



# City Research Online

## City St George's, University of London

**Citation:** Li, X., Guo, Z., Zhu, R., Ma, Z., Guo, J. & Xue, J-H. (2024). A simple scheme to amplify inter-class discrepancy for improving few-shot fine-grained image classification. *Pattern Recognition*, 156, 110736. doi: 10.1016/j.patcog.2024.110736

This is the accepted version of the paper.

This version of the publication may differ from the final published version. To cite this item please consult the publisher's version.

**Permanent repository link:** <https://openaccess.city.ac.uk/id/eprint/33245/>

**Link to published version:** <https://doi.org/10.1016/j.patcog.2024.110736>

**Copyright and Reuse:** Copyright and Moral Rights remain with the author(s) and/or copyright holders. Copies of full items can be used for personal research or study, educational, or not-for-profit purposes without prior permission or charge, unless otherwise indicated, provided that the authors, title and full bibliographic details are credited, a hyperlink and/or URL is given for the original metadata page and the content is not changed in any way. For full details of reuse please refer to [City Research Online policy](#).

## Journal Pre-proof

A simple scheme to amplify inter-class discrepancy for improving few-shot fine-grained image classification

Xiaoxu Li, Zijie Guo, Rui Zhu, Zhayu Ma, Jun Guo, Jing-Hao Xue



PII: S0031-3203(24)00487-4  
DOI: <https://doi.org/10.1016/j.patcog.2024.110736>  
Reference: PR 110736

To appear in: *Pattern Recognition*

Received date: 22 October 2023  
Revised date: 12 May 2024  
Accepted date: 26 June 2024

Please cite this article as: X. Li, Z. Guo, R. Zhu et al., A simple scheme to amplify inter-class discrepancy for improving few-shot fine-grained image classification, *Pattern Recognition* (2024), doi: <https://doi.org/10.1016/j.patcog.2024.110736>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2024 The Author(s). Published by Elsevier Ltd. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

# A simple scheme to amplify inter-class discrepancy for improving few-shot fine-grained image classification

Xiaoxu Li<sup>a</sup>, Zijie Guo<sup>a</sup>, Rui Zhu<sup>b,\*</sup>, Zhayu Ma<sup>c</sup>, Jun Guo<sup>c</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>*Faculty of Actuarial Science and Insurance, Bayes Business School, City, University of London, London EC1Y 8TZ, UK*

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

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

---

## Abstract

Few-shot image classification is a challenging topic in pattern recognition and computer vision. Few-shot fine-grained image classification is even more challenging, due to not only the few shots of labelled samples but also the subtle differences to distinguish subcategories in fine-grained images. A recent method called task discrepancy maximisation (TDM) can be embedded into the feature map reconstruction network (FRN) to generate discriminative features, by preserving the appearance details through reconstructing the query image and then assigning higher weights to more discriminative channels, producing the state-of-the-art performance for few-shot fine-grained image classification. However, due to the small inter-class discrepancy in fine-grained images and the small training set in few-shot learn-

---

\*Corresponding author. Tel.: +44(0)1227 827008

*Email addresses:* [lixiaoxu@lut.edu.cn](mailto:lixiaoxu@lut.edu.cn) (Xiaoxu Li), [gzej18801586376@163.com](mailto:gzej18801586376@163.com) (Zijie Guo), [ruizhu@city.ac.uk](mailto:ruizhu@city.ac.uk) (Rui Zhu), [mazhanyu@bupt.edu.cn](mailto:mazhanyu@bupt.edu.cn) (Zhayu Ma), [guojun@bupt.edu.cn](mailto:guojun@bupt.edu.cn) (Jun Guo), [jinghao.xue@ucl.ac.uk](mailto:jinghao.xue@ucl.ac.uk) (Jing-Hao Xue)

ing, the training of FRN+TDM can result in excessively flexible boundaries between subcategories and hence overfitting. To resolve this problem, we propose a simple scheme to amplify inter-class discrepancy and thus improve FRN+TDM. To achieve this aim, instead of developing new modules, our scheme only involves two simple amendments to FRN+TDM: relaxing the inter-class score in TDM, and adding a centre loss to FRN. Extensive experiments on five benchmark datasets showcase that, although embarrassingly simple, our scheme is quite effective to improve the performance of few-shot fine-grained image classification. The code is available at <https://github.com/Airgods/AFRN.git>.

*Keywords:* Few-shot learning, fine-grained image classification, metric-based methods.

---

## 1. Introduction

Few-shot fine-grained image classification is a challenging task that draws wide attention in the pattern recognition and computer vision communities. Although deep neural networks learnt from a large amount of labelled training data can provide impressive image classification performances, few-shot learning that trains a model with little labelled data for each class remains difficult. Moreover, the fine-grained setting brings further challenges, as each class is divided to a large number of subcategories, which makes the inter-class discrepancy even smaller and the classification task much harder.

Metric-based methods are effective for few-shot learning [1]. They aim to learn a metric function to measure the similarities/dissimilarities between different classes and assign the test image to the class with the highest similarity

13 or lowest dissimilarity. For example, the prototypical networks (ProtoNet)  
14 proposed by Snell et al. [2] adopt the average of features of all images from  
15 the same class in the support set as the prototype of that class, and as-  
16 sign the query image to the class with the shortest Euclidean distances from  
17 the class prototypes. Recent works enhance ProtoNet by generating more  
18 representative prototypes [3]. The matching networks (MatchingNet) [4]  
19 utilise a bidirectional LSTM network to map the support set and an at-  
20 tention mechanism-based LSTM to map the query set, and adopt the cosine  
21 similarity as the metric function. In addition to the common metric func-  
22 tions, Zhang et al. [5] propose a new metric function EMD, which assigns  
23 different weights to different positions of the image and calculates the best  
24 matching between the image blocks of the support set and the query set to  
25 represent their similarities. To maintain feature discriminability, Nguyen et  
26 al. [6] propose the square root of the sum of the Euclidean distance and the  
27 norm distance as the metric function. Similarities between images can also  
28 be measured via a properly structured neural network [7].

29 However, when the high similarities between subclasses are not carefully  
30 considered, metric-based methods can fail to classify fine-grained images.  
31 Thus it is crucial to extract features with strong discriminative power to  
32 distinguish the ultra-fine differences between subclasses. Li et al. intro-  
33 duce the bi-similarity network (BSNet) with two similarity metrics to learn  
34 such discriminative features [8]. Huang et al. propose the low-rank pairwise  
35 aligned bilinear network (LRPABN), which utilises bilinear pooling opera-  
36 tions to distinguish support and query images [9]. Huang et al. also propose  
37 the targeted alignment network (TOAN), which can increase the inter-class

38 variation by extracting discriminative fine-grained features while reducing  
 39 intra-class variation by matching support and query features [10].

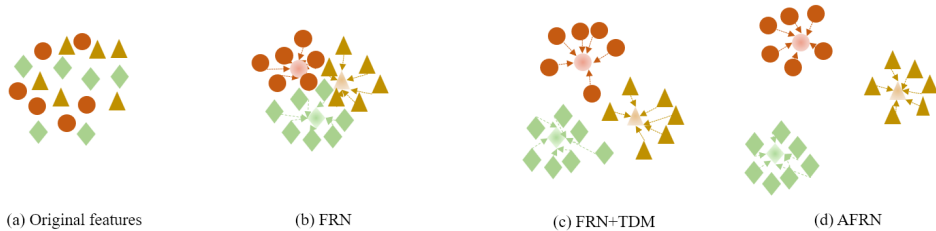


Figure 1: An illustration of the motivation of the adaptive feature map reconstruction network (AFRN). The solid circles, triangles and diamonds represent the instances from three classes, respectively, and the transparent circle, triangle and diamond represent the corresponding prototypes of the three classes, respectively. In (a), we depict a challenging classification task, with severe overlapping between the three classes in the original features space. This challenge is partially resolved by FRN in (b), because the appearance details of images are well preserved by reconstruction, which potentially makes the embedded features more discriminative. In (c), TDM is incorporated to FRN to assign high weights to channels with strong discriminative abilities, and thus the classes become more separable. Finally, in (d), AFRN further improves FRN+TDM by amplifying the inter-class discrepancy, and thus the three classes can be more easily distinguished.

40 There is a problem in many previous metric-based learning algorithms  
 41 that the input to the metric function has to be reshaped to vectors, resulting  
 42 in deficient spatial information. To resolve this problem, Wertheimer et  
 43 al. [11] propose a novel metric-based classification mechanism, feature map  
 44 reconstruction networks (FRN), for few-shot learning. FRN predicts the  
 45 membership of the query image by reconstructing the query feature map  
 46 via the pooled support features of each class. The idea behind FRN is that  
 47 the query feature map is expected to be well reconstructed by the support

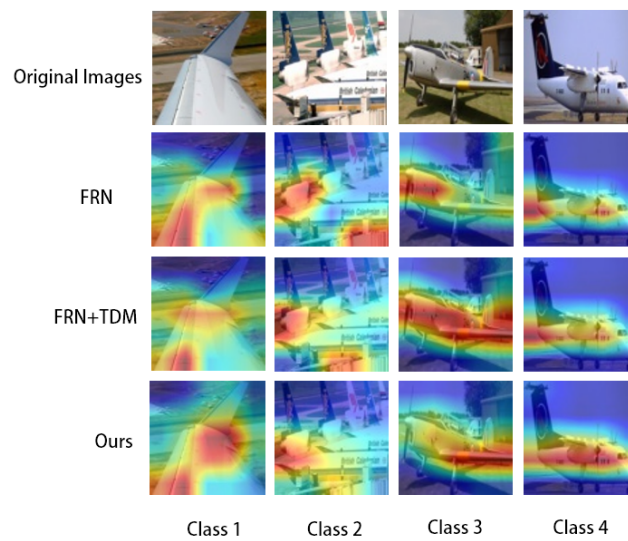


Figure 2: Examples of the features captured by FRN, FRN+TDM and AFRN on four subclasses of airplanes. Apparently, FRN focuses on the objects as well as the nuisance background. Involving TDM in FRN makes the features more discriminative and the focus on background is reduced slightly. In comparison, our AFRN can identify the most discriminative features to distinguish the subclasses with the least focus on the background.

48 features from the correct class with the smallest reconstruction error. Hence,  
49 through the reconstruction process, FRN can well preserve the appearance  
50 details of the images.

51 However, in FRN, all channels are treated equally with the same weights,  
52 without stressing the different importance of different channels. Hence, Lee  
53 et al. [12] propose the task discrepancy maximisation (TDM) module to  
54 identify channels with high discriminative power and assign higher weights  
55 to these channels to improve the classification results of few-shot methods,  
56 such as FRN, for fine-grained images. TDM produces channel weights for  
57 both support and query sets via the support attention module (SAM) and  
58 the query attention module (QAM), respectively. SAM provides class-wise  
59 channel weights to highlight the discriminative channels to distinguish be-  
60 tween classes, while QAM provides object-wise channel weights to focus more  
61 on the object-relevant channels. Lee et al. [12] demonstrate that by incorpo-  
62 rating TDM to FRN, namely FRN+TDM, a state-of-the-art performance of  
63 few-shot fine-grained image classification can be achieved.

64 However, due to the small inter-class discrepancy omnipresent in fine-  
65 grained images and the small training set in the setting of few-shot learning,  
66 FRN+TDM can produce excessively flexible boundaries between subcate-  
67 gories and hence overfitting. To resolve this problem, we propose a simple  
68 scheme to amplify inter-class discrepancy and thus improve FRN+TDM. To  
69 this end, instead of developing new modules to further enhance the extraction  
70 of discriminative features, our scheme only involves two simple amendments  
71 to FRN+TDM: relaxing the inter-class score in TDM, and adding a centre  
72 loss to FRN. We name the network incorporating our scheme to FRN+TDM

73 the adaptive feature map reconstruction network (AFRN).

74 The centre loss [13] aims to achieve intra-class compactness by penalising  
 75 the distance between the learnt features and their corresponding class  
 76 centres, which is vital to distinguish subclasses with high similarity normally  
 77 occurring in fine-grained image classification. In Figure 1, we illustrate  
 78 the motivation of AFRN by a challenging classification of three overlapping  
 79 classes, which is typical in fine-grained image classification with small inter-  
 80 class discrepancy. By involving the centre loss in AFRN, we expect that  
 81 the three classes can be intra-class more compact and thus inter-class more  
 82 separated to make the classification easier. Moreover, in Figure 2, we demon-  
 83 strate one real-data example of the discriminative features extracted by FRN,  
 84 FRN+TDM and AFRN on four subclasses of airplanes. The original FRN  
 85 focuses on the airplanes as well as the nuisance backgrounds; incorporating  
 86 TDM can improve this situation with less focus on the backgrounds; while,  
 87 in comparison, AFRN can identify the most discriminative features with the  
 88 least focus on the backgrounds. For instance, in class 2, the background in  
 89 the lower right corner is least highlighted in our method.

90 More importantly, we observe that FRN+TDM can produce excessively  
 91 flexible boundaries between subcategories and thus overfitting, as the inter-  
 92 class score in TDM to measure the discrepancy between classes is the Eu-  
 93 clidean distance between one class and its *nearest* class. Such an inter-class  
 94 score can result in extremely flexible classification boundaries for fine-grained  
 95 images and thus overfitting to the seen classes in the training set. In few-shot  
 96 fine-grained learning, this problem is severer, because in the test phase, few-  
 97 shot learning aims to classify the novel set with completely different classes

98 from those in the training set. Thus we propose to relax the inter-class score  
 99 in TDM simply to the Euclidean distance between one class and its *furthest*  
 100 class, to mitigate the potential overfitting to a large extent. This amendment  
 101 makes the original TDM module a relaxed TDM module.

102 In summary, the main contributions of our work are as follows.

- 103 • We propose AFRN, a simple scheme to amplify inter-class discrepancy  
 104 and thus improve the few-shot fine-grained image classification. Our  
 105 scheme only involves two simple amendments to FRN+TDM: relaxing  
 106 the inter-class score in TDM, and adding a centre loss to FRN.
- 107 • By relaxing the inter-class score in TDM, we are able to remarkably  
 108 mitigate the negative impact, from the overfitting to the seen training  
 109 set of fine-grained subclasses, on the inference of unseen novel classes  
 110 in the few-shot learning setting.
- 111 • By incorporating the guidance of the centre loss to FRN, we are able to  
 112 enhance the discriminative power of the learnt features for fine-grained  
 113 image classification, through enlarging the omnipresent subtle distances  
 114 between fine-grained subclasses.
- 115 • The experiments on five benchmark fine-grained datasets demonstrate  
 116 that our scheme, although very simple, is quite effective to improve the  
 117 performance of few-shot fine-grained image classification.

118 The rest of the paper is organised as follows. In section 2, we discuss  
 119 the literature that is closely related to our work. The technical details of  
 120 FRN+TDM and AFRN are presented in section 3. In section 4, we demon-  
 121 strate the superior classification performances of AFRN through extensive

122 experimental results and ablation studies. Lastly, we draw conclusions in  
123 section 5.

## 124 2. Related Work

### 125 2.1. Metric-based few-shot methods for image classification

126 Metric-based few-shot methods aim to learn discriminative feature em-  
127 beddings that can be well generalized to new classes based on a predefined or  
128 a learnt distance metric, such as Euclidean distance [2], cosine distance [14],  
129 hyperbolic distance [15], or distance parameterized by neural networks [16].  
130 MatchingNet [4] adopts the cosine similarity to assign the label of the query  
131 image. ProtoNet [2] calculates prototypes as the average features of each  
132 class in the support set and assign the query image to the nearest class  
133 prototype by Euclidean distance. Instead of using a predefined metric, Rela-  
134 tionNet [16, 17] utilises a neural network to compute the nonlinear similar-  
135 ities between different samples. Moreover, Satorras and Estrach propose to  
136 utilise graph neural networks to measure the similarities between images [18].  
137 A large amount of work has also been done to extend the metric-based meth-  
138 ods for fine-grained images. For example, BSNet involves two similarity met-  
139 rics to learn discriminative features [8] and LRPABN adopts bilinear pooling  
140 operations [9].

### 141 2.2. Feature alignment-based few-shot methods for image classification

142 Feature alignment methods usually aim to align the object positions  
143 between the support and query sets to improve the classification perfor-  
144 mance [19]. CrossTransformers (CTX) [20] utilises the transformer-based

145 network to explore the spatially-correlated features and calculate the sim-  
 146 ilarity between two images. A more recent transformer-based method is  
 147 QSFormer [21], which effectively learns consistent representations of the sup-  
 148 port and query sets via the global sample transformer and the local patch  
 149 transformer. Dynamic meta-filer (DMF) [22] considers both channel-wise  
 150 and spatial-wise alignments by neural ordinary differential equation. Re-  
 151 lational embedding network (RENet) utilises the self-correlational repre-  
 152 sentation (SCR) module and the cross-correlational attention (CCA) mod-  
 153 ule, where the SCR module transforms the basic feature maps into self-  
 154 correlational tensors and extracts structural patterns, while the CCA module  
 155 calculates the cross-correlations between images and generates common at-  
 156 tention between them. FRN [11] aligns the features maps of the query image  
 157 and the support set via reconstructing the query image based on the pooled  
 158 support features, where the ridge regression-based reconstruction with close-  
 159 form solutions makes the process efficient. Besides the  $L_2$  norm adopted  
 160 in FRN, Sun et al. [23] propose to utilise the  $L_{2,1}$  norm for feature recon-  
 161 struction. To alleviate overfitting of the reconstruction-based methods, Li et  
 162 al. [24] propose the self-reconstruction network that can diversify the query  
 163 features by reconstructing the query features by themselves.

### 164 3. Methodology

#### 165 3.1. Problem definition

166 Few-shot learning aims to learn discriminative knowledge from a small  
 167 amount of labelled data to classify test instances from new tasks. In few-  
 168 shot learning, the dataset is usually divided into a base set  $\mathcal{D}_B$ , a validation

169 set  $\mathcal{D}_V$  and a novel set  $\mathcal{D}_N$ , where the classes of the three subsets do not  
 170 intersect. Few-shot learning learns from the tasks on  $\mathcal{D}_B$  to classify instances  
 171 of new tasks on  $\mathcal{D}_N$ . The instances in  $\mathcal{D}_V$  assist to find the best model during  
 172 the training process. In this paper, we follow the classic setting of  $N$ -way  
 173  $K$ -shot, i.e. the model is trained by the support set,  $\mathcal{S} = \{\mathbf{x}_i, y_i\}_{i=1}^{N \times K}$ , of  $N$   
 174 classes with  $K$  instances each class, and evaluated on the query set of the  
 175 same classes in  $\mathcal{S}$ ,  $\mathcal{Q} = \{\mathbf{x}_j, y_j\}_j^{N \times q}$ , of  $N$  classes with  $q$  instances each class.  
 176 The classification performance of the trained model is finally tested on  $\mathcal{D}_N$   
 177 with its average classification accuracy as the performance measure.

### 178 3.2. FRN+TDM

179 In metric-based few-shot learning methods, reshaping feature maps to  
 180 feature vectors as input to metric function can lead to loss of spatial details.  
 181 FRN [11] aims to resolve this problem by reconstructing every location of the  
 182 query feature map by the pooled support features from each class through  
 183 ridge regression. The class membership of the query instance is then assigned  
 184 based on the reconstruction error. However, in FRN, all channels are treated  
 185 equally with the same weights, which cannot stress the regions with high  
 186 discriminative abilities. To identify the discriminative regions, the TDM  
 187 module can be embedded in the FRN framework.

188 Specifically, TDM [12] takes the features extracted from the embedding  
 189 module to calculate the task-wise channel weight vector  $\beta_n$  of the  $n$ th class  
 190 as a linear combination of the support weight vector  $\beta_n^S$  and the query weight  
 191 vector  $\beta^Q$ :

$$\beta_n = \alpha \beta_n^S + (1 - \alpha) \beta^Q \in \mathbb{R}^C, \quad (1)$$

192 where  $\alpha \in [0, 1]$  is a hyper-parameter.  $\beta_n^S$  and  $\beta^Q$  are obtained from the  
 193 support attention module (SAM) and the query attention module (QAM),  
 194 respectively, based on the task-wise intra-class scores  $\mathbf{r}_n^{\text{intra}}$  and inter-class  
 195 scores  $\mathbf{r}_n^{\text{inter}}$ .

196 The input to SAM is the prototype of each class  $\mathcal{P}_n \in \mathbb{R}^{H \times W \times C}$ , i.e. the  
 197 average of all support set instances in the  $n$ th class. The  $c$ th element of  $\mathbf{r}_n^{\text{intra}}$   
 198 is then calculated as

$$r_{n,c}^{\text{intra}} = \frac{1}{H \times W} \|\mathcal{P}_{n,c} - \mathbf{M}_n\|_2^2, \quad (2)$$

199 where  $H$  and  $W$  are the height and width of the feature maps,  $C$  is the  
 200 number of channels,  $\mathcal{P}_{n,c} \in \mathbb{R}^{H \times W}$  is the  $c$ th channel of the  $n$ th prototype and  
 201  $\mathbf{M}_n \in \mathbb{R}^{H \times W}$  is the average of the channels in  $\mathcal{P}_n$ , i.e.  $\mathbf{M}_n = \frac{1}{C} \sum_{c=1}^C \mathcal{P}_{n,c}$ .  
 202 Thus  $\mathbf{r}_n^{\text{intra}}$  measures the dispersion of the channels in the prototype of each  
 203 class. On the contrary, the  $c$ th element of  $\mathbf{r}_n^{\text{inter}}$  involves information from  
 204 different classes:

$$r_{n,c}^{\text{inter}} = \frac{1}{H \times W} \min_{1 \leq l \leq N, n \neq l} \|\mathcal{P}_{n,c} - \mathbf{M}_l\|_2^2, \quad (3)$$

205 where  $\mathbf{M}_l$  denotes the mean spatial features of the  $l$ th class. It is clear  
 206 that  $r_{n,c}^{\text{inter}}$  measures the difference between each channel and its closest mean  
 207 spatial features of a different class. Finally, we obtain  $\beta_n^S$  as

$$\beta_n^S = \eta(g^{\text{inter}}(\mathbf{r}_n^{\text{inter}})) + (1 - \eta)(g^{\text{intra}}(\mathbf{r}_n^{\text{intra}})), \quad (4)$$

208 where  $g^{\text{inter}}$  and  $g^{\text{intra}}$  are fully-connected blocks and  $\eta \in [0, 1]$ . We adopt the  
 209 same structure for  $g$  as in [12].

210 Since the labels of query images are unknown, only the intra-class score

211 is involved in QAM:

$$r_{Q,c}^{\text{intra}} = \frac{1}{H \times W} \|\mathcal{P}_{Q,c} - \mathbf{M}_Q\|_2^2, \quad (5)$$

212 where  $\mathcal{P}_{Q,c}$  is the  $c$ th channel of the query feature maps and  $\mathbf{M}_Q$  is the mean  
213 of all channels of  $\mathcal{P}_Q$ . Then,  $\beta^Q$  is calculated as

$$\beta^Q = g^Q(\mathbf{r}_Q^{\text{intra}}), \quad (6)$$

214 where  $g^Q$  is a fully-connected block with the same structure as  $g^{\text{inter}}$  and  $g^{\text{intra}}$ .  
215 By substituting equation(4) and equation(6) to equation(1), we obtain the  
216 task-wise weights  $\beta_n$ .

217 In FRN+TDM, suppose the pooled support features of the  $n$ th class is  
218  $\mathbf{S}_n \in \mathbb{R}^{(K \times H \times W) \times C}$  while the the query features are  $\mathbf{Q} \in \mathbb{R}^{(H \times W) \times C}$ .  $\mathbf{Q}$  is  
219 reconstructed by each  $\mathbf{S}_n$  via ridge regression:

$$\hat{\mathbf{W}} = \underset{\mathbf{W}}{\operatorname{argmin}} \|\mathbf{Q} - \mathbf{W}\mathbf{S}_n\|_2^2 + \lambda \|\mathbf{W}\|_2^2, \quad (7)$$

220 where  $\mathbf{W} \in \mathbb{R}^{(H \times W) \times (K \times H \times W)}$  is the weight matrix and  $\lambda$  is a non-negative  
221 value that controls the contribution of the ridge penalty. The reconstructed  
222 query image by the  $n$ th class is calculated as

$$\hat{\mathbf{Q}}_n = \hat{\mathbf{W}}\mathbf{S}_n. \quad (8)$$

Then, the task-wise weight vector  $\beta_n$  is applied to the original and the reconstructed query feature maps to re-weight the channels:

$$\begin{aligned} \mathbf{Q}_n^r &= (\mathbf{1}_{H \times W} \beta_n^T) \odot \mathbf{Q}, \\ \hat{\mathbf{Q}}_n^r &= (\mathbf{1}_{H \times W} \beta_n^T) \odot \hat{\mathbf{Q}}_n, \end{aligned} \quad (9)$$

223 where  $\mathbf{1}_{H \times W}$  is a vector of  $H \times W$  1s and  $\odot$  is the element-wised Hadamard  
 224 product.

225 Lastly, to assign the class membership of the  $j$ th query image, we calculate  
 226 its probability of belonging to the  $n$ th class as

$$P(\hat{y}_j = n | \mathbf{x}_j) = \frac{e^{-\gamma d(\mathbf{Q}_{n,j}^r, \hat{\mathbf{Q}}_{n,j}^r)}}{\sum_{n' \in [1, N]} e^{-\gamma d(\mathbf{Q}_{n',j}^r, \hat{\mathbf{Q}}_{n',j}^r)}}, \quad (10)$$

227 where  $d(\mathbf{Q}_n^r, \hat{\mathbf{Q}}_n^r) = \frac{1}{H \times W} \|\mathbf{Q}_n^r - \hat{\mathbf{Q}}_n^r\|_2^2$  and  $\gamma$  is a non-negative parameter.

The training process of FRN+TDM is guided by the cross-entropy loss and the auxiliary loss in FRN:

$$\begin{aligned} L_{FRN} &= L_{CE} + L_{AUX} \\ &= - \sum_{j=1}^{N_q} \log(P(\hat{y}_j = y_j | \mathbf{x}_j)) \\ &\quad + \sum_{n \in [1, N]} \sum_{n' \in [1, N], n' \neq n} \|\hat{\mathbf{S}}_n (\hat{\mathbf{S}}_{n'})^T\|^2, \end{aligned} \quad (11)$$

228 where  $\hat{\mathbf{S}}_n$  is the row-normalised  $\mathbf{S}_n$ .

### 229 3.3. Adaptive feature map reconstruction network (AFRN)

230 Although FRN+TDM has achieved a state-of-the-art performance in few-  
 231 shot fine-grained image classification, due to the small inter-class discrepancy  
 232 omnipresent in fine-grained images and the small training set in the setting  
 233 of few-shot learning, the training of FRN+TDM can still result in excessively  
 234 flexible boundaries between subcategories and hence overfitting to the seen  
 235 subclasses in the training set. To mitigate this issue, we propose a simple  
 236 scheme to amplify inter-class discrepancy and thus improve FRN+TDM.  
 237 Our scheme only involves two simple amendments to FRN+TDM: relaxing

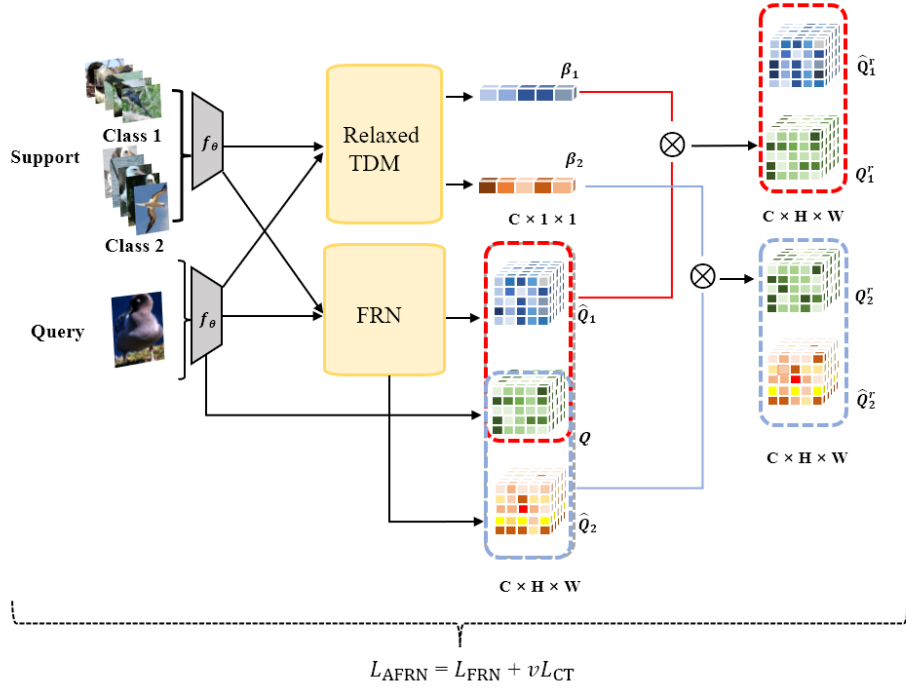


Figure 3: The structure of AFRN with an example of 2-way 5-shot classification. The embedded features of the support set and the query set are input to the FRN and the relaxed TDM modules. The FRN module reconstructs the query feature map by the pooled support features of each class and output the reconstructed query feature maps  $\hat{Q}_1$  and  $\hat{Q}_2$ . The relaxed TDM module produce the task-wise channel weights  $\beta_1$  and  $\beta_2$ . Then, the original query feature map  $Q$  and the reconstructed  $\hat{Q}_1$  are re-weighted by  $\beta_1$  to obtain  $Q_1^r$  and  $\hat{Q}_1^r$ . Similarly,  $Q$  and  $\hat{Q}_2$  are re-weighted by  $\beta_2$  to obtain  $Q_2^r$  and  $\hat{Q}_2^r$ . Lastly, the two pairs of re-weighted query features are used to obtain probabilities in equation(10) to assign the membership of the query image.

238 the inter-class score in TDM, and adding a centre loss to FRN. We call the  
 239 network incorporating our scheme to FRN+TDM the adaptive feature map  
 240 reconstruction network (AFRN). The structure of AFRN is illustrated in  
 241 Figure 3.

### 242 3.3.1. Relaxing inter-class score in TDM

243 In equation(3),  $r_{n,c}^{\text{inter}}$  measures the minimum distance between each chan-  
 244 nel and its closest mean spatial features of a different class. Therefore, the  
 245 classes that are mostly difficult to distinguish are specifically considered.  
 246 However, this may lead to extremely flexible classification boundaries in the  
 247 setting of fine-grained image classification, which is even severer in the few-  
 248 shot setting where the classes in the base set and the novel set are not the  
 249 same, due to the overfitting to the seen subclasses in the base set. To miti-  
 250 gate this problem, we propose the relaxed TDM by revising the calculation  
 251 of  $r_{n,c}^{\text{inter}}$  in equation (3) as

$$r_{n,c}^{\text{inter}} = \frac{1}{H \times W} \max_{1 \leq l \leq N, n \neq l} \|\mathcal{P}_{n,c} - \mathbf{M}_l\|_2^2. \quad (12)$$

252 In this way,  $r_{n,c}^{\text{inter}}$  measures the differences between classes that are less dif-  
 253 ficult to distinguish, which makes the classification boundaries less flexible  
 254 and thus mitigates the overfitting to a large extent.

### 255 3.3.2. Adding centre loss to FRN

256 The centre loss  $L_{CT}$  measures the intra-class variation of each class, which  
 257 is calculated as

$$L_{CT} = \sum_{j=1}^{Nq} \|\mathbf{Q}_j - \mathbf{C}_{y_j}\|_2^2, \quad (13)$$

258 where  $\mathbf{C}_{y_j}$  denotes the centre of the  $y_j$ th class, and  $\mathbf{Q}_j$  represents the feature  
 259 of the  $j$ th query. To effectively update the centre, we compute the centre as  
 260 the average of the query samples in one task.

261 Hence, the total loss function of AFRN is a simple amendment to that of  
 262 FRN in equation (11):

$$L_{AFRN} = L_{FRN} + \nu L_{CT}. \quad (14)$$

## 263 4. Experiments

264 In this section, we empirically demonstrate the superior classification per-  
 265 formance of AFRN on five fine-grained image datasets, by comparing it with  
 266 eight state-of-the-art methods: MatchingNet [4], ProtoNet [2], CTX [20],  
 267 DeepEMD [5], RENet [25], MixFSL [26], FRN [11] and FRN+TDM [12].

### 268 4.1. Datasets

269 We choose five publicly-available benchmark datasets for few-shot image  
 270 classification, namely CUB-200-2011 [27], aircraft [28], Oxford flowers [29],  
 271 Stanford cars [30] and Stanford dogs [31]. We name these datasets CUB,  
 272 aircraft, flowers, cars and dogs for short hereafter.

273 The CUB dataset contains 200 species of birds, with a total of 11,788  
 274 images. We randomly divide the 200 categories into the training, validation  
 275 and test sets, each consisting of 100, 50 and 50 categories, respectively.

276 The aircraft dataset has 100 classes of aircrafts, with a total of 10,000  
 277 images. We randomly divide the dataset into the training set with 50 classes,  
 278 the validation set with 25 classes and the test set with 25 classes.

279 The flowers dataset consists of 102 categories of flowers with 8,189 images.  
280 Each type of flower consists of 40 to 258 images, mainly featuring common  
281 British flowers. We randomly select 51 classes as the training set, 26 classes  
282 as the validation set, and 25 classes as the test set.

283 The cars dataset includes 196 classes of cars, with a total of 16,185 images.  
284 We randomly divide the dataset into the training set with 130 classes, the  
285 validation set with 17 classes and the test set with 49 classes.

286 The dogs dataset contains 120 breeds of dogs, with a total of 20,580  
287 images. We randomly divide the 120 categories into the training set with 60  
288 categories, the validation set with 30 categories and the testing set with 30  
289 categories.

#### 290 *4.2. Implementation details*

291 We adopt ResNet-12 as the backbone with the same implementation de-  
292 tails as in [28, 32, 33]. The ResNet-12 backbone consists of 4 residual blocks,  
293 and each residual block has 3 convolutional layers. We adopt the leaky ReLU  
294 with  $\alpha = 0.1$  and  $2 \times 2$  max pooling. We also adopt the deep block from the  
295 original implementation [32, 28, 33], so the output size of each residual block  
296 is 64, 160, 320 and 640. Therefore, the shape of the output feature map of  
297 an input image of size  $84 \times 84$  is  $640 \times 5 \times 5$ . During the training process, we  
298 implement the standard data augmentation step, including random cropping,  
299 horizontal flipping and color jittering, as in [28, 5, 34, 35].

300 Following [14, 33], we train ResNet-12 for 1,200 epochs and reduce the  
301 learning rate proportionally at the 400th and 800th epochs. We use the  
302 validation set to select the best performing model during the training process  
303 and validate every 20 epochs. We train the models with the 10-way 5-shot

304 setting and test the models with the 5-way 1-shot and 5-way 5-shot setting.

305 For AFRN, we follow TDM [12] to set  $\alpha = \eta = 0.5$ , and set  $\nu = 0.05$ . In  
306 section 4.5, we will show the robustness of  $\nu$ .

307 AFRN and FRN+TDM have the same amount of parameters and they  
308 have the same FLOPs. For the 5-way 1-shot task with 16 query images, their  
309 FLOPs is 299.6G per task while for the 5-way 5-shot setting with 16 query  
310 images, their FLOPs is 370G per task.

Table 1: 5-way few-shot classification accuracies on the CUB, aircraft, flowers, cars and dogs datasets with the ResNet-12 backbone. Methods labeled by † denote our implementations. The best classification accuracies are labelled in bold fonts.

Method	CUB		Aircraft		Flowers		Cars		Dogs	
	1-shot	5-shot	1-shot	5-shot	1-shot	5-shot	1-shot	5-shot	1-shot	5-shot
MatchingNet[4] †	71.87±0.24	85.08±0.24	56.74±0.87	73.75±0.69	71.89±0.90	85.46±0.59	45.29±0.82	64.00±0.74	66.48±0.88	79.57±0.63
ProtoNet[2] †	81.02±0.20	91.93±0.11	46.68±0.81	71.27±0.27	75.41±0.22	89.46±0.14	82.29±0.20	93.11±0.10	73.81±0.21	87.39±0.12
CTX[20] †	80.39±0.20	91.01±0.11	65.60±0.25	80.20±0.25	-	-	85.03±0.19	92.63±0.11	73.22±0.22	85.90±0.13
DeepEMD[5] †	75.59±0.30	88.23±0.18	-	-	70.00±0.35	83.63±0.26	73.30±0.29	88.37±0.17	70.38±0.30	85.24±0.18
RENet[25] †	77.45±0.45	90.50±0.26	59.16±0.47	76.48±0.37	79.91±0.42	92.33±0.22	79.66±0.44	91.95±0.22	71.69±0.47	85.60±0.30
MixFSL[26] †	64.53±0.92	80.67±0.64	60.55±0.86	77.57±0.69	72.60±0.91	86.52±0.65	58.15±0.87	80.54±0.63	67.26±0.90	82.05±0.56
FRN[11] †	82.33±0.19	92.02±0.11	70.26±0.22	83.58±0.14	81.68±0.20	92.61±0.11	86.59±0.18	95.01±0.08	76.49±0.21	88.22±0.12
FRN+TDM[12] †	83.31±0.19	92.70±0.10	70.61±0.21	84.53±0.13	82.95±0.19	93.61±0.10	<b>89.38±0.16</b>	<b>96.98±0.06</b>	76.67±0.21	88.53±0.12
<b>Ours</b>	<b>83.95±0.18</b>	<b>93.17±0.10</b>	<b>72.19±0.21</b>	<b>85.59±0.13</b>	<b>83.59±0.19</b>	<b>94.05±0.09</b>	89.27±0.16	96.89±0.06	<b>77.01±0.21</b>	<b>88.60±0.12</b>

Table 2: The results of the one-sided paired  $t$ -test of comparing the classification accuracies of our method with those of the state-of-the-art methods in Table 1. The null hypothesis  $H_0$  is  $\mu_{\text{AFRN}} < \mu_m$ , where  $\mu$  is the mean classification accuracy and  $m \in \{\text{MatchingNet, ProtoNet, CTX, DeepEMD, RENet, MixFSL, FRN, FRN+TDM}\}$ .

Ours vs.	MatchingNet	ProtoNet	CTX	DeepEMD	RENet	MixFSL	FRN	FRN+TDM
$p$ value	$1 \times 10^{-3}$	$7 \times 10^{-3}$	$3.9 \times 10^{-5}$	$2.8 \times 10^{-4}$	$2.8 \times 10^{-4}$	$1.4 \times 10^{-4}$	$3.3 \times 10^{-5}$	$7 \times 10^{-3}$
Reject at 1% level	✓	✓	✓	✓	✓	✓	✓	✓

### 311 4.3. Comparison with the state-of-the-art methods

312 We report the classification accuracies of AFRN and the eight state-of-  
313 the-art methods on five fine-grained image datasets in Table 1. Obviously,

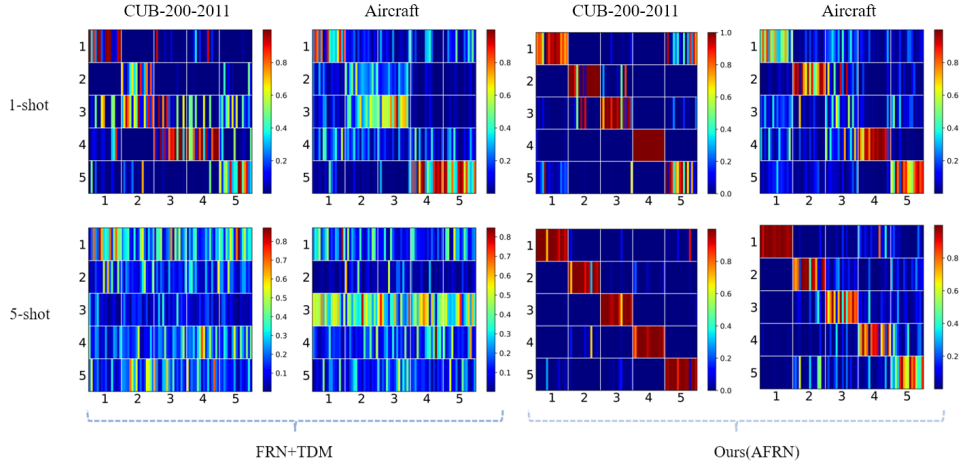


Figure 4: The visualisations of the confusion matrices of AFRN and FRN+TDM on the CUB and aircraft datasets under the 5-way 1-shot and 5-way 5-shot settings.

Table 3: The ablation study on the relaxed TDM module and the centre loss.

	Relaxed TDM	Centre loss	CUB		Aircraft	
			1-shot	5-shot	1-shot	5-shot
(a)	-	-	83.31±0.19	92.70±0.10	70.61±0.21	84.53±0.13
(b)	✓	-	83.73±0.12	92.86±0.10	71.59±0.22	85.06±0.13
(c)	-	✓	83.77±0.18	93.09±0.10	71.05±0.21	84.58±0.13
(d)	✓	✓	<b>83.95±0.18</b>	<b>93.17±0.10</b>	<b>72.19±0.21</b>	<b>85.59±0.13</b>

314 our method can beat all state-of-the-art methods on the CUB, aircraft, flow-  
 315 ers and dogs dataset, while providing competitive classification results with  
 316 FRN+TDM on the cars dataset. This demonstrates the effectiveness of in-  
 317 volving the centre loss and the relaxed TDM module. To have a deep insight  
 318 to the results, we compare the visualisations of the confusion matrices of  
 319 AFRN and FRN+TDM in Figure 4 on the CUB and aircraft datasets. It is  
 320 clear that AFRN is better than FRN+TDM on the two datasets with more  
 321 deep red stripes or higher values on the diagonals. To confirm that AFRN is  
 322 significantly better than the state-of-the-art methods, we perform one-sided  
 323 paired  $t$ -test to compare the classification accuracies of AFRN and those of  
 324 other methods in Table 1, with a null hypothesis  $H_0$  of  $\mu_{AFRN} < \mu_m$ , where  $\mu$   
 325 is the mean classification accuracy and  $m \in \{\text{MatchingNet, ProtoNet, CTX,}$   
 326  $\text{DeepEMD, RENet, MixFSL, FRN, FRN+TDM}\}$ .  $H_0$  can be rejected at 1%  
 327 level for all methods compared, suggesting that the classification accuracy of  
 328 AFRN is significantly better than those of other methods.

#### 329 4.4. Ablation studies

330 Here we explore the impacts of the relaxed TDM module and the centre  
 331 loss on the classification performance and report the results on the CUB and  
 332 aircraft datasets in Table 3. For the relaxed TDM column, ‘-’ represents  
 333 adopting the original TDM module while ‘✓’ is for the proposed relaxed  
 334 TDM module. For the centre loss column, ‘-’ is to train the model by the  
 335 original FRN loss in (11) while ‘✓’ represents training the model by the  
 336 AFRN loss in (14). Thus, scenario-(a) corresponds to FRN+TDM while  
 337 scenario-(d) represents AFRN. Clearly, the classification accuracy of TDM  
 338 can be raised by only modifying the inter-class score via the relaxed TDM

339 in scenario-(b). It is worth noting that, for the 1-shot classification of the  
 340 aircraft dataset, the accuracy is improved greatly by almost 1%, suggesting  
 341 that the subcategories of aircraft are highly similar and the relaxed score is  
 342 required to reduce potential overfitting. In scenario-(c), when we only involve  
 343 the additional centre loss, the improvement is more substantial for the CUB  
 344 dataset, suggesting that the variation within each subcategory of the CUB  
 345 dataset is relatively large and thus making intra-class variation smaller via  
 346 centre loss is beneficial. Finally, utilising the relaxed TDM module as well  
 347 as the centre loss can provide the best classification accuracies.

Table 4: The effect of  $\nu$  in (14) of the AFRN loss.

$\nu$	CUB		Flowers	
	1-shot	5-shot	1-shot	5-shot
0.5	83.22±0.19	92.75±0.10	82.75±0.19	93.46±0.10
0.05	<b>83.95±0.18</b>	<b>93.17±0.10</b>	<b>83.59±0.19</b>	<b>94.05±0.09</b>
0.005	83.69±0.18	93.07±0.10	82.35±0.20	93.22±0.10

#### 348 4.5. The effect of $\nu$ in (14)

349 In this section, we present the effect of  $\nu$  in (14), i.e. the parameter con-  
 350 trolling the contribution of the centre loss, on the classification performance.  
 351 The classification accuracies of the CUB and flowers datasets for three values  
 352 of  $\nu$ , 0.5, 0.05 and 0.005, are summarised in Table 4. It shows that 0.05 is a  
 353 proper choice. In addition, the accuracies of using the three values of  $\nu$  are  
 354 all higher than or competitive with FRN+TDM.

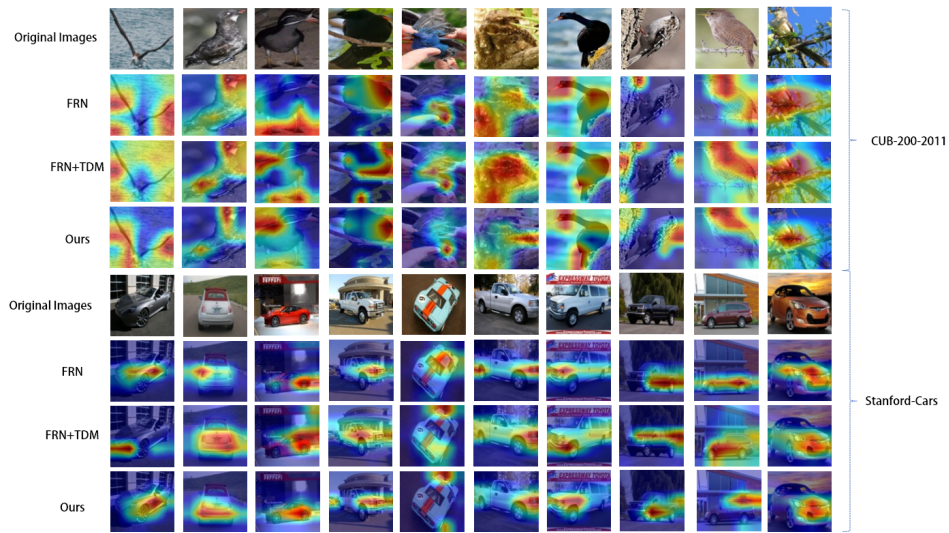


Figure 5: The visualisation of the discriminative features extracted by FRN, FRN+TDM and AFRN ('Ours') on the CUB and cars datasets. AFRN focuses on the most discriminative regions compared with FRN and FRN+TDM.

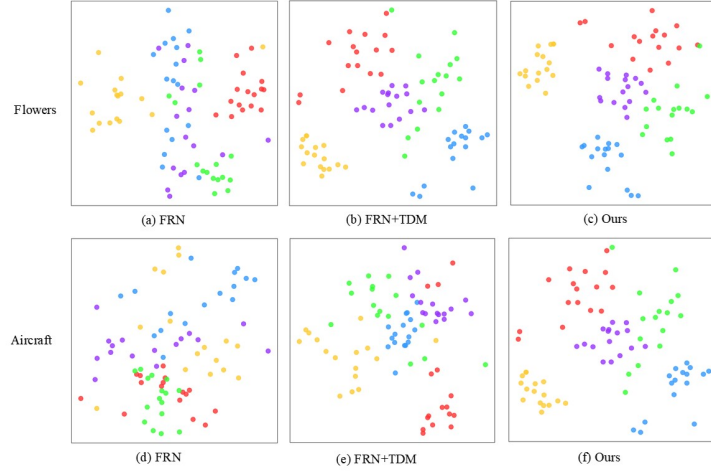


Figure 6: The visualisation of the feature embeddings of FRN, FRN+TDM and AFRN ('Ours') on the flowers and aircraft datasets. AFRN can provide the best separation of different classes.

#### 355 4.6. The visual comparisons of FRN, FRN+TDM and AFRN

##### 356 4.6.1. Visualisation of discriminative features

357 To demonstrate that AFRN can focus on the most discriminative regions  
 358 for classification, we visually compare the discriminative regions identified by  
 359 FRN, FRN+TDM and AFRN, following the Grad-CAM technology [36] in  
 360 Figure 5. For the CUB and cars datasets, we randomly select 10 images for  
 361 visualisation. We can observe that FRN tends to focus on both the objects  
 362 and irrelevant backgrounds. FRN+TDM can improve this by identifying  
 363 smaller discriminative regions, while AFRN can usually make the areas even  
 364 smaller by focusing on the highly discriminative ones.

365 *4.6.2. Visualisation of feature embeddings*

366 To further show that AFRN can amplify the inter-class discrepancy, we  
 367 visualise the feature embeddings learnt by FRN, FRN+TDM and AFRN  
 368 via  $t$ -distributed stochastic neighbour embedding ( $t$ -SNE) [37] in Figure 6.  
 369 The results of the flowers and aircraft datasets are presented in the first and  
 370 second rows in Figure 6, respectively. For each dataset, we randomly select  
 371 five classes with 16 test samples each and label them by different colours. The  
 372 five classes are severely mixed in FRN while better separated in FRN+TDM.  
 373 Obviously, the best separation of the classes is achieved by FRN: the inter-  
 374 class discrepancy is amplified, which also supports our motivation in Figure 1.

375 *4.7. Discussion*

Table 5: The classification accuracies of FRN, FRN+TDM and AFRN (‘Ours’) on two coarse-grained datasets, mini-ImageNet and FC100, with the ResNet-12 backbone. The best classification accuracies are labelled in bold fonts.

	mini-ImageNet		FC100	
	1-shot	5-shot	1-shot	5-shot
FRN	<b>63.26±0.21</b>	77.68±0.15	<b>40.31±0.17</b>	<b>55.34±0.17</b>
FRN+TDM	62.18±0.20	78.41±0.15	39.84±0.17	54.16±0.17
Ours	62.78±0.20	<b>78.60±0.15</b>	40.09±0.18	54.38±0.18

376 In this section, we further test the ability of AFRN to classify coarse-  
 377 grained data, where larger categories or super-categories with large intra-  
 378 class variations are considered. We adopt two benchmark coarse-grained  
 379 datasets, the mini-ImageNet dataset [4] and the FC100 dataset [38]. The  
 380 mini-ImageNet dataset contains 60,000 images distributed evenly over 100  
 381 classes. We randomly divide the dataset to a training set with 64 classes,

382 a validation set with 16 classes and a test set with 20 classes. The FC100  
 383 dataset has 100 object categories which are merged to 20 super-categories.  
 384 We randomly divide it to a training set with 12 super-categories containing  
 385 60 object categories, a validation set with 4 super-categories containing 20  
 386 object categories and a test set with 4 super-categories containing 20 object  
 387 categories.

388 The classification accuracies of FRN, FRN+TDM and AFRN on coarse-  
 389 grained datasets are reported in Table 5. Clearly, the original FRN dominates  
 390 FRN+TDM and AFRN in most scenarios, except for the classification of 5-  
 391 shot mini-ImageNet. However, we note that AFRN performs slightly better  
 392 than FRN+TDM in all cases, which demonstrate that the two amendments  
 393 also work on coarse-grained data, but not effective enough to beat the original  
 394 FRN. One explanation to this result is that TDM or relaxed TDM put too  
 395 much attention on few channels while ignore information from other channels  
 396 that may be valuable for coarse-grained data. Thus, they perform worse than  
 397 the original FRN when all channels are considered equally.

## 398 5. Conclusions

399 In this paper, we propose AFRN, a simple scheme to amplify the inter-  
 400 class discrepancy and thus improve the classification performance of FRN+TDM  
 401 on few-shot fine-grained images. To mitigate the potential overfitting to the  
 402 seen subclasses, we propose to relax the inter-class score in TDM. To enlarge  
 403 the subtle differences between the subclasses of fine-grained images, we pro-  
 404 pose to incorporate the centre loss to FRN. Extensive experiments on five  
 405 fine-grained datasets showcase that our scheme can produce the state-of-the-

406 art performance, verified by statistical tests. Results in ablation study also  
407 reveal the effectiveness of each amendment. Moreover, we note one limitation  
408 of our method on classifying coarse-grained data, which we identify as our  
409 future work.

#### 410 **Acknowledgement**

411 This research was supported in part by the National Natural Science  
412 Foundation of China (NSFC) under Grant 62176110, the Key Research and  
413 Development Program of Gansu Province under Grant 22YF7GA130, S&T  
414 Program of Hebei under grant SZX2020034, Hong-liu Distinguished Young  
415 Talents Foundation of Lanzhou University of Technology and the Royal So-  
416 ciety under International Exchanges Award IEC\NSFC\201071.

#### 417 **References**

- 418 [1] X. Li, X. Yang, Z. Ma, J.-H. Xue, Deep metric learning for few-shot im-  
419 age classification: A review of recent developments, Pattern Recognition  
420 (2023) 109381.
- 421 [2] J. Snell, K. Swersky, R. Zemel, Prototypical networks for few-shot learn-  
422 ing, in: Advances in Neural Information Processing Systems, Vol. 30,  
423 2017.
- 424 [3] X. Huang, S. H. Choi, SAPENet: Self-attention based prototype en-  
425 hancement network for few-shot learning, Pattern Recognition 135  
426 (2023) 109170.

- 427 [4] O. Vinyals, C. Blundell, T. Lillicrap, k. kavukcuoglu, D. Wierstra,  
428 Matching networks for one shot learning, in: Advances in Neural In-  
429 formation Processing Systems, Vol. 29, 2016.
- 430 [5] C. Zhang, Y. Cai, G. Lin, C. Shen, DeepEMD: Few-shot image classifi-  
431 cation with differentiable earth mover's distance and structured classi-  
432 fiers, in: 2020 IEEE/CVF Conference on Computer Vision and Pattern  
433 Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020, 2020,  
434 pp. 12200–12210.
- 435 [6] V. N. Nguyen, S. Løkse, K. Wickstrøm, M. Kampffmeyer, D. Roverso,  
436 R. Jenssen, SEN: A novel feature normalization dissimilarity measure for  
437 prototypical few-shot learning networks, in: Computer Vision–ECCV  
438 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020,  
439 Proceedings, Part XXIII 16, Springer, 2020, pp. 118–134.
- 440 [7] C. Chen, K. Li, W. Wei, J. T. Zhou, Z. Zeng, Hierarchical graph neu-  
441 ral networks for few-shot learning, IEEE Transactions on Circuits and  
442 Systems for Video Technology 32 (1) (2021) 240–252.
- 443 [8] X. Li, J. Wu, Z. Sun, Z. Ma, J. Cao, J.-H. Xue, BSNet: Bi-similarity net-  
444 work for few-shot fine-grained image classification, IEEE Transactions  
445 on Image Processing 30 (2020) 1318–1331.
- 446 [9] H. Huang, J. Zhang, J. Zhang, J. Xu, Q. Wu, Low-rank pairwise align-  
447 ment bilinear network for few-shot fine-grained image classification,  
448 IEEE Transactions on Multimedia 23 (2020) 1666–1680.

- 449 [10] H. Huang, J. Zhang, L. Yu, J. Zhang, Q. Wu, C. Xu, TOAN: Target-  
450 oriented alignment network for fine-grained image categorization with  
451 few labeled samples, *IEEE Transactions on Circuits and Systems for*  
452 *Video Technology* 32 (2) (2021) 853–866.
- 453 [11] D. Wertheimer, L. Tang, B. Hariharan, Few-shot classification with fea-  
454 ture map reconstruction networks, in: *IEEE Conference on Computer*  
455 *Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021,*  
456 *2021, pp. 8012–8021.*
- 457 [12] S. B. Lee, W. Moon, J. Heo, Task discrepancy maximization for fine-  
458 grained few-shot classification, in: *IEEE/CVF Conference on Computer*  
459 *Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA,*  
460 *June 18-24, 2022, 2022, pp. 5321–5330.*
- 461 [13] Y. Wen, K. Zhang, Z. Li, Y. Qiao, A discriminative feature learning  
462 approach for deep face recognition, in: *Computer Vision–ECCV 2016:*  
463 *14th European Conference, Amsterdam, The Netherlands, October 11–*  
464 *14, 2016, Proceedings, Part VII 14, Springer, 2016, pp. 499–515.*
- 465 [14] S. Gidaris, N. Komodakis, Dynamic few-shot visual learning without  
466 forgetting, in: *2018 IEEE Conference on Computer Vision and Pattern*  
467 *Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018,*  
468 *2018, pp. 4367–4375.*
- 469 [15] V. Khrulkov, L. Mirvakhabova, E. Ustinova, I. V. Oseledets, V. S. Lem-  
470 pitsky, Hyperbolic image embeddings, in: *2020 IEEE/CVF Conference*

- 471 on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA,  
472 USA, June 13-19, 2020, 2020, pp. 6417–6427.
- 473 [16] F. Sung, Y. Yang, L. Zhang, T. Xiang, P. H. S. Torr, T. M. Hospedales,  
474 Learning to compare: Relation network for few-shot learning, in: 2018  
475 IEEE Conference on Computer Vision and Pattern Recognition, CVPR  
476 2018, Salt Lake City, UT, USA, June 18-22, 2018, 2018, pp. 1199–1208.
- 477 [17] Z. Li, Z. Hu, W. Luo, X. Hu, SaberNet: Self-attention based effective  
478 relation network for few-shot learning, *Pattern Recognition* 133 (2023)  
479 109024.
- 480 [18] V. G. Satorras, J. B. Estrach, Few-shot learning with graph neural net-  
481 works, in: 6th International Conference on Learning Representations,  
482 ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Confer-  
483 ence Track Proceedings, 2018.
- 484 [19] F. Hao, F. He, J. Cheng, D. Tao, Global-local interplay in semantic align-  
485 ment for few-shot learning, *IEEE Transactions on Circuits and Systems*  
486 *for Video Technology* 32 (7) (2021) 4351–4363.
- 487 [20] C. Doersch, A. Gupta, A. Zisserman, CrossTransformers: Spatially-  
488 aware few-shot transfer, in: *Advances in Neural Information Processing*  
489 *Systems*, Vol. 33, 2020, pp. 21981–21993.
- 490 [21] X. Wang, X. Wang, B. Jiang, B. Luo, Few-shot learning meets trans-  
491 former: Unified query-support transformers for few-shot classification,  
492 *IEEE Transactions on Circuits and Systems for Video Technology*  
493 33 (12) (2023) 7789–7802.

- 494 [22] C. Xu, Y. Fu, C. Liu, C. Wang, J. Li, F. Huang, L. Zhang, X. Xue,  
495 Learning dynamic alignment via meta-filter for few-shot learning, in:  
496 IEEE Conference on Computer Vision and Pattern Recognition, CVPR  
497 2021, virtual, June 19-25, 2021, 2021, pp. 5182–5191.
- 498 [23] J. Sun, X. Shen, Q. Sun, Efficient feature reconstruction via  $l_{2,1}$ -norm  
499 regularization for few-shot classification, IEEE Transactions on Circuits  
500 and Systems for Video Technology 33 (12) (2023) 7452–7465.
- 501 [24] X. Li, Z. Li, J. Xie, X. Yang, J.-H. Xue, Z. Ma, Self-reconstruction  
502 network for fine-grained few-shot classification, Pattern Recognition 153  
503 (2024) 110485.
- 504 [25] D. Kang, H. Kwon, J. Min, M. Cho, Relational embedding for few-shot  
505 classification, in: 2021 IEEE/CVF International Conference on Com-  
506 puter Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021,  
507 2021, pp. 8802–8813.
- 508 [26] A. Afrasiyabi, J. Lalonde, C. Gagné, Mixture-based feature space learn-  
509 ing for few-shot image classification, in: 2021 IEEE/CVF International  
510 Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada,  
511 October 10-17, 2021, 2021, pp. 9021–9031.
- 512 [27] C. Wang, H. K. Galoogahi, C. Lin, S. Lucey, Deep-LK for efficient  
513 adaptive object tracking, in: 2018 IEEE International Conference on  
514 Robotics and Automation, ICRA 2018, Brisbane, Australia, May 21-25,  
515 2018, 2018, pp. 627–634.

- 516 [28] H. Ye, H. Hu, D. Zhan, F. Sha, Few-shot learning via embedding adap-  
517 tation with set-to-set functions, in: 2020 IEEE/CVF Conference on  
518 Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA,  
519 USA, June 13-19, 2020, 2020, pp. 8805–8814.
- 520 [29] M. Nilsback, A. Zisserman, Automated flower classification over a large  
521 number of classes, in: Sixth Indian Conference on Computer Vision,  
522 Graphics & Image Processing, ICVGIP 2008, Bhubaneswar, India, 16-  
523 19 December 2008, 2008, pp. 722–729.
- 524 [30] J. Krause, M. Stark, J. Deng, L. Fei-Fei, 3D object representations for  
525 fine-grained categorization, in: 2013 IEEE International Conference on  
526 Computer Vision Workshops, ICCV Workshops 2013, Sydney, Australia,  
527 December 1-8, 2013, 2013, pp. 554–561.
- 528 [31] A. Khosla, N. Jayadevaprakash, B. Yao, F.-F. Li, Novel dataset for fine-  
529 grained image categorization: Stanford dogs, in: Proc. CVPR workshop  
530 on fine-grained visual categorization (FGVC), Vol. 2, Citeseer, 2011.
- 531 [32] Y. Tian, Y. Wang, D. Krishnan, J. B. Tenenbaum, P. Isola, Rethinking  
532 few-shot image classification: a good embedding is all you need?, in:  
533 Computer Vision–ECCV 2020: 16th European Conference, Glasgow,  
534 UK, August 23–28, 2020, Proceedings, Part XIV 16, Springer, 2020, pp.  
535 266–282.
- 536 [33] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez,  
537 L. u. Kaiser, I. Polosukhin, Attention is all you need, in: Advances in  
538 Neural Information Processing Systems, Vol. 30, 2017.

- 539 [34] Y. Wang, W.-L. Chao, K. Q. Weinberger, L. Van Der Maaten, Sim-  
540 pleShot: Revisiting nearest-neighbor classification for few-shot learning,  
541 arXiv preprint arXiv:1911.04623.
- 542 [35] Z. Lin, M. Feng, C. N. d. Santos, M. Yu, B. Xiang, B. Zhou, Y. Ben-  
543 gio, A structured self-attentive sentence embedding, arXiv preprint  
544 arXiv:1703.03130.
- 545 [36] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, D. Batra,  
546 Grad-CAM: Visual explanations from deep networks via gradient-based  
547 localization, in: IEEE International Conference on Computer Vision,  
548 ICCV 2017, Venice, Italy, October 22-29, 2017, 2017, pp. 618–626.
- 549 [37] G. Hinton, L. Van Der Maaten, Visualizing data using t-sne, Journal of  
550 Machine Learning Research 9 (2008) 2579–2605.
- 551 [38] B. Oreshkin, P. Rodríguez López, A. Lacoste, TADAM: Task dependent  
552 adaptive metric for improved few-shot learning, in: Advances in Neural  
553 Information Processing Systems, Vol. 31, 2018.

- We propose AFRN, a simple scheme to amplify inter-class discrepancy.
- We relax the inter-class score in TDM to mitigate the negative impact of overfitting.
- We incorporate the guidance of the centre loss to FRN to enhance the discriminative power of learnt features.
- The experimental results demonstrate that our scheme is simple yet effective to improve the few-shot classification performance.

Journal Pre-proof

**Xiaoxu Li** received the Ph.D. degree from Beijing University of Posts and Telecommunications in 2012. She is currently an Associate Professor with the School of Computer and Communication, Lanzhou University of Technology. Her research interests include machine learning fundamentals with a focus on applications in image and video understanding. She is also a member of the China Computer Federation.

**Zijie Guo** received the B.E. degree in Management from Hankou University in 2021. He is a postgraduate student in Lanzhou University of Technology. His research interests include computer vision and few-shot learning.

**Rui Zhu** received the Ph.D. degree in statistics from University College London in 2017. She is a Senior Lecturer in the Faculty of Actuarial Science and Insurance, City, University of London. Her research interests include machine learning and its applications in image quality assessment, hyperspectral image analysis and actuarial science. She is an Associate Editor of Neurocomputing.

**Zhanyu Ma** is currently a Professor at Beijing University of Posts and Telecommunications, Beijing, China, since 2019. He received the Ph.D. degree in electrical engineering from KTH Royal Institute of Technology, Sweden, in 2011. From 2012 to 2013, he was a Postdoctoral Research Fellow with the School of Electrical Engineering, KTH. He has been an Associate Professor with the Beijing University of Posts and Telecommunications, Beijing, China, from 2014 to 2019. His research interests include pattern recognition and machine learning fundamentals with a focus on applications in computer vision, multimedia signal processing. He is a Senior Member of IEEE.

**Jun Guo** received B.E. and M.E. degrees from Beijing University of Posts and Telecommunications (BUPT), China in 1982 and 1985, respectively, Ph.D. degree from the Tohoku-Gakuin University, Japan in 1993. At present he is a professor of BUPT. His research interests include pattern recognition theory and application, information retrieval, content based information security, and network management.

**Jing-Hao Xue** received the Dr.Eng. degree in signal and information processing from Tsinghua University in 1998 and the Ph.D. degree in statistics from the University of Glasgow in 2008. He is a Professor in the Department of Statistical Science at University College London. His research interests include statistical pattern recognition, machine learning and computer vision. He is an Associate Editor of the IEEE Transactions on Circuits and Systems for Video Technology, the IEEE Transactions on Cybernetics, and the IEEE Transactions on Neural Networks and Learning Systems.

**Declaration of interests**

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests:

Journal Pre