color) 2D embedding space (a) before and (b) after fine-tuning ResNetS. The black circles indicate base categories of *mini*ImageNet and novel categories are expressed as colored squares. (c)-(g) are enlarged figures to corresponding parts of (b). We observe that the features of novel categories are well placed between base categories.



(a) Original    (b) After fine-tuning    (c) Before fine-tuning    (d) Original    (e) After fine-tuning    (f) Before fine-tuning

Fig. 5. Feature maps for novel categories. (a),(d): original images, (b),(e): feature maps after fine-tuning and (c),(f): feature maps before fine-tuning. Our algorithm focuses on important parts of objects and removes noisy parts.

(a) Cosine similarity of ResNetS          (b) Cosine similarity of C64F          (c) Cosine similarity of ResNet10
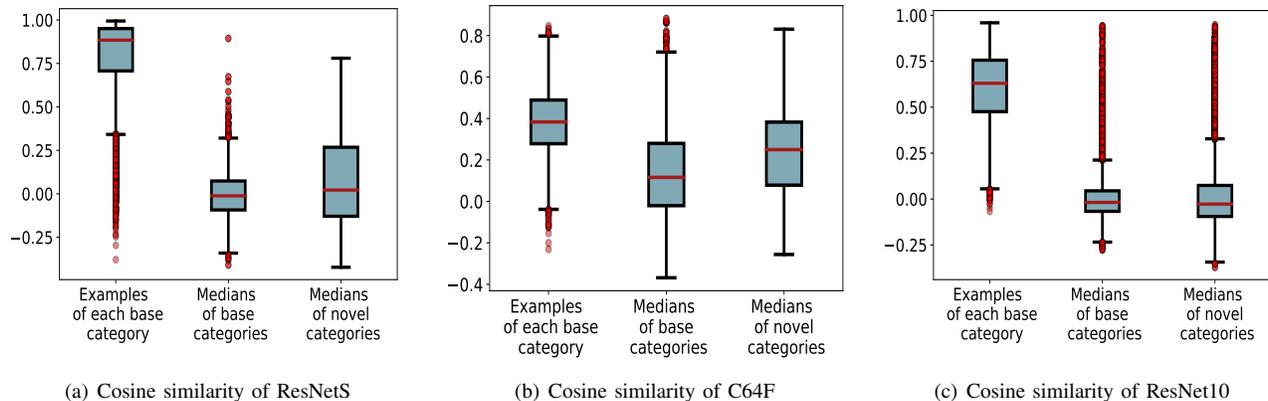
Fig. 6. Box plots of cosine similarity. Given a network trained on base categories, we first extracted features of the examples for base and novel categories. Based on the feature space, we obtained the medians of the features of each category. For figures (a)-(c), we calculated the cosine similarity among (left) examples belonging to each base category, (middle) medians of base categories and (right) medians of novel categories.

TABLE IV
COMPARISON BETWEEN THE PROPOSED METHOD AND THE SINGLE-STAGE FINE-TUNING APPROACH ON THE IMAGENET WITH RESNET10-DROPOUT.
FOR BASE CATEGORIES, THE AVERAGE ACCURACY OF ALL SHOTS IS REPORTED.

| Models | Novel | | | | | Both | | | | | Both with prior | | | | | Base |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | k=1 | 2 | 5 | 10 | 20 | k=1 | 2 | 5 | 10 | 20 | k=1 | 2 | 5 | 10 | 20 | |
| Proposed-Dropout | 49.57 | 60.89 | 71.07 | 76.11 | 78.84 | 60.39 | 67.44 | 74.22 | 77.32 | 79.05 | 58.91 | 65.84 | 72.71 | 75.96 | 77.88 | **93.55** |
| | ±.24 | ±.16 | ±.09 | ±.07 | ±.05 | ±.14 | ±.10 | ±.06 | ±.05 | ±.03 | ±.14 | ±.10 | ±.06 | ±.05 | ±.04 | |
| Single-stage fine-tuning | 45.08 | 57.01 | 69.98 | 75.1 | 78.04 | 55.23 | 63.34 | 71.08 | 73.82 | 75.24 | 54.44 | 61.55 | 69.76 | 73.27 | 75.34 | **90.45** |
| | ±.43 | ±.28 | ±.16 | ±.13 | ±.09 | ±.34 | ±.20 | ±.13 | ±.08 | ±.05 | ±.29 | ±.20 | ±.14 | ±.07 | ±.05 | |

ResNetS. We used T-SNE [38] for dimensionality reduction. As shown in Fig. 4, the features of novel categories are well placed between base categories. Without fine-tuning, features of a category are spread out. This spreadability implies that it is crucial to fine-tune a feature extractor for some novel categories.

Furthermore, we visualize feature maps for novel categories using Grad-CAM [39] that is a visualization method to analyse what deep networks learn. As shown in Fig. 5, our algorithm focuses on important parts of objects and removes noisy parts, which helps improve the classification accuracy.

## VI. DISCUSSION

We first discuss why the proposed method significantly outperforms the state of the art on *mini*ImageNet and is slightly better on the ImageNet. Then, we discuss the following questions: Does a single-stage fine-tuning approach work? How does the performance vary according to the possible combinations of the geometric constraints in Eq. 8? How sensitive is the performance to the margin value in Eq. 7? Finally, we discuss the limitation of the proposed method.

### A. Network Characteristic

Since performance gains over other methods vary depending on the networks, we analyse the reason in terms of network characteristics. Given a network trained on base categories, we first extracted features of the examples for base and novel categories. Based on the feature space, we obtained

TABLE V
COMPARISON BETWEEN THE PROPOSED METHOD AND THE SINGLE-STAGE FINE-TUNING APPROACH ON *mini*IMAGENET WITH RESNETS.

| Models | 5-way 5-shot | | | 5-way 1-shot | | |
|---|---|---|---|---|---|---|
| | Novel | Both | Base | Novel | Both | Base |
| Proposed | 78.00±0.61 | 68.16 | **79.78** | 58.52±0.82 | 56.05 | **79.78** |
| Single-stage fine-tuning | 76.78±0.59 | 58.62 | **64.58** | 55.11±0.78 | 52.95 | **72.32** |

the medians of the features of each category. For example, if we have 64 categories, this process produces 64 medians corresponding to each category. Then, we calculated the cosine similarity among (a) features of the examples belonging to each base category, (b) the medians of base categories and (c) the medians of novel categories. Specifically, (a) is to show how well the features cluster for a category they belong to. (b) and (c) are to investigate how sufficiently the categories are separated from each other. For statistical analysis, we used box plots. As shown in Fig. 6, features from ResNetS well cluster together for each category and categories are separated sufficiently. However, a lot of categories are close together on C64F and ResNet10.

Based on the above characteristics, we compare the proposed method with the novel weight generator by Gidaris and Komodakis [19] that shows comparable results to ours. To train the novel weight generator [19], the algorithm samples some base categories and regards them as novel categories.

This is to mimic the environment when real novel categories are given. For this training procedure to be effective and generalized, the distribution of features for base categories should resemble that of novel categories. However, this is not the case with ResNetS. As a result, despite the much higher capacity, the performance of ResNetS is worse than that of C64F. On the contrary, our proposed method relies on the feature discriminability trained on base categories. Thus, the proposed method produces valuable results with ResNetS but is less effective with C64F and ResNet10. Even though performance gains over other methods vary depending on network characteristic, we achieve the best performance on all datasets.

### B. Single-Stage Fine-Tuning Approach

To validate that a simple fine-tuning approach with only a few examples destroys a pre-trained network, we present the following experimental setting: after training a network using base categories, we added novel weights to the classifier. Given a few training examples, we then fine-tuned the novel weights of the classifier and the last convolution block of the feature extractor with cross-entropy loss. For this experiment, the feature extractor was shared for base and novel categories. We used SGD with learning rates $1e$-4 for ResNetS and $1e$-5 for ResNet10. As shown in Tables V and IV, the single-stage approach destroys the pre-trained network. Especially, the performance on base categories cannot be preserved unlike the proposed method.

### C. Effectiveness of each Loss Function

Since we use the three loss functions in Eq. 8, the effect of each loss function can be discussed. The accuracy for all possible combinations of the loss functions is shown in Fig. 7. We can observe the following trends. (a) For both categories, the angular weight separation plays a crucial role in the performance. (b) For novel categories, the use of all loss functions does not always result in the best performance. However, it stably produce comparable results for all experiments and the importance increases as the number of training examples increases.

### D. Sensitivity of the Margin Value

For the experiments, we have shown the results with the margin value that was empirically set for the angular weight separation. The purpose of the margin is to maintain the classification weights far from each other. Thus, the value should be determined according to the geometry of the feature space. Nonetheless, Fig. 8 shows that the performance is not too sensitive to the variation of the margin value. In other words, the performance does not vary irregularly and this enable us to find a peak point through validation examples.

### E. Limitation

Although the proposed method achieves the state-of-the art performance, it relies on the network capacity. Since we have only a few training examples, it is difficult to adjust

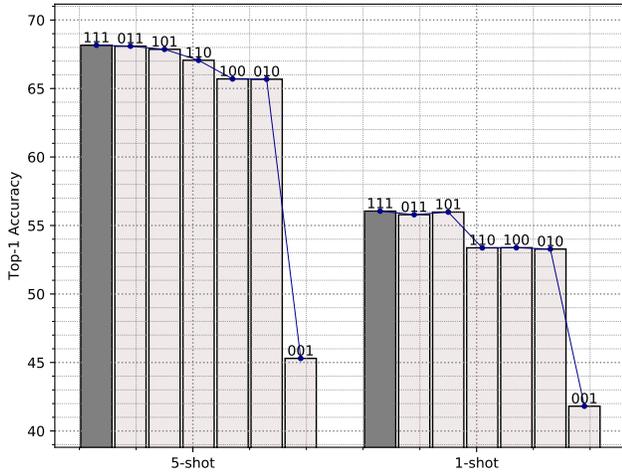| Model | k-Shot | Optimizer | Learning Rate | |
| --- | --- | --- | --- | --- |
| | | | Feature Extractor | Classifier |
| C64F | 1 | | 1E-03 | 1E-03 |
| | 5 | | 1E-05 | 1E-02 |
| C64F-Dropout | 1 | | 1E-04 | 5E-05 |
| | 5 | | 1E-05 | 1E-03 |
| ResNetS | 1 | Adam | 1E-04 | 1E-02 |
| | 5 | | 1E-04 | 1E-02 |
| ResNet10 & ResNet10 -Dropout | 1 | | 1E-05 | 1E-03 |
| | 2 | | 1E-05 | 1E-03 |
| | 5 | | 1E-05 | 1E-03 |
| | 10 | | 1E-05 | 1E-03 |
| | 20 | | 1E-05 | 1E-03 |

features for both base and novel categories. This is why we fine-tune only the parameters for novel categories and freeze the parameters for base categories. In this sense, to boost the performance, we need a network well trained on base categories where the features belonging to a category cluster together and features for different categories are sufficiently separated. To overcome the limitation, although inefficient, training examples for base categories can be used together to adjust the features of both base and novel categories. However, we have only a few training examples for novel categories as ever, which causes imbalanced categories. Thus, we believe that in any case we need a data generation technique for few-shot learning. We would like to leave this for the future work.
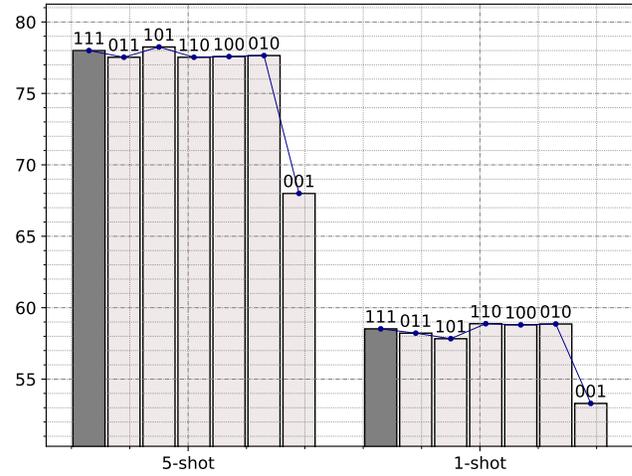
### VII. CONCLUSION AND FUTURE WORK

In this paper, we presented a fine-tuning strategy for few-shot learning. We considered a network trained for base categories with a large number of training examples and we aimed to add novel categories to it that had only a few training examples. We proposed two geometric constraints to extract discriminative features for the novel categories while preserving the feature space learned for the base categories. The first constraint enabled features of the novel categories to cluster near the category weights when combined with cross-entropy loss. The second maintained the weights of the novel categories far from the weights of the base categories. Applying the constraints, we showed that the accuracy on base categories was not affected even after fine-tuning the network and comparable performances to the base categories could be achieved on novel categories. The proposed method gained state-of-the art performance especially using a network with Dropout. Since we verified a simple way to augment training examples improves the performance, future work involves developing example generation techniques that are advantageous to our idea.
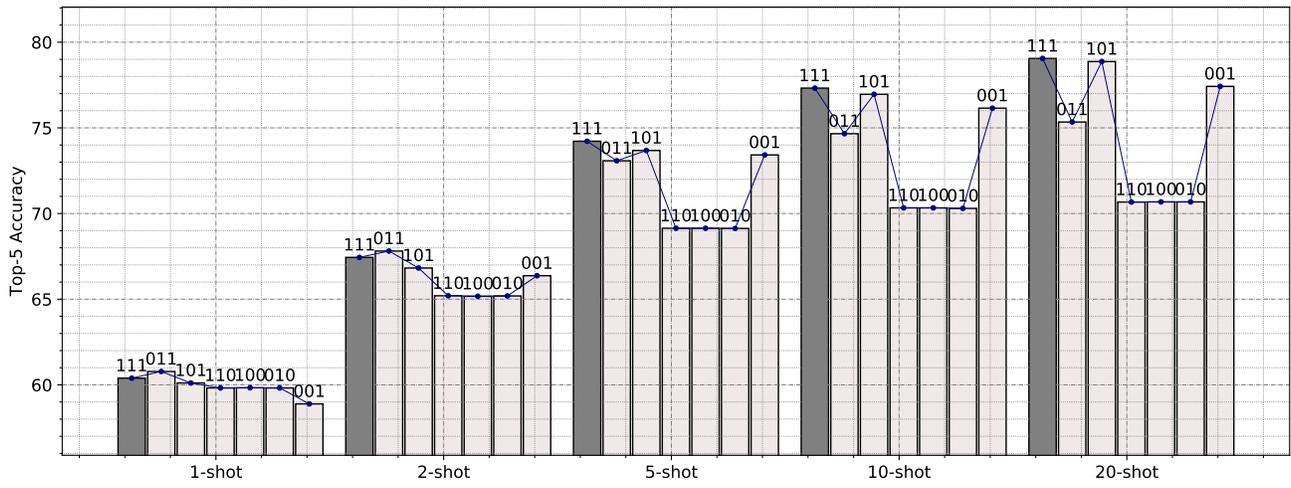
### APPENDIX A

Since the performance of neural networks depends on the choice of optimizers and hyperparameters, it is common to search for proper values from validation examples of a dataset. Thus, in Table VI, we provide the optimizers and the hyperparameters where we used for our experiments. For all experiments, we used the Adam [40] optimizer.
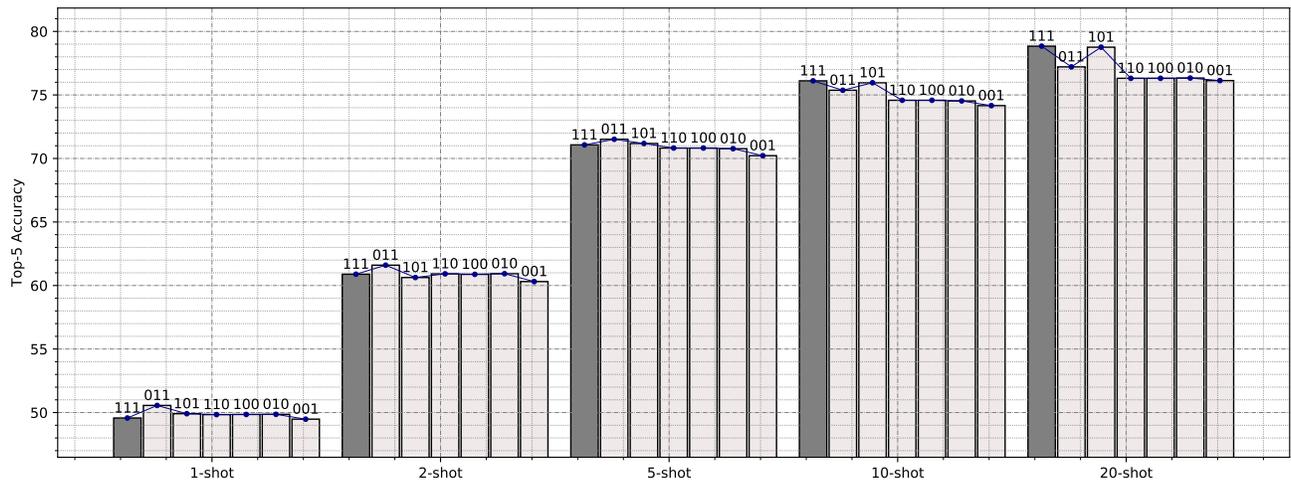
(a) Both categories with ResNetS on *mini*ImageNet

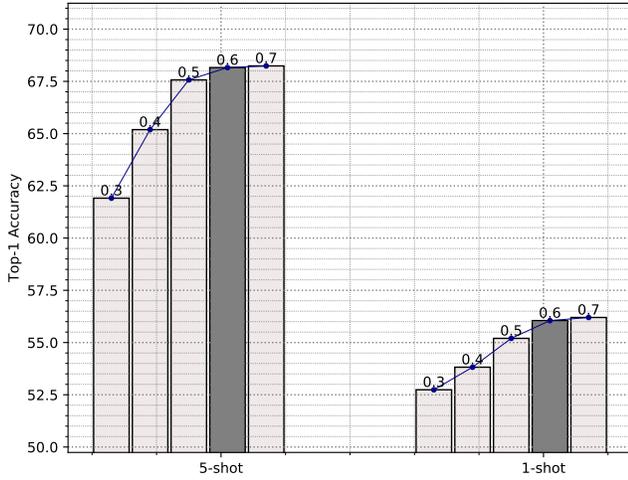(b) Novel categories with ResNetS on *mini*ImageNet

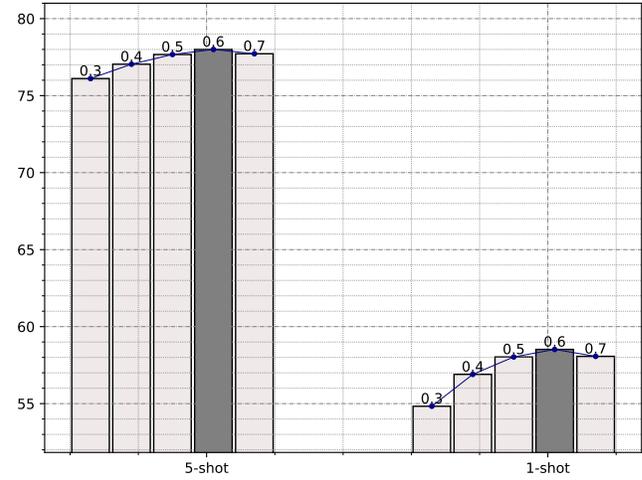(c) Both categories with ResNet10-Dropout on the ImageNet

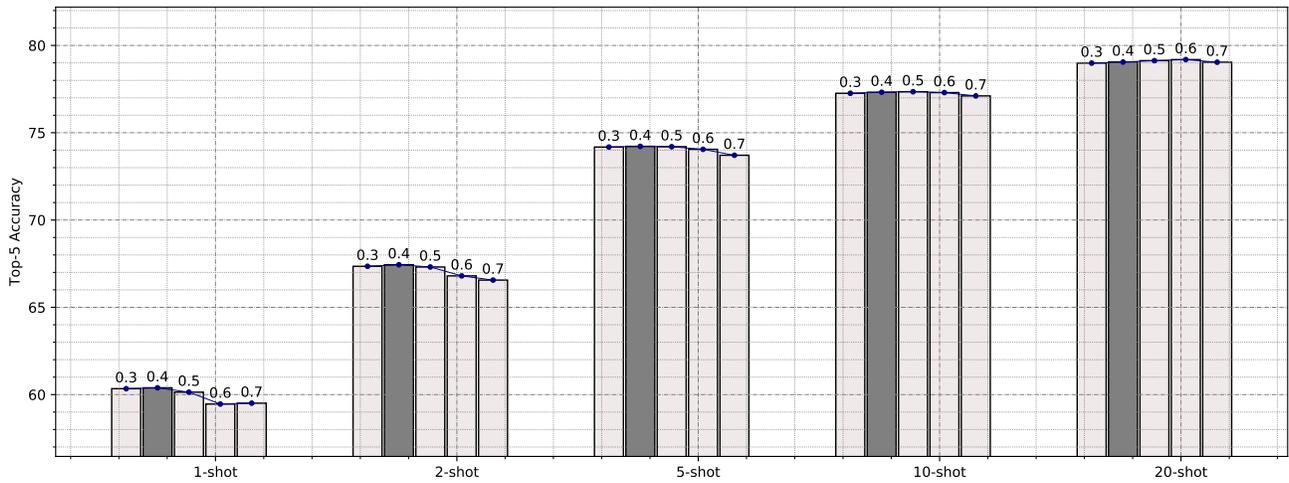(d) Novel categories with ResNet10-Dropout on the ImageNet

Fig. 7. Classification accuracy according to the combinations of the three loss functions. The binary numbers correspond to $\gamma\alpha\beta$ in Eq. 8. The gray bar indicates 111, which means we use cross entropy loss, the weight-centric feature clustering and the angular weight separation.
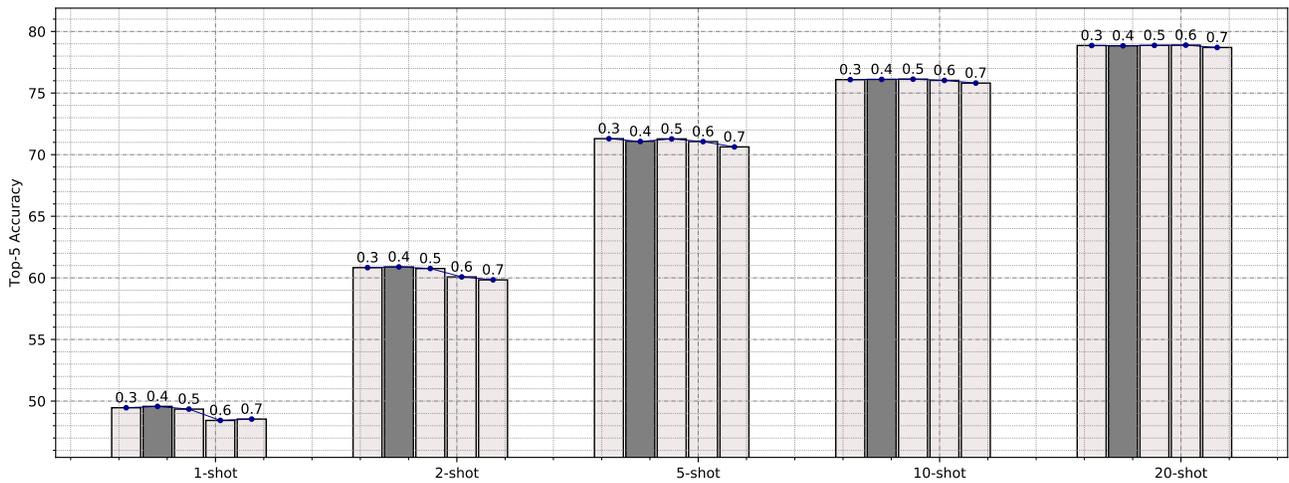
(a) Both categories with ResNetS on *mini*ImageNet

(b) Novel categories with ResNetS on *mini*ImageNet

(c) Both categories with ResNet10-Dropout on the ImageNet

(d) Novel categories with ResNet10-Dropout on the ImageNet

Fig. 8. Classification accuracy according to various margin values. This margin value is used for the angular weight separation in Eq. 7. The gray bar indicates the accuracy we reported in Tables I and II.

## References

[1] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2016, pp. 770–778.

[2] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2017, pp. 4700–4708.

[3] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *Proc. Int. Conf. Comput. Vis. (ICCV)*, 2015, pp. 1–9.

[4] X. Chen, J. Weng, W. Lu, J. Xu, and J. Weng, "Deep manifold learning combined with convolutional neural networks for action recognition," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 29, no. 9, pp. 3938–3952, 2017.

[5] B. Cao, N. Wang, J. Li, and X. Gao, "Data augmentation-based joint learning for heterogeneous face recognition," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 30, no. 6, pp. 1731–1743, 2018.

[6] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, 2015, pp. 91–99.

[7] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2015, pp. 3431–3440.

[8] P. Arbeláez, B. Hariharan, C. Gu, S. Gupta, L. Bourdev, and J. Malik, "Semantic segmentation using regions and parts," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2012, pp. 3378–3385.

[9] J. Sánchez, F. Perronnin, T. Mensink, and J. Verbeek, "Image classification with the fisher vector: Theory and practice," *Int. J. Comput. Vis.*, vol. 105, no. 3, pp. 222–245, 2013.

[10] U. Park, H.-C. Choi, A. K. Jain, and S.-W. Lee, "Face tracking and recognition at a distance: A coaxial and concentric ptz camera system," *IEEE Trans. Inf. Forensics Secur.*, vol. 8, no. 10, pp. 1665–1677, 2013.

[11] M.-C. Roh, H.-K. Shin, and S.-W. Lee, "View-independent human action recognition with volume motion template on single stereo camera," *Pattern Recognit. Lett.*, vol. 31, no. 7, pp. 639–647, 2010.

[12] B.-W. Hwang, V. Blanz, T. Vetter, and S.-W. Lee, "Face reconstruction from a small number of feature points," in *Proc. 15th Int. Conf. Pattern Recognit. (ICPR)*, 2000, pp. 838–841.

[13] L. A. Schmidt, "Meaning and compositionality as statistical induction of categories and constraints," Ph.D. dissertation, Massachusetts Institute of Technology, 2009.

[14] M. D. Zeiler and R. Fergus, "Visualizing and understanding convolutional networks," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*, 2014, pp. 818–833.

[15] J. Kirkpatrick, R. Pascanu, N. Rabinowitz, J. Veness, G. Desjardins, A. A. Rusu, K. Milan, J. Quan, T. Ramalho, A. Grabska-Barwinska, D. Hassabis, C. Clopath, D. Kumaran, and R. Hadsell, "Overcoming catastrophic forgetting in neural networks," *Proc. Natl. Acad. Sci. U.S.A. (PNAS)*, vol. 114, no. 13, pp. 3521–3526, 2017.

[16] A. Krizhevsky and G. Hinton, "Learning multiple layers of features from tiny images," University of Toronto, Tech. Rep., 2009.

[17] J. Mao, X. Wei, Y. Yang, J. Wang, Z. Huang, and A. L. Yuille, "Learning like a child: Fast novel visual concept learning from sentence descriptions of images," in *Proc. Int. Conf. Comput. Vis. (ICCV)*, 2015, pp. 2533–2541.

[18] S. Qiao, C. Liu, W. Shen, and A. Yuille, "Few-shot image recognition by predicting parameters from activations," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2018, pp. 7229–7238.

[19] S. Gidaris and N. Komodakis, "Dynamic few-shot visual learning without forgetting," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2018, pp. 4367–4375.

[20] T. Munkhdalai and H. Yu, "Meta networks," in *Proc. Int. Conf. Mach. Learn. (ICML)*, 2017, pp. 2554–2563.

[21] H. Qi, M. Brown, and D. G. Lowe, "Low-shot learning with imprinted weights," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2018, pp. 5822–5830.

[22] J. Snell, K. Swersky, and R. Zemel, "Prototypical networks for few-shot learning," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, 2017, pp. 4077–4087.

[23] O. Vinyals, C. Blundell, T. Lillicrap, D. Wierstra *et al.*, "Matching networks for one shot learning," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, 2016, pp. 3630–3638.

[24] F. S. Y. Yang, L. Zhang, T. Xiang, P. H. Torr, and T. M. Hospedales, "Learning to compare: relation network for few-shot learning," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2018, pp. 1199–1208.

[25] Q. Cai, Y. Pan, T. Yao, C. Yan, and T. Mei, "Memory matching networks for one-shot image recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2018, pp. 4080–4088.

[26] Y.-X. Wang, R. Girshick, M. Hebert, and B. Hariharan, "Low-shot learning from imaginary data," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2018, pp. 7278–7286.

[27] B. Hariharan and R. B. Girshick, "Low-shot visual recognition by shrinking and hallucinating features," in *Proc. Int. Conf. Comput. Vis. (ICCV)*, 2017, pp. 3018–3027.

[28] Y. Wang, X. Wu, Q. Li, J. Gu, W. Xiang, L. Zhang, and V. O. K. Li, "Large margin few-shot learning," *arXiv preprint arXiv:1807.02872*, 2018.

[29] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, "Imagenet large scale visual recognition challenge," *Int. J. Comput. Vis.*, vol. 115, no. 3, pp. 211–252, 2015.

[30] F. Wang, X. Xiang, J. Cheng, and A. L. Yuille, "Normface: l2 hypersphere embedding for face verification," in *Proc. ACM Multimedia (ACM MM)*, 2017, pp. 1041–1049.

[31] S. Ravi and H. Larochelle, "Optimization as a model for few-shot learning," in *Proc. Int. Conf. on Learn. Representations (ICLR)*, 2017.

[32] C. Finn, P. Abbeel, and S. Levine, "Model-agnostic meta-learning for fast adaptation of deep networks," in *Proc. Int. Conf. Mach. Learn. (ICML)*, 2017, pp. 1126–1135.

[33] N. Mishra, M. Rohaninejad, X. Chen, and P. Abbeel, "A simple neural attentive meta-learner," in *Proc. Int. Conf. Learn. Representations (ICLR)*, 2018.

[34] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *Proc. Int. Conf. Mach. Learn. (ICML)*, 2015, pp. 448–456.

[35] V. Nair and G. E. Hinton, "Rectified linear units improve restricted boltzmann machines," in *Proc. Int. Conf. Mach. Learn. (ICML)*, 2010, pp. 807–814.

[36] A. Radford, L. Metz, and S. Chintala, "Unsupervised representation learning with deep convolutional generative adversarial networks," *arXiv preprint arXiv:1511.06434*, 2015.

[37] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: A simple way to prevent neural networks from overfitting," *J. Mach. Learn. Res.*, vol. 15, no. 1, pp. 1929–1958, 2014.

[38] L. v. d. Maaten and G. Hinton, "Visualizing data using t-sne," *J. Mach. Learn. Res.*, vol. 9, pp. 2579–2605, 2008.

[39] 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 *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2017, pp. 618–626.

[40] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

**Hong-Gyu Jung** received a B.S. degree in Electronic Engineering and a M.S. degree in Information Telecommunication Engineering from Soongsil University, Seoul, Korea, in 2012 and 2014, respectively. He is currently a Ph.D. student in the Department of Brain and Cognitive Engineering at Korea University. His current research interests include artificial intelligence and pattern recognition.

**Seong-Whan Lee** (S84M89SM96-F10) received a B.S. degree in computer science and statistics from Seoul National University, Korea, in 1984, and M.S. and Ph.D. degrees in computer science from the Korea Advanced Institute of Science and Technology, Korea, in 1986 and 1989, respectively. He is currently the head of the Department of Artificial Intelligence, Korea University. His current research interests include artificial intelligence, pattern recognition, and brain engineering. He is a fellow of the IAPR and the Korea Academy of Science and Technology.