DeepEMD: Differentiable Earth Mover's Distance for Few-Shot Learning

Chi Zhang, Yujun Cai, Guosheng Lin, Chunhua Shen

Abstract—In this work, we develop methods for few-shot image classification from a new perspective of optimal matching between image regions. We employ the Earth Mover's Distance (EMD) as a metric to compute a structural distance between dense image representations to determine image relevance. The EMD generates the optimal matching flows between structural elements that have the minimum matching cost, which is used to calculate the image distance for classification. To generate the important weights of elements in the EMD formulation, we design a cross-reference mechanism, which can effectively alleviate the adverse impact caused by the cluttered background and large intra-class appearance variations. To implement *k*-shot classification, we propose to learn a structured fully connected layer that can directly classify dense image representations with the EMD. Based on the implicit function theorem, the EMD can be inserted as a layer into the network for end-to-end training. Our extensive experiments validate the effectiveness of our algorithm which outperforms state-of-the-art methods by a significant margin on five widely used few-shot classification benchmarks, namely, minilmageNet, tieredImageNet, Fewshot-CIFAR100 (FC100), Caltech-UCSD Birds-200-2011 (CUB), and CIFAR-FewShot (CIFAR-FS). We also demonstrate the effectiveness of our method on the image retrieval task in our experiments.

Index Terms—few-shot classification, meta learning, metric learning.

1 Introduction

DEEP neural networks have achieved great success in many vision tasks, typically requiring a large amount of labeled data. A notorious drawback of deep learning methods is that they suffer from poor sample efficiency. This is in sharp contrast to how we humans learn. In machine learning, few-shot learning is the task that addresses this issue, which is often solved as a special case of the broader meta learning. Meta learning attempts to learn a model that can generalize to new tasks with minimum adaption effort. One of the most well-studied test-beds for meta-learning algorithms is few-shot image classification, which aims to perform classification on new image categories with only a limited amount of labeled training data. This is the focus of our work here.

To tackle this problem, a line of previous work in literature adopts metric-based methods [1], [2], [3], [4], [5], [6] that learn to represent image data in an appropriate feature space and use a distance function to predict image labels. Following the formulation of the standard image classification networks [7], [8], metric-based methods often employ a convolution neural network to learn image feature representations and replace the fully connected layer with a distance function, *e.g.*, the *cosine* distance and Euclidean

- C. Zhang is with School of Computer Science and Engineering, Nanyang Technological University, Singapore, 639798.
 E-mail: chi007@e.ntu.edu.sg
- Y. Cai is with Institute for Media Innovation, Nanyang Technological University, Singapore, 639798.
 E-mail: yujun001@e.ntu.edu.sg
- G. Lin is with School of Computer Science and Engineering, Nanyang Technological University, Singapore, 639798.
 E-mail: gslin@ntu.edu.sg
- C. Shen is with Zhejiang University, China. E-mail: chunhua@me.com

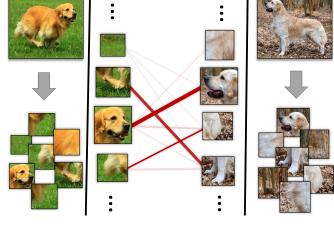


Fig. 1 – Illustration of using the Earth Mover's Distance for oneshot image classification. Our method uses an optimal matching cost between image regions to represent the image distance more faithfully.

distance. Such distance functions directly compute the distances between the embeddings of the test images and training images for classification, which bypasses the difficult optimization problem in learning a classifier in the fewshot setting. The network is usually trained by sampling from a distribution of tasks, in the hope of acquiring a good generalization ability to unseen but similar tasks.

Although these methods have achieved some degree of success, we observe that the cluttered background and large intra-class appearance variations may drive the imagelevel embeddings from the same category far apart in a given metric space. This issue can be largely alleviated by deep neural networks under the setting of fully supervised

learning, thanks to the large capacity of deep models and sufficient training images. However, it is almost inevitably amplified in low-data regimes and thus adversely impacts the image classification accuracy. Moreover, a mixed global representation would struggle to well capture image structures and is likely to lose useful local feature characteristics. Local features can provide discriminative and transferable information across categories, which can be important cues for image classification in the few-shot scenario. Therefore, a desirable metric-based algorithm should have the ability to exploit the local discriminative information and minimize the distraction caused by irrelevant regions.

A natural approach to determine the similarity of two complex structured representations is to compare their building blocks. The difficulty lies in that we do not have their correspondence supervision for training and not all building elements can always find their counterparts in the other structures. To solve the problems above, in this paper, we formalize the few-shot classification as an instance of optimal matching, and we propose to use the optimal matching cost between two structures to represent their dissimilarity. Given the local feature representation sets generated by two images, we use the Earth Mover's Distance (EMD) [9] to compute their structural similarity. The EMD is a metric for computing distance between structured representations, which was originally proposed for image retrieval. Given the distance between all element pairs, the EMD can acquire the optimal matching flows between two structures that have the minimum overall distance. It can also be interpreted as the minimum cost to reconstruct a structured representation against the other one. An illustration of our motivation is shown in Fig. 1. The EMD has the formulation of the transportation problem [10] and the global minimum can be attained by solving a Linear Programming problem. To embed the optimization problem into the model for end-to-end training, we apply the implicit function theorem [11], [12], [13] to form the Jacobian matrix of the optimal optimization variables with respect to the problem parameters [13]. We explore multiple ways to extract local representations from an image, including fully convolutional networks, image grids, and image region sampling. We also investigate pyramid structures at both the feature level and the image level to capture local representations at different scales.

An important problem-specific parameter in the EMD formulation is the weight of each element. Elements with large weights generate more matching flows and thus contribute more to the overall distance. Ideally, the algorithm should accommodate the flexibility to assign less weight to irrelevant regions such that they contribute less to the overall distance no matter which elements they match with. To achieve this goal, we propose a cross-reference mechanism to determine the importance of the elements. In the proposed cross-reference mechanism, we determine the weight of each node by comparing it with the global statistics of the other structure. Intuitively, the image region that shows greater relevance to the other image is more likely to be the object region and should be assigned with a larger weight, while the weights of high-variance background regions and the object parts that are not co-occurrent in two images should be eliminated as much as possible when computing

the matching cost.

In the k-shot setting where multiple support images are presented, we propose to learn a structured fully connected (FC) layer as the classifier for classification to make use of the increasing number of training images. The structured FC layer includes a group of learnable vectors for each class. At inference time, we use the EMD to compute the distance between the image embeddings and the learnable vector set in each class for classification. The structured FC is an extension of the standard fully connected layer in that it replaces dot product operations between vectors with the EMD function between vector sets such that the structured FC layer can directly classify feature maps. The structured FC layer can also be interpreted as learning the prototype embeddings generated by a dummy image for each category such that the test images can be matched with each of them for classification.

To validate our algorithm, we conduct extensive experiments on multiple datasets to demonstrate the effectiveness of our algorithm. Our main contributions are summarized as follows:

- We propose to formalize the few-shot image classification as an optimal matching problem and employ the Earth Mover's Distance as the distance metric between structured representations. The EMD layer can be embedded into the network for end-to-end training.
- We propose a cross-reference mechanism to generate the weights of elements in the EMD formulation, which can effectively reduce the noise introduced by the irrelevant background regions in images.
- We propose to learn a structured fully connected layer in the *k*-shot settings, which is able to directly classify the structured representations of an image using the Earth Mover's Distance.
- Experiments on five popular few-shot classification benchmark datasets—miniImagenet, tieredImagenet, FC100, CUB, and CIFAR-FS show that our algorithm on both 1-shot and 5-shot classification tasks significantly outperforms the baseline methods and achieves new state-of-the-art performance. We also demonstrate that our method can effectively improve many deep metric learning methods on the image retrieval task.

Our preliminary result was published in [14]. To facilitate future research, the source code and trained models are made available at this link¹.

2 RELATED WORK

Few-Shot Learning. The research literature on few-shot learning shows great diversity [15], [16], [17], [18], [19], [20], [21], [22], [23], [24], [25]. There are two main streams in the few-shot classification literature, metric-based approaches and optimization-based approaches. Optimization-based methods, *e.g.*, [26], [27], [28], [29], [30], [31], [32], [33], [34], [35], [36], [37], [38], [39], [40], [41], [42], [43], [44], [45], [46], [47], [48], [49], [50], [51], [52], [53], target at effectively

adapting model parameters to new tasks in the low-shot regime. For example, MAML [54] and many of its variants [55], [56], [57], [58] aim to learn a good model initialization that can rapidly adapt to novel tasks with limited optimization steps. Our design is more related to the metricbased methods [1], [2], [3], [4], [5], [6], [59], [60], [61], [62], [63], [64], which aim to represent samples in an appropriate feature space where data from different categories can be distinguished with distance metrics. To achieve this goal, most previous methods represent the whole image as a data point in the feature space. There are also some works utilizing local features to make predictions. For example, Lifchitz et al. [38] directly make predictions with each local feature and fuse their results. Li et al. [6] adopt k-NN to fuse local distances. Cross Attention Networks [59] use attention mechanisms to highlight the target object regions and generate discriminative features for few-shot classification. CrossTransformers [65] compute the distances between spatially-corresponding features in the query and labeled samples for few-shot classifications. Our solution to the kshot problem also draws connections to optimization-based methods since we learn a classifier that can directly classify structured representations with the Earth Mover's Distance, which can benefit from the increasing number of support samples.

Besides the two popular approaches, many other promising methods have also been proposed to tackle the few-shot classification problem, such as works based on graph theories [66], [67], [68], [69], reinforcement learning [70], differentiable SVM [71], generative models [72], [73], [74], [75], [76], [77], [78], [79], [80], transductive learning [81], [82], [83], [84], [85], recurrent models [86], [87], self-supervised learning [88], [89], the recent capsule network [90], and temporal convolutions [91]. Few-shot learning has also been investigated for other computer vision tasks, such as image segmentation [92], [93], [94] and object detection [95].

Earth Mover's Distance. Earth Mover's Distance (EMD) was originally proposed in [9] as a metric for color and texture based image retrieval. EMD has the formulation of the well studied transportation problem in linear programming, and thus the global optimal matching can be found by solving a linear program. EMD has several desirable properties that make it a popular method to compare structured representations. First, EMD can generate a structural similarity without explicit alignment information. It extends the distance between single elements to the distance between sets or distributions. Second, the number of elements in the sets can vary and EMD allows for partial matching when the total weights of two sets are not equal. EMD has been widely applied to many areas. For example, Kusner et al. [96] use EMD to measure the similarity between two documents, which calculates the minimal cost to transfer the word embeddings in a document to the other for document classification. Wang and Chan [97] propose to represent the hand shapes and textures with superpixels and employ EMD to measure the dissimilarity between the hand gestures for gesture recognition. In [98], Nikolentzos et al. represent graph data as a set of vectors corresponding to the vertices and use EMD to determine the similarity of two graphs for graph comparison. Schulter et al. [99] solve the multiobject tracking problem with a network flow formulation that learns features for network-flow-based data association. Zhao *et al.* [100] propose to use the differential EMD to tackle the visual tracking problem based on the sensitivity analysis of the simplex method. Li [101] uses a tensor-SIFT based EMD to tackle the contour tracking problem.

Parameterized optimization. Parameterized optimization problems have a parameterized object function and constraints that depend on input data. Many previous works have investigated differentiation through the argmin operators. In [102], Gould et al. present the methods for differentiation through optimization problems with only equality constraints. Agrawal et al. [103] propose a method that can compute the gradient of the solution with respect to the coefficients in the convex cone program, which can scale to large problems. Barratt [13] describes the general case of using the implicit function theorem and interior point methods to compute the Jacobian of the solution with respect to the problem parameters. With the same theory, Amos and Kolter [104] design a batched Quadratic Programming solver as a layer that can be integrated into a neural network for end-to-end training. In [105], Agrawal et al. propose a differentiable convex optimization layer that can differentiate through disciplined convex programs and allow users to define problems in a natural syntax without converting problems to canonical forms. Vlastelica et al. [106] introduce the combinatorial building blocks into neural networks and the end-to-end trainable network can generate informative backward gradients through any black-box implementations of combinatorial solvers. Based on such building blocks, Rolínek et al. [107] design an endto-end network that incorporates a combinatorial solver to solve the graph matching problem.

3 PRELIMINARY

Before presenting our algorithm in detail, we first introduce some preliminary concepts in the few-shot classification literature. The general meta-learning algorithm aims to learn transferable knowledge across tasks, where knowledge learned on training tasks can be used to solve novel tasks with only a small amount of training data. In the fewshot classification scenario, a task T_i is to undertake classification over a set of sampled classes, characterized by scarce training images. Specifically, an *N*-way *K*-shot task denotes classification over N classes with K training samples in each class. To acquire generalization ability across tasks, the training and testing of the model are often aligned with the episodic paradigm [1] where batched tasks are sampled for training or evaluation. For each sampled task, the training set $S = \{(x_1^s, y_1^s), ..., (x_{NK}^s, y_{NK}^s)\}$ is called the support set and the testing set $\mathbb{Q}=\{(x_1^q,y_1^q),...,(x_{NK_Q}^q,y_{NK_Q}^q)\}$ is called the query set, where x_i is an image, y_i is its corresponding label, $y_i \in \{1,...,N\}$, and K_Q is the number of testing images per class. At training time, the ground-truth label of the query sets provides learning supervision, and at inference time, we repeatedly sample tasks for evaluation and record their mean accuracy.

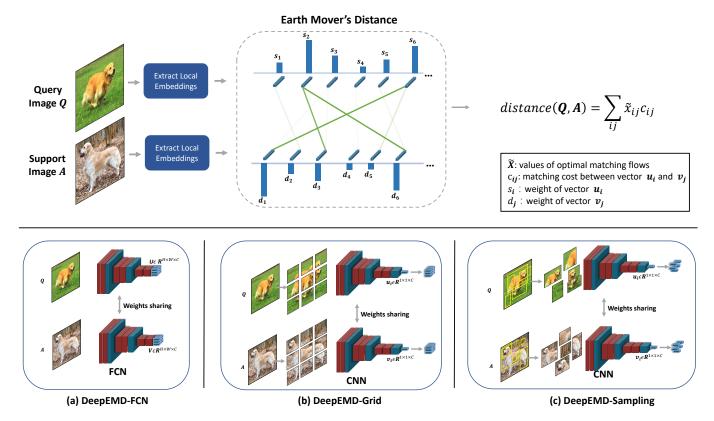


Fig. 2 – Our framework for 1-shot image classification. Given a pair of images, we first extract their local embeddings, which are two sets of feature vectors. Then we use the Earth Mover's Distance to generate the optimal matching flows between two sets, which have the minimum overall matching cost. Finally, based on the optimal matching flows and matching costs, we can compute the distance between two images, which is used for classification. We explore three methods to extract local embeddings: (a) fully convolutional networks, (b) cropping image patches based on grids; and (c) random sampling of image patches. The details of the three methods are provided in Section 4.2.

4 OUR METHOD

In this section, we first present a brief review of the Earth Mover's Distance and describe how we formulate the one-shot classification as an optimal matching problem that can be trained end-to-end. Then, we describe our cross-reference mechanism to generate the weight of each node, which is an important parameter in the EMD formulation. Finally, we demonstrate how to use the EMD to tackle k-shot learning with our proposed structured fully connected layer. The overview of our framework for one-shot classification is shown in Fig. 2.

4.1 Revisiting the Earth Mover's Distance

The Earth Mover's Distance is a distance measure between two sets of weighted objects or distributions, which is built upon the basic distance between individual objects and the weight of each element. It has the form of the well-studied transportation problem from Linear Programming. Specifically, suppose that a set of sources or suppliers $S = \{s_i \mid i=1,2,...m\}$ are required to transport goods to a set of destinations or demanders $\mathcal{D} = \{d_j \mid j=1,2,...k\}$, where s_i denotes the supply units of supplier i and d_j represents the demand of j-th demander. The cost per unit transported from supplier i to demander j is denoted by c_{ij} , and the number of units transported is denoted by x_{ij} . The goal of the transportation problem is then to find a least-

expensive flow of goods $\tilde{\mathcal{X}} = \{\tilde{x}_{ij} \mid i = 1,...m, j = 1,...k\}$ from the suppliers to the demanders:

minimize
$$\sum_{i=1}^{m} \sum_{j=1}^{k} c_{ij} x_{ij}$$
 subject to
$$x_{ij} \geqslant 0, i = 1, ..., m, j = 1, ..., k$$

$$\sum_{j=1}^{k} x_{ij} = s_i, \quad i = 1, ...m$$

$$\sum_{i=1}^{m} x_{ij} = d_j, \quad j = 1, ...k$$
 (1)

The roles of suppliers and demanders can be switched without affecting the total transportation cost. Here s_i and d_j are also called the weights of the nodes, which controls the total matching flows generated by each node. EMD seeks an optimal matching $\tilde{\mathfrak{X}}$ between suppliers and demanders such that the overall matching cost can be minimized. The global optimal matching flows $\tilde{\mathfrak{X}}$ can be achieved by solving a Linear Programming problem.

4.2 EMD for Few-Shot Classification

In the few-shot classification task, metric-based methods aim to find a good distance metric and data representations to compute the distance between images, which are used to compare images for classification. Different from the previous methods [1], [2] that perform distance computation between the image-level embeddings, our approach advocates the use of discriminative local information. The intuition is that as the goal of few-shot learning is to undertake the

classification task on novel categories, directly generating a category-level embedding that corresponds to a new class is difficult. On the other hand, we can decompose an object into a set of object parts that may have been seen in the training process. For example, wheel can be a shared building element across vehicle categories, and if such representation is learned during training, it can be useful to classify unseen vehicle categories. Therefore, local discriminative representations are likely to provide more transferable information across categories. In our framework, we decompose images into a set of local representations, and by assigning appropriate weights to local embeddings in two images, we can use the optimal matching cost between them to represent their dissimilarity. We explore three strategies to generate local representations from an image, as illustrated in Fig. 2: 1) Fully Convolutional Networks. We can deploy a fully convolutional network (FCN) [108] to generate the dense representation $\mathbf{U} \in \mathbb{R}^{H \times W \times C}$ of an image, where H and W denote the spatial size of the feature map and C is the feature dimension. Each image representation contains a collection of local feature vectors $\{\mathbf{u}_1, \mathbf{u}_2, ... \mathbf{u}_{HW}\}$, and each vector \mathbf{u}_i can be seen as a node in the set. Thus, the dissimilarity of two images can be represented as the optimal matching cost between two sets of vectors.

- 2) Dividing the input image into grids. We crop the image evenly into an $H \times W$ grid before feeding it to the CNN, and each image patch in the grid cell is encoded by the CNN individually and generates a feature vector. The feature vectors generated by all the patches constitute the embedding set of an image.
- 3) Random sampling of image patches. Instead of generating the image patches by grids, here we randomly sample M patches in the images with different sizes and aspect ratios. The randomly sampled patches are then re-scaled to the same input size and are encoded by CNNs. The embeddings of these sampled patches make up the embedding set of an image.

We denote our networks adopting the three strategies above by DeepEMD-FCN, DeepEMD-Grid, and DeepEMD-Sampling, respectively. As the size of the grid in DeepEMD-Grid and the number of patches in DeepEMD-Sampling are hyperparameters in our network, we conduct various experiments to investigate the influence of these parameters in Section 5.3. We also investigate pyramid structures on the image level and the feature level to capture local representations at multiple scales, illustrated in Fig. 3. Specifically, we add a feature pyramid structure to DeepEMD-FCN and an image pyramid structure to DeepEMD-Grid. The feature pyramid applies the RoI pooling to the feature maps generated by the FCN, and the resulting feature vectors together with the raw feature vectors constitute the embedding set. For the image pyramid, we simply crop the patches according to different grid sizes and send all patches to the CNN to generate the embedding set.

After acquiring the embedding sets of two images, we can follow the original EMD formulation in Equation (1) to compute the distance. Concretely, assuming there are $H \times W$ vectors in each set, the cost per unit is obtained by computing the pairwise distance between embedding nodes

 \mathbf{u}_i , \mathbf{v}_j from two image features:

$$c_{ij} = 1 - \frac{\mathbf{u}_i^T \mathbf{v}_j}{\|\mathbf{u}_i\| \|\mathbf{v}_i\|},\tag{2}$$

where nodes with similar representations tend to generate small matching costs between each other. As to the generation of weights s_i and d_j , we leave the detailed elaborations in Section 4.4. Once acquiring the optimal matching flows $\tilde{\mathcal{X}}$, we can compute the similarity score s between image representations with:

$$s(\mathbf{U}, \mathbf{V}) = \sum_{i=1}^{HW} \sum_{j=1}^{HW} (1 - c_{ij}) \tilde{x}_{ij}.$$
 (3)

4.3 End-to-End Training

In order to embed the optimal matching problem into a neural network for end-to-end training, it is important to make the solution of the optimal matching \tilde{X} differentiable with respect to the problem parameter θ . As is indicated by [13], we can apply the implicit function theorem [11], [12], [13] on the optimality (KKT) conditions to obtain the Jacobian. For the sake of completeness, we transform the optimization in Equation (1) to a compact matrix form:

minimize
$$c(\theta)^T x$$

subject to $G(\theta)x \leq h(\theta)$, (4)
 $A(\theta)x = b(\theta)$.

Here $x \in \mathbb{R}^n$ is our optimization variable, with $n=m \times k$ representing the total number of matching flows in \mathcal{X} . θ is the problem parameter that relates to the earlier layers in a differentiable way. Ax=b represents the equality constraints and $Gx \leqslant h$ denotes the inequality constraint in Equation (1). Specifically, to construct the compact matrix form of the original optimization, we can build up the sparse matrix below for the equality constraints:

and the inequality constraint can be written as:

Accordingly, the Lagrangian of the LP problem in Equation (4) is given by:

$$L(\theta, x, \nu, \lambda) = c^T x + \lambda^T (Gx - h) + \nu^T (Ax - b), \quad (7)$$

where ν denotes the dual variables on the equality constraints and $\lambda\geqslant 0$ denotes the dual variables on the inequality constraints.

Following the KKT conditions with notational convenience, we can obtain the optimum $(\tilde{x}, \tilde{\nu}, \tilde{\lambda})$ of the objective function by solving $g(\theta, \tilde{x}, \tilde{\nu}, \tilde{\lambda}) = 0$ with primal-dual interior point methods, where

$$g(\theta, x, \nu, \lambda) = \begin{bmatrix} \nabla_{\theta} L(\theta, x, \nu, \lambda) \\ \mathbf{diag}(\lambda) (G(\theta) x - h(\theta)) \\ A(\theta) x - b(\theta) \end{bmatrix}. \tag{8}$$

Then, the following theorem holds to help us derive the gradients of the LP parameters.

Theorem 1 (From Barratt [13]) Suppose $g(\theta, \tilde{\lambda}, \tilde{\nu}, \tilde{x}) = 0$. Then, when all derivatives exist, the partial Jacobian of \tilde{x} with respect to θ at the optimal solution $(\tilde{\lambda}, \tilde{\nu}, \tilde{x})$, namely $J_{\theta}\tilde{x}$, can be obtained by satisfying:

$$J_{\theta}\tilde{x} = -J_{x}g(\theta, \tilde{\lambda}, \tilde{\nu}, \tilde{x})^{-1}J_{\theta}g(\theta, \tilde{x}, \tilde{\nu}, \tilde{\lambda}). \tag{9}$$

Here the formula for the Jacobian of the solution mapping is obtained by applying the implicit function theorem to the KKT conditions. For instance, the (partial) Jacobian with respect to θ can be defined as

$$J_{\theta}g(\theta, \tilde{\lambda}, \tilde{\nu}, \tilde{x}) = \begin{bmatrix} J_{\theta} \nabla_{x} L(\theta, \tilde{x}, \tilde{\nu}, \tilde{\lambda}) \\ \operatorname{diag}(\tilde{\lambda}) J_{\theta}(G(\theta) x - h(\theta)) \\ J_{\theta}(A(\theta) \tilde{x} - b(\theta)) \end{bmatrix}. \tag{10}$$

Therefore, once getting the optimal solution \tilde{x} for the LP problem, we can obtain a closed-form expression for the gradient of \tilde{x} with respect to the input LP parameters θ . This helps us achieve an efficient backpropagation through the entire optimization process without perturbation of the initialization and optimization trajectory.

4.4 Weight Generation

As can be observed in the EMD formulation, an important problem parameter is the weight of each node, *e.g.*, s_i , which controls the total matching flows $\sum_{j=1}^k x_{ij}$ from it. Intuitively, the node with a larger weight plays a more important role in the comparison of two sets, while a node with a very

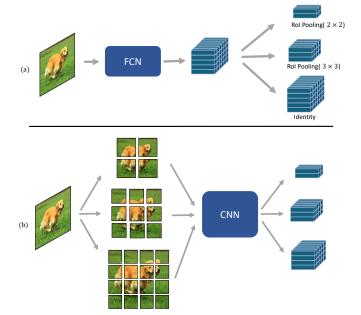


Fig. 3 – Pyramid structures applied on DeepEMD-FCN and DeepEMD-Grid to extract local embeddings. The feature pyramid structure (a) adopts RoI poolings with different output sizes on the feature maps to generate local embeddings at multiple scales while the image pyramid structure (b) crops the input image into patches according to different grid sizes, and all patches are sent to the CNN to generate local embeddings.

small weight can hardly influence the overall distance no matter which nodes it connects with. In the pioneering work that adopts EMD for color-based image retrieval [9], they use the histogram as the elementary feature and perform feature clustering over all pixels to generate the nodes. The weight of each node is set as the size of the corresponding cluster. It makes sense because, for color-based image retrieval, large weights should be assigned to the dominant colors with more pixels, such that the retrieved images can be visually close to the query images. However, for few-shot image classification tasks where features for classification often contain high-level semantic information, the number of pixels does not necessarily reflect the importance. It is common to find image data with greater background regions than the target objects in classification datasets, e.g., ImageNet. Therefore, large weights should be given to the foreground object region in a matching algorithm. However, it may be difficult to define what is the foreground region, particularly when there exist multiple object categories in a single image. An object can be both foreground and background in different cases. Instead of determining the weights by inspecting individual images alone, we argue that for the few-shot classification task, the co-occurrent regions in two images are more likely to be the foreground and the weights of node features should be generated by comparing the nodes on both sides. To achieve this goal, we propose a cross-reference mechanism that uses dot product between a node feature and the average node feature in the other structure to generate a relevance score as the weight value:

$$s_i = \max \left\{ \mathbf{u}_i^T \cdot \frac{\sum_{j=1}^{HW} \mathbf{v}_j}{HW}, 0 \right\}, \tag{11}$$

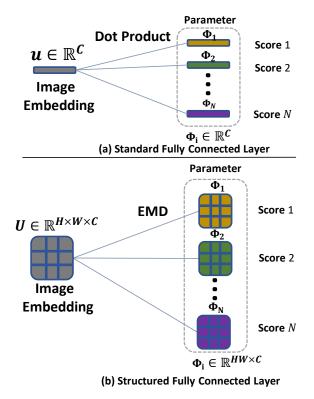


Fig. 4 – Comparison of the standard fully connected layer (a) and our proposed structured fully connected layer (SFC) (b). The SFC learns a group of vectors as the prototype for each class such that we can use the EMD to generate category scores.

where \mathbf{u}_i and \mathbf{v}_j denote the vectors from two feature maps, and function $max(\cdot)$ ensures the weights are always nonnegative. Finally, we normalize all the weights in the structure to make both sides have the same total weights for matching:

$$\hat{s}_i = s_i \frac{HW}{\sum_{j=1}^{HW} s_j}.$$
 (12)

For simplicity, here we take s_i as an example and d_i can be obtained in the same manner. The cross-reference mechanism aims to give less weight to the high-variance background regions and more to the co-occurrent object regions in two images. This can also put less weight on the object parts that do not co-occur in two images and thus allows partial matching to some extent. As a result, the proposed distance metric is based only on confident regions and confident features that have high responses.

4.5 Structured Fully Connected Layer

Thus far we have discussed using the Earth Mover's Distance as the metric to generate the distance value between paired images, *i.e.*, the one-shot case. A question is then raised—how do we tackle the k-shot setting where multiple support images are available? Before presenting our design in detail, let us revisit how the standard fully connected layer classifies an image embedding extracted by CNNs. An FC layer, parameterized by $[\Phi_1,...,\Phi_N] \in \mathbb{R}^{C \times N}$ contains a set of learnable vectors $\Phi_i \in \mathbb{R}^C$ corresponding to each category. Given an image embedding $\mathbf{u} \in \mathbb{R}^C$ generated by the convolutional layer, the FC layer generates the score of class i by computing the dot product between the image

Algorithm 1: A testing episode for an N-way K-shot task. N_{iter} is the number of optimization iterations.

Input: a trained feature extractor Θ , the support set ${\mathbb S}$, and the query set ${\mathbb Q}$

Output: Testing accuracy *Acc* of query set.

- 1 Initialize SFC layer: $\Phi = \Phi'$;
- 2 for i from 1 to N_{iter} do
- 3 | Sample a mini-batch \mathcal{B}_i from S;
- 4 make predictions for \mathcal{B}_i with our model $[\Theta, \Phi]$;
- 5 calculate cross entropy loss $\mathcal{L}_{\mathfrak{I}}$;
- 6 optimize Φ with SGD;
- 7 end
- 8 make predictions for images in the query set Ω with the model $[\Theta, \Phi]$;
- 9 calculate accuracy Acc;
- 10 return Acc.

vector \mathbf{u} and the parameter vector $\mathbf{\Phi}_i$, and this process is applied to all the categories in parallel by matrix multiplication. There are also some previous works replacing the dot product operation in the FC layer with the *cosine* function for computing the category scores [28], [109]. The learning of the FC layer can be seen as finding a prototype vector for each class such that we can use distance metrics to classify an image. An illustration of the standard FC layer is shown in Fig. 4 (a).

With the same formulation, we can learn a structured fully connected layer that adopts EMD as the distance function to directly classify a structured feature representation. The learnable parameter for each class becomes a group of vectors, rather than one vector, such that we can use the structural distance function EMD to undertake image classification. This can also be interpreted as learning a prototype feature map generated by a dummy image for each class. The comparison of the structured FC and the standard FC can be found in Fig. 4. Algorithm 1 provides the pseudocode of a testing episode at inference time. We fix the network backbone and use SGD to learn the parameters in the structured fully connected layer by sampling data from the support set. After several iterations of optimization, we can generate the category scores by computing the EMD between the query images and each of the prototypes in the SFC.

5 EXPERIMENTS

To evaluate the performance of our proposed algorithm for few-shot classification, we conduct extensive experiments on multiple datasets. In this section, we first present dataset information and some important implementation details in our network design. Then we conduct various ablative experiments to validate the effectiveness of each component in our network and compare our model with the state-of-the-art methods on popular benchmark datasets. Finally, we validate the effectiveness of our model on the image retrieval task.

TABLE 1 – Comparison with different baseline methods for 1-shot classification. Our model with EMD as the distance metric significantly outperforms baseline models based on image-level representations and local representations.

Model	Embedding	Metric	5-way	10-way
ProtoNet [2]	global	Euclidean	60.37	44.34
MatchingNet [1]	global	cosine	63.08	47.09
FC [28]	global	dot	59.41	44.08
FC [28]	global	cosine	55.43	40.42
KNN [6]	local	cosine	62.52	47.08
Prediction Fusion [38]	local	cosine	62.38	47.04
DeepEMD-FCN (our)	local	EMD	65.91	49.66

5.1 Implementation Details

Network. For a fair comparison with previous works, we employ a 12-layer ResNet (ResNet12) as our model backbone, which is widely used in the few-shot classification literature. For DeepEMD-Grid and DeepEMD-Sampling, we remove the fully connected layer in ResNet, such that the network generates a vector for each input image patch. For DeepEMD-FCN, we further remove the global average pooling layer, such that the network is transformed into a fully convolutional network. Specifically, given an image of size 84×84 , the model generates a feature map of size $5 \times 5 \times 512$, *i.e.*, 25 feature vectors. For DeepEMD-Grid, we slightly enlarge the region of the local patches in the grid by a factor of 2 to incorporate context information, which is found helpful to generate local representations.

Training. At training time, we use the GPU accelerated convex optimization solver QPTH [104] to solve the Linear Programming problem in our network and compute gradients for back-propagation. As is commonly implemented in the state-of-the-art literature [3], we adopt a feature pretraining step followed by episodic meta-training [1] to learn our network. At the pre-training stage, we train a standard classification model with all training classes and we use DeepEMD-FCN for validation with the validation set in the training process. After pre-training, the model with the highest validation accuracy is further optimized by episodic training for 5,000 episodes. In each training episode, we randomly sample a 5-way 1-shot task with 16 query images, which is aligned with the testing episodes. For the *k*-shot classification task, we re-use the trained 1-shot model as the network backbone to extract features and fix it during training and testing. We initialize the parameters in the structured FC layer with the average local representations of all support data in each class and sample a mini-batch of 5 images from the support set to finetune the structured FC layer for 100 iterations.

5.2 Dataset Description

We conduct few-shot classification experiments on five popular benchmark datasets, namely, *mini*ImageNet [1], *tiered*ImageNet [61], Fewshot-CIFAR100 (FC100) [4], Caltech-UCSD Birds-200-2011 (CUB) [110], and CIFAR-FewShot (CIFAR-FS) [111].

*mini*ImageNet. *mini*ImageNet was first proposed in [1] and becomes the most popular benchmark in the few-shot classification literature. It contains 100 classes with 600 images in each class, which are built upon the ImageNet dataset [112]. The 100 classes are divided into 64, 16, and

TABLE 2 – Different methods for setting the weights in the EMD. We report the 1-shot performance with only the feature pre-training step. DC denotes dense connections. K denotes the number of clusters in K-means; BD denotes the bandwidth in Mean-shift; mK denotes the average number of clusters generated by Mean-shift. EMD with our cross-reference (**CR**) mechanism yields the best result. The model variant that is solely based on the cross-reference mechanism as attention without EMD causes a significant performance drop.

Method	Operation	5-way	10-way
DC.	Average	55.16	40.88
DC	CR	55.41	41.60
EMD	Equal	56.95	42.89
EMD	K-means ($K = 25$)	56.95	42.89
EMD	K-means $(K=10)$	56.25	41.85
EMD	K-means $(K=5)$	55.92	41.57
EMD	K-means $(K=2)$	56.02	41.75
EMD	K-means $(K=1)$	58.65	44.13
EMD	Mean-shift (BD $\approx 0, mK = 25$)	56.95	42.89
EMD	Mean-shift (BD = $2.5, mK = 24.7$)	56.93	42.85
EMD	Mean-shift (BD = $5, mK = 22.8$)	56.66	42.63
EMD	Mean-shift (BD = $7.5, mK = 18.1$)	56.28	42.13
EMD	Mean-shift (BD = $10, mK = 10.6$)	54.99	40.55
EMD	Mean-shift (BD = $12.5, mK = 4.3$)	54.18	40.21
EMD	Mean-shift (BD = $15, mK = 1.8$)	55.48	41.80
EMD	Mean-shift (BD $\approx \infty, mK = 1$)	58.65	44.13
EMD	Mean-shift (auto BD, $mK = 3.4$)	53.56	39.70
EMD	CR	61.13	46.92

20 for meta-training, meta-validation, and meta-testing, respectively.

tieredImageNet. tieredImageNet is also a subset of ImageNet, which includes 608 classes from 34 superclasses. Compared with miniImageNet, the splits of metatraining(20), meta-validation(6), and meta-testing(8) are set according to the super-classes to enlarge the domain difference between training and testing phases. The dataset also includes more images for training and evaluation (779,165 images in total).

Fewshot-CIFAR100. FC100 is a few-shot classification dataset built on CIFAR100 [113]. We follow the split division proposed in [4], where 36 super-classes were divided into 12 (including 60 classes), 4 (including 20 classes), and 4 (including 20 classes), for meta-training, meta-validation, and meta-testing, respectively, and each class contains 100 images.

CIFAR-FewShot. CIFAR-FS [111] is also a few-shot classification dataset built on CIFAR100 [113]. It contains 64, 15, and 20 classes for training, validation, and testing, respectively.

Caltech-UCSD Birds-200-2011. CUB was originally proposed for fine-grained bird classification, which contains 11,788 images from 200 classes. We follow the splits in [3] where 200 classes are divided into 100, 50, and 50 for metatraining, meta-validation, and meta-testing, respectively.

5.3 Analysis

In this section, we implement various experiments to evaluate the effectiveness of our algorithm. We also explore multiple design variants of our network and compare them with baseline solutions. All the experiments are conducted on the *mini*ImageNet dataset.

Comparison with methods based on image-level representations. In the beginning, we first compare our method with a set of methods that utilize image-level vector representations on the 1-shot task. These methods maintain

the global average pooling operation in ResNet to generate vector representations for images and use various distance metrics for classification. We select the representative methods in the literature for comparison: 1) Prototypical Network [2] with Euclidean distance. 2) Matching Network [1] with cosine distance. 3) Finetuning a FC classifier. In [28], Chen et al. propose to fix the pretrained feature extractor and finetune the FC layer with the support images. For fair comparisons, we adopt the same backbones and training schemes for all these baseline methods. We use FCN to extract local features in this experiment such that the only difference in the backbone is the lack of global average pooling operation in our model. The experiment result is shown in Table 1. As we can see, our algorithm significantly outperforms baseline methods that rely on imagelevel vector representations under both 1-shot 5-way and 1shot 10-way settings, which validates the effectiveness of the optimal matching based method that relies on local features.

Comparison with methods based on local representations. There are also a few methods in the literature focusing on local representations to solve few-shot classification. They all remove the global average pooling in the CNN to obtain dense representations of images. In [6], Li et al. use the top k nearest vectors (KNN) between two feature maps to represent the image-level distance. Lifchitz et al. [38] propose to make predictions with each local representation and average their output probabilities. We replace our EMD head with their methods for comparison. The result is shown in Table 1. Our optimal matching based algorithm outperforms all other model variants. Compared with other methods based on local features, the advantage of our method is that, although the basic ground distance in the EMD is based on local features, our algorithm compares the two structures in a global way. Predictions solely based on nearest local features in two images may not extract sufficient information to differentiate between images. For example, eyes can be the nearest feature between animal images, but such a feature can hardly be used to differentiate between animal species.

Weights in the EMD. We next investigate the influence of weights in the EMD. We use the FCN to extract local embeddings for all the methods in this experiment. The first baseline is to set equal weight values for all local embeddings, which is denoted by Equal. The vanilla EMD [9] for image retrieval uses the pixel color as the feature and clusters pixels to generate nodes. The weight of the node is set with the portion of pixels in this cluster. We examine two clustering algorithms as baselines to generate weights: Kmeans [115] and Mean-shift [116] which are implemented by the Scikit-learn [117] library. For both clustering algorithms, we use the cluster mean as the local embedding and the cluster size, *i.e.*, the number of feature vectors in each cluster, as the weight, for computing EMD. It is important to note that there are two special cases that deviate from the goal of using clustering algorithms to set the weights. The first case is that each individual local embedding is a cluster, which means no clustering is applied. This corresponds to K-means with K = 25, and Mean-shift with bandwidth (**BD**) ≈ 0 . As a result, there are totally 25 clusters of an image with equal weights assigned to each cluster. This amounts to our first baseline that sets equal weights in EMD. The second

TABLE 3 – Comparison of different local embedding extractors described in Section 4.2 on 1-shot tasks. $P_{\mathbf{feat}}$ denotes the feature pyramid structure applied on DeepEMD-FCN and $P_{\mathbf{grid}}$ denotes the image pyramid structure applied on DeepEMD-Grid. The parameters in $P_{\mathbf{feat}}$ and $P_{\mathbf{grid}}$ are the RoI pooling size and the grid size, respectively. DeepEMD-Sampling outperforms the plain version of other two methods and both pyramid structures in feature level and image level can effectively boost the performance.

Method	Embedding	5-way	10-way
	5×5	65.91	49.66
DaggEMD ECN	$P_{feat}(5,3)$	66.27	49.85
DeepEMD-FCN	$P_{feat}(5,2)$	66.42	50.02
	$\mathbf{P_{feat}}(5,2,1)$	66.50	50.09
	6×6	63.61	48.72
	5×5	65.38	50.08
	4×4	66.57	51.42
	3×3	67.31	52.26
DeepEMD-Grid	2×2	66.09	50.94
	$P_{grid}(3,2)$	67.74	52.76
	$\mathbf{P_{grid}}(5,3)$	67.01	51.89
	$\mathbf{P_{grid}}(5,2)$	67.39	52.20
	$\mathbf{P_{grid}}(5,3,2)$	67.83	52.85
	9 patches	68.09	53.04
DeepEMD-Sampling	16 patches	68.54	53.61
	25 patches	68.77	53.83

TABLE 4 – Model pre-training with self-supervised auxiliary tasks. The self-supervised auxiliary task at the pre-training stage can effectively boost the 1-shot performance of our models, while applying the random rotation as data augmentation alone degrades the performance. Please refer to Section 5.3 for analysis.

Model	Rotation	Self-Supervision	5-way	10-way
-			65.91	49.66
DoomEMD ECN	\checkmark		65.45	49.70
DeepEMD-FCN		✓	66.64	50.99
			67.31	52.26
DoomEMD Crid	\checkmark		66.96	52.09
DeepEMD-Grid		✓	67.87	53.08
			68.77	53.83
DEMD Clin-	\checkmark		67.98	53.42
DeepEMD-Sampling		✓	69.17	54.52

case is that all local embeddings of an image are clustered into one cluster. This corresponds to K-means with K = 1, and Mean-shift with BD $\approx \infty$. In this case, an image is represented by a single vector, which is the cluster mean, and the EMD is no longer a structured distance function and becomes the cosine distance between two global vector embeddings, which amounts to the baseline in Table 1. Therefore, for K-means, we choose the number of clusters from $\{25, 10, 5, 2, 1\}$. For Mean-shift, we manually select multiple values to set the hyper-parameter, bandwidth, from $\{0, 2.5, 5, 7.5, 10, 12.5, 15, \infty\}$. We also use the automatically estimated bandwidth value for each image provided in the Scikit-learn [117] library. To better observe the influence of the bandwidth values, we additionally list the average number of clusters generated by Mean-shift in each image during testing. As the clustering process of the aforementioned algorithms is non-differentiable, for a fair comparison, we fix the parameters in the backbone after pre-training to evaluate all *methods*. To further test whether our performance advantage is solely brought by the cross-reference mechanism, we also compare our network with a model variant that is solely based on the cross-reference mechanism without EMD. Concretely, we compute the *cosine* distance between all vector pairs, denoted by **Dense Connections**, and compute a weighted sum of these distances with the node weights

TABLE 5 – Cross-domain experiments (*mini*Imagenet \rightarrow CUB). We report the performance with 95% confidence intervals on 1-shot 5-way and 5-shot 5-way tasks. We use FCN to extract local features for KNN and our method. Our proposed algorithm outperforms baseline methods with a large margin.

Model	1-shot	5-shot
ProtoNet [2]	50.01 ± 0.82	72.02 ± 0.67
MatchingNet [1]	51.65 ± 0.84	69.14 ± 0.72
cosine classifier [28]	44.17 ± 0.78	69.01 ± 0.74
linear classifier [28]	50.37 ± 0.79	73.30 ± 0.69
KNN [6]	50.84 ± 0.81	71.25 ± 0.69
DeepEMD-FCN	54.24 ± 0.86	78.86 ± 0.65

TABLE 6 – Computation time of the EMD layer in DeepEMD-FCN. In a 5-way 1-shot task with 10 query images, we vary the spatial size and the feature dimension of two input feature maps and record their computation time.

Solver	Spatial Size	Dimension	Time
	3×3	256	0.294 s
	3×3	512	0.294 s
	3×3	1024	0.295 s
	3×3	2048	0.295 s
	4×4	256	1.489 s
QPTH [104]	4×4	512	1.493 s
	4×4	1024	1.498 s
	4×4	2048	1.499 s
	5×5	256	8.595 s
	5×5	512	$8.638 \mathrm{s}$
	5×5	1024	$8.638 \mathrm{s}$
	5×5	2048	$8.650 \mathrm{\ s}$
	5×5	256	0.018 s
OpenCV [114]	5×5	512	0.019 s
	5×5	1024	0.021 s
	5×5	2048	$0.025 \mathrm{s}$

generated by the cross-reference mechanism. This baseline draws connections with CrossTransformers [65] and Cross Attention Networks [59], where spatial attentions based on local correspondence between samples are used to refine the data embeddings for better classification. Here we use the cross-reference mechanism to weight the distances between local embeddings, which plays a similar role with attentions.

As we can see from the results in Table 2, our crossreference mechanism can bring an improvement of up to 4.2% over the baseline with equal weights, while using the clustering-based methods to set the weights can not improve the performance, which validates our hypothesis that the number of pixels does not necessarily correspond to the importance in few-shot classification. The optimal results of K-means and Mean-shift are both obtained in the second special case, where all local embeddings of an image are clustered into one cluster. This indicates that EMD with weights set by clustering algorithms is inferior to cosine distance with global representations. We find that for clustering-based models, better results are always obtained when the clustering results are close to the two special cases described above. Using the auto-estimated bandwidth for Mean-shift results in the poorest results, although the average number of generated clusters looks normal, which indicates that a good clustering result does not necessarily lead to better classification performance. For the model variant solely based on the cross-reference mechanism as an attention, it can only slightly improve the result of the simple average operation, while a combination of the

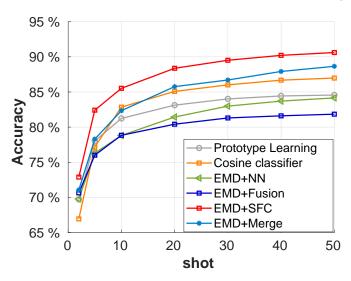


Fig. 5 – Experiment on 5-way k-shot classification. The proposed structured FC layer significantly outperforms baselines and previous k-shot solutions.

cross-reference mechanism and the EMD can yield a significant performance improvement, which again validates the advantages of using the EMD as the metric and the effectiveness of the cross-reference mechanism.

Local embedding extractor. In Table 3, we compare the three methods to extract local embeddings of the input image described in Section 4.2. We also investigate some key parameters in respective methods, *e.g.* the size of grids in DeepEMD-Grid and the number of patches in DeepEMD-Sampling. From the experiment result, we have the following findings:

- 1) DeepEMD-Sampling shows distinct advantages over the plain versions of the other two methods. For DeepEMD-Sampling, increasing the number of sampled patches can consistently boost the performance, even if more low-quality patches, *e.g.* patches from the background region, are likely to be sampled. We observe the weights of these patches and find that they are usually assigned with small weights, thanks to our proposed cross-reference mechanism. Therefore, these noisy patches contribute less to the overall distance, and increasing the number of sampled patches does not negatively influence the prediction.
- 2) Increasing the number of cells in DeepEMD-Grid results in small input image patches and the performance degrades. A possible explanation is that small image patches decompose the object in the image into small pieces, which may lose context information and raises the difficulty in generating high-level representations.
- 3) We speculate that the advantage of DeepEMD-Sampling in the performance over other two methods is gained by multi-scale information implicitly captured by random sampling, as the embedding set contains randomly sampled local descriptors that cover different sizes of local regions. Moreover, as all sampled patches are resized to the input size, an object may appear at an appropriate size in the input image, which can have a stronger feature response in the CNN. To validate our hypothesis, we add pyramid structures to the other two methods to incorporate multi-scale information, as illustrated in Fig. 3. Then the new

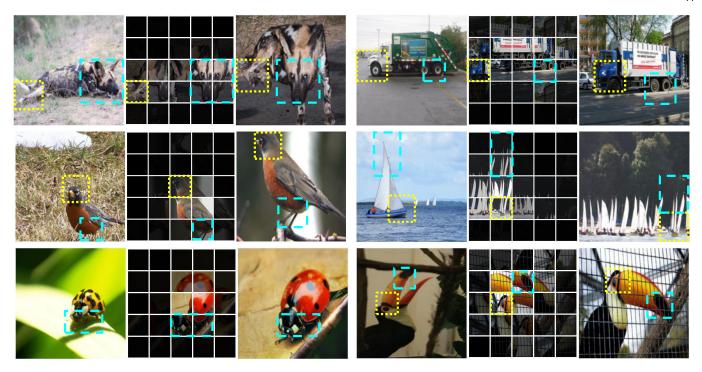


Fig. 6 – Visual reconstruction experiment by DeepEMD-FCN. Given two images (left and right), we plot in the middle the best matched patch (the one with the highest flow value) of each local region in the left image. It can be seen as the reconstruction of the left image using patches from the right one. The weight controls the brightness of the corresponding region. Our algorithm can effectively establish semantic correspondence between local regions and assign small weights to the background regions.

embedding sets generated by the pyramid structures are used for computing the EMD. As is shown in Table 3, the pyramid structures can effectively boost the performance of DeepEMD-FCN and DeepEMD-Grid, which indicates that multi-scale information is useful for reasoning the relations between images in the low-data regime, and pyramid structures can work at both the image level and the feature level.

Comparison with other k**-shot methods.** We next compare our proposed SFC layer (denoted by **EMD** + **SFC**) with some baseline methods and other k-shot solutions in the literature.

- EMD + NN. As the EMD is a paired function for two structures, the first baseline model for *k*-shot experiment is the nearest neighbor (NN) method. We compute the Earth Mover's distance between the query image and all support images and then classify the query image as the category of the nearest support sample.
- EMD + Fusion. Instead of taking the label of the nearest sample, we can fuse the distances of all support images belonging to the same class and classify the query image to the class with the minimum overall distance.
- EMD + Merge. As the number of elements in the compared sets can differ, we can also merge the local embeddings of all support samples in each class into a big embedding set to compute EMD, and there will be $k \times 5 \times 5$ embedding vectors in the set. The problem in this baseline is that as the size of the set grows linearly with respect to the number of shots, the time complexity of solving the transportation problem increases cubically with the QPTH solver,

- which makes the inference of many-shot tasks very slow. Please refer to Section 5.6 for the analysis of time complexity.
- Prototype Learning. In [2], they average the feature embeddings of support images belonging to the same class as the category prototype and apply the nearest neighbor method for classification.
- Finetuning a cosine classifier [28]. Similar to our proposed SFC layer, the network backbone is fixed and the support images are used to learn a fully connected layer as the classifier. The difference is the parameters of the prototypes and the distance metrics for classification, illustrated earlier in Fig. 4.

We compare these models on the k-shot 5-way tasks with multiple k values, and the results are shown in Fig. 5. We can find that the performance of non-optimization based methods often gets saturated quickly and the performance increases slowly with more support samples available. Our structured FC layer can consistently outperform baseline models, and with the number of support images increasing, our network shows even more advantages. The comparison between our method and the finetuned cosine classifier further shows the advantages of structured prototypes and EMD as the metric in the proposed SFC layer.

Pre-training with self-supervised auxiliary tasks. As the weight generation and the computation of optimal matching flows rely on the representations encoded by the backbone, the feature pre-training plays an important role in our framework. We next investigate the use of self-supervision at the pre-training stage. Self-supervised learning [120], [122] has been recently employed as a pre-text task to learn generic representations, which benefits

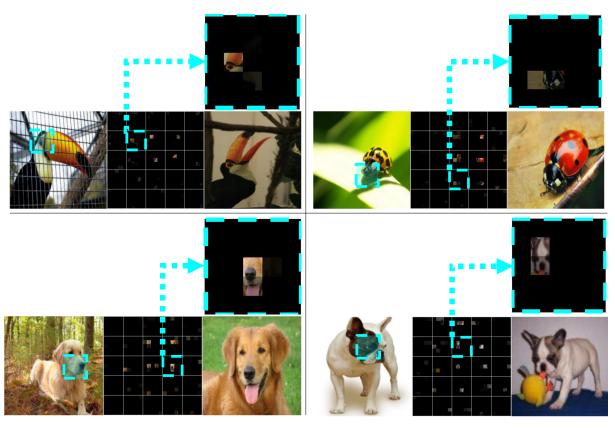


Fig. 7 – Visualization of full matching flows in DeepEMD-FCN. Each grid cell in the middle contains the matched patches (from the right image) of the corresponding region in the left image. The brightness is controlled by the flow values and the weights. *Please zoom in for details*.

many downstream vision tasks, such as object detection and instance segmentation. Following [120], we add a selfsupervised auxiliary learning task during backbone pretraining and observe its influence on our method. Specifically, we randomly rotate the input images by r degrees, and $r \in \{0^{\circ}, 90^{\circ}, 180^{\circ}, 270^{\circ}\}$. The self-supervised learning task is to predict which rotation angle is applied to each image. More concretely, we add a 4-way linear classifier after the network backbone to predict the selections, as is done in [120]. The auxiliary cross-entropy loss is weighted by a factor of $\gamma = 0.5$, which is searched based on the validation accuracy. As can be seen from the result in Table 4, the self-supervised auxiliary task at the feature pre-training stage can effectively boost the performance of our proposed method by up to 1.3%, while applying the random rotation as data augmentation alone degrades the performance in many tasks. This indicates that that classagnostic self-supervised learning task is beneficial to the representation learning in our framework.

5.4 Cross-Domain Experiments

Following the experimental setups in [28], we perform a cross-domain experiment where models are trained on *mini*Imagenet and evaluated on the CUB dataset. Due to the large domain gap, we can better evaluate the models' ability to tackle large domain difference between tasks. We compare our proposed method with baseline models in the Table 5. As we can see, our algorithm *outperforms the baseline models with a large margin*. This shows that local features can provide more transferable and discriminative

information across domains. Moreover, due to our cross-reference mechanism, the optimal matching can be restricted within the co-occurrent object regions that have high feature response, such that the final distance is based on confident regions and representations, and thus has the ability to filter noise when there is a huge domain shift.

5.5 Visualization of Matching Flows and Weights

It is interesting to visualize the optimal matching flows and node weights in the network inference process. We conduct two visualization experiments: First, since we have the correspondence information between the regions in two images, we can reconstruct one image with the local patches from the other image. In Fig. 6, we paste the best-matched patches from the right image to the corresponding position in the left image, and the weights of patches control the brightness of the regions. As we can see, our algorithm can effectively establish semantic correspondence between local regions, and the background regions in the left image are assigned with small weights, thus contributing less to the overall distance. Second, as the correspondence between regions is not strictly 1-to-1, we then plot the full optimal matching flows in Fig. 7. As is shown, one patch can be related to multiple regions in the other image with different weights, which is a useful property when the sizes of the same object are different in two images.

5.6 Time Complexity

Compared with the methods with a closed-form distance metric, the training and inference of DeepEMD come with

TABLE 7 - Comparison with the state-of-the-art 1-shot 5-way and 5-shot 5-way performance (%) with 95% confidence intervals on miniImageNet (a), tieredImageNet (a), CIFAR-FewShot (a) Fewshot-CIFAR100 (b), and Caltech-UCSD Birds-200-2011 (c) datasets. Our model achieves new state-of-the-art performance on all datasets and even outperforms methods with deeper backbones.

Mathad	Backbone	miniIn	nagenet	tieredIı	tieredImagenet		CIFAR-FS	
Method	backbone	1-shot	5-shot	1-shot	5-shot	1-shot	5-shot	
cosine classifier [28]	ResNet12	55.43 ± 0.81	77.18 ± 0.61	61.49 ± 0.91	$82.37. \pm 0.67$	-	-	
TADAM [4]	ResNet12	58.50 ± 0.30	76.70 ± 0.30	-	-	-	-	
ECM [118]	ResNet12	$59.00 \pm -$	$77.46 \pm -$	$63.99 \pm -$	$81.97 \pm -$	$69.15 \pm -$	$84.70 \pm -$	
TPN [119]	ResNet12	$59.46 \pm -$	$75.65 \pm -$	59.91 ± 0.94	73.30 ± 0.75	-	-	
PPA [31]	WRN-28-10 [†]	59.60 ± 0.41	73.74 ± 0.19	65.65 ± 0.92	$83.40. \pm 0.65$	-	-	
Dhillon et al. [82]	WRN-28-10 [†]	57.73 ± 0.62	78.17 ± 0.49	66.58 ± 0.70	85.55 ± 0.48	68.72 ± 0.67	86.11 ± 0.47	
ProtoNet [2]	ResNet12	60.37 ± 0.83	78.02 ± 0.57	65.65 ± 0.92	$83.40. \pm 0.65$	-	-	
wDAE-GNN [67]	$WRN-28-10^{\dagger}$	61.07 ± 0.15	76.75 ± 0.11	68.18 ± 0.16	83.09 ± 0.12	-	-	
MTL [57]	ResNet12	61.20 ± 1.80	75.50 ± 0.80	-	-	-	-	
LEO [56]	WRN-28-10 [†]	61.76 ± 0.08	77.59 ± 0.12	66.33 ± 0.05	81.44 ± 0.09	-	-	
DC [38]	ResNet12	62.53 ± 0.19	79.77 ± 0.19	-	-	-	-	
MetaOptNet [71]	ResNet12	62.64 ± 0.82	78.63 ± 0.46	65.99 ± 0.72	81.56 ± 0.53	72.00 ± 0.70	84.20 ± 0.50	
FEAT [3]	ResNet24 [†]	62.96 ± 0.20	78.49 ± 0.15	-	-	-	-	
MatchNet [1]	ResNet12	63.08 ± 0.80	75.99 ± 0.60	68.50 ± 0.92	80.60 ± 0.71	-	-	
CAN [59]	ResNet12	63.85 ± 0.48	79.44 ± 0.34	69.89 ± 0.51	84.23 ± 0.37	-	-	
DSN [47]	ResNet12	62.64 ± 0.66	78.83 ± 0.45	66.22 ± 0.75	82.79 ± 0.48	72.30 ± 0.80	85.10 ± 0.60	
Tian et al. [89]	ResNet12	64.82 ± 0.60	82.14 ± 0.43	71.52 ± 0.69	86.03 ± 0.49	73.90 ± 0.80	86.90 ± 0.50	
CTM [60]	ResNet18 [†]	64.12 ± 0.82	80.51 ± 0.13	68.41 ± 0.39	84.28 ± 1.73	-	-	
S2M2-R [120]	ResNet34 [†]	63.74 ± 0.18	79.45 ± 0.12	-	-	62.77 ± 0.23	75.75 ± 0.13	
Negative Margin [22]	ResNet12 [†]	63.85 ± 0.81	81.57 ± 0.56	-	-	-	-	
Kim et al. [48]	ResNet12	65.08 ± 0.86	82.70 ± 0.54	-	-	73.51 ± 0.92	85.49 ± 0.68	
Centroid [21]	ResNet18 [†]	59.88 ± 0.67	80.35 ± 0.73	69.29 ± 0.56	85.97 ± 0.49	-	-	
AM3-TADAM [63]	ResNet12	65.30 ± 0.49	78.10 ± 0.36	69.08 ± 0.47	82.58 ± 0.31	-	-	
E ³ BM [121]	ResNet25 [†]	$64.30 \pm$ -	$81.00\pm$ -	$70.00\pm$ -	$85.00\pm$ -	-	-	
DeepEMD-FCN	ResNet12	66.50 ± 0.80	82.41 ± 0.56	72.65 ± 0.31	86.03 ± 0.58	$\textbf{74.58} \pm \textbf{0.29}$	86.92 ± 0.41	
DeepEMD-Grid	ResNet12	67.83 ± 0.29	$\textbf{83.14} \pm \textbf{0.57}$	$\textbf{73.13} \pm \textbf{0.32}$	87.08 ± 0.60	$\textbf{73.31} \pm \textbf{0.29}$	85.43 ± 0.37	
DeepEMD-Sampling		68.77 ± 0.29	84.13 ± 0.53	74.29 ± 0.32	86.98 ± 0.60	$\textbf{74.48} \pm \textbf{0.29}$	86.37 ± 0.36	

(a) Results on miniImageNet, tieredImageNet, and CIFAR-FewShot datasets.

Method	Backbone	1-shot	5-shot
cosine classifier [28]	ResNet12	38.47 ± 0.70	57.67 ± 0.77
TADAM [4]	ResNet12	40.10 ± 0.40	56.10 ± 0.40
MetaOptNet [71]	ResNet12	41.10 ± 0.60	55.5 ± 0.60
ProtoNet [2]	ResNet12	41.54 ± 0.76	57.08 ± 0.76
DC [38]	ResNet12	42.04 ± 0.17	57.05 ± 0.16
MatchNet [1]	ResNet12	43.88 ± 0.75	57.05 ± 0.71
MTL [57]	ResNet12	45.10 ± 1.8	57.6 ± 0.9
Centroid [21]	ResNet18 [†]	45.83 ± 0.48	59.74 ± 0.56
Tian et al. [89]	ResNet12	44.60 ± 0.70	60.90 ± 0.6
E ³ BM [121]	ResNet25 [†]	45.00± -	$60.50 \pm$ -
DeepEMD-FCN	ResNet12	46.60 ± 0.26	63.22 ± 0.71
DeepEMD-Grid	ResNet12	$\textbf{45.23} \pm \textbf{0.26}$	61.39 ± 0.76
DeepEMD-Sampling	ResNet12	$\textbf{45.37} \pm \textbf{0.25}$	$\textbf{61.51} \pm \textbf{0.70}$

(b) I	Results on	Fewshot-CIFAR100	dataset.

Backbone	1-shot	5-shot
ResNet12	66.09 ± 0.92	82.50 ± 0.58
ResNet34 [†]	66.20 ± 0.99	82.30 ± 0.58
ResNet50†	66.95 ± 1.06	77.11 ± 0.78
ResNet34 [†]	67.28 ± 1.08	83.47 ± 0.59
ResNet12	67.30 ± 0.86	84.75 ± 0.60
ResNet12	71.87 ± 0.85	85.08 ± 0.57
ResNet34 [†]	72.92 ± 0.83	86.55 ± 0.51
ResNet18†	74.22 ± 1.09	88.65 ± 0.55
ResNet18 [†]	73.87 ± 0.76	84.95 ± 0.59
ResNet12	77.14 ± 0.29	88.98 ± 0.49
ResNet12	77.64 ± 0.29	89.25 ± 0.53
ResNet12	79.27 \pm 0.29	89.80 ± 0.51
	ResNet12 ResNet34 [†] ResNet50 [†] ResNet34 [†] ResNet12 ResNet12 ResNet18 [†] ResNet18 [†] ResNet18 [†] ResNet12 ResNet12	$\begin{array}{lll} ResNet12 & 66.09 \pm 0.92 \\ ResNet34^{\dagger} & 66.20 \pm 0.99 \\ ResNet50^{\dagger} & 66.95 \pm 1.06 \\ ResNet34^{\dagger} & 67.28 \pm 1.08 \\ ResNet12 & 67.30 \pm 0.86 \\ ResNet12 & 71.87 \pm 0.85 \\ ResNet34^{\dagger} & 72.92 \pm 0.83 \\ ResNet18^{\dagger} & 74.22 \pm 1.09 \\ ResNet18^{\dagger} & 73.87 \pm 0.76 \\ ResNet12 & 77.14 \pm 0.29 \\ ResNet12 & 77.64 \pm 0.29 \\ \end{array}$

(c) Results on the Caltech-UCSD Birds-200-2011 dataset.

more computation cost, as an LP problem must be solved for each forward pass. As is discussed in [104], the main computation lies in the factorization of the KKT matrix as well as back-substitution when using the interior point method to solve the LP problem, which have cubic and quadratic time complexity respectively with respect to the number of optimization variables. The depth of the backbone network has negligible influence on the computation time of the EMD layer, as only the weight of each vector and ground distance matrix are needed for computing EMD, and they can be efficiently computed in parallel in modern tensor processing libraries, e.g. PyTorch, TensorFlow, and NumPy. Therefore, for DeepEMD-FCN, poolings are helpful in reducing computation time when a feature map with a large spatial size is generated by the backbone. In Table 6, we compare the computation time of the EMD layer in a forward pass of a 5-way 1-shot task, given different sizes of feature maps and dimensions. At training time, we use

the QPTH library which adopts the interior point method to solve the LP problem. As we can see, the feature dimension has little influence on the computation time, which indicates that we can replace the backbone with deeper ones, e.g., ResNet-101 (2048 channels), without significantly increasing the inference time. As the interior point method is only necessary for computing gradients at training time, after the model is trained, we can replace the solver with other solvers that can make faster inference, e.g., Sinkhorn [123] and Simplex. At inference time, we deploy the model with the OpenCV [114] library, which adopts a modified Simplex algorithm to solve the LP problem. As is shown in Table 6, the Simplex solver is much faster than interior point method. Therefore, we use QPTH to train the network and use OpenCV to deploy a trained model for inference.

5.7 Comparison with State-of-the-art Methods

Finally, we compare our algorithm with the state-of-theart methods. We report 1-shot 5-way and 5-shot 5-way performance on 5 popular benchmarks: miniImageNet, tieredImageNet, FC100, CUB and CIFAR-FS. For the 1-shot experiment, we repeatedly sample 5,000 testing episodes and record their average accuracy, and for 5-shot experiments, we sample 600 episodes. We reproduce the methods in some earlier works [1], [2], [28] with our network backbone and training strategies, and report the higher performance between our results and their reported ones. The results are shown in Table 7. Our algorithm achieves new stateof-the-art performance on all datasets without seeking any extra data. In particular, our results outperform the state-of-the-art performance by a significant margin on multiple tasks, e.g., 1-shot (3.47%) and 5-shot (1.43%) on the miniImageNet dataset; 1-shot (2.77%) and 5-shot (1.05%) on the tieredImageNet dataset. We observe that on the FC100 and CIFAR-FS datasets, DeepEMD-FCN outperforms DeepEMD-Grid and DeepEMD-Sampling, which is different from the observations on other datasets. The possible reason is that these two datasets are built upon CIFAR100, where the size of images is 32×32 . Further cropping patches from the image results in very small input images, which makes the CNNs difficult to generate useful representations.

5.8 Experiments on Image Retrieval

We next validate the effectiveness of our method on the image retrieval task, which is also evaluated based on the pairwise image similarity. Different from most existing works in the deep metric learning literature that focus on the strategy of optimization and losses, our method provides a generic metric that is complementary to many existing deep metric learning algorithms for images. We observe the gain on the performance when we combine some high-performing baseline methods with our DeepEMD. To do so, we simply replace the globally pooled vector representations with local representations and replace the metric with DeepEMD.

Details. Our experiment mainly follows a recent comprehensive empirical study in [124], where existing deep metric learning algorithms are fairly evaluated by three metrics based on image similarities, including, Recall@1 (P@1), R-Precision (RP), and Mean Average Precision at R (MAP@R). The experiment is conducted on CUB dataset, where the first 100 classes are used for training, and the rest classes are used for evaluation. We use ImageNet pretrained ResNet-50 as the backbone, and a linear layer is added to project the output embedding to 128 channels. The training and evaluation configurations are kept same for the baselines and our method for fair comparisons. For DeepEMD-FCN, we use the feature pyramid of $\{3, 2, 1\}$; for DeepEMD-Grid, we use the image pyramid of {3, 2}; for DeepEMD-Sampling, we randomly sample 25 patches. Please refer to [124] for more details about the evaluation metrics and baselines.

Results. The comparisons between baselines and our methods are shown in Table 8. As we can see, our method can effectively improve the performance of baselines under three evaluation metrics.

TABLE 8 – Image retrieval experiment on Caltech-UCSD Birds-200-2011 dataset. Our method can effectively improve the performance of deep metric learning methods. The relative performance improvements over the baselines are indicated (↑).

Method	Method	P@1 ↑	RP ↑	MAP@R↑
	Raw	61.83	33.18	22.18
N. Coftman, [105]	+DeepEMD-F	63.40 ↑ 1.57	34.41 ↑ 1.23	$23.42 \uparrow 1.24$
N. Softmax [125]	+DeepEMD-G	63.93 ↑ 2.10	34.27 ↑ 1.09	$23.38 \uparrow 1.20$
	+DeepEMD-S	$64.12 \uparrow 2.29$	35.40 ± 2.22	24.34 ↑ 2.16
	Raw	62.51	34.68	23.58
Contrastive [126]	+DeepEMD-F	64.74 + 2.23	35.77 ↑ 1.09	24.78 + 1.20
Contrastive [126]	+DeepEMD-G	66.58 + 4.07	36.85 + 2.17	25.79 + 2.21
	+DeepEMD-S	$67.74 \uparrow 5.23$	37.91 ↑ 3.23	26.78 ↑ 3.20
	Raw	62.05	33.98	23.34
SoftTriple [127]	+DeepEMD-F	63.15 ↑ 1.10	34.82 ± 0.84	24.19 ± 0.85
301111ple [127]	+DeepEMD-G	63.57 ↑ 1.52	35.08 ↑ 1.10	$24.56 \uparrow 1.22$
	+DeepEMD-S	65.19 ↑ 3.14	35.90 ↑ 1.92	25.33 ↑ 1.99
	Raw	63.63	34.48	23.45
SNR [128]	+DeepEMD-F	65.21 ↑ 1.58	35.71 ↑ 1.23	24.80 ↑ 1.35
31VK [120]	+DeepEMD-G	67.74 ↑ 4.11	36.74 + 2.26	25.78 + 2.33
	+DeepEMD-S	68.52 + 4.89	37.90 ↑ 3.42	26.91 ↑ 3.46
	Raw	62.74	34.01	22.85
ArcFace [129]	+DeepEMD-F	63.92 ↑ 1.45	35.21 ↑ 1.20	$24.10 \uparrow 1.25$
Aicrace [129]	+DeepEMD-G	65.29 ↑ 2.55	35.56 ↑ 1.55	24.58 ↑ 1.73
	+DeepEMD-S	66.49 ↑ 3.75	$36.82 \uparrow 2.81$	25.67 ↑ 2.82
	Raw	63.99	33.71	22.73
MS [130]	+DeepEMD-F	66.39 ↑ 2.40	35.82 ↑ 2.11	24.89 ↑ 2.16
1013 [130]	+DeepEMD-G	67.94 ↑ 3.95	36.60 ↑ 2.89	25.65 ↑ 2.92
	+DeepEMD-S	$69.51 \uparrow 5.52$	38.29 + 4.58	27.38 + 4.65

6 CONCLUSION

We have proposed a few-shot classification framework that employs the Earth Mover's Distance as the distance metric. The implicit function theorem allows our network to be end-to-end trainable. Our proposed cross-reference mechanism for setting the weights of nodes turns out crucial in the EMD formulation and can effectively minimize the negative impact caused by irrelevant regions. The learnable structured fully connected layer can directly classify dense representations of images in the k-shot settings. Our algorithm achieves new state-of-the-art performance on multiple datasets.

ACKNOWLEDGMENTS

This research was supported by the National Research Foundation Singapore under its AI Singapore Programme (Award Number: AISG-RP-2018-003) and the MOE Tier-1 research grants: RG28/18 (S), RG22/19 (S) and RG95/20.

REFERENCES

- [1] O. Vinyals, C. Blundell, T. Lillicrap, K. Kavukcuoglu, and D. Wierstra, "Matching networks for one shot learning," in *Proc. Advances in Neural Inf. Process. Syst.*, 2016.
- [2] J. Snell, K. Swersky, and R. S. Zemel, "Prototypical networks for few-shot learning," in *Proc. Advances in Neural Inf. Process. Syst.*, 2017.
- [3] H. Ye, H. Hu, D. Zhan, and F. Sha, "Learning embedding adaptation for few-shot learning," arXiv, vol. 1812.03664, 2018.
- [4] B. N. Oreshkin, P. Rodríguez, and A. Lacoste, "TADAM: task dependent adaptive metric for improved few-shot learning," in Proc. Advances in Neural Inf. Process. Syst., 2018.
- [5] F. Sung, Y. Yang, L. Zhang, T. Xiang, P. H. S. Torr, and T. M. Hospedales, "Learning to compare: Relation network for fewshot learning," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2018.
- [6] W. Li, L. Wang, J. Xu, J. Huo, Y. Gao, and J. Luo, "Revisiting local descriptor based image-to-class measure for few-shot learning," in Proc. IEEE Conf. Comp. Vis. Patt. Recogn., June 2019.
- [7] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Proc. Advances in Neural Inf. Process. Syst.*, 2012, pp. 1097–1105.

- [8] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2016, pp. 770–778.
- [9] Y. Rubner, C. Tomasi, and L. J. Guibas, "The earth mover's distance as a metric for image retrieval," *Int. J. Comput. Vision*, vol. 40, no. 2, pp. 99–121, 2000.
- [10] F. L. Hitchcock, "The distribution of a product from several sources to numerous localities," J. Mathematics and Physics, vol. 20, no. 1-4, pp. 224–230, 1941.
- [11] S. G. Krantz and H. R. Parks, The implicit function theorem: history, theory, and applications. Springer Science & Business Media, 2012.
- [12] A. L. Dontchev and R. T. Rockafellar, "Implicit functions and solution mappings," Springer Monographs in Mathematics. Springer, vol. 208, 2009.
- [13] S. Barratt, "On the differentiability of the solution to convex optimization problems," arXiv preprint arXiv:1804.05098, 2018.
- [14] C. Zhang, Y. Cai, G. Lin, and C. Shen, "DeepEMD: Few-shot image classification with differentiable earth mover's distance and structured classifiers," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2020.
- [15] R. Keshari, M. Vatsa, R. Singh, and A. Noore, "Learning structure and strength of CNN filters for small sample size training," in Proc. IEEE Conf. Comp. Vis. Patt. Recogn., 2018.
- [16] F. Li, R. Fergus, and P. Perona, "One-shot learning of object categories," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 28, no. 4, pp. 594–611, 2006.
- [17] S. Gidaris and N. Komodakis, "Dynamic few-shot visual learning without forgetting," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2018
- [18] S. Yan, S. Zhang, and X. He, "A dual attention network with semantic embedding for few-shot learning," in *Proc. AAAI Conf.* Artificial Intell., 2019.
- [19] D. Wertheimer and B. Hariharan, "Few-shot learning with localization in realistic settings," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, June 2019.
- [20] Y. Wang, C. Xu, C. Liu, L. Zhang, and Y. Fu, "Instance credibility inference for few-shot learning," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2020, pp. 12836–12845.
- [21] A. Afrasiyabi, J.-F. Lalonde, and C. Gagné, "Associative alignment for few-shot image classification," in *Proc. Eur. Conf. Comp. Vis.* Springer, 2020, pp. 18–35.
- [22] B. Liu, Y. Cao, Y. Lin, Q. Li, Z. Zhang, M. Long, and H. Hu, "Negative margin matters: Understanding margin in few-shot classification," in *Proc. Eur. Conf. Comp. Vis.* Springer, 2020, pp. 438–455.
- [23] M. Lichtenstein, P. Sattigeri, R. Feris, R. Giryes, and L. Karlinsky, "Tafssl: Task-adaptive feature sub-space learning for few-shot classification," in *Proc. Eur. Conf. Comp. Vis.* Springer, 2020, pp. 522–539.
- [24] Z. Yu, L. Chen, Z. Cheng, and J. Luo, "Transmatch: A transfer-learning scheme for semi-supervised few-shot learning," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2020, pp. 12856–12864.
- [25] W. Xu, Y. Xu, H. Wang, and Z. Tu, "Attentional constellation nets for few-shot learning," 2021.
- [26] T. Munkhdalai, X. Yuan, S. Mehri, and A. Trischler, "Rapid adaptation with conditionally shifted neurons," in *Proc. Int. Conf. Mach. Learn.*, 2018.
- [27] A. Santoro, S. Bartunov, M. Botvinick, D. Wierstra, and T. P. Lillicrap, "Meta-learning with memory-augmented neural networks," in *Proc. Int. Conf. Mach. Learn.*, 2016.
- [28] W.-Y. Chen, Y.-C. Liu, Z. Kira, Y.-C. Wang, and J.-B. Huang, "A closer look at few-shot classification," in *Proc. Int. Conf. Learn. Representations*, 2019.
- [29] T. Munkhdalai and H. Yu, "Meta networks," in Proc. Int. Conf. Mach. Learn., 2017.
- [30] T. R. Scott, K. Ridgeway, and M. C. Mozer, "Adapted deep embeddings: A synthesis of methods for k-shot inductive transfer learning," in *Proc. Advances in Neural Inf. Process. Syst.*, 2018.
- [31] S. Qiao, C. Liu, W. Shen, and A. L. Yuille, "Few-shot image recognition by predicting parameters from activations," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2018.
- [32] F. Zhou, B. Wu, and Z. Li, "Deep meta-learning: Learning to learn in the concept space," arXiv, vol. 1802.03596, 2018.
- [33] L. Metz, N. Maheswaranathan, B. Cheung, and J. Sohl-Dickstein, "Meta-learning update rules for unsupervised representation learning," in *Proc. Int. Conf. Learn. Representations*, 2019.

- [34] Y. Lee and S. Choi, "Gradient-based meta-learning with learned layerwise metric and subspace," in Proc. Int. Conf. Mach. Learn., 2018.
- [35] Z. Li, F. Zhou, F. Chen, and H. Li, "Meta-sgd: Learning to learn quickly for few shot learning," in Proc. Int. Conf. Mach. Learn., 2018
- [36] J. Luketina, T. Raiko, M. Berglund, and K. Greff, "Scalable gradient-based tuning of continuous regularization hyperparameters," in *Proc. Int. Conf. Mach. Learn.*, 2016.
- [37] D. K. Naik and R. Mammone, "Meta-neural networks that learn by learning," in Proc. Int. Joint Conf. Neural Networks, 1992.
- [38] Y. Lifchitz, Y. Avrithis, S. Picard, and A. Bursuc, "Dense classification and implanting for few-shot learning," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, June 2019.
- [39] E. Schonfeld, S. Ebrahimi, S. Sinha, T. Darrell, and Z. Akata, "Generalized zero- and few-shot learning via aligned variational autoencoders," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, June 2019.
- [40] S. Flennerhag, A. A. Rusu, R. Pascanu, H. Yin, and R. Hadsell, "Meta-learning with warped gradient descent," arXiv preprint arXiv:1909.00025, 2019.
- [41] X. Li, Q. Sun, Y. Liu, Q. Zhou, S. Zheng, T.-S. Chua, and B. Schiele, "Learning to self-train for semi-supervised few-shot classification," in *Proc. Advances in Neural Inf. Process. Syst.*, 2019, pp. 10276–10286.
- [42] E. Park and J. B. Oliva, "Meta-curvature," in Proc. Advances in Neural Inf. Process. Syst., 2019, pp. 3309–3319.
- [43] L. Franceschi, P. Frasconi, S. Salzo, R. Grazzi, and M. Pontil, "Bilevel programming for hyperparameter optimization and meta-learning," arXiv preprint arXiv:1806.04910, 2018.
- [44] A. Rajeswaran, C. Finn, S. M. Kakade, and S. Levine, "Meta-learning with implicit gradients," in *Proc. Advances in Neural Inf. Process. Syst.*, 2019, pp. 113–124.
- [45] L. Franceschi, P. Frasconi, S. Salzo, R. Grazzi, and M. Pontil, "Bilevel programming for hyperparameter optimization and meta-learning," in *Proc. Int. Conf. Mach. Learn.*, 2018.
- [46] H. Qi, M. Brown, and D. G. Lowe, "Low-shot learning with imprinted weights," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2018, pp. 5822–5830.
- [47] C. Simon, P. Koniusz, R. Nock, and M. Harandi, "Adaptive subspaces for few-shot learning," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2020, pp. 4136–4145.
- [48] J. Kim, H. Kim, and G. Kim, "Model-agnostic boundary-adversarial sampling for test-time generalization in few-shot learning," *Proc. Eur. Conf. Comp. Vis.*, pp. 599–617, 2020.
- [49] M. Zhang, J. Zhang, Z. Lu, T. Xiang, M. Ding, and S. Huang, "IEPT: Instance-level and episode-level pre-text tasks for fewshot learning," 2021.
- [50] N. Fei, Z. Lu, T. Xiang, and S. Huang, "Melr: Meta-learning via modeling episode-level relationships for few-shot learning," in Proc. Int. Conf. Learn. Representations, 2021.
- [51] J. Oh, H. Yoo, C. Kim, and S.-Y. Yun, "BOIL: Towards representation change for few-shot learning," in *Proc. Int. Conf. Learn. Representations*, 2021.
- [52] J. Snell and R. Zemel, "Bayesian few-shot classification with onevs-each polya-gamma augmented Gaussian processes," in *Proc. Int. Conf. Learn. Representations*, 2021.
- [53] M. Patacchiola, J. Turner, E. J. Crowley, M. O'Boyle, and A. Storkey, "Bayesian meta-learning for the few-shot setting via deep kernels," 2020.
- [54] C. Finn, P. Abbeel, and S. Levine, "Model-agnostic meta-learning for fast adaptation of deep networks," in *Proc. Int. Conf. Mach. Learn.*, 2017.
- [55] A. Antoniou, H. Edwards, and A. Storkey, "How to train your maml," in Proc. Int. Conf. Learn. Representations, 2019.
- [56] A. A. Rusu, D. Rao, J. Sygnowski, O. Vinyals, R. Pascanu, S. Osindero, and R. Hadsell, "Meta-learning with latent embedding optimization," in *Proc. Int. Conf. Learn. Representations*, 2019.
- [57] Q. Sun, Y. Liu, T.-S. Chua, and B. Schiele, "Meta-transfer learning for few-shot learning," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2019.
- [58] M. A. Jamal and G.-J. Qi, "Task agnostic meta-learning for fewshot learning," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, June 2019.
- [59] R. Hou, H. Chang, B. Ma, S. Shan, and X. Chen, "Cross attention network for few-shot classification," in *Proc. Advances in Neural Inf. Process. Syst.*, 2019.

- [60] H. Li, D. Eigen, S. Dodge, M. Zeiler, and X. Wang, "Finding task-relevant features for few-shot learning by category traversal," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2019.
- [61] M. Ren, E. Triantafillou, S. Ravi, J. Snell, K. Swersky, J. B. Tenenbaum, H. Larochelle, and R. S. Zemel, "Meta-learning for semi-supervised few-shot classification," in *Proc. Int. Conf. Learn. Representations*, 2018.
- [62] J. Liu, L. Song, and Y. Qin, "Prototype rectification for few-shot learning," arXiv preprint arXiv:1911.10713, 2019.
- [63] C. Xing, N. Rostamzadeh, B. N. Oreshkin, and P. O. Pinheiro, "Adaptive cross-modal few-shot learning," arXiv preprint arXiv:1902.07104, 2019.
- [64] C. Zhang, N. Song, G. Lin, Y. Zheng, P. Pan, and Y. Xu, "Few-shot incremental learning with continually evolved classifiers," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, June 2021.
- [65] C. Doersch, A. Gupta, and A. Zisserman, "CrossTransformers: spatially-aware few-shot transfer," in *Proc. Advances in Neural Inf. Process. Syst.*, 2020.
- [66] J. Kim, T. Kim, S. Kim, and C. D. Yoo, "Edge-labeling graph neural network for few-shot learning," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, June 2019.
- [67] S. Gidaris and N. Komodakis, "Generating classification weights with gnn denoising autoencoders for few-shot learning," in Proc. IEEE Conf. Comp. Vis. Patt. Recogn., June 2019.
- [68] V. Garcia and J. Bruna, "Few-shot learning with graph neural networks," arXiv preprint arXiv:1711.04043, 2017.
- [69] L. Yang, L. Li, Z. Zhang, X. Zhou, E. Zhou, and Y. Liu, "Dpgn: Distribution propagation graph network for few-shot learning," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2020, pp. 13390– 13399.
- [70] W.-H. Chu, Y.-J. Li, J.-C. Chang, and Y.-C. F. Wang, "Spot and learn: A maximum-entropy patch sampler for few-shot image classification," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, June 2019
- [71] K. Lee, S. Maji, A. Ravichandran, and S. Soatto, "Meta-learning with differentiable convex optimization," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2019.
- [72] S. Bartunov and D. P. Vetrov, "Few-shot generative modelling with generative matching networks," in *Proc. Int. Conf. Artificial Intell. & Stat.*, 2018.
- [73] Y. Wang, R. B. Girshick, M. Hebert, and B. Hariharan, "Low-shot learning from imaginary data," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2018.
- [74] E. Schwartz, L. Karlinsky, J. Shtok, S. Harary, M. Marder, R. S. Feris, A. Kumar, R. Giryes, and A. M. Bronstein, "Delta-encoder: an effective sample synthesis method for few-shot object recognition," in *Proc. Advances in Neural Inf. Process. Syst.*, 2018.
- [75] A. Mehrotra and A. Dukkipati, "Generative adversarial residual pairwise networks for one shot learning," arXiv, vol. 1703.08033, 2017.
- [76] R. Zhang, T. Che, Z. Grahahramani, Y. Bengio, and Y. Song, "Metagan: An adversarial approach to few-shot learning," in Proc. Advances in Neural Inf. Process. Syst., 2018.
- [77] H. Zhang, J. Zhang, and P. Koniusz, "Few-shot learning via saliency-guided hallucination of samples," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, June 2019.
- [78] W. Shen, Z. Shi, and J. Sun, "Learning from adversarial features for few-shot classification," arXiv preprint arXiv:1903.10225, 2019.
- [79] K. Li, Y. Zhang, K. Li, and Y. Fu, "Adversarial feature hallucination networks for few-shot learning," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2020, pp. 13470–13479.
- [80] S. Yang, L. Liu, and M. Xu, "Free lunch for few-shot learning: Distribution calibration," arXiv: Comp. Res. Repository, 2021.
- [81] P. Rodríguez, I. Laradji, A. Drouin, and A. Lacoste, "Embedding propagation: Smoother manifold for few-shot classification," in Proc. Eur. Conf. Comp. Vis. Springer, 2020, pp. 121–138.
- [82] G. S. Dhillon, P. Chaudhari, A. Ravichandran, and S. Soatto, "A baseline for few-shot image classification," arXiv: Comp. Res. Repository, 2019.
- [83] S. X. Hu, P. G. Moreno, Y. Xiao, X. Shen, G. Obozinski, N. D. Lawrence, and A. Damianou, "Empirical Bayes transductive meta-learning with synthetic gradients," arXiv: Comp. Res. Repository, 2020.
- [84] Z. Yue, H. Zhang, Q. Sun, and X.-S. Hua, "Interventional few-shot learning," arXiv preprint arXiv:2009.13000, 2020.

- [85] M. Boudiaf, Z. I. Masud, J. Rony, J. Dolz, P. Piantanida, and I. B. Ayed, "Transductive information maximization for few-shot learning," arXiv preprint arXiv:2008.11297, 2020.
- [86] P. Shyam, S. Gupta, and A. Dukkipati, "Attentive recurrent comparators," in *Proc. Int. Conf. Mach. Learn*. JMLR. org, 2017, pp. 3173–3181.
- [87] S. Ravi and H. Larochelle, "Optimization as a model for few-shot learning," in *Proc. Int. Conf. Learn. Representations*, 2017.
- [88] J.-C. Su, S. Maji, and B. Hariharan, "When does self-supervision improve few-shot learning?" in *Proc. Eur. Conf. Comp. Vis.* Springer, 2020, pp. 645–666.
- [89] Y. Tian, Y. Wang, D. Krishnan, J. B. Tenenbaum, and P. Isola, "Rethinking few-shot image classification: a good embedding is all you need?" arXiv preprint arXiv:2003.11539, 2020.
- [90] F. Wu, J. S. Smith, W. Lu, C. Pang, and B. Zhang, "Attentive prototype few-shot learning with capsule network-based embedding," in *Proc. Eur. Conf. Comp. Vis.* Springer, 2020, pp. 237–253.
- [91] N. Mishra, M. Rohaninejad, X. Chen, and P. Abbeel, "Snail: A simple neural attentive meta-learner," in *Proc. Int. Conf. Learn.* Representations, 2018.
- [92] C. Zhang, G. Lin, F. Liu, R. Yao, and C. Shen, "Canet: Class-agnostic segmentation networks with iterative refinement and attentive few-shot learning," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2019, pp. 5217–5226.
- [93] C. Zhang, G. Lin, F. Liu, J. Guo, Q. Wu, and R. Yao, "Pyramid graph networks with connection attentions for region-based oneshot semantic segmentation," in *Proc. IEEE Int. Conf. Comp. Vis.*, 2019, pp. 9587–9595.
- [94] W. Liu, C. Zhang, G. Lin, and F. Liu, "CRNet: Cross-reference networks for few-shot segmentation," arXiv: Comp. Res. Repository, 2020.
- [95] Z. Yang, Y. Wang, X. Chen, J. Liu, and Y. Qiao, "Context-transformer: Tackling object confusion for few-shot detection," arXiv: Comp. Res. Repository, 2020.
- [96] M. Kusner, Y. Sun, N. Kolkin, and K. Weinberger, "From word embeddings to document distances," in *Proc. Int. Conf. Mach. Learn.*, 2015, pp. 957–966.
- [97] C. Wang and S. Chan, "A new hand gesture recognition algorithm based on joint color-depth superpixel earth mover's distance," in *Proc. Int. Workshop Cognitive Information Processing*. IEEE, 2014, pp. 1–6.
- [98] G. Nikolentzos, P. Meladianos, and M. Vazirgiannis, "Matching node embeddings for graph similarity," in *Proc. AAAI Conf. Artificial Intell.*, 2017.
- [99] S. Schulter, P. Vernaza, W. Choi, and M. Chandraker, "Deep network flow for multi-object tracking," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2017, pp. 6951–6960.
- [100] Q. Zhao, Z. Yang, and H. Tao, "Differential earth mover's distance with its applications to visual tracking," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 32, no. 2, pp. 274–287, 2008.
- [101] P. Li, "Tensor-sift based earth mover's distance for contour tracking," J. Mathematical Imaging & Vision, vol. 46, no. 1, pp. 44–65, 2013.
- [102] S. Gould, B. Fernando, A. Cherian, P. Anderson, R. S. Cruz, and E. Guo, "On differentiating parameterized argmin and argmax problems with application to bi-level optimization," arXiv preprint arXiv:1607.05447, 2016.
- [103] A. Agrawal, S. Barratt, S. Boyd, E. Busseti, and W. M. Moursi, "Differentiating through a conic program," arXiv preprint arXiv:1904.09043, 2019.
- [104] B. Amos and J. Z. Kolter, "OptNet: Differentiable optimization as a layer in neural networks," in *Proc. Int. Conf. Mach. Learn.*, vol. 70, 2017, pp. 136–145.
- [105] A. Agrawal, B. Amos, S. Barratt, S. Boyd, S. Diamond, and J. Z. Kolter, "Differentiable convex optimization layers," in *Proc. Advances in Neural Inf. Process. Syst.*, 2019, pp. 9558–9570.
- [106] M. Vlastelica, A. Paulus, V. Musil, G. Martius, and M. Rolínek, "Differentiation of blackbox combinatorial solvers," arXiv preprint arXiv:1912.02175, 2019.
- [107] M. Rolínek, P. Swoboda, D. Zietlow, A. Paulus, V. Musil, and G. Martius, "Deep graph matching via blackbox differentiation of combinatorial solvers," arXiv preprint arXiv:2003.11657, 2020.
- [108] E. Shelhamer, J. Long, and T. Darrell, "Fully convolutional networks for semantic segmentation," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 39, no. 4, pp. 640–651, 2017.
- [109] T. Salimans and D. P. Kingma, "Weight normalization: A simple reparameterization to accelerate training of deep neural net-

- works," in Proc. Advances in Neural Inf. Process. Syst., 2016, pp. 901–909.
- [110] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie, "The caltech-ucsd birds-200-2011 dataset," 2011.
- [111] L. Bertinetto, J. F. Henriques, P. H. Torr, and A. Vedaldi, "Metalearning with differentiable closed-form solvers," arXiv preprint arXiv:1805.08136, 2018.
- [112] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. S. Bernstein, A. C. Berg, and F. Li, "Imagenet large scale visual recognition challenge," *Int. J. Comput. Vision*, vol. 115, no. 3, pp. 211–252, 2015.
- [113] A. Krizhevsky, "Learning multiple layers of features from tiny images," *University of Toronto*, 2009.
- [114] G. Bradski and A. Kaehler, Learning OpenCV: Computer vision with the OpenCV library. "O'Reilly Media, Inc.", 2008.
 [115] X. Jin and J. Han, K-Means Clustering. Boston, MA:
- [115] X. Jin and J. Han, K-Means Clustering. Boston, MA: Springer US, 2010, pp. 563–564. [Online]. Available: https://doi.org/10.1007/978-0-387-30164-8_425
- [116] D. Comaniciu and P. Meer, "Mean shift: A robust approach toward feature space analysis," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 24, no. 5, pp. 603–619, May 2002.
- [117] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, "Scikit-learn: Machine learning in Python," J. Machine Learning Research, vol. 12, pp. 2825–2830, 2011.
- [118] A. Ravichandran, R. Bhotika, and S. Soatto, "Few-shot learning with embedded class models and shot-free meta training," arXiv preprint arXiv:1905.04398, 2019.
- [119] Y. Liu, J. Lee, M. Park, S. Kim, E. Yang, S. J. Hwang, and Y. Yang, "Learning to propagate labels: Transductive propagation network for few-shot learning," arXiv preprint arXiv:1805.10002, 2018.
- [120] P. Mangla, N. Kumari, A. Sinha, M. Singh, B. Krishnamurthy, and V. N. Balasubramanian, "Charting the right manifold: Manifold mixup for few-shot learning," in *The IEEE Winter Conference on Applications of Computer Vision*, 2020, pp. 2218–2227.
- [121] Y. Liu, B. Schiele, and Q. Sun, "An ensemble of epoch-wise empirical bayes for few-shot learning," in *Proc. Eur. Conf. Comp. Vis.* Springer, 2020, pp. 404–421.
- [122] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, "Momentum contrast for unsupervised visual representation learning," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2020.
- [123] M. Cuturi, "Sinkhorn distances: Lightspeed computation of optimal transport," Proc. Advances in Neural Inf. Process. Syst., vol. 26, pp. 2292–2300, 2013.
- [124] K. Musgrave, S. Belongie, and S.-N. Lim, "A metric learning reality check," in *Proc. Eur. Conf. Comp. Vis.* Springer, 2020, pp. 681–699
- [125] A. Zhai and H.-Y. Wu, "Classification is a strong baseline for deep metric learning," arXiv preprint arXiv:1811.12649, 2018.
- [126] R. Hadsell, S. Chopra, and Y. LeCun, "Dimensionality reduction by learning an invariant mapping," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2006, pp. 1735–1742.
- [127] Q. Qian, L. Shang, B. Sun, J. Hu, H. Li, and R. Jin, "Softtriple loss: Deep metric learning without triplet sampling," in *Proc. IEEE Int. Conf. Comp. Vis.*, 2019, pp. 6450–6458.
- [128] T. Yuan, W. Deng, J. Tang, Y. Tang, and B. Chen, "Signal-to-noise ratio: A robust distance metric for deep metric learning," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2019, pp. 4815–4824.
- [129] J. Deng, J. Guo, N. Xue, and S. Zafeiriou, "Arcface: Additive angular margin loss for deep face recognition," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2019, pp. 4690–4699.
- [130] X. Wang, X. Han, W. Huang, D. Dong, and M. R. Scott, "Multi-similarity loss with general pair weighting for deep metric learning," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2019, pp. 5022–5030



Chi Zhang is a PhD candidate with the School of Computer Science and Engineering, Nanyang Technological University, Singapore. He received the B.S. degree from China University of Mining and Technology in 2017. His research interests are in computer vision and machine learning.



Yujun Cai received the B.Eng. degree in Information Science and Engineering from Southeast University in 2017. She is now a PhD candidate with the Institute for Media Innovation, Interdisciplinary Graduate School, Nanyang Technological University, Singapore. Her research interests mainly include computer vision, machine learning and human-computer interaction.



Guosheng Lin is an Assistant Professor at School of Computer Science and Engineering, Nanyang Technological University, Singapore. His research interests are in computer vision and machine learning.

Chunhua Shen is a Professor of Computer Science at Zhejiang University, China.