

Li, X., Yang, X. , Ma, Z. and Xue, J.-H. (2023) Deep metric learning for few-shot image classification: a review of recent developments. *Pattern Recognition*, 138, 109381. (doi: [10.1016/j.patcog.2023.109381](https://doi.org/10.1016/j.patcog.2023.109381))

The material cannot be used for any other purpose without further permission of the publisher and is for private use only.

There may be differences between this version and the published version. You are advised to consult the publisher's version if you wish to cite from it.

<https://eprints.gla.ac.uk/291768/>

Deposited on 10 February 2023

Enlighten – Research publications by members of the University of  
Glasgow

<http://eprints.gla.ac.uk>

# Deep metric learning for few-shot image classification: A Review of recent developments

Xiaoxu Li<sup>a,b,1</sup>, Xiaochen Yang<sup>c,1</sup>, Zhanyu Ma<sup>b,\*</sup>, Jing-Hao Xue<sup>d</sup>

*<sup>a</sup>School of Computer and Communication, Lanzhou University of  
Technology, Lanzhou, 730050, China*

*<sup>b</sup>Pattern Recognition and Intelligent System Laboratory, School of Artificial Intelligence,  
Beijing University of Posts and Telecommunications, Beijing, 100876, China*

*<sup>c</sup>School of Mathematics and Statistics, University of Glasgow, Glasgow, G12 8QQ, UK*

*<sup>d</sup>Department of Statistical Science, University College London, London, WC1E 6BT, UK*

---

## Abstract

Few-shot image classification is a challenging problem that aims to achieve the human level of recognition based only on a small number of training images. One main solution to few-shot image classification is deep metric learning. These methods, by classifying unseen samples according to their distances to few seen samples in an embedding space learned by powerful deep neural networks, can avoid overfitting to few training images in few-shot image classification and have achieved the state-of-the-art performance. In this paper, we provide an up-to-date review of deep metric learning methods for few-shot image classification from 2018 to 2022 and categorize them into three groups according to three stages of metric learning, namely learning feature embeddings, learning class representations, and learning distance measures. Under this taxonomy, we identify the trends of transitioning from

---

\*Corresponding author

Email address: [mazhanyu@bupt.edu.cn](mailto:mazhanyu@bupt.edu.cn). (Zhanyu Ma)

<sup>1</sup>X. Li and X. Yang contribute equally.

learning task-agnostic features to task-specific features, from simple computation of prototypes to computing task-dependent prototypes or learning prototypes, from using analytical distance or similarity measures to learning similarities through convolutional or graph neural networks. Finally, we discuss the current challenges and future directions of few-shot deep metric learning from the perspectives of effectiveness, optimization and applicability, and summarize their applications to real-world computer vision tasks.

*Keywords:* Few-shot learning, Metric learning, Image classification, Deep neural networks

---

## 1. Introduction

Image classification is an important task in machine learning and computer vision. With the rapid development of deep learning, recent years have witnessed breakthroughs in this area [1, 2, 3, 4]. Such progress, however, hinges on collecting and labeling a vast amount of data (in the order of millions), which can be difficult and costly. More severely, this learning mechanism is in stark contrast with that of humans, where one or few examples suffice for learning a new concept [5]. Therefore, to reduce the data requirement and imitate human intelligence, many researchers started to focus on few-shot classification [6, 7, 8], i.e., learning a classification rule from few (typically 1-5) labeled examples.

The biggest challenge in few-shot classification is a high risk of model overfitting to the few labeled training samples. To alleviate this problem, researchers have proposed various approaches, such as meta-learning methods, transfer learning methods, and metric learning methods. Meta-learning

16 methods train a meta-learner on many different classification tasks to extract  
 17 generalizable knowledge, which enables rapid learning on a new related task  
 18 with few examples [7, 9]. Transfer learning methods presume shared knowl-  
 19 edge between the source and target domains, and fine-tune the model trained  
 20 on abundant source data to fit few labeled target samples [10, 11]. Metric  
 21 learning methods learn feature embeddings [6] and/or distance measures (or  
 22 inversely, similarity measures) [12] and classify an unseen sample based on  
 23 its distance to labeled samples or class representations; samples of the same  
 24 class are expected to locate close together in the embedding space and sam-  
 25 ples of different classes should be far apart. Note that the above methods  
 26 can be applied simultaneously, for example learning feature embeddings of  
 27 metric learning methods by using a meta-learning strategy [7].

28 In this paper, we present a review of recent deep metric learning methods  
 29 for few-shot image classification. Metric learning methods deserve special at-  
 30 tention as they do not require learning additional parameters for new classes  
 31 once the metric is learned, and thus able to avoid overfitting to the few labeled  
 32 samples of new classes in few-shot learning. They have also demonstrated  
 33 impressive classification performance on benchmark datasets. Moreover, in  
 34 this review we decouple metric learning into three learning stages, namely  
 35 learning feature embeddings, learning class representations, and learning dis-  
 36 tance measures. Such decomposition facilitates exchange of ideas between  
 37 researchers from two underpinning communities: few-shot image classifica-  
 38 tion and deep metric learning. For example, latest developments in learning  
 39 generalizable feature embeddings can be adopted for few-shot image classifi-  
 40 cation, and the idea of learning prototypes, one type of class representations,

41 can be extended for long-tailed visual recognition [13].

42 A number of surveys on few-shot learning (FSL) have been published or  
43 preprinted. [14] is the first survey on small sample learning, summarizing  
44 methods for different small sample learning scenarios, including zero-shot  
45 learning and FSL, and for various tasks, such as image classification, object  
46 detection, visual question answering, and neural machine translation. Since  
47 the survey was conducted early in 2018, it includes relatively limited work  
48 on few-shot classification, particularly metric learning methods. [15] provides  
49 the first comprehensive review on FSL. In addition to defining FSL and dis-  
50 tinguishing it from related machine learning problems, the authors discuss  
51 FSL from the fundamental perspective of error decomposition in supervised  
52 learning and classify all methods in terms of augmenting the training data  
53 for reducing the estimation error, learning models from prior knowledge for  
54 constraining the hypothesis space and reducing the approximation error, and  
55 learning initializations or optimizers which improve the search for the optimal  
56 hypothesis within the hypothesis space. The survey has limited coverage on  
57 metric learning methods and categorize them all under learning embedding  
58 models, which does not fully describe the merits of these methods. [16] is an-  
59 other comprehensive survey, reviewing literature over a long period from the  
60 2000s to 2020 as well as summarizing applications of FSL in various fields. It  
61 includes early, non-deep approaches of metric learning methods and, since the  
62 survey emphasizes on meta-learning methods, categorizes most recent, deep  
63 approaches under meta-learning as learning-to-measure. Compared with [16]  
64 which links different meta-learning metric learning methods to three classi-  
65 cal methods, our review provides a deeper insight into how metric learning

| Conferences   | Journals   |
|---|--|
| AAAI Conference on Artificial Intelligence (AAAI)   | IEEE Trans. on Circuits and Systems for Video Technology (TCSVT) |
| Int. Conference on Artificial Intelligence and Statistics (AISTATS)   | IEEE Trans. on Image Processing (TIP)                            |
| Conference on Computer Vision and Pattern Recognition (CVPR)  | IEEE Trans. on Multimedia (TMM)                                  |
| European Conference on Computer Vision (ECCV)   | IEEE Trans. on Neural Networks and Learning Systems (TNNLS)      |
| Int. Conference on Computer Vision (ICCV)   | IEEE Trans. on Pattern Analysis and Machine Intelligence (TPAMI) |
| Int. Conference on Learning Representations (ICLR)  | Pattern Recognition (PR)   |
| Int. Joint Conference on Artificial Intelligence (IJCAI)  |  |
| Conference on Neural Information Processing Systems (NeurIPS)   |  |
| Keywords: few-shot/one-shot learning, few-shot/one-shot classification, few-shot/one-shot image recognition |  |

Table 1: Selected conferences and journals (listed in alphabetical order of their abbreviations). Papers that include at least one of the keywords were considered for further investigation.

66 methods evolve in order to generalize better and be more applicable in the  
67 settings that mimic the reality more closely. Moreover, the rapid develop-  
68 ment of FSL leads to a considerable amount of methods proposed since the  
69 publications of [15] and [16]. These new approaches have been discussed in  
70 this review. [17] is the latest review on FSL published in 2021, but it is en-  
71 tirely devoted to meta-learning approaches and has very little overlap with  
72 our work. In short, this paper provides an up-to-date review of deep metric  
73 learning methods for few-shot image classification and a careful examination  
74 of different components of these methods to understand their strengths and  
75 limitations. The conferences and journals being surveyed are listed in Ta-  
76 ble 1. Papers that include at least one of the keywords are considered for  
77 further investigation on their relevance and contribution.

78 The rest of this review is organized as follows. Firstly for completeness,  
79 in Section 2 we give the definition of few-shot classification and introduce the  
80 evaluation procedure and commonly used datasets. Secondly, in Section 3

we review classical few-shot metric learning algorithms and recent influential works published from 2018 to 2022. In the light of the procedure of metric learning, these methods are classified into learning feature embeddings, learning class representations, and learning distance or similarity measures. Finally, we discuss some remaining challenges, future directions, and real-world applications in Section 4 and conclude this review in Section 5.

## 2. The framework of few-shot image classification

### 2.1. Notation and definitions

We first establish the notation and give a unified definition of various types of few-shot classification by generalizing the definition of few-shot learning [12].

Few-shot classification involves two datasets, **base dataset** and **novel dataset**. The novel dataset is the dataset on which the classification task is performed. The base dataset is an auxiliary dataset used to facilitate the learning of the classifier by transferring knowledge. We use  $\mathbb{D}_{base} = \{(X_i, Y_i); X_i \in \mathcal{X}_{base}, Y_i \in \mathcal{Y}_{base}\}_{i=1}^{N_{base}}$  to denote the base dataset, where  $Y_i$  is the class label of instance  $X_i$ ; in the case of image classification,  $X_i$  denotes the feature vector of the  $i$ th image. The novel dataset is denoted similarly by  $\mathbb{D}_{novel} = \{(\tilde{X}_j, \tilde{Y}_j); \tilde{X}_j \in \mathcal{X}_{novel}, \tilde{Y}_j \in \mathcal{Y}_{novel}\}_{j=1}^{N_{novel}}$ .  $\mathbb{D}_{base}$  and  $\mathbb{D}_{novel}$  have no overlap in the label space, i.e.,  $\mathcal{Y}_{base} \cap \mathcal{Y}_{novel} = \emptyset$ . To train and test the classifier, we split  $\mathbb{D}_{novel}$  into the support set  $\mathbb{D}_S$  and the query set  $\mathbb{D}_Q$ .

**Definition 1.** Suppose the support set  $\mathbb{D}_S$  is available, and the sample size of each class in  $\mathbb{D}_S$  is very small (e.g., from 1 to 5). The **few-shot classification** task aims to learn from  $\mathbb{D}_S$  a classifier  $f : \mathcal{X}_{novel} \rightarrow \mathcal{Y}_{novel}$  that can

correctly classify instances in the query set  $\mathbb{D}_Q$ . In particular, if  $\mathbb{D}_S$  contains  $C$  classes and  $K$  labeled examples per class, the task is called  **$C$ -way  $K$ -shot classification**; if the sample size of each class in  $\mathbb{D}_S$  is one, then the task is called **one-shot classification**.

Before presenting the next definition, we introduce the concept of domain. A *domain* consists of two components, namely a feature space  $\mathcal{X}$  and a marginal distribution  $P(X)$  over  $\mathcal{X}$  [18].

**Definition 2.** A few-shot classification task is called **cross-domain few-shot classification** if the base dataset and the novel dataset come from two different domains, i.e.,  $\mathcal{X}_{base} \neq \mathcal{X}_{novel}$  or  $P(X) \neq P(\tilde{X})$ , where  $X \in \mathcal{X}_{base}$  and  $\tilde{X} \in \mathcal{X}_{novel}$ .

**Definition 3.** The **generalized few-shot classification** task aims to learn a classifier  $f : \mathcal{X}_{novel} \cup \mathcal{X}_{base} \rightarrow \mathcal{Y}_{novel} \cup \mathcal{Y}_{base}$  that can correctly classify instances in the query set  $\mathbb{D}_Q$ , where  $\mathbb{D}_Q$  includes instance-label pairs from  $\mathbb{D}_{base}$  in addition to existing pairs from  $\mathbb{D}_{novel}$ .

## 2.2. Evaluation procedure of few-shot classification

We provide a general procedure to evaluate the performance of a classifier for  $C$ -way  $K$ -shot classification in Algorithm 1. The evaluation procedure includes many episodes (i.e., tasks). In each episode, we first randomly select  $C$  classes from the novel label set, and then randomly select  $K$  samples from each of the  $C$  classes to form a support set and  $M$  samples from the remaining samples of those  $C$  classes to form a query set. Let  $\mathbb{X}^{(e)}$  and  $\mathbb{Y}^{(e)}$  denote the set of instances and the set of labels in the query set at the  $e$ th episode, respectively. A learning algorithm returns a classifier  $f(\cdot | \mathbb{D}_{base}, \mathbb{D}_S^{(e)})$  upon



---

**Algorithm 1** Evaluation procedure of  $C$ -way  $K$ -shot classification

---

**Input:**  $\mathbb{D}_{base} = \{(X_i, Y_i); X_i \in \mathcal{X}_{base}, Y_i \in \mathcal{Y}_{base}\}_{i=1}^{N_{base}}; \mathbb{D}_{novel} = \{(\tilde{X}_j, \tilde{Y}_j); \tilde{X}_j \in \mathcal{X}_{novel}, \tilde{Y}_j \in \mathcal{Y}_{novel}\}_{j=1}^{N_{novel}}; \text{ number of episodes } E.$

1: **for**  $e = 1, \dots, E$  **do**

2:   Randomly select  $C$  classes from  $\mathcal{Y}_{novel}$ .

3:   Randomly select  $K$  samples from each class as the support set  $\mathbb{D}_S^{(e)}$ .

129 4:   Randomly select  $M$  samples from the remaining samples of  $C$  classes as the query set  $\{(\mathbb{X}^{(e)}, \mathbb{Y}^{(e)})\}$ .

5:   Record predicted labels  $\hat{\mathbb{Y}}^{(e)} = f(\mathbb{X}^{(e)} | \mathbb{D}_{base}, \mathbb{D}_S^{(e)})$ .

6:   Compute accuracy  $a^{(e)} = \frac{1}{M} \sum_{j=1}^M \mathbb{1}[\hat{\mathbb{Y}}^{(e)} = \mathbb{Y}^{(e)}]^a$ .

7: **end for**

8: **return** mean accuracy  $\frac{1}{E} \sum_{e=1}^E a^{(e)}$ .

---

130   <sup>a</sup> $\mathbb{1}$  denotes the element-wise indicator function.

131 receiving the base dataset and the  $e$ th support set, which predicts labels of  
132 query instances as  $\hat{\mathbb{Y}}^{(e)} = f(\mathbb{X}^{(e)} | \mathbb{D}_{base}, \mathbb{D}_S^{(e)})$ . Let  $a^{(e)}$  denote the classification  
133 accuracy on the  $e$ th episode. The performance of a learning algorithm is  
134 measured by the classification accuracy averaged over all episodes.

### 135 2.3. Datasets for few-shot image classification

136   In this section, we briefly introduce benchmark datasets for few-shot im-  
137 age classification. Statistics of the datasets and commonly used experimental  
138 settings are listed below, and sample images are shown in Figure 1.

139 *Omniglot* [19]: one of the most widely used datasets for evaluating few-shot  
140 classification algorithms. It contains 1623 characters from 50 languages. The  
141 dataset is often augmented by rotations of 90, 180, 270 degrees, resulting in

142 6492 classes, which are split into 4112 base, 688 validation, and 1692 novel  
143 classes. The validation classes are used for model selection. The dataset is  
144 used less often in the latest studies as many methods can attain over 99%  
145 accuracy on the 5-way 1-shot classification task.

146 *Mini-ImageNet and Tiered-ImageNet*: another two widely used datasets de-  
147 rived from the ImageNet dataset [20]. Mini-ImageNet consists of 100 selected  
148 classes with 600 images for each class. This dataset was first proposed by  
149 Vinyals et al. [7], but recent studies follow the experimental setting provided  
150 by Ravi and Larochelle [21], which splits 100 classes into 64 base, 16 val-  
151 idation, and 20 novel classes. Tiered-ImageNet is a larger dataset with a  
152 hierarchical structure [22]. It is constructed from 34 super-classes with 608  
153 classes in total and include 779,165 images. These super-classes are split  
154 into 20 base, 6 validation, and 8 novel super-classes, which correspond to 351  
155 base, 97 validation, and 160 novel classes, respectively.

156 *CIFAR-FS and FC100*: two datasets derived from CIFAR-100 [23]. CIFAR-  
157 FS [24] contains 100 classes with 600 images per class, and it is split into 64  
158 base, 16 validation, and 20 novel classes. FC100 [25] divides 100 classes into  
159 20 super-classes, with five classes in each super-class. The dataset is split  
160 into 12 base, 4 validation, and 4 novel super-classes.

161 *Stanford Dogs [26]*: one of the benchmark datasets for fine-grained classifi-  
162 cation task, which contains 120 breeds (classes) of dogs with a total number  
163 of 20,580 images. These classes are divided into 70 base, 20 validation, and  
164 30 novel classes.

165 *CUB-200-2010/2011*: another fine-grained dataset of 200 bird species. The



Figure 1: Sample images of some benchmark datasets for few-shot image classification. Datasets include Onimiglot, Mini-ImageNet, Fewshot-CIFAR100, Stanford Dogs, and CUB-200-2011.

166 initial version in 2010 collects 6033 images [27] and is extended in 2011 to  
 167 11,788 images [28]. The CUB-200-2010 dataset is commonly split into 130  
 168 base, 20 validation, and 50 novel classes [29], while the CUB-200-2011 dataset  
 169 is commonly split into 100 base, 50 validation, and 50 novel classes [30].

170 *Mini-ImageNet*  $\rightarrow$  *CUB*: a dataset used for cross-domain few-shot classifica-  
 171 tion. Mini-ImageNet serves as the base dataset, 50 classes of CUB-200-2011  
 172 serve as the validation classes, and the remaining 50 classes serve as novel  
 173 classes.

174 *Meta-Dataset*: a new, large-scale dataset for evaluating few-shot classifica-  
 175 tion methods, particularly cross-domain methods. It initially consists of 10  
 176 diverse image datasets [31], e.g., ImageNet, CUB, and MS COCO [32], and  
 177 later expanded with three additional datasets [33]. There are two train-

ing procedures and two evaluation protocols. In the more commonly used setting of training on all datasets (multi-domain learning) [33, 34, 35], the methods are trained on the official training splits of the first eight datasets, and they are evaluated on the test splits of the same datasets for in-domain performance and the remaining five datasets for out-of-domain performance. The other setting is training only on the Meta-Dataset version of ImageNet (single-domain learning), and evaluating on the test split of ImageNet for in-domain performance and the rest 12 datasets for out-of-domain performance.

### 3. Few-shot deep metric learning methods

The goal of supervised metric learning is to learn a distance metric to measure the similarity among samples such that it is optimal for the subsequent learning tasks. For example, for classification, samples from the same (different, resp.) class should be assigned with a small (large, resp.) distance. In the case of few-shot classification, the metric is learned on the base dataset; query images of the novel class are classified by computing their distances to novel support images with respect to the learned measure, followed by applying a distance-based classifier such as the  $k$ -nearest neighbor ( $k$ NN) algorithm. Traditional metric learning methods learn a Mahalanobis distance, which is equivalent to learning a linear transformation of original features [36]. However, in deep metric learning, the distance measure and feature embeddings are often learned separately so as to capture the nonlinear data structure and generate more discriminative feature representations. Moreover, instead of comparing with individual samples, many few-shot metric learning methods compare query samples with class repre-

202 sentations such as prototypes and subspaces. In the remainder of this section,  
203 we provide a review of representative approaches, which are categorized into  
204 three groups according to the aspect they are improving on, namely 1) learn-  
205 ing feature embeddings, 2) learning class representations, and 3) learning  
206 distance or similarity measures. A summary of these methods is provided in  
207 Figure 2.

### 208 3.1. *Learning feature embeddings*

209 Methods of learning feature embeddings implicitly assume that the net-  
210 work is powerful to extract discriminative features and can generalize well  
211 to novel classes. Early approaches aim at a task-agnostic embedding model  
212 that is effective for any task. More recently, endeavors are made to learn a  
213 task-specific embedding model for better distinguishing the classes at hand.  
214 Furthermore, techniques for data augmentation and multi-task learning are  
215 leveraged to address the issues of data scarcity and overfitting.

#### 216 3.1.1. *Learning task-agnostic features*

217 The Siamese Convolutional Neural Network [6] is the first deep metric  
218 learning method for one-shot image classification. The Siamese Network, first  
219 introduced in [37], consists of two sub-networks with identical architectures  
220 and shared weights. [6] adopted the VGG-styled convolutional layers as the  
221 sub-network to extract high-level features from two images and employed  
222 the weighted  $L_1$  distance as the distance between the two feature vectors.  
223 Weights of the network, as well as those of component-wise distance, are  
224 trained using the conventional technique of mini-batch gradient descent.

225 The Matching Network [7] encoded support and query images using

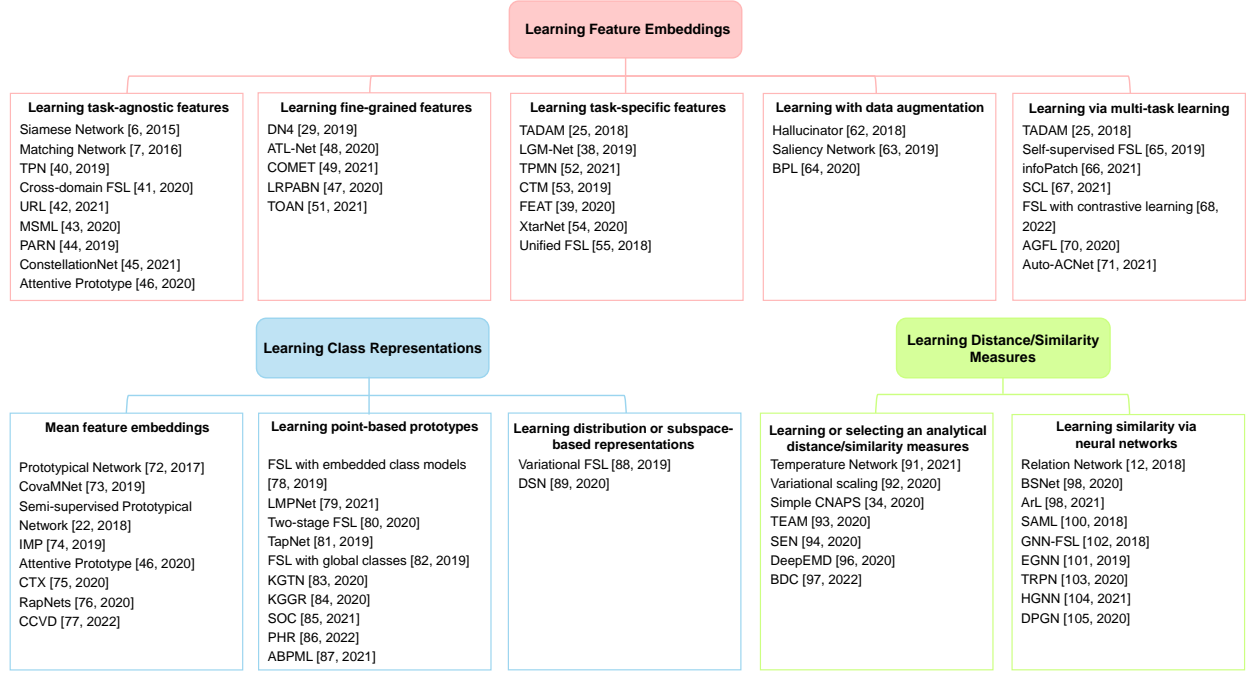


Figure 2: Taxonomy of few-shot deep metric learning methods reviewed in this paper. Some methods contribute to two aspects of metric learning and thus appear twice.

different networks in the context of the entire support set, and it first introduced episodic training to few-shot classification. A support image is embedded via a bidirectional LSTM network, which takes account of not only the image itself but also other images in the set; a query image is embedded via an LSTM with an attention mechanism to enable dependency on the support set. However, the sequential nature of bidirectional LSTM results in feature embeddings that will change with different ordering of samples in the support set. This issue can be sidestepped, such as by applying a pooling operation [38] or using self-attention [39]. The classification mechanism of Matching Network is suitable for few-shot learning. The network outputs a label distribution by computing a convex combination of one-hot label vectors of all support samples, with coefficients defined by using a softmax over

cosine similarities; the class with the highest probability is selected as the predicted class. Another valuable contribution of [7] is the episode-based training strategy, which has been adopted by many subsequent works. Following meta-learning, the training phase on the base dataset should mimic the prediction phase where only few support samples are available. That is, gradient updates should be performed on episodes with  $C$  classes randomly sampled from the base label set and  $K$  examples for each class.

The episodic training strategy closes the gap between training and test distributions and thus alleviates the issue of overfitting to few labeled training images. The overfitting issue can be further addressed by utilizing query instances (i.e., excluding query labels) via transductive inference. Transductive Propagation Network (TPN) [40] is the first work adopting transductive inference for few-shot learning and introduced the idea of label propagation. Concretely, the network contains a feature embedding module and a graph construction module. The graph construction module, taking feature embeddings as inputs, learns a label propagation graph to exploit the manifold structure of support and query samples. Based on the learned  $k$ NN graph, labels are propagated from the support set to the query set; a closed-form solution of label propagation is used to speed up the prediction procedure. While transductive learning takes advantage of query instances, it is unsuitable for online learning where data arrive sequentially.

The aforementioned methods, designed for classifying novel data from the same domain, degrade when novel data comes from different domains [30]. Tseng et al. [41] noticed that this is caused by the large discrepancy between the feature distributions in different domains and proposed to simulate var-

263 ious feature distributions in the training stage as a general solution to en-  
 264 hance the domain generalization ability of metric learning methods. This  
 265 is achieved by inserting multiple feature-wise transformation layers into the  
 266 feature extractor; each transformation simulates one distribution, and the  
 267 hyperparameters of affine transformations can be tuned via a meta-learning  
 268 approach so that they are optimal to a particular metric learning method  
 269 and capture the complex variation in feature distributions. Li et al. [42] pro-  
 270 posed to learn a universal feature representation that works well for multiple  
 271 domains. The technique of knowledge distillation is applied, where a multi-  
 272 domain network is learned to generate universal features which align with  
 273 features from multiple single-domain networks up to a linear transformation.

274 Motivated by the observation that the interested object may locate only  
 275 in a region of an image and at different positions across images, a series of  
 276 improvements on feature embedding have been proposed, such as by learning  
 277 local features [29] and multi-scale features [43] and encoding the position in-  
 278 formation [44]. Local feature-based methods, while can be applied to generic  
 279 few-shot image classification, are particularly effective for fine-grained im-  
 280 age classification and thus will be discussed separately in the next subsec-  
 281 tion. Jiang et al. [43] proposed the Multi-Scale Metric Learning (MSML)  
 282 network, which constructs multiple feature embeddings corresponding to dif-  
 283 ferent scales of the image. The similarity between support and query features  
 284 at each scale is computed using the Relation Network (which will be intro-  
 285 duced in Section 3.3.2). Wu et al. [44] proposed the Position-Aware Relation  
 286 Network (PARN) to reduce the sensitivity of Relation Network to the spatial  
 287 position of semantic objects. PARN adopts deformable convolutional layers



288 to extract more effective features which filter out unrelated information like  
 289 the background, and a dual correlation attention module to incorporate each  
 290 spatial position of an image with the global information about the compared  
 291 image and the image itself, so that the subsequent convolution operations,  
 292 even subject to local connectivity, can perceive and compare semantic fea-  
 293 tures in different positions. Compared with standard ways of overcoming  
 294 position sensitivity, such as by using larger kernels or more layers, PARN  
 295 is more parameter efficient. Xu et al. [45] proposed the ConstellationNet  
 296 which extracts part-based features and encodes the spatial relationship be-  
 297 tween these representations by using self-attention with an explicit, learnable  
 298 positional encoding. The spatial relationship between different parts of the  
 299 image has also been encoded in [46] by using a capsule network.

### 300 3.1.2. *Learning task-agnostic features for fine-grained image classification*

301 Fine-grained image classification aims to distinguish different sub-categories  
 302 under the same basic-level category. It is particularly challenging due to the  
 303 subtle differences between different sub-categories and large variance in the  
 304 same sub-category which may result from variations in the object’s pose,  
 305 scale, rotation, etc. Therefore, for effective classification, several methods  
 306 have been proposed to extract local features and second-order features.

307 In deep nearest neighbor neural network (DN4) [29], the feature embed-  
 308 ding module extracts multiple local descriptors from an image, which are  
 309 essentially the feature maps learned via CNNs prior to adding the final image-  
 310 level pooling layer. The classification is performed at an image-to-class level,  
 311 meaning that the local descriptors from support images of the same class  
 312 are put into one pool,  $k$ NNs in each class pool are searched for each query

313 local descriptor, and the total distance over all local descriptors and  $k$ NNs  
 314 is the distance between the query image and the corresponding class. The  
 315 method is shown to be particularly effective on fine-grained datasets, and  
 316 the idea of learning local descriptors has been adopted in other fine-grained  
 317 classification methods [47]. The Adaptive Task-aware Local representations  
 318 Network (ATL-Net) [48] improved DN4 by selecting local descriptors with  
 319 learned thresholds and assigning them different weights based on episodic  
 320 attention, which brings more flexibility than using  $k$ NNs and adjusts for the  
 321 discriminability between classes, respectively. In contrast to learning one  
 322 feature embedding over spatially local patches, COMET [49] learns multiple  
 323 embedding functions over various parts of input features. A set of fixed bi-  
 324 nary masks, termed concepts, are applied to input features to separate an  
 325 image into human-interpretable segments. For each concept, a feature em-  
 326 bedding is learned to map masked features into a new discriminative feature  
 327 space. The query image is classified according to the distances aggregated  
 328 from all concept-specific spaces.

329 Huang et al. [47] proposed the Low-Rank Pairwise Alignment Bilinear  
 330 Network (LRPABN) which aligns features spatially and extracts discrimi-  
 331 native, second-order features. After learning first-order features from base  
 332 images, the method trains a two-layer multi-layer perceptron network with  
 333 two designed feature alignment losses to transform the positions of image  
 334 features of a query image to match those of a support image, and designs a  
 335 low-rank pairwise bilinear pooling layer which adapts the self-bilinear pool-  
 336 ing [50] to extract second-order features from a pair of support and query  
 337 images. The classification is performed as in the Relation Network. In the

338 follow-up work, [51] improves the spatial alignment part by using the cross-  
 339 channel attention to generate spatially matched support and query features  
 340 and groups features in the convolutional channel dimension before the pool-  
 341 ing layer as each group corresponds to a semantic concept.

### 342 3.1.3. *Learning task-specific features*

343 Methods reviewed in the preceding sections generate the same feature em-  
 344 bedding for an image, regardless of the subsequent classification task. While  
 345 this avoids the risk of overfitting, these generic features may not be suffi-  
 346 ciently discriminative to distinguish novel classes. To this end, task-specific  
 347 embedding models have been proposed to adapt features to the particular  
 348 task; it should be noted that the adaptation is learned on the base dataset  
 349 and does not involve any re-training on the novel dataset.

350 TADAM [25] is the first metric learning method which explicitly performs  
 351 task adaptation. Exploiting the technique of conditional batch normaliza-  
 352 tion, it applies a task-specific affine transformation to each convolutional  
 353 layer of a task-agnostic feature extractor. The task is represented by the  
 354 mean of class prototypes, and the scale and shift parameters of the affine  
 355 transformation are generated from a separate network, called the Task Em-  
 356 bedding Network (TEN). As TEN introduces more parameters and causes  
 357 difficulty in optimization, the training scheme is revised to add the standard  
 358 training, i.e., to distinguish all classes in the base dataset, as an auxiliary  
 359 task to the episodic training.

360 Li et al. [38] proposed a meta-learning approach that can adapt weights  
 361 of Matching Network to novel data. The proposed LGM-Net consists of a  
 362 meta-learner termed MetaNet and a task-specific learner termed TargetNet.

363 The MetaNet module learns to produce a representation of each task from  
 364 the support set and construct a mapping from the representation to weights  
 365 of TargetNet. The TargetNet module, set as the Matching Network, em-  
 366 beds support and query images and performs classification. The proposed  
 367 meta-learning strategy can be potentially implemented to adapt network pa-  
 368 rameters of other metric learning methods. Wu et al. [52] also proposed to  
 369 learn task-specific parameters, but they combined the idea with local fea-  
 370 tures. The proposed Task-aware Part Mining Network (TPMN) learns to  
 371 generate parameters of filters used for extracting part-based features.

372 Different from the above two works which generate parameters for task-  
 373 specific embedding layers, Li et al. [53] proposed to modify the generic fea-  
 374 tures output from the task-agnostic embedding layers. A task-specific fea-  
 375 ture mask is generated from the Category Traversal Module (CTM), which  
 376 includes a concentrator unit and a projector unit to extract features for intra-  
 377 class commonality and inter-class uniqueness, respectively. It is noted that  
 378 CTM can be easily embedded into most few-shot metric learning methods,  
 379 such as Matching Network, Prototypical Network, and Relation Network; the  
 380 latter two methods will be introduced in the following sections. Ye et al. [39]  
 381 also proposed to adjust features directly, but instead of applying a mask,  
 382 set-to-set functions are used to transform a set of task-agnostic features into  
 383 a set of task-specific ones. These functions can model interactions between  
 384 images in a set and hence enable co-adaptation of each image. Four set-to-set  
 385 function approximators are presented in [39], and the one with Transformer,  
 386 termed FEAT, is shown to be most effective.

387 Yoon et al. [54] proposed XtarNet to learn task-specific features for a new

388 setting of generalized few-shot learning, where the model is trained on the  
 389 base dataset, adapted given the support set of the novel dataset, and used to  
 390 classify instances from both base and novel classes. XtarNet contains three  
 391 meta-learners. The MetaCNN module adapts feature embeddings for each  
 392 task. The MergeNet module produces weights for mixing pre-trained features  
 393 and meta-learned features. As the classification is performed by comparing  
 394 the mixed features with class prototypes, the TconNet module adapts pro-  
 395 totypes of base and novel classes to improve discriminability. Rahman et  
 396 al. [55] proposed a unified approach for zero-shot learning, generalized zero-  
 397 shot learning and few-shot learning, which classifies a query image based on  
 398 the similarity between its semantic representation and the textual features  
 399 of each class. The semantic representation is a combination of two parts –  
 400 one is a linear combination of base samples’ semantic features, and the other  
 401 one is based on the linear mapping learned from support images.

#### 402 3.1.4. *Feature learning with data augmentation*

403 Data augmentation is a strategy that expands the support set in an ar-  
 404 tificial or model-based way with label preserving transformations, and thus  
 405 is well-suited when the support samples are limited. One commonly used  
 406 method is deformation [56, 57, 58], such as cropping, padding, and hori-  
 407 zontal flipping. Besides this, generating more training samples [59, 60] and  
 408 pseudo labels [61] are also popular techniques to augment data.

409 In few-shot learning, there is one class of works which places the data  
 410 augmentation process into a model, that is, they embed a generator that can  
 411 generate the augmented data to learn or imagine the diversity of data. Wang  
 412 et al. [62] constructed an end-to-end few-shot learning method, in which the

413 training data goes through two streams to output – one is from the original  
 414 data to the classifier directly, and the other one is from the original data to a  
 415 ‘hallucination’ network to augment data and then from the augmented data  
 416 to classifier. Zhang et al. [63] developed a saliency-based data generation  
 417 strategy. The Saliency Network obtains foregrounds and backgrounds of an  
 418 image, which are used to achieve the hallucination for the image. In [64],  
 419 a much simpler feature synthesis strategy was proposed, which synthesizes  
 420 novel features by perturbing the semantic representations (i.e., word vectors  
 421 of class labels) and projecting them into the visual feature space. In addition,  
 422 when learning the projection function, a competitive learning formulation is  
 423 adopted to push the synthesized sample towards the center of the most likely  
 424 unseen class and away from that of the second best class.

### 425 3.1.5. *Multi-task feature learning*

426 Besides generating more training data, some works tried to exploit auxil-  
 427 iary information of samples to perform multi-task learning, which creates a  
 428 regularization effect and helps learn discriminative features.

429 As briefly discussed above, TADAM [25] used an auxiliary task of training  
 430 a normal global classifier on the base dataset to co-train the few-shot classi-  
 431 fier; the task is sampled with a probability during the training process. An  
 432 alternative auxiliary task is to exploit generative [65] or contrastive [66] self-  
 433 supervised learning, which adopts self-defined pseudo labels as supervision  
 434 to learn generalizable feature embeddings. In [65], support samples are arti-  
 435 ficially rotated to different number of degrees. A shared feature embedding is  
 436 learned through two branches of networks, one for the original classification  
 437 task and the other for identifying the rotation degree. In [66], infoPatch was

438 proposed, which trains the embedding network episodically according to the  
 439 standard classification loss and an auxiliary contrastive loss. The contrastive  
 440 pairs are constructed for each query image, with the positive pair using sup-  
 441 port images of the same class and the negative pair using supports of different  
 442 classes. To generate hard pairs, random blocks are applied to support im-  
 443 ages to mask parts of the image, and a query image is split into patches with  
 444 one of them exchanged with a patch of another image. Not only in episodic  
 445 training, contrastive learning can also be introduced in pre-training [67] or  
 446 in both training stages [68]. In particular, in the episodic training stage  
 447 of [68], the entire episode is regarded as the shared context, and two data  
 448 augmentation strategies are applied to construct contrastive episodes. How-  
 449 ever, as noted by Xiao et al. [69], these contrastive learning methods require  
 450 hand selecting augmentations and carefully tuning the hyperparameters to  
 451 control the strength of augmentation. More severely, they implicitly assume  
 452 invariance to particular transformations, e.g., rotation and color, which may  
 453 be beneficial to some downstream tasks but harmful to others. One solution  
 454 proposed in [69] is to use a multi-head network with a shared backbone to  
 455 learn several embedding spaces, one for invariance to all augmentations and  
 456 the others for invariance to all but one augmentation. The downstream task  
 457 can flexibly utilize the optimal set of invariant features. The solution was  
 458 proposed in a transfer learning setting; more research is needed for metric  
 459 learning.

460 Zhu et al. [70] suggested that base and novel classes, despite being dis-  
 461 joint, can be connected by some visual attributes. Based on this insight,  
 462 they used attribute learning as an auxiliary task. Visual attributes are pro-

463 vided as additional information during training, and the embedding network  
 464 is learned to correctly predict both attribute labels and class labels. [71] also  
 465 utilized attribute information, but in a richer way which requires an addi-  
 466 tional prediction of common and different attributes between an image pair.  
 467 Moreover, the neural architecture search was first introduced to few-shot  
 468 learning for automatically identifying the optimal operation from max pool-  
 469 ing, convolution, identity mapping, etc for layers in the feature embedding  
 470 network and attribute learning network.

### 471 *3.2. Learning class representations*

472 Early few-shot metric learning methods such as Siamese Network and  
 473 Matching Network classify a query sample by measuring and comparing its  
 474 distance to support samples. However, since support samples are scarce, they  
 475 have limited capacity in representing the novel class. To alleviate this issue,  
 476 some researchers propose to use class prototypes, which serve as reference  
 477 vectors for each class. Prototypes can be constructed by taking simple or  
 478 weighted average of feature embeddings, or learned in an end-to-end manner  
 479 so as to further improve their representation ability. Besides point-based pro-  
 480 totypes, some works consider the distribution of each class or use subspaces  
 481 as class representations.

#### 482 *3.2.1. Feature embeddings-based prototypes*

483 Prototypical Network [72] is a classical method that performs classifica-  
 484 tion by calculating the Euclidean distance to class prototypes in the learned  
 485 embedding space. It builds on the hypothesis that there exists an embedding  
 486 space in which each class can be represented by a single prototype and all



instances cluster around the prototype of their corresponding classes. In [72], the prototype of each class is set as the mean of feature embeddings of support samples in the class. Feature embeddings, and thus class prototypes, are learned using episodic training with the objective of minimizing the cross entropy loss. In [73], the class prototype is represented using the covariance matrix of feature embeddings. A covariance-based metric is also proposed to measure the similarity between the query and the class.

To make use of both labeled support samples and unlabeled samples, Ren et al. [22] proposed semi-supervised Prototypical Network, which is the first work of semi-supervised few-shot learning. The method adopts soft  $k$ -means to compute assignment score of unlabeled samples and computes the prototype as the mean of weighted samples based on assignment scores.

Considering that the dataset may exhibit multi-modality and multiple prototypes would be more suitable in this scenario, Infinite Mixture Prototypes (IMP) [74] was proposed to model multiple clusters in each class, and each cluster is modeled as a Gaussian distribution. Concretely, the probability that a sample follows the Gaussian distribution of each cluster determines which cluster the sample is assigned into. Moreover, the cluster variance of the Gaussian distributions, which needs to be learned, can affect the number of class prototype and performance of IMP.

Wu et al. [46] proposed to compute query-dependent prototypes. An attentive prototype is computed for each query as the weighted average of support samples and the weights are given by the Gaussian kernel with the Euclidean distance between the query and the support samples. As support samples that are more relevant to the query have greater influence on

the classification, the method is more robust to outliers in support samples. Query-dependent prototypes have also been studied in CrossTransformers (CTX) [75], but they are computed separately for each spatial location. In other words, a local region of a query image is compared with an attentive prototype specific to this query and region, and the overall distance between the query and the prototype is the averaged distances over all local regions. Moreover, self-supervised episodes are constructed to train CTX.

Lu et al. [76] proposed the Robust attentive profile Networks (RapNets) to enhance the robustness of prototypes against outliers and label noises. The network transforms raw feature embeddings into correlation features in a nonparametric way and then inputs these features into a parametric bidirectional LSTM and fully-connected network to generate attention scores which serve as weights to combine support images. Moreover, training episodes are revised to include noisy data, and a new evaluation metric is proposed to evaluate the robustness of few-shot classification methods.

Ma et al. [77] provided a geometric interpretation of Prototypical Network, regarding it as a Voronoi diagram. In addition, the authors extended this perspective and proposed the Cluster-to-Cluster Voronoi Diagram (CCVD), which can ensemble models learned with different data augmentation, built on single or multiple feature transformations, and using linear or nearest neighbor classifier.

### 3.2.2. Point-based learnable prototypes

Ravichaandran et al. [78] adopted an implicit way to learn class representation instead of determining class prototypes as in the aforementioned methods. The prototype is modeled as a learnable and parameterized func-

537 tion of feature embedding of labeled samples in the class and is obtained by  
 538 minimizing a loss which measures the distance between the feature embed-  
 539 ding of a sample and the class prototype. Meanwhile, the function is shot  
 540 free, that is, it allows sample sizes of classes in novel data to be unbalanced.  
 541 In [79], prototypes are represented as weighted averages of feature embed-  
 542 dings, but different from [22, 46] discussed in the previous section, weights are  
 543 learned end-to-end via episodic training. Moreover, instead of using image-  
 544 level features, [79] combines local descriptors of one class following the idea of  
 545 DN4 and learns multiple weight vectors to generate multiple prototypes per  
 546 class. Das and Lee proposed a two-stage approach for generating class pro-  
 547 totypes [80]. In the first stage, feature embeddings are learned, from which  
 548 coarse prototypes of base and novel classes can be obtained from mean em-  
 549 beddings. In the second stage, the novel class prototype is refined through a  
 550 meta-learnable function of its own prototype and related base prototypes.

551 Besides the above methods, TapNet [81] explicitly modeled class pro-  
 552 totypes as learnable parameters. Prototypes and feature embeddings are  
 553 learned simultaneously on the base dataset following the training procedure  
 554 of Prototypical Network. In addition, to make prototypes and feature em-  
 555 beddings more specific to the current task, both of them are projected into  
 556 a new classification space via a linear projection matrix. The projection  
 557 matrix is obtained by using a linear nulling operation and does not include  
 558 any learnable parameter. Luo et al. [82] proposed to learn prototypes of  
 559 base and novel classes simultaneously by including the support set of novel  
 560 classes in the training process. In each episode, local prototypes are gener-  
 561 ated from the sample synthesis module, which aims to increase the diversity

562 of novel classes. They are then used in the registration module to update  
 563 the global prototypes towards better separability. The query image is clas-  
 564 sified by searching for the nearest neighbor among global prototypes. As  
 565 both base and novel prototypes are learned, the method can be readily ap-  
 566 plied to the generalized few-shot learning setting. Chen et al. [83] shared the  
 567 same aim of learning base and novel prototypes, but additionally took advan-  
 568 tage of the semantic correlations among these classes. A Knowledge Graph  
 569 Transfer Network (KGTN) is proposed, which employs a gated graph neural  
 570 network to represent class prototypes and correlations as nodes and edges,  
 571 respectively. By propagating through the graph, information from correlated  
 572 base classes is used to guide the learning of novel prototypes. This work is  
 573 extended in [84] to the multi-label classification setting, which employs the  
 574 attention mechanism and an additional graph for learning class-specific fea-  
 575 ture vectors. In [85], the Shared Object Concentrator (SOC) algorithm was  
 576 proposed to learn a series of prototypes for each novel class from local fea-  
 577 tures of support images. The first prototype is learned to have the largest  
 578 cosine similarity with one of the local features, the second prototype has the  
 579 second largest value, and so forth. The query image is classified according to  
 580 the weighted sum of similarities between its local features and all prototypes,  
 581 with weights decaying exponentially to account for the decreasing influence  
 582 of prototypes. Zhou et al. [86] proposed the Progressive Hierarchical Refine-  
 583 ment (PHR) method to update prototypes iteratively using all novel data.  
 584 In each iteration, support images and a random subset of query images are  
 585 embedded into features at local, global and semantic levels, and a loss func-  
 586 tion defined over these hierarchical features is used to refine prototypes for

587 better inter-class separability. As each update is based on a random subset  
588 of queries, the method is less likely to overfit to noisy query samples, though  
589 it implicitly assumes the availability of a large number of queries.

590 Sun et al. [87] proposed to treat prototypes as random variables. The  
591 posterior distributions of latent class prototypes are learned by using amor-  
592 tized variational inference, a technique which enables prototype learning to  
593 be formulated as a probabilistic generative model without encountering se-  
594 vere computational and inferential difficulties.

### 595 3.2.3. *Distribution or subspace-based representations*

596 Considering that single point-based metric learning is sensitive to noise,  
597 Zhang et al. [88] proposed a variational Bayesian framework for few-shot  
598 learning and used the Kullback-Leibler divergence to measure the distance  
599 of samples. The framework can compute the confidence that a query image  
600 is assigned into each class by estimating the distribution of each class based  
601 on a neural network.

602 Simon et al. [89] proposed Deep Subspace Network (DSN) to represent  
603 each class using a low-dimensional subspace, constructed from support sam-  
604 ples via singular value decomposition. Query samples are classified according  
605 to the nearest subspace classifier, that is to assign the query to the class which  
606 has the shortest Euclidean distance between the query and its projection onto  
607 the class-specific subspace. The method is shown to be more robust to noises  
608 and outliers than Prototypical Network.

### 609 3.3. *Learning distance or similarity measures*

610 Methods reviewed in Sections 3.1 and 3.2 focus on learning a discrimi-

native feature embedding or obtaining an accurate class representation. For classification, they mostly adopt a fixed distance or similarity measure, such as the Euclidean distance [72] and the cosine similarity [7]. More recently, researchers propose to learn parameters in these fixed measures or define novel measures so as to further improve the classification accuracy. Moreover, considerable effort has been made to learn similarity scores by using fully-connected neural networks or Graph Neural Networks (GNNs).

### 3.3.1. *Learning or selecting an analytical distance or similarity measure*

In TADAM [25], Oreshkin et al. mathematically analyzed the effect of metric scaling on the loss function. Since then, many works tune the scaling parameter via cross-validation [48, 90]. Zhu et al. [91] proposed to use two different scaling parameters for the ground-truth class and other classes to enforce the same-class distance is much smaller than the different-classes distance. Moreover, the scaling parameters are gradually tuned every few episodes, which implements the idea of self-paced learning to learn from easy to hard. Chen et al. [92] proposed to learn the scaling parameter in a Bayesian framework. By assuming a univariate or multivariate Gaussian prior and applying the stochastic variational inference technique for approximating the posterior distribution, a scaling parameter or a scaling vector can be learned respectively, which scales the distance equally over all dimensions or differently for each dimension. Task-specific scaling vectors can also be learned by learning a neural network from the task to variational parameters.

The traditional Mahalanobis distance decorrelates and scales features using the inverse of the covariance matrix. In Simple CNAPS [34], after extracting features using the architecture of Conditional Neural Adaptive Processes

(CNAPS) [33], the classification is performed based on the Mahalanobis distance between query instances and class prototypes. Task-specific class-specific covariance matrices are estimated as convex combinations of sample covariance matrices estimated from instances of the task and instances of the class and regularized toward an identity matrix. Transductive Episodic-wise Adaptive Metric (TEAM) [93] learned task-specific metric from support and query samples. TEAM contains three modules, namely a feature extractor, a task-specific metric module, and a similarity computation module. The task-specific metric module learns a Mahalanobis distance to shrink the distance between similar pairs and enlarge the distance between dissimilar pairs, following the objective function of the pioneering metric learning method [36].

Nguyen et al. [94] proposed a dissimilarity measure termed SEN, which combines the Euclidean distance and the difference in the  $L_2$ -norm. Minimizing this measure will encourage feature normalization and consequently benefit the classification performance [95]. DeepEMD [96] combined a structural distance over dense image representations, Earth Mover’s Distance (EMD) and convolutional feature embedding to conduct few-shot learning. The optimal matching flow parameters in EMD and the parameters in the feature embedding are trained in an end-to-end fashion. Xie et al. [97] introduced the Brownian Distance Covariance (BDC) metric, a new distance measure founded on the characteristic function of random vectors. The metric has a closed-form expression for discrete feature vectors and can be computed easily by first computing the BDC matrix for every image and then calculating the inner product between two BDC matrices. The computation of BDC matrices also only involves standard matrix operations and can be formu-

lated as a pooling layer, thus endowing the method with high computational efficiency and ease of integrating with other few-shot classification methods.

### 3.3.2. *Learning similarity scores via neural networks*

The Relation Network [12] is the first work of introducing a neural network to model the similarity of feature embeddings in few-shot learning. It consists of an embedding module and a relation module. The embedding module builds on convolutional blocks for mapping original images into an embedding space, and the relation module consists of two convolutional blocks and two fully-connected layers for computing the similarity between each pair of support and query images. The learnable similarity measure enhances the model flexibility. Li et al. [98] pointed out that a single similarity measure may not be sufficient to learn discriminative features for fine-grained image classification and thus proposed the Bi-Similarity Network (BSNet), which integrates the proposed cosine module with existing similarity measures such as the relation module, forcing features to adapt to two similarity measures of diverse characteristics and consequently generating a more compact feature space. In principle, the method can be further developed to ensemble multiple metrics, and more importantly, an elegant way to determine the optimal set of metrics to be combined is needed. Relation Network and subsequent methods all use class labels to form binary supervision, indicating whether the image pair comes from the same class. Zhang et al. [99] argued that such binary relations are not sufficient to capture the similarity nuance in the real-world setting and therefore proposed a new method termed Absolute-relative Learning (ArL) which, in addition to binary relations, constructs continuous-valued relations from attributes of images, such as colors



686 and textures.

687 Different from Relation Network, Semantic Alignment Metric Learning  
688 (SAML) [100] adopted the Multi-Layer Perceptron (MLP) network for com-  
689 puting the similarity score. Specifically, SAML contains a feature embedding  
690 module and a semantic alignment module. In the semantic alignment mod-  
691 ule, a relation matrix at the level of local features is first computed by using  
692 fixed similarity measures and an attention mechanism, and then fed into a  
693 MLP network which outputs the similarity score between the query and the  
694 support class. Due to the use of relation matrix as the input, the MLP net-  
695 work has more parameters than Relation Network, thus increasing the risk  
696 of overfitting.

697 Recently, some researchers adopt Graph Neural Networks (GNNs) to im-  
698 plement few-shot classification. Like the above reviewed works, GNN-based  
699 methods also use a neural network to model the similarity measure, while  
700 its advantage lies in the rich relational structure on samples [101]. Garcia  
701 et al. [102] proposed the first GNN-based neural network for few-shot learn-  
702 ing, short for GNN-FSL here. It contains two modules, a feature embedding  
703 module and a GNN module. In the GNN module, a node represents a sam-  
704 ple, and more specifically, equals the concatenation of features of the sample  
705 and its label. For a query sample, its initial label in the first GNN layer  
706 uses uniform distribution over  $K$ -simplex ( $K$  is number of classes), and its  
707 predicted label in the last GNN layer is used for computing the loss func-  
708 tion. Like GNN-FSL, Edge-labeling Graph Neural Network (EGNN) [101]  
709 also contains a feature embedding module and a GNN module with three  
710 layers. However, rather than labeling nodes, EGNN learns to label edges in

711 GNN layers so that it can cluster samples explicitly by employing the intra-  
 712 cluster similarity and inter-cluster dissimilarity. In EGNN, each GNN layer  
 713 has its own loss function that is computed based on edge values in the layer,  
 714 and the total loss is the weighted sum of loss functions of all GNN layers.  
 715 The Transductive Relation-Propagation graph neural Network (TRPN) [103]  
 716 explicitly modeled the relation of support-query pairs by treating them as  
 717 graph nodes. After relation propagation, a similarity function is learned to  
 718 map the updated node to a similarity score, which represents the probability  
 719 that the support and query samples are of the same class. The class with the  
 720 highest sum of scores is the predicted class. The Hierarchical Graph Neural  
 721 Network (HGNN) [104], aimed at modeling the hierarchical structure within  
 722 classes, first down-samples support nodes to build a hierarchy of graphs and  
 723 then performs up-sampling to reconstruct all support nodes for prediction.

724 The previous GNN-based methods focus simply on the relation between a  
 725 pair of samples. In Distribution Propagation Graph Network (DPGN) [105],  
 726 the global relation between a sample and all support samples is considered by  
 727 generating a distribution feature from the similarity vector. A dual complete  
 728 graph is built to proceed sample-level and distribution-level features inde-  
 729 pendently, and a cyclic update policy is used to propagate between the two  
 730 graphs. Information from the distribution graph refines sample-level node  
 731 features and hence improves the classification based on edge similarities.

732 Table 2 summarizes few-shot deep metric learning methods, listing the  
 733 backbone network for feature embedding, classification mechanism, similarity  
 734 measure, training strategy, datasets studied in the experiment, and classifica-  
 735 tion performance. As the methods were implemented with different backbone

networks and tested on different datasets, for a fair comparison, we select Conv-4 and ResNet-12 backbones whenever possible and report the 5-way 1-shot and 5-way 5-shot classification accuracy on Mini-ImageNet. Moreover, we notice that some methods were trained with higher ways or higher shots, which may lead to better performance, and thus this information is included under training strategy. Nevertheless, there are other factors which may affect the performance, such as the use of data augmentation techniques, optimization strategy, and the number of test episodes. Table 3 is a summary for few-shot fine-grained image classification. Here we note that the CUB dataset was split into training, validation and test sets in multiple ways.

## 4. Further research

Even though few-shot metric learning methods have achieved the promising performance, there remains several important challenges that need to be dealt with in the future. In this section, we will discuss issues related to generalization and robustness of few-shot learning methods, training strategy, and applicability, as well as listing some promising applications of few-shot metric learning methods.

### 4.1. Challenges and future directions

*1. Improving generalized feature learning on few samples.* In the existing few-shot metric learning methods or even the entire few-shot learning methods, researchers mostly try to learn discriminative feature based on the attention mechanism, data augmentation, multi-task learning, and so on. To learn feature with good generalization ability from few labeled examples, new ways of evaluation and feature learning need to be developed.

Table 2: Summary of deep metric learning methods for few-shot image classification.

| Method                                    | Classification mechanism       | Similarity measure  | Training strategies                       | Mini-ImageNet (Conv-4) | 1-shot           | 5-shot                       | 1-shot                       | 5-shot           | Additional datasets or embedding architectures      |
|---|--------------------------------|---|---|------------------------|------------------|------------------------------|------------------------------|------------------|---|
| Siamese Network [6]                       | w.r.t. instances               | weighted $L_1$ distance                                     | minibatch training                        | -                      | -                | -                            | -                            | -                | Omniglot  |
| Matching Network [7]                      | w.r.t. instances               | cosine similarity   | episodic training                         | 46.60                  | 60.00            | -                            | -                            | -                | Omniglot  |
| TPN [40]                                  | w.r.t. instances               | weighted Euclidean distance (learnable weights)             | episodic training (higher-shot training)  | 55.51                  | 69.86 (T)        | -                            | -                            | -                | Tiered-ImageNet                                     |
| Cross-domain FSL [41]                     | w.r.t. instances               | learned distance  | pre-train + episodic training             | -                      | -                | -                            | 66.32 $\pm$ 0.80 (ResNet-10) | 81.98 $\pm$ 0.55 | -   |
| URL [42]                                  | w.r.t. prototypes              | cosine similarity   | episodic training                         | -                      | -                | -                            | -                            | -                | Meta-Dataset  |
| MSML [43]                                 | w.r.t. prototypes              | learned distance  | pre-train + episodic training             | -                      | -                | -                            | 72.41 $\pm$ 1.72 (ResNet-50) | 84.33 $\pm$ 1.14 | Tiered-ImageNet                                     |
| PARN [44]                                 | w.r.t. instances               | learned distance  | episodic training                         | 55.22 $\pm$ 0.84       | 71.55 $\pm$ 0.66 | -                            | -                            | -                | Omniglot  |
| ConstellationNet [45]                     | w.r.t. prototypes              | cosine similarity   | episodic training                         | 58.82 $\pm$ 0.23       | 75.00 $\pm$ 0.18 | 64.89 $\pm$ 0.23             | 79.95 $\pm$ 0.17             | 79.95 $\pm$ 0.17 | CIFAR-FS, FC100                                     |
| TADAM [25]                                | w.r.t. prototypes              | Euclidean distance  | episodic training + co-training           | -                      | -                | 58.5 $\pm$ 0.3               | 76.7 $\pm$ 0.3               | 76.7 $\pm$ 0.3   | FC100   |
| LGM-Net [38]                              | w.r.t. instances               | cosine similarity   | episodic training                         | 69.13 $\pm$ 0.35       | 71.18 $\pm$ 0.68 | -                            | -                            | -                | Omniglot  |
| TPMN [52]                                 | w.r.t. prototypes              | weighted sum of dot products                                | pre-train + episodic training             | -                      | -                | -                            | 67.64 $\pm$ 0.63             | 83.44 $\pm$ 0.43 | Tiered-ImageNet, CIFAR-FS, FC100                    |
| CTM [53]                                  | w.r.t. instances or prototypes | Any, e.g., cosine similarity<br>Euclidean, learned distance | pre-train (opt.) + episodic training      | -                      | -                | -                            | 64.12 $\pm$ 0.82 (ResNet-18) | 80.51 $\pm$ 0.13 | Tiered-ImageNet                                     |
| FEAT [39]                                 | w.r.t. instances               | cosine similarity   | pre-train + fine-tune temperature scaling | 55.15 $\pm$ 0.20       | 71.61 $\pm$ 0.16 | 66.78 $\pm$ 0.20             | 82.05 $\pm$ 0.14             | 82.05 $\pm$ 0.14 | Tiered-ImageNet, OfficeHome; WRN                    |
| Hallucinator [62]                         | w.r.t. prototypes              | cosine similarity   | episodic training                         | -                      | -                | -                            | -                            | -                | ImageNet; ResNet-10, ResNet-50                      |
| Saliency Network [63]                     | w.r.t. instances               | learned distance  | episodic training                         | 57.45 $\pm$ 0.88       | 72.01 $\pm$ 0.67 | -                            | -                            | -                | Open MIC  |
| BPL [64]                                  | w.r.t. prototypes              | Euclidean distance with learned projection matrix           | pre-train + episodic training             | 54.20 $\pm$ 0.58       | 65.28 $\pm$ 0.59 | 59.57 $\pm$ 0.63             | 76.86 $\pm$ 0.49             | 76.86 $\pm$ 0.49 | on ZSL, GZSL; WRN                                   |
| Self-supervised FSL [65]                  | w.r.t. prototypes              | Euclidean distance  | episodic training / minibatch training    | 54.83 $\pm$ 0.43       | 71.86 $\pm$ 0.33 | 62.93 $\pm$ 0.46 (WRN)       | 79.87 $\pm$ 0.33             | 79.87 $\pm$ 0.33 | Tiered-ImageNet, ImageNet-FS; Conv-4-512, ResNet-10 |
| infoPatch [66]                            | w.r.t. instances               | cosine similarity   | episodic training                         | -                      | -                | 67.67 $\pm$ 0.45             | 82.44 $\pm$ 0.31             | 82.44 $\pm$ 0.31 | Tiered-ImageNet, FC100                              |
| SCL [67]                                  | w.r.t. prototypes              | Euclidean distance  | pre-train / episodic training             | -                      | -                | -                            | 77.60 (ResNet-18)            | 77.60            | Tiered-ImageNet; CIFAR-FS, FC100                    |
| FSL with contrastive learning [68]        | w.r.t. prototypes              | Euclidean distance  | pre-train + episodic training             | -                      | -                | 70.19 $\pm$ 0.46             | 84.66 $\pm$ 0.29             | 84.66 $\pm$ 0.29 | Tiered-ImageNet, CIFAR-FS                           |
| AGFL [70]                                 | w.r.t. instances or prototypes | Any, e.g., cosine similarity, Euclidean, learned distance   | episodic training                         | -                      | -                | 56.59 $\pm$ 0.64 (ResNet-50) | 73.58 $\pm$ 0.48             | 73.58 $\pm$ 0.48 | CUB-200-2011, AWA                                   |
| Prototypical Network [72]                 | w.r.t. prototypes              | Euclidean distance  | episodic training (higher-way training)   | 49.42 $\pm$ 0.78       | 68.20 $\pm$ 0.66 | -                            | -                            | -                | Omniglot; CUB-200-2011 (for ZSL)                    |
| Semi-supervised Prototypical Network [22] | w.r.t. prototypes              | Euclidean distance  | episodic training                         | 50.41 $\pm$ 0.31       | 64.39 $\pm$ 0.24 | -                            | -                            | -                | Omniglot, Tiered-ImageNet                           |
| IMP [74]                                  | w.r.t. prototypes              | Euclidean distance  | episodic training                         | 49.60 $\pm$ 0.80       | 68.10 $\pm$ 0.80 | -                            | -                            | -                | Omniglot, Tiered-ImageNet                           |
| Attentive Prototype [46]                  | w.r.t. prototypes              | Euclidean distance  | episodic training                         | -                      | -                | 66.43 $\pm$ 0.26 (DeepCaps)  | 82.13 $\pm$ 0.21             | 82.13 $\pm$ 0.21 | Tiered-ImageNet, FC100                              |
| CTX [75]                                  | w.r.t. prototypes              | Euclidean distance  | episodic training                         | -                      | -                | -                            | -                            | -                | Meta-Dataset; ResNet-34                             |
| RapNets [76]                              | w.r.t. prototypes              | Euclidean distance  | episodic training (higher-way training)   | -                      | 70.89 $\pm$ 0.64 | -                            | -                            | -                | Omniglot, FC100, CUB-200-2011                       |

Table 2 (cont.)

| Method                                     | Classification mechanism | Similarity measure  | Training strategies   | Mini-ImageNet (Conv-4) | 1-shot           | 5-shot           | 1-shot            | 5-shot | ResNet-12) Additional architectures or datasets                         |
|--|--------------------------|---|---|------------------------|------------------|------------------|-------------------|--------|---|
| CCVD [77]                                  | w.r.t. prototypes        | Euclidean distance  | episodic training   | 48.47 $\pm$ 0.86       | 65.86 $\pm$ 0.73 | 69.48 $\pm$ 0.45 | 86.75 $\pm$ 0.28  | (WRN)  | Tiered-ImageNet, CUB-200-2011; MobileNet, ResNet-10/18/34, DenseNet-121 |
| FSL with embedded class models [78]        | w.r.t. prototypes        | Euclidean distance with learned projection matrix               | episodic training   | 49.07 $\pm$ 0.43       | 65.73 $\pm$ 0.36 | 59.00            | 77.46             |        | Tiered-ImageNet, CIFAR-FS   |
| LMPNet [79]                                | w.r.t. prototypes        | cosine similarity   | episodic training   | 49.87 $\pm$ 0.20       | 68.81 $\pm$ 0.16 | 62.74 $\pm$ 0.11 | 80.23 $\pm$ 0.52  |        | Tiered-ImageNet, CUB-200-2010, Stanford Dogs, Stanford Cars             |
| Two-Stage FSL [80]                         | w.r.t. prototypes        | Mahalanobis distance  | episodic training (higher-way training in the first stage)    | 52.68 $\pm$ 0.51       | 70.91 $\pm$ 0.85 | -                | -                 |        | Omniglot, CIFAR-FS, CUB-200-2011  |
| TapNet [81]                                | w.r.t. prototypes        | Mahalanobis distance  | episodic training (higher-way training)                       | 50.68 $\pm$ 0.11       | 69.00 $\pm$ 0.09 | 61.65 $\pm$ 0.15 | 76.36 $\pm$ 0.10  |        | Omniglot, Tiered-ImageNet   |
| FSL with global class representations [82] | w.r.t. prototypes        | Euclidean distance  | pre-train + episodic training                                 | 53.21 $\pm$ 0.40       | 72.34 $\pm$ 0.32 | -                | -                 |        | Omniglot  |
| KGTN [83]                                  | w.r.t. prototypes        | dot product, cosine similarity, Pearson correlation coefficient | pre-train + minibatch training                                | -                      | -                | -                | -                 |        | ImageNet-FS, ImageNet-6K; ResNet-50                                     |
| SOC [85]                                   | w.r.t. prototypes        | cosine similarity   | pre-train + fine-tune/episodic training of feature embeddings | -                      | -                | 69.28 $\pm$ 0.49 | 85.16 $\pm$ 0.42  |        | Tiered-ImageNet   |
| PHR [86]                                   | w.r.t. prototypes        | Euclidean distance  | pre-train + episodic training                                 | 65.10 $\pm$ 0.70       | 78.10 $\pm$ 0.40 | 74.90 $\pm$ 0.60 | 84.50 $\pm$ 0.30  |        | CIFAR-FS, FC100, CUB-200-2011; ResNet-18                                |
| ABPML [87]                                 | w.r.t. prototypes        | probability from Gaussian distribution                          | episodic training   | 53.28 $\pm$ 0.91       | 70.44 $\pm$ 0.72 | -                | -                 |        | Omniglot, CUB-200-2011, Stanford Dogs                                   |
| Variational FSL [88]                       | w.r.t. distributions     | probability from Gaussian distribution                          | pre-train + episodic training                                 | 57.15 $\pm$ 0.31       | 71.54 $\pm$ 0.23 | 61.23 $\pm$ 0.26 | 77.69 $\pm$ 0.17  |        | Omniglot; cluttered Omniglot (for segmentation)                         |
| DSN [89]                                   | w.r.t. subspaces         | Euclidean distance  | episodic training   | 55.88 $\pm$ 0.90       | 70.50 $\pm$ 0.68 | 64.60 $\pm$ 0.72 | 79.51 $\pm$ 0.50  |        | Tiered-ImageNet, CIFAR-FS, Open MIC                                     |
| Temperature Network [91]                   | w.r.t. prototypes        | Euclidean distance  | episode training  | 52.39                  | 67.89            | -                | -                 |        | Stanford Dogs, Stanford Cars, Dermnet skin disease                      |
| Variational scaling [92]                   | w.r.t. prototypes        | Euclidean distance or cosine similarity                         | episodic training   | 49.34 $\pm$ 0.29       | 67.83 $\pm$ 0.16 | 56.09 $\pm$ 0.19 | 74.46 $\pm$ 0.17  |        |   |
| Simple CNAPS [34]                          | w.r.t. prototypes        | Mahalanobis distance  | pre-train   | -                      | -                | 82.16            | 89.80 (ResNet-18) |        | Tiered-ImageNet, Meta-Dataset   |
| TEAM [93]                                  | w.r.t. prototypes        | Mahalanobis distance  | pre-train (for Conv-4 only) + episodic training               | 56.57                  | 72.04 (T)        | 60.07            | 75.90 (T)         |        | CIFAR-FS, CUB-200-2011  |
| SEN [94]                                   | w.r.t. prototypes        | SEN dissimilarity measure                                       | episodic training   | -                      | 69.80            | -                | 72.3 (WRN-16-6)   |        | Omniglot, FC100   |
| DeepEMD [96]                               | w.r.t. prototypes        | Earth Mover's Distance  | pre-train + episodic training                                 | -                      | -                | 65.91 $\pm$ 0.82 | 82.41 $\pm$ 0.56  |        | Tiered-ImageNet, FC100, CUB-200-2011                                    |
| DeepBDC [97]                               | w.r.t. prototypes        | Brownian Distance Covariance metric                             | pre-train / episodic training                                 | -                      | -                | 67.83 $\pm$ 0.43 | 85.45 $\pm$ 0.29  |        | Tiered-ImageNet, CUB-200-2011   |
| Relation Network [12]                      | w.r.t. prototypes        | learned distance  | episodic training   | 50.44 $\pm$ 0.82       | 65.32 $\pm$ 0.70 | -                | -                 |        | Omniglot; AwA, CUB-200-2011 (for ZSL)                                   |
| ArL [99]                                   | w.r.t. prototypes        | learned distance  | episodic training   | 57.48 $\pm$ 0.65       | 72.64 $\pm$ 0.45 | 65.21 $\pm$ 0.58 | 80.41 $\pm$ 0.49  |        | CUB-200-2011, Flowers   |
| SAML [100]                                 | w.r.t. prototypes        | learned distance  | episodic training   | 57.69 $\pm$ 0.2        | 73.03 $\pm$ 0.16 | -                | -                 |        | CUB-200-2011  |

Table 2 (cont.)

| Method        | Classification mechanism | Similarity measure | Training strategies | Mini-ImageNet (Conv-4)  | Mini-ImageNet (ResNet-12) | Additional architectures or datasets                    |
|---------------|--------------------------|--------------------|---------------------|---|---------------------------|---|
| GNN-FSL [102] | w.r.t. instances         | learned distance   | episodic training   | 1-shot 50.33 $\pm$ 0.36 5-shot 66.41 $\pm$ 0.63                             | 1-shot - 5-shot -         | Omniglot  |
| EGNN [101]    | w.r.t. instances         | learned distance   | episodic training   | - 76.37 (T)   | - -                       | Tiered-ImageNet   |
| TRPN [103]    | w.r.t. instances         | learned distance   | episodic training   | 57.84 $\pm$ 0.51 78.57 $\pm$ 0.44 (T) 68.25 $\pm$ 0.50 85.40 $\pm$ 0.39 (T) | (WRN)                     | Tiered-ImageNet   |
| HGNN [104]    | w.r.t. instances         | learned distance   | episodic training   | 60.03 $\pm$ 0.51 79.64 $\pm$ 0.36 (T)                                       |                           | Tiered-ImageNet, CUB-200-2011                           |
| DPGN [105]    | w.r.t. instances         | learned distance   | episodic training   | 66.01 $\pm$ 0.36 82.83 $\pm$ 0.41 (T) 67.77 $\pm$ 0.32 84.60 $\pm$ 0.43 (T) |                           | Tiered-ImageNet, CIFAR-FS, CUB-200-2011; ResNet-18, WRN |

All experimental results are reported for 5-way classification. (T) denotes transductive setting. Unless specified otherwise, Conv-4 uses 4 convolutional layer with 64 filters, and WRN uses 28 convolutional layers with a widening factor of 10.

Table 3: Summary of deep metric learning methods for few-shot fine-grained image classification.

| Method          | Classification mechanism      | Similarity measure                   | Training strategies | CUB-200-2011   | Stanford Dogs                                      | Additional datasets or embedding architectures |
|-----------------|-------------------------------|--------------------------------------|---------------------|--|--|--|
| DN4 [29]        | w.r.t. bags of local features | cosine similarity                    | episodic training   | 1-shot 53.15 $\pm$ 0.84 5-shot 81.90 $\pm$ 0.60 (CUB-2010) | 1-shot 45.73 $\pm$ 0.76 5-shot 66.33 $\pm$ 0.66    | Mini-ImageNet, Stanford Cars                   |
| ATL-Net [48]    | w.r.t. bags of local features | cosine similarity & learned distance | episodic training   | 60.91 $\pm$ 0.91 (CUB-2010)                                | 54.49 $\pm$ 0.92 73.20 $\pm$ 0.69                  | Mini-ImageNet, Stanford Cars                   |
| CovaMNet [73]   | w.r.t. bags of local features | covariance metric                    | episodic training   | 52.42 $\pm$ 0.76 (CUB-2010)                                | 49.10 $\pm$ 0.76 63.04 $\pm$ 0.65                  | Mini-ImageNet, Stanford Cars                   |
| COMET [49]      | w.r.t. prototypes             | Euclidean distance                   | episodic training   | 67.9 $\pm$ 0.9 85.3 $\pm$ 0.5                              | - -  | Flowers; Conv-6                                |
| LRPABN [47]     | w.r.t. prototypes             | learned distance                     | episodic training   | 63.63 $\pm$ 0.77 (120/30/50)                               | 45.72 $\pm$ 0.75 60.94 $\pm$ 0.66                  | Stanford Cars                                  |
| TOAN [51]       | w.r.t. prototypes             | learned distance                     | episodic training   | 65.34 $\pm$ 0.75 (120/30/50)                               | 51.83 $\pm$ 0.80 69.83 $\pm$ 0.66                  | Stanford Cars; ResNet-12                       |
| Auto-ACNet [71] | w.r.t. instances              | learned distance                     | minibatch training  | 57.93 $\pm$ 0.54 (80/40/80; ACNet)                         | - -  | AwA2   |
| BSNet [98]      | w.r.t. prototypes             | cosine similarity + learned distance | episodic training   | 62.84 $\pm$ 0.95   | 85.39 $\pm$ 0.56 43.42 $\pm$ 0.86 71.90 $\pm$ 0.68 | Stanford Cars                                  |

All experimental results are reported for 5-way classification with Conv-4-64 as the embedding architecture. Unless specified otherwise, CUB-200-2011 [28] is split into 100/50 training/validation/test sets. CUB-2010 refers to CUB-200-2010 [27] with the split of 120/30/50. For Stanford Dogs, the dataset is split into 70/20/30 training/validation/test sets.

761 *2. Enhancing stability to support samples and robustness to adversarial per-*  
762 *turbations and distribution shifts.* Despite the continuous improvement in  
763 classification accuracy, few-shot classification methods are vulnerable in var-  
764 ious scenarios, hindering their usage in safety-critical applications such as  
765 medical image analysis. Prior works show that existing methods are non-  
766 robust to input or label outliers [76], adversarial perturbations (i.e., small, vi-  
767 sually imperceptible changes of data that fool the classifier to make incorrect  
768 predictions) added to support [106] or query images [107], and distribution  
769 shift between support and query datasets [108]. In [109], it is demonstrated  
770 that even non-perturbed and in-distribution support images can significantly  
771 deteriorate the classification accuracy of several popular methods. Further  
772 exploration of vulnerability in existing approaches and design of robust and  
773 stable models will be very valuable.

774 *3. Rethinking the use of episodic training strategy.* While episodic training is  
775 a common practice to train metric learning methods in the few-shot learning  
776 setting, it is rigid to require each training episode to have the same number of  
777 classes and images as the evaluation episode; in fact, [72] observed the benefit  
778 of training with a larger number of classes. Moreover, the model gets updated  
779 after receiving an episode without regard to its quality and thus is prone to  
780 poorly sampled images like outliers. [110] is the first attempt to alleviate  
781 this problem by exploiting the relationship between episodes; more solutions  
782 are needed to identify episodes that are high-quality and useful to the novel  
783 task. Furthermore, we notice that episodic training can result in models  
784 that underfit the base dataset. One possible reason is that, by using episodic  
785 training, methods adopt continual learning on plenty of tasks sampled from

786 the base dataset and suffer from catastrophic forgetting [111, 112], i.e., the  
787 model learned from previous tasks is supplanted after learning on a new task.  
788 Therefore, how to avoid this problem and enhance the model fitting ability of  
789 metric learning methods on both base and novel datasets remains a challenge.

790 *4. Developing metric learning methods for cross-domain few-shot classifi-*  
791 *cation.* While base and novel datasets may come from different domains in  
792 practice, currently only few works focus on cross-domain few-shot classi-  
793 fication. More recently and severely, [113] reported that all meta-trained  
794 methods, including the reviewed work [41], are outperformed by the simple  
795 transductive fine-tuning in the presence of a large domain shift, specifically,  
796 when training on natural images and evaluating beyond them, such as on  
797 agriculture and satellite images. The difficulty is that the base data and the  
798 novel data usually have different metric spaces. Therefore, how to alleviate  
799 domain shift between the training and evaluation phases needs to be explored  
800 in the future.

#### 801 *4.2. Applications*

802 The superior performance of deep metric learning methods for few-shot  
803 image classification motivates researchers to extend these methods to non-  
804 natural images from various disciplines. For example, the methods have been  
805 developed for diagnosing and classifying diseases based on dermoscopic [114]  
806 images and computerised tomography (CT) images [115], classifying plant  
807 diseases based on leaf images collected in the field [116], scene classification  
808 in aerial images [117] and remote sensing images [118], and hyperspectral  
809 image classification [119, 120].



810 Deep metric learning has also been applied beyond image classification, to  
811 more challenging computer vision applications. A notable example is person  
812 re-identification (Re-ID), whose aim is to retrieve a person of interest across  
813 multiple non-overlapping cameras [121, 122]. Metric learning is particularly  
814 effective for Re-ID, as this is an open-set classification task with different  
815 people in the training and test classes and often there is only one image  
816 available for the query person [123]. Metric learning also shows impressive  
817 results on face recognition, in both closed-set [124] and open-set [125] set-  
818 tings, and content-based image retrieval [126, 127], which can be formulated  
819 as a ranking problem.

## 820 5. Conclusions

821 This paper presents a review of recent few-shot deep metric learning meth-  
822 ods. After providing the definitions and a general evaluation framework for  
823 few-shot learning and expounding on the widely used datasets and their set-  
824 tings, we review the novelty and limitations of existing methods. In partic-  
825 ular, there is a pattern of progressing towards learning task-specific feature  
826 embeddings, task-dependent prototypes, and more flexible similarity mea-  
827 sures. In addition, we list applications where few-shot deep metric learning  
828 prevails and suggest future research on improving feature generalizability,  
829 method robustness, training strategy, and applicability to cross-domain set-  
830 tings.

## 831 Acknowledgements

832 This work was supported in part by the Beijing Natural Science Foun-  
833 dation under Grant Z200002, the Royal Society under International Ex-  
834 changes Award IEC\NSFC\201071, the National Natural Science Foundation  
835 of China (NSFC) under Grant 62111530146, 62176110, 61906080, 61922015,  
836 U19B2036, 62225601, and Young Doctoral Fund of Education Department  
837 of Gansu Province under Grant 2021QB-038, and Hong-liu Distinguished  
838 Young Talents Foundation of Lanzhou University of Technology.

## 839 References

- 840 [1] A. Krizhevsky, I. Sutskever, G. E. Hinton, ImageNet classification with  
841 deep convolutional neural networks, in: Advances in Neural Informa-  
842 tion Processing Systems, 2012, pp. 1097–1105.
- 843 [2] K. Simonyan, A. Zisserman, Very deep convolutional networks for  
844 large-scale image recognition, arXiv preprint arXiv:1409.1556 (2014).
- 845 [3] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Er-  
846 han, V. Vanhoucke, A. Rabinovich, Going deeper with convolutions,  
847 in: IEEE Conference on Computer Vision and Pattern Recognition,  
848 2015, pp. 1–9.
- 849 [4] J. Gu, Z. Wang, J. Kuen, L. Ma, A. Shahroudy, B. Shuai, T. Liu,  
850 X. Wang, G. Wang, Recent advances in convolutional neural networks,  
851 arXiv preprint arXiv:1512.07108 (2015).

- 852 [5] F.-F. Li, R. Fergus, P. Perona, One-shot learning of object categories,  
853 IEEE Transactions on Pattern Analysis and Machine Intelligence 28 (4)  
854 (2006) 594–611.
- 855 [6] G. Koch, R. Zemel, R. Salakhutdinov, Siamese neural networks for  
856 one-shot image recognition, in: International Conference on Machine  
857 Learning deep learning workshop, Vol. 2, 2015.
- 858 [7] O. Vinyals, C. Blundell, T. Lillicrap, D. Wierstra, et al., Matching  
859 networks for one shot learning, in: Advances in Neural Information  
860 Processing Systems, 2016, pp. 3630–3638.
- 861 [8] A. Santoro, S. Bartunov, M. Botvinick, D. Wierstra, T. Lillicrap, One-  
862 shot learning with memory-augmented neural networks. arxiv preprint,  
863 arXiv preprint arXiv:1605.06065 (2016).
- 864 [9] C. Finn, P. Abbeel, S. Levine, Model-agnostic meta-learning for fast  
865 adaptation of deep networks, in: International Conference on Machine  
866 Learning, JMLR. org, 2017, pp. 1126–1135.
- 867 [10] M. Rohrbach, S. Ebert, B. Schiele, Transfer learning in a transductive  
868 setting, in: Advances in Neural Information Processing Systems, 2013,  
869 pp. 46–54.
- 870 [11] Q. Sun, Y. Liu, T.-S. Chua, B. Schiele, Meta-transfer learning for few-  
871 shot learning, in: IEEE Conference on Computer Vision and Pattern  
872 Recognition, 2019, pp. 403–412.
- 873 [12] F. Sung, Y. Yang, L. Zhang, T. Xiang, P. H. Torr, T. M. Hospedales,  
874 Learning to compare: Relation network for few-shot learning, in: IEEE

- 875 Conference on Computer Vision and Pattern Recognition, 2018, pp.  
876 1199–1208.
- 877 [13] Z. Liu, Z. Miao, X. Zhan, J. Wang, B. Gong, S. X. Yu, Large-scale  
878 long-tailed recognition in an open world, in: IEEE/CVF Conference  
879 on Computer Vision and Pattern Recognition, 2019, pp. 2537–2546.
- 880 [14] J. Shu, Z. Xu, D. Meng, Small sample learning in big data era, arXiv  
881 preprint arXiv:1808.04572 (2018).
- 882 [15] Y. Wang, Q. Yao, J. T. Kwok, L. M. Ni, Generalizing from a few  
883 examples: A survey on few-shot learning, ACM Computing Surveys  
884 53 (3) (2020) 1–34.
- 885 [16] J. Lu, P. Gong, J. Ye, C. Zhang, Learning from very few samples: A  
886 survey, arXiv preprint arXiv:2009.02653 (2020).
- 887 [17] X. Li, Z. Sun, J.-H. Xue, Z. Ma, A concise review of recent few-shot  
888 meta-learning methods, Neurocomputing 456 (2021) 463–468.
- 889 [18] S. J. Pan, Q. Yang, A survey on transfer learning, IEEE Transactions  
890 on Knowledge and Data Engineering 22 (10) (2009) 1345–1359.
- 891 [19] B. Lake, R. Salakhutdinov, J. Gross, J. Tenenbaum, One shot learning  
892 of simple visual concepts, in: Proceedings of the annual meeting of the  
893 cognitive science society, Vol. 33, 2011.
- 894 [20] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma,  
895 Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al., ImageNet

- 896 large scale visual recognition challenge, *International Journal of Com-*  
897 *puter Vision* 115 (3) (2015) 211–252.
- 898 [21] S. Ravi, H. Larochelle, Optimization as a model for few-shot learning,  
899 *International Conference on Learning Representations* (2017).
- 900 [22] M. Ren, E. Triantafillou, S. Ravi, J. Snell, K. Swersky, J. B. Tenen-  
901 baum, H. Larochelle, R. S. Zemel, Meta-learning for semi-supervised  
902 few-shot classification, *International Conference on Learning Represen-*  
903 *tations* (2018).
- 904 [23] A. Krizhevsky, Learning multiple layers of features from tiny images,  
905 *University of Toronto* (2009).
- 906 [24] L. Bertinetto, J. F. Henriques, P. H. Torr, A. Vedaldi, Meta-learning  
907 with differentiable closed-form solvers, *International Conference on*  
908 *Learning Representations* (2019).
- 909 [25] B. Oreshkin, P. R. López, A. Lacoste, TADAM: Task dependent adap-  
910 tive metric for improved few-shot learning, in: *Advances in Neural*  
911 *Information Processing Systems*, 2018, pp. 721–731.
- 912 [26] A. Khosla, N. Jayadevaprakash, B. Yao, F.-F. Li, Novel dataset for  
913 fine-grained image categorization: Stanford dogs, in: *CVPR Workshop*  
914 *on Fine-Grained Visual Categorization*, Vol. 2, 2011.
- 915 [27] P. Welinder, S. Branson, T. Mita, C. Wah, F. Schroff, S. Belongie,  
916 P. Perona, *Caltech-UCSD birds 200* (2010).

- [28] C. Wah, S. Branson, P. Welinder, P. Perona, S. Belongie, The Caltech-UCSD birds-200-2011 dataset (2011).
- [29] W. Li, L. Wang, J. Xu, J. Huo, Y. Gao, J. Luo, Revisiting local descriptor based image-to-class measure for few-shot learning, in: IEEE Conference on Computer Vision and Pattern Recognition, 2019.
- [30] W.-Y. Chen, Y.-C. Liu, Z. Kira, Y.-C. F. Wang, J.-B. Huang, A closer look at few-shot classification, in: International Conference on Learning Representations, 2019.
- [31] E. Triantafillou, T. Zhu, V. Dumoulin, P. Lamblin, U. Evci, K. Xu, R. Goroshin, C. Gelada, K. Swersky, P. Manzagol, H. Larochelle, Meta-dataset: A dataset of datasets for learning to learn from few examples, in: International Conference on Learning Representations, 2020.
- [32] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, C. L. Zitnick, Microsoft COCO: common objects in context, in: European conference on computer vision, Springer, 2014, pp. 740–755.
- [33] J. Requeima, J. Gordon, J. Bronskill, S. Nowozin, R. E. Turner, Fast and flexible multi-task classification using conditional neural adaptive processes, in: Advances in Neural Information Processing Systems, 2019.
- [34] P. Bateni, R. Goyal, V. Masrani, F. Wood, L. Sigal, Improved few-shot visual classification, in: IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2020, pp. 14493–14502.

- [35] W. Li, X. Liu, H. Bilen, Cross-domain few-shot learning with task-specific adapters, in: IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2022.
- [36] E. Xing, M. Jordan, S. J. Russell, A. Ng, Distance metric learning with application to clustering with side-information, *Advances in Neural Information Processing Systems* 15 (2002) 521–528.
- [37] J. Bromley, J. W. Bentz, L. Bottou, I. Guyon, Y. LeCun, C. Moore, E. Säckinger, R. Shah, Signature verification using a “siamese” time delay neural network, *International Journal of Pattern Recognition and Artificial Intelligence* 7 (04) (1993) 669–688.
- [38] H. Li, W. Dong, X. Mei, C. Ma, F. Huang, B.-G. Hu, LGM-Net: Learning to generate matching networks for few-shot learning, in: *International Conference on Machine Learning*, 2019, pp. 3825–3834.
- [39] H.-J. Ye, H. Hu, D.-C. Zhan, F. Sha, Few-shot learning via embedding adaptation with set-to-set functions, in: *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 8808–8817.
- [40] Y. Liu, J. Lee, M. Park, S. Kim, E. Yang, S. J. Hwang, Y. Yang, Learning to propagate labels: Transductive propagation network for few-shot learning, in: *International Conference on Learning Representations*, 2019.
- [41] H.-Y. Tseng, H.-Y. Lee, J.-B. Huang, M.-H. Yang, Cross-domain few-shot classification via learned feature-wise transformation, in: *International Conference on Learning Representations*, 2020.

- 963 [42] W.-H. Li, X. Liu, H. Bilen, Universal representation learning from mul-  
 964 tiple domains for few-shot classification, in: IEEE/CVF International  
 965 Conference on Computer Vision, 2021, pp. 9526–9535.
- 966 [43] W. Jiang, K. Huang, J. Geng, X. Deng, Multi-scale metric learning  
 967 for few-shot learning, IEEE Transactions on Circuits and Systems for  
 968 Video Technology (2020).
- 969 [44] Z. Wu, Y. Li, L. Guo, K. Jia, PARN: Position-aware relation networks  
 970 for few-shot learning, in: IEEE International Conference on Computer  
 971 Vision, 2019.
- 972 [45] W. Xu, Y. Xu, H. Wang, Z. Tu, Attentional constellation nets for  
 973 few-shot learning, in: International Conference on Learning Represen-  
 974 tations, 2021.
- 975 [46] F. Wu, J. S. Smith, W. Lu, C. Pang, B. Zhang, Attentive prototype  
 976 few-shot learning with capsule network-based embedding, in: European  
 977 Conference on Computer Vision, Springer, 2020, pp. 237–253.
- 978 [47] H. Huang, J. Zhang, J. Zhang, J. Xu, Q. Wu, Low-rank pairwise align-  
 979 ment bilinear network for few-shot fine-grained image classification,  
 980 IEEE Transactions on Multimedia (2020).
- 981 [48] C. Dong, W. Li, J. Huo, Z. Gu, Y. Gao, Learning task-aware local rep-  
 982 resentations for few-shot learning, in: International Joint Conference  
 983 on Artificial Intelligence, 2020.
- 984 [49] K. Cao, M. Brbic, J. Leskovec, Concept learners for few-shot learning,  
 985 in: International Conference on Learning Representations, 2021.



- 986 [50] X.-S. Wei, P. Wang, L. Liu, C. Shen, J. Wu, Piecewise classifier map-  
 987 pings: Learning fine-grained learners for novel categories with few ex-  
 988 amples, *IEEE Transactions on Image Processing* 28 (12) (2019) 6116–  
 989 6125.
- 990 [51] H. Huang, J. Zhang, L. Yu, J. Zhang, Q. Wu, C. Xu, TOAN: Target-  
 991 oriented alignment network for fine-grained image categorization with  
 992 few labeled samples, *IEEE Transactions on Circuits and Systems for*  
 993 *Video Technology* (2021).
- 994 [52] J. Wu, T. Zhang, Y. Zhang, F. Wu, Task-aware part mining network for  
 995 few-shot learning, in: *IEEE/CVF International Conference on Com-*  
 996 *puter Vision*, 2021, pp. 8433–8442.
- 997 [53] H. Li, D. Eigen, S. Dodge, M. Zeiler, X. Wang, Finding task-relevant  
 998 features for few-shot learning by category traversal, in: *IEEE Confer-*  
 999 *ence on Computer Vision and Pattern Recognition*, 2019, pp. 1–10.
- 1000 [54] S. W. Yoon, D.-Y. Kim, J. Seo, J. Moon, XtarNet: Learning to extract  
 1001 task-adaptive representation for incremental few-shot learning, in: *In-*  
 1002 *ternational Conference on Machine Learning*, 2020, pp. 10852–10860.
- 1003 [55] S. Rahman, S. Khan, F. Porikli, A unified approach for conventional  
 1004 zero-shot, generalized zero-shot, and few-shot learning, *IEEE Transac-*  
 1005 *tions on Image Processing* 27 (11) (2018) 5652–5667.
- 1006 [56] T. D. Kulkarni, W. F. Whitney, P. Kohli, J. Tenenbaum, Deep convo-  
 1007 lutional inverse graphics network, in: *Advances in Neural Information*  
 1008 *Processing Systems*, 2015, pp. 2539–2547.

- [57] A. J. Ratner, H. Ehrenberg, Z. Hussain, J. Dunnmon, C. Ré, Learning to compose domain-specific transformations for data augmentation, in: Advances in Neural Information Processing Systems, 2017, pp. 3236–3246.
- [58] L. Perez, J. Wang, The effectiveness of data augmentation in image classification using deep learning, arXiv preprint arXiv:1712.04621 (2017).
- [59] A. Shrivastava, T. Pfister, O. Tuzel, J. Susskind, W. Wang, R. Webb, Learning from simulated and unsupervised images through adversarial training, in: IEEE Conference on Computer Vision and Pattern Recognition, 2017, pp. 2242–2251.
- [60] A. Antoniou, A. Storkey, H. Edwards, Data augmentation generative adversarial networks, International Conference on Learning Representations Workshop (2018).
- [61] A. J. Ratner, C. M. De Sa, S. Wu, D. Selsam, C. Ré, Data programming: Creating large training sets, quickly, in: Advances in Neural Information Processing Systems, 2016, pp. 3567–3575.
- [62] Y.-X. Wang, R. Girshick, M. Hebert, B. Hariharan, Low-shot learning from imaginary data, in: IEEE Conference on Computer Vision and Pattern Recognition, 2018, pp. 7278–7286.
- [63] H. Zhang, J. Zhang, P. Koniusz, Few-shot learning via saliency-guided hallucination of samples, in: IEEE Conference on Computer Vision and Pattern Recognition, 2019, pp. 2770–2779.

- 1032 [64] J. Guan, Z. Lu, T. Xiang, A. Li, A. Zhao, J.-R. Wen, Zero and few  
1033 shot learning with semantic feature synthesis and competitive learning,  
1034 IEEE Transactions on Pattern Analysis and Machine Intelligence 43 (7)  
1035 (2020) 2510–2523.
- 1036 [65] S. Gidaris, A. Bursuc, N. Komodakis, P. Perez, M. Cord, Boosting  
1037 few-shot visual learning with self-supervision, in: IEEE International  
1038 Conference on Computer Vision, 2019.
- 1039 [66] C. Liu, Y. Fu, C. Xu, S. Yang, J. Li, C. Wang, L. Zhang, Learning a few-  
1040 shot embedding model with contrastive learning, in: AAAI Conference  
1041 on Artificial Intelligence, Vol. 35, 2021, pp. 8635–8643.
- 1042 [67] Y. Ouali, C. Hudelot, M. Tami, Spatial contrastive learning for few-  
1043 shot classification, in: Joint European Conference on Machine Learning  
1044 and Knowledge Discovery in Databases, Springer, 2021, pp. 671–686.
- 1045 [68] Z. Yang, J. Wang, Y. Zhu, Few-shot classification with contrastive  
1046 learning, in: European Conference on Computer Vision, Springer, 2022,  
1047 pp. 293–309.
- 1048 [69] T. Xiao, X. Wang, A. A. Efros, T. Darrell, What should not be con-  
1049 trastive in contrastive learning, in: International Conference on Learn-  
1050 ing Representations, 2021.
- 1051 [70] Y. Zhu, W. Min, S. Jiang, Attribute-guided feature learning for few-  
1052 shot image recognition, IEEE Transactions on Multimedia 23 (2020)  
1053 1200–1209.

- 1054 [71] L. Zhang, S. Wang, X. Chang, J. Liu, Z. Ge, Q. Zheng, Auto-FSL:  
1055 Searching the attribute consistent network for few-shot learning, IEEE  
1056 Transactions on Circuits and Systems for Video Technology (2021).
- 1057 [72] J. Snell, K. Swersky, R. Zemel, Prototypical networks for few-shot  
1058 learning, in: Advances in Neural Information Processing Systems, 2017,  
1059 pp. 4077–4087.
- 1060 [73] W. Li, J. Xu, J. Huo, L. Wang, Y. Gao, J. Luo, Distribution consis-  
1061 tency based covariance metric networks for few-shot learning, in: AAAI  
1062 Conference on Artificial Intelligence, Vol. 33, 2019, pp. 8642–8649.
- 1063 [74] K. Allen, E. Shelhamer, H. Shin, J. Tenenbaum, Infinite mixture pro-  
1064 totypes for few-shot learning, in: International Conference on Machine  
1065 Learning, 2019, pp. 232–241.
- 1066 [75] C. Doersch, A. Gupta, A. Zisserman, CrossTransformers: spatially-  
1067 aware few-shot transfer, in: Advances in Neural Information Processing  
1068 Systems, 2020.
- 1069 [76] J. Lu, S. Jin, J. Liang, C. Zhang, Robust few-shot learning for user-  
1070 provided data, IEEE Transactions on Neural Networks and Learning  
1071 Systems 32 (4) (2020) 1433–1447.
- 1072 [77] C. Ma, Z. Huang, M. Gao, J. Xu, Few-shot learning via dirichlet tessel-  
1073 ation ensemble, in: International Conference on Learning Representa-  
1074 tions, 2022.

- 1075 [78] A. Ravichandran, R. Bhotika, S. Soatto, Few-shot learning with embed-  
1076 ded class models and shot-free meta training, in: IEEE International  
1077 Conference on Computer Vision, 2019.
- 1078 [79] H. Huang, Z. Wu, W. Li, J. Huo, Y. Gao, Local descriptor-based  
1079 multi-prototype network for few-shot learning, Pattern Recognition 116  
1080 (2021) 107935.
- 1081 [80] D. Das, C. G. Lee, A two-stage approach to few-shot learning for image  
1082 recognition, IEEE Transactions on Image Processing 29 (2020) 3336–  
1083 3350.
- 1084 [81] S. W. Yoon, J. Seo, J. Moon, TapNet: Neural network augmented  
1085 with task-adaptive projection for few-shot learning, in: International  
1086 Conference on Machine Learning, 2019, pp. 7115–7123.
- 1087 [82] A. Li, T. Luo, T. Xiang, W. Huang, L. Wang, Few-shot learning  
1088 with global class representations, in: IEEE International Conference  
1089 on Computer Vision, 2019.
- 1090 [83] R. Chen, T. Chen, X. Hui, H. Wu, G. Li, L. Lin, Knowledge graph  
1091 transfer network for few-shot recognition, in: AAAI Conference on  
1092 Artificial Intelligence, Vol. 34, 2020, pp. 10575–10582.
- 1093 [84] T. Chen, L. Lin, X. Hui, R. Chen, H. Wu, Knowledge-guided multi-  
1094 label few-shot learning for general image recognition, IEEE Transac-  
1095 tions on Pattern Analysis and Machine Intelligence (2020).
- 1096 [85] X. Luo, L. Wei, L. Wen, J. Yang, L. Xie, Z. Xu, Q. Tian, Rectifying the

- 1097 shortcut learning of background for few-shot learning, in: Advances in  
1098 Neural Information Processing Systems, 2021, pp. 13073–13085.
- 1099 [86] Y. Zhou, Y. Guo, S. Hao, R. Hong, Hierarchical prototype refinement  
1100 with progressive inter-categorical discrimination maximization for few-  
1101 shot learning, IEEE Transactions on Image Processing (2022).
- 1102 [87] Z. Sun, J. Wu, X. Li, W. Yang, J.-H. Xue, Amortized bayesian pro-  
1103 totype meta-learning: A new probabilistic meta-learning approach to  
1104 few-shot image classification, in: International Conference on Artificial  
1105 Intelligence and Statistics, 2021, pp. 1414–1422.
- 1106 [88] J. Zhang, C. Zhao, B. Ni, M. Xu, X. Yang, Variational few-shot learn-  
1107 ing, in: IEEE International Conference on Computer Vision, 2019.
- 1108 [89] C. Simon, P. Koniusz, R. Nock, M. Harandi, Adaptive subspaces for  
1109 few-shot learning, in: IEEE/CVF Conference on Computer Vision and  
1110 Pattern Recognition, 2020, pp. 4136–4145.
- 1111 [90] B. Liu, Y. Cao, Y. Lin, Q. Li, Z. Zhang, M. Long, H. Hu, Negative  
1112 margin matters: Understanding margin in few-shot classification, in:  
1113 European Conference on Computer Vision, Springer, 2020, pp. 438–  
1114 455.
- 1115 [91] W. Zhu, W. Li, H. Liao, J. Luo, Temperature network for few-shot  
1116 learning with distribution-aware large-margin metric, Pattern Recog-  
1117 nition 112 (2021) 107797.

- 1118 [92] J. Chen, L.-M. Zhan, X.-M. Wu, F.-l. Chung, Variational metric scal-  
1119 ing for metric-based meta-learning, in: AAAI Conference on Artificial  
1120 Intelligence, Vol. 34, 2020, pp. 3478–3485.
- 1121 [93] L. Qiao, Y. Shi, J. Li, Y. Wang, T. Huang, Y. Tian, Transductive  
1122 episodic-wise adaptive metric for few-shot learning, in: IEEE Interna-  
1123 tional Conference on Computer Vision, 2019.
- 1124 [94] V. N. Nguyen, S. Løkse, K. Wickstrøm, M. Kampffmeyer, D. Roverso,  
1125 R. Jenssen, SEN: A novel feature normalization dissimilarity measure  
1126 for prototypical few-shot learning networks, in: European Conference  
1127 on Computer Vision, Vol. 12368, Springer, 2020, pp. 118–134.
- 1128 [95] Y. Zheng, D. K. Pal, M. Savvides, Ring loss: Convex feature normal-  
1129 ization for face recognition, in: IEEE conference on computer vision  
1130 and pattern recognition, 2018, pp. 5089–5097.
- 1131 [96] C. Zhang, Y. Cai, G. Lin, C. Shen, DeepEMD: Few-shot image classi-  
1132 fication with differentiable earth mover’s distance and structured clas-  
1133 sifiers, in: IEEE/CVF Conference on Computer Vision and Pattern  
1134 Recognition, 2020, pp. 12203–12213.
- 1135 [97] J. Xie, F. Long, J. Lv, Q. Wang, P. Li, Joint distribution matters: Deep  
1136 brownian distance covariance for few-shot classification, in: IEEE/CVF  
1137 Conference on Computer Vision and Pattern Recognition, 2022.
- 1138 [98] X. Li, J. Wu, Z. Sun, Z. Ma, J. Cao, J.-H. Xue, BSNet: Bi-similarity  
1139 network for few-shot fine-grained image classification, IEEE Transac-  
1140 tions on Image Processing 30 (2020) 1318–1331.

- 1141 [99] H. Zhang, P. Koniusz, S. Jian, H. Li, P. H. Torr, Rethinking class rela-  
1142 tions: Absolute-relative supervised and unsupervised few-shot learning,  
1143 in: IEEE/CVF Conference on Computer Vision and Pattern Recogni-  
1144 tion, 2021, pp. 9432–9441.
- 1145 [100] F. Hao, F. He, J. Cheng, L. Wang, J. Cao, D. Tao, Collect and select:  
1146 Semantic alignment metric learning for few-shot learning, in: IEEE  
1147 International Conference on Computer Vision, 2019.
- 1148 [101] J. Kim, T. Kim, S. Kim, C. D. Yoo, Edge-labeling graph neural network  
1149 for few-shot learning, in: IEEE Conference on Computer Vision and  
1150 Pattern Recognition, 2019, pp. 11–20.
- 1151 [102] V. Garcia, J. Bruna, Few-shot learning with graph neural networks,  
1152 International Conference on Learning Representations (2018).
- 1153 [103] Y. Ma, S. Bai, S. An, W. Liu, A. Liu, X. Zhen, X. Liu, Transductive  
1154 relation-propagation network for few-shot learning, in: International  
1155 Joint Conference on Artificial Intelligence, 2020, pp. 804–810.
- 1156 [104] C. Chen, K. Li, W. Wei, J. T. Zhou, Z. Zeng, Hierarchical graph neu-  
1157 ral networks for few-shot learning, IEEE Transactions on Circuits and  
1158 Systems for Video Technology (2021).
- 1159 [105] L. Yang, L. Li, Z. Zhang, X. Zhou, E. Zhou, Y. Liu, DPGN: Distribu-  
1160 tion propagation graph network for few-shot learning, in: IEEE/CVF  
1161 Conference on Computer Vision and Pattern Recognition, 2020, pp.  
1162 13390–13399.



- 1163 [106] E. T. Oldewage, J. F. Bronskill, R. E. Turner, Attacking few-shot clas-  
1164 sifiers with adversarial support poisoning, in: ICML Workshop on Ad-  
1165 versarial Machine Learning, 2021.
- 1166 [107] M. Goldblum, L. Fowl, T. Goldstein, Adversarially robust few-shot  
1167 learning: A meta-learning approach, *Advances in Neural Information*  
1168 *Processing Systems* 33 (2020) 17886–17895.
- 1169 [108] E. Bennequin, V. Bouvier, M. Tami, A. Toubhans, C. Hudelot, Bridg-  
1170 ing few-shot learning and adaptation: new challenges of support-query  
1171 shift, in: *Joint European Conference on Machine Learning and Knowl-*  
1172 *edge Discovery in Databases*, Springer, 2021, pp. 554–569.
- 1173 [109] M. Agarwal, M. Yurochkin, Y. Sun, On sensitivity of meta-learning to  
1174 support data, *Advances in Neural Information Processing Systems* 34  
1175 (2021) 20447–20460.
- 1176 [110] N. Fei, Z. Lu, T. Xiang, S. Huang, MELR: Meta-learning via model-  
1177 ing episode-level relationships for few-shot learning, in: *International*  
1178 *Conference on Learning Representations*, 2021.
- 1179 [111] M. McCloskey, N. J. Cohen, Catastrophic interference in connectionist  
1180 networks: The sequential learning problem, in: *Psychology of Learning*  
1181 *and Motivation*, Vol. 24, Elsevier, 1989, pp. 109–165.
- 1182 [112] S. Gidaris, N. Komodakis, Dynamic few-shot visual learning with-  
1183 out forgetting, in: *IEEE Conference on Computer Vision and Pattern*  
1184 *Recognition*, 2018, pp. 4367–4375.

- 1185 [113] Y. Guo, N. C. Codella, L. Karlinsky, J. V. Codella, J. R. Smith,  
1186 K. Saenko, T. Rosing, R. Feris, A broader study of cross-domain few-  
1187 shot learning, in: European Conference on Computer Vision, Springer,  
1188 2020, pp. 124–141.
- 1189 [114] K. Mahajan, M. Sharma, L. Vig, Meta-dermdiagnosis: Few-shot skin  
1190 disease identification using meta-learning, in: IEEE/CVF Conference  
1191 on Computer Vision and Pattern Recognition Workshops, 2020, pp.  
1192 730–731.
- 1193 [115] X. Chen, L. Yao, T. Zhou, J. Dong, Y. Zhang, Momentum contrastive  
1194 learning for few-shot covid-19 diagnosis from chest ct images, Pattern  
1195 Recognition 113 (2021) 107826.
- 1196 [116] D. Argüeso, A. Picon, U. Irusta, A. Medela, M. G. San-Emeterio,  
1197 A. Bereciartua, A. Alvarez-Gila, Few-shot learning approach for plant  
1198 disease classification using images taken in the field, Computers and  
1199 Electronics in Agriculture 175 (2020) 105542.
- 1200 [117] L. Li, X. Yao, G. Cheng, J. Han, Aifs-dataset for few-shot aerial im-  
1201 age scene classification, IEEE Transactions on Geoscience and Remote  
1202 Sensing 60 (2022) 1–11.
- 1203 [118] X. Li, D. Shi, X. Diao, H. Xu, Scl-mlnet: Boosting few-shot re-  
1204 mote sensing scene classification via self-supervised contrastive learn-  
1205 ing, IEEE Transactions on Geoscience and Remote Sensing 60 (2021)  
1206 1–12.

- 1207 [119] Z. Li, M. Liu, Y. Chen, Y. Xu, W. Li, Q. Du, Deep cross-domain few-  
1208 shot learning for hyperspectral image classification, *IEEE Transactions*  
1209 *on Geoscience and Remote Sensing* 60 (2021) 1–18.
- 1210 [120] Z. Xue, Y. Zhou, P. Du, S3net: Spectral-spatial siamese network for  
1211 few-shot hyperspectral image classification, *IEEE Transactions on Geo-*  
1212 *science and Remote Sensing* (2022).
- 1213 [121] M. Ye, J. Shen, G. Lin, T. Xiang, L. Shao, S. C. Hoi, Deep learning for  
1214 person re-identification: A survey and outlook, *IEEE Transactions on*  
1215 *Pattern Analysis and Machine Intelligence* 44 (6) (2021) 2872–2893.
- 1216 [122] G. Zou, G. Fu, X. Peng, Y. Liu, M. Gao, Z. Liu, Person re-identification  
1217 based on metric learning: A survey, *Multimedia Tools and Applications*  
1218 80 (17) (2021) 26855–26888.
- 1219 [123] W.-S. Zheng, S. Gong, T. Xiang, Towards open-world person re-  
1220 identification by one-shot group-based verification, *IEEE Transactions*  
1221 *on Pattern Analysis and Machine Intelligence* 38 (3) (2015) 591–606.
- 1222 [124] Y. Wu, H. Liu, Y. Fu, Low-shot face recognition with hybrid classifiers,  
1223 in: *IEEE International Conference on Computer Vision Workshops*,  
1224 2017, pp. 1933–1939.
- 1225 [125] H. Du, H. Shi, Y. Liu, J. Wang, Z. Lei, D. Zeng, T. Mei, Semi-siamese  
1226 training for shallow face learning, in: *European Conference on Com-*  
1227 *puter Vision*, Springer, 2020, pp. 36–53.
- 1228 [126] F. Cakir, K. He, X. Xia, B. Kulis, S. Sclaroff, Deep metric learning

- 1229 to rank, in: IEEE/CVF conference on computer vision and pattern  
1230 recognition, 2019, pp. 1861–1870.
- 1231 [127] X. Shen, Y. Xiao, S. X. Hu, O. Sbai, M. Aubry, Re-ranking for image  
1232 retrieval and transductive few-shot classification, Advances in Neural  
1233 Information Processing Systems 34 (2021) 25932–25943.