

Journal Pre-proof

Self-reconstruction network for fine-grained few-shot classification

Xiaoxu Li, Zhen Li, Jiyang Xie, Xiaochen Yang, Jing-Hao Xue,
Zhanyu Ma



PII: S0031-3203(24)00236-X
DOI: <https://doi.org/10.1016/j.patcog.2024.110485>
Reference: PR 110485

To appear in: *Pattern Recognition*

Received date: 16 March 2023
Revised date: 10 October 2023
Accepted date: 5 April 2024

Please cite this article as: X. Li, Z. Li, J. Xie et al., Self-reconstruction network for fine-grained few-shot classification, *Pattern Recognition* (2024), doi: <https://doi.org/10.1016/j.patcog.2024.110485>.

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 Published by Elsevier Ltd.

Self-Reconstruction Network for Fine-Grained Few-Shot Classification

Xiaoxu Li^{a,b}, Zhen Li^a, Jiyang Xie^b, Xiaochen Yang^{c,*}, Jing-Hao Xue^d,
Zhanyu Ma^b

^a*School of Computer and Communication, Lanzhou University of
Technology, Lanzhou, 730050, China*

^b*Pattern Recognition and Intelligent System Laboratory, School of Artificial Intelligence,
Beijing University of Posts and Telecommunications, Beijing, 100876, China*

^c*School of Mathematics and Statistics, University of Glasgow, Glasgow, G12 8QQ, UK*

^d*Department of Statistical Science, University College London, London, WC1E 6BT, UK*

Abstract

Metric-based methods are one of the most common methods to solve the problem of few-shot image classification. However, traditional metric-based few-shot methods suffer from overfitting and local feature misalignment. The recently proposed feature reconstruction-based approach, which reconstructs query image features from the support set features of a given class and compares the distance between the original query features and the reconstructed query features as the classification criterion, effectively solves the feature misalignment problem. However, the issue of overfitting still has not been considered. To this end, we propose a self-reconstruction metric module for diversifying query features and a restrained cross-entropy loss for avoiding over-confident predictions. By introducing them, the proposed self-reconstruction network can effectively alleviate overfitting. Ex-

*Corresponding author

Email address: xiaochen.yang@glasgow.ac.uk (Xiaochen Yang)

tensive experiments on five benchmark fine-grained datasets demonstrate that our proposed method achieves state-of-the-art performance on both 5-way 1-shot and 5-way 5-shot classification tasks. Code is available at <https://github.com/liz-lut/SRM-main>.

Keywords: Few-shot learning, Fine-grained image classification, Deep neural network, Self-reconstruction network

1. Introduction

Deep learning has achieved impressive performance in computer vision. However, deep networks usually demand numerous labeled data for training, which is impractical in many tasks where data acquisition is costly and time-consuming. For this reason, researchers in the computer vision community turn considerable attention to few-shot learning in recent years, especially for few-shot classification [1, 2, 3].

The goal of few-shot classification is to recognize unseen query sample with very limited (often less than 10) labeled support samples. Existing methods usually can be categorized into three classes, metric-based methods [4], optimization-based methods [5], and transfer learning-based methods [6, 7]. Among these methods, metric-based methods are relatively simple but effective, achieving state-of-the-art performance in many few-shot tasks. Metric-based methods usually adopt episodic training strategy to train a feature extractor with a fixed distance metric or a parameterized distance metric, and then fix them, i.e., without fine-tuning, to classify unseen novel query samples. However, early works, e.g., prototypical network [8] and relation network [9], mainly build on global features, which suffer from inaccurate

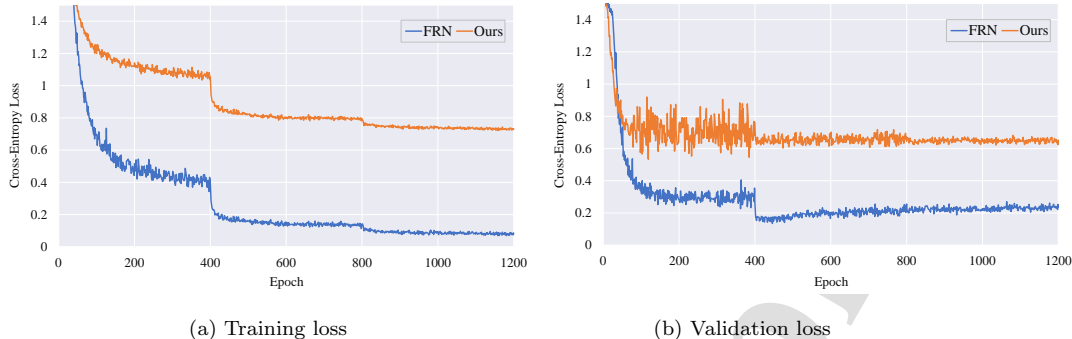


Figure 1: Motivation of the self-reconstruction network. FRN [13] encounters the overfitting problem as its training loss keeps decreasing but its validation loss starts to increase after around 400 epochs. In contrast, the validation loss of our method stays stable after 400 epochs, demonstrating its effectiveness in addressing overfitting.

19 similarity measure between two samples due to the mismatch of key infor-
 20 mation in images. This is particularly detrimental to few-shot fine-grained
 21 image classification, as sub-categories have subtle differences and the valu-
 22 able and discriminative information is likely to locate in different regions. To
 23 address such an issue, some subsequent works start to focus on learning a
 24 metric on local features [10] or aligning local features [11, 12].

25 Recently, some metric-based approaches introducing new alignment [12]
 26 or reconstruction [13] techniques have achieved impressive performance in
 27 fine-grained few-shot image classification. However, in the experiment, we no-
 28 ticed that the state-of-the-art method, feature reconstruction network (FRN) [13],
 29 suffered from overfitting during episodic training. As shown in Figure 1, while
 30 the loss function of FRN keeps decreasing on the training set, it increases on
 31 the validation set after 400 epochs. This overfitting phenomenon may occur
 32 as, comparing with that of ImageNet, CIFAR, etc., the numbers of images

33 and classes in fine-grained datasets are relatively small and thus the task
34 diversity in episodic training is limited.

35 To alleviate overfitting, we propose a self-reconstruction network for few-
36 shot fine-grained classification, which introduces a new self-reconstruction
37 metric module and a restrained cross-entropy loss. The self-reconstruction
38 metric module not only reconstructs query features from support features as
39 in FRN, but more importantly, it also reconstructs query features from them-
40 selves. Such self-reconstruction can effectively augment and diversify query
41 features without introducing artifacts, and as shown in the experiment, it
42 avoids over-reliance on one discriminative feature. After feature reconstruc-
43 tion, both the distance between the original query features and its support-
44 reconstructed features and the distance between the support-reconstructed
45 and self-reconstructed query features are calculated, whose sum is used to
46 match a query sample and the support class. The self-reconstruction network
47 is trained according to the classical cross-entropy loss and a new restrained
48 cross-entropy loss. The latter can prevent the model from producing over-
49 confident predictions which were due to overfitting the training data. By
50 utilizing self-reconstruction and the restrained cross-entropy loss, the pro-
51 posed method can avoid overfitting, as shown in Figure 1.

52 To summarize, the contributions of our work are three-fold:

- 53 1. We are the first to propose using self-reconstruction to increase the
54 diversity of query features, which can effectively expand the represen-
55 tation capability of the learned query feature space and mitigate the
56 overfitting problem.
- 57 2. We further propose a restrained cross-entropy loss, which can be easily

58 equipped with existing metric-based few-shot classification models.

- 59 3. Experiments on five fine-grained datasets demonstrate the superiority
60 of the proposed method, with detailed ablation studies showing that
61 both self-reconstruction and restrained loss are effective in alleviating
62 overfitting.

63 2. Related Work

64 In this section, we first provide a concise review on fine-grained few-shot
65 learning methods. Next, we review two types of methods that are most rele-
66 vant to this work, namely metric-based methods and alignment-based meth-
67 ods. For a more comprehensive review on fine-grained few-shot classification,
68 we refer readers to [14].

69 2.1. Fine-Grained Few-Shot Classification

70 Fine-grained few-shot classification faces the dual challenges of scarce la-
71 beled data and the subtle distinctions between different sub-categories, e.g.,
72 distinguishing between beagle and pug within the category of dog. Global
73 features, which capture image-level concepts, are insufficient to discriminate
74 between fine-grained categories. Therefore, a line of research focus on lo-
75 cal features. DN4 [10], a notable method for fine-grained classification, first
76 utilized intermediate levels of CNN as local features and performed classifica-
77 tion according to the aggregated distances calculated over local features and
78 its k -nearest neighbors. As an improvement of DN4, LSA Net [15] allowed
79 for different scales of local patches to better capture the structure informa-
80 tion and assigned different weights to query patches to suppress the back-
81 ground and highlight the targets. MCL-Katz [16] aggregated local features

82 into global features using weights set as the stationary distribution of local
83 features. AGPF-FSFG [17] constructed multi-scale features and reweighted
84 them via multi-level attention. Our method also makes use of local fea-
85 tures to build the support set or query set, which are then used to generate
86 support-reconstructed query features and self-reconstructed query features.

87 In addition, our method shares the idea of data augmentation. Methods
88 such as Hallucinator [18] and FOT [19] assume that variations in illumination,
89 backgrounds or poses can be shared across classes and thus can be utilized
90 to diversify the limited support samples. In contrast, our work focuses on
91 diversifying the query samples, and we employ the ridge regression technique
92 for efficient self-reconstruction, eliminating the necessity to model intra-class
93 variations.

94 2.2. Metric-Based Methods

95 Metric-based methods constitute one mainstream approach in few shot
96 learning. These methods learn a transferable feature embedding network
97 such that queries can be classified according to the similarity between query
98 features and support features. The similarity can be pre-defined such as
99 by using cosine similarity [20] or learned via a neural network [9]. A pi-
100 oneering metric-based method is the prototypical network [8], which first
101 constructs prototypes as the average of support features and then compares
102 queries with these class representations. To adapt to fine-grained classifica-
103 tion, LMPNet [21] used multiple prototypes per class and constructed pro-
104 totypes as weighted averages of feature embeddings with learnable weights.
105 COMET [22] learned multiple embedding functions, one for each image seg-
106 ment or concept, and accordingly constructed multiple concept prototypes.

107 PHR [23] learned feature embeddings at local, global, and semantic levels
108 and updated prototypes according to novel data. SAPENet [24] obtained
109 more representative prototypes by emphasizing discriminative local features
110 and channels using self-attention and the proposed intra-class attention, re-
111 spectively.

112 Another direction of development in metric-based methods is on the dis-
113 tance measure itself. To list a few, DeepBDC [25] proposed the Brownian
114 distance covariance metric to exploit the joint distributions between support
115 and query features. BSNet [26] combined cosine similarity and relation score
116 for learning more discriminative features in fine-grained images. Tempera-
117 ture Network [27], despite using a single similarity measure, gradually tuned
118 the temperature scaling parameter in the measure, which acts similarly to
119 enforcing a large-margin metric. Different from these methods, our method
120 combines two Euclidean distances – one is between the original query and the
121 support-reconstructed query, and the other is between the self-reconstructed
122 query and the same set of support-reconstructed query.

123 *2.3. Alignment-Based Methods*

124 One issue with metric-based methods is that position information of the
125 embedded features of labeled samples may not correspond to that of unseen
126 samples and therefore the distance calculated directly over these features
127 can be very large, even for samples from the same category. To this end,
128 alignment-based methods have been proposed [12, 28, 13]. LRPABN [28]
129 trained a position transformation matrix to re-arrange the position of sup-
130 port local features to match the query ones. DeepEMD [12] addressed the
131 spatial inconsistency by adopting the earth mover’s distance, which can be

132 interpreted as the optimal matching cost of aligning two sets of local features
133 extracted from a support and a query image. FRN [13] reconstructed query
134 features from support features based on ridge regression, which avoids intro-
135 ducing many parameters as in aforementioned methods and admits a closed-
136 form solution. Building on FRN, LCCRN [29] improved the features by
137 utilizing information from neighborhood pixels. The resulting features were
138 used to construct four cross-reconstruction tasks, whose reconstruction errors
139 were combined using learnable weights. Besides spatial alignment, channel
140 alignment has also been considered [30, 31, 32]. TDM [30] performed channel
141 alignment, which used attention on the support set to highlight class-wise
142 discriminative channels and on a query instance to highlight object-relevant
143 channels. SaberNet [31] adopted Swin Transformer as the feature extrac-
144 tor to capture spatial long-range dependencies between local features and
145 aligned query features and refined prototype features at both spatial and
146 channel levels.

147 In this work, we perform feature reconstruction in a similar manner
148 to FRN [13] by adopting the ridge regression. However, our method self-
149 reconstructs the query features from themselves, which diversifies the query
150 features without introducing artifacts. Consequently, the representation ca-
151 pability of the learned query feature space is expanded, thus alleviating over-
152 fitting and leading to better generalization.

153 3. Self-Reconstruction Network

154 3.1. Problem Definition

155 The few-shot classification problem usually divides a given dataset into
 156 three sub-datasets according to different stages. D_{train} is used for model
 157 training, D_{val} is used for model evaluation in the training stage, and D_{test} is
 158 used for the final test of the trained model. The three sub-datasets contain
 159 different image categories.

160 N -way K -shot is a common setting for few-shot classification, which
 161 means that a classification task consists of N classes and each class has K
 162 labeled samples. At different stages, each sub-dataset (D_{train} , D_{val} and D_{test})
 163 is used to construct a series of tasks and each task contains a support set $\mathcal{S} =$
 164 $\{(x_i, y_i)\}_{i=1}^n$ ($n = N \times K$) and a query set $\mathcal{Q} = \{(x_q, y_q)\}_{q=1}^m$ ($m = N \times M$),
 165 where x denotes the image and y denotes the class label. The support set
 166 is formed by first randomly selecting N classes and then randomly selecting
 167 K images for each of these N classes. The query set is formed by randomly
 168 selecting M images for the same N classes. Features and/or metrics are
 169 learned from the labeled support set and used to perform classification on
 170 the unlabeled query set.

171 3.2. Feature Reconstruction by Ridge Regression

172 Ridge regression is one of the most widely-used penalized regression meth-
 173 ods for analyzing multivariate data with multicollinearity. FRN [13] adopted
 174 this technique to reconstruct the features to solve the few-shot image clas-
 175 sification task and achieved state-of-the-art results. For this reason, we also
 176 propose our model based on the strategy of ridge regression.

177 Ridge regression shrinks the coefficient estimates by adding a penalty on
 178 squared coefficient values, i.e., by minimizing the following penalized residual
 179 sum of squares:

$$\hat{\boldsymbol{\beta}}^{\text{ridge}} = \arg \min_{\boldsymbol{\beta}} \{ \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \lambda \|\boldsymbol{\beta}\|^2 \}, \quad (1)$$

180 where \mathbf{y} is the response vector, \mathbf{X} is the design matrix, $\boldsymbol{\beta}$ is the coefficient
 181 vector, and λ is a parameter controlling the magnitude of the penalty.

182 In the case of feature-map reconstruction here, the response is a matrix.
 183 Therefore, the coefficient should also be a matrix and, following the idea of
 184 ridge regression, the objective function is revised as follows:

$$\hat{\mathbf{A}} = \arg \min_{\mathbf{A}} \{ \|\mathbf{Y} - \mathbf{A}\mathbf{X}\|_F^2 + \lambda \|\mathbf{A}\|_F^2 \}, \quad (2)$$

185 where \mathbf{Y} denotes the response matrix, \mathbf{A} denotes the coefficient matrix, and
 186 $\|\cdot\|_F$ denotes the Frobenius norm. The optimal solution is given by

$$\hat{\mathbf{A}}(\lambda) = \mathbf{Y}\mathbf{X}^\top (\mathbf{X}\mathbf{X}^\top + \lambda\mathbf{I})^{-1}. \quad (3)$$

187 The most expensive cost in calculating Eq. 3 is the inverse operation,
 188 which is $\mathcal{O}(q^3)$ for an $q \times p$ matrix \mathbf{X} . When $p < q$, it is computationally
 189 more efficient to calculate the following equivalent solution, which is obtained
 190 by applying the Woodbury matrix identity [33]:

$$\hat{\mathbf{A}}(\lambda) = \mathbf{Y} (\mathbf{X}^\top \mathbf{X} + \lambda\mathbf{I})^{-1} \mathbf{X}^\top. \quad (4)$$

191 The computational cost of Eq. 4 is $\mathcal{O}(p^3)$, which is smaller than that of Eq. 3
 192 when $p < q$.

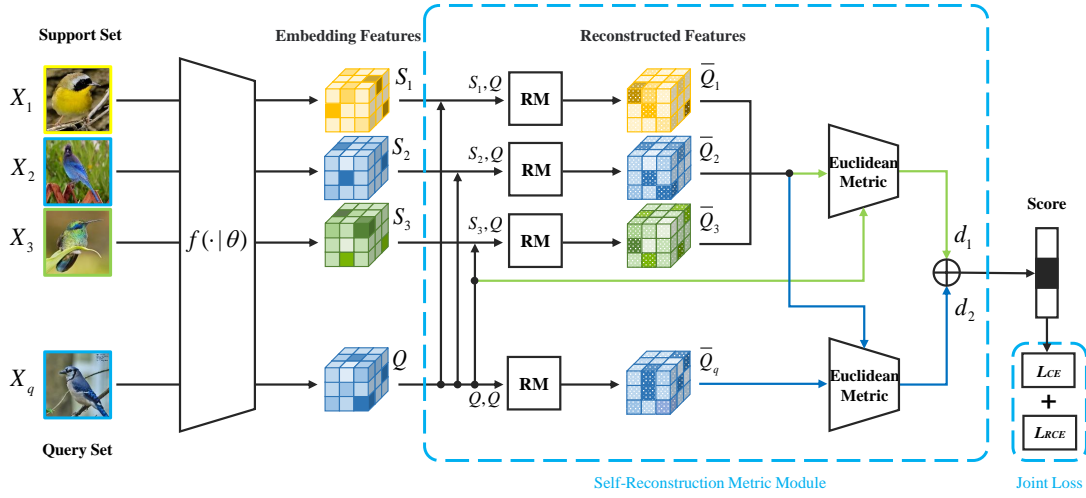


Figure 2: The model architecture of the self-reconstruction network. Support and query images are mapped to the feature space using $f(\cdot|\theta)$. Next, these features are sent to the proposed self-reconstruction metric module, which generates reconstructed features through the reconstruction module RM and calculates the distances. The network is trained according to the proposed joint loss, which combines the cross entropy loss (L_{CE}) and restricted cross entropy loss (L_{RCE}).

193 3.3. Architecture Overview

194 Our network structure, as shown in Figure 2, mainly consists of three
 195 modules. The first module is the feature extraction module, which maps the
 196 original image to the embedding feature. The second is the self-reconstruction
 197 metric module (SRM), which feeds the support features and a query feature
 198 into the feature reconstruction modules (RMs) to obtain two types of re-
 199 constructed query features, one reconstructed by the support features and
 200 the other self-reconstructed by the query feature. Then the distance be-
 201 tween the support-reconstructed query features and the original feature and

202 the distance between the support-reconstructed query features and the self-
203 reconstructed query feature are measured separately. The third is a joint
204 loss module based on the output distances. In the training phase, the net-
205 work is optimized according to the joint loss, which includes the proposed
206 restrained cross-entropy loss L_{RCE} and the standard cross-entropy loss L_{CE} .
207 The pseudo algorithm of our method is provided in Algorithm 1.

Algorithm 1 Training procedure of self-reconstruction network for N -way K -shot classification

Input: training data D_{train} , number of classes N , number of support images per class K , number of query images per class M , number of episodes t , optimizer (SGD).

Output: model parameter θ .

- 1: **for** $e=1$ to t **do**
 - 2: Sample a task \mathcal{T} from D_{train} ;
 - 3: Split \mathcal{T} into the support set \mathcal{S} and the query set \mathcal{Q} ;
 - 4: Use $f(\cdot|\theta)$ to extract class-specific features $\{\mathbf{S}_c\}_{c=1}^N$ and query features $\{\mathbf{Q}_q\}_{q=1}^{NM}$;
 - 5: **for** \mathbf{Q}_q in $\{\mathbf{Q}_q\}_{q=1}^{NM}$ **do**
 - 6: Compute the class-specific reconstructed query feature $\bar{\mathbf{Q}}_c$ using Eq. 5;
 - 7: Compute the self-reconstructed query feature $\bar{\mathbf{Q}}_q$ using Eq. 6;
 - 8: Compute the squared Euclidean distance $d_{c,1}$ between the $\bar{\mathbf{Q}}_c$ and \mathbf{Q}_q using Eq. 9;
 - 9: Compute the squared Euclidean distance $d_{c,2}$ between the $\bar{\mathbf{Q}}_c$ and $\bar{\mathbf{Q}}_q$ using Eq. 10;
 - 10: Obtain the distance between the query and the class c : $d_{c,q} = d_{c,1} + d_{c,2}$;
 - 11: Obtain the final classification probability P using Eq. 12;
 - 12: Compute the loss of model $Loss$ using Eq. 15;
 - 13: Update model parameter θ to minimize $Loss$ using the optimizer.
-

208 *3.4. Self-Reconstruction Metric Module (SRM)*

209 In this paper, we propose a self-reconstruction metric module, which re-
 210 constructs query features from both support features and from the query fea-
 211 tures themselves while enabling both sides to be metricized. Reconstructing

212 query features from support features could achieve spatial alignment between
 213 query and support, and reconstructing query features from themselves could
 214 augment query features and increase task diversity.

215 In the N -way K -shot setting, the feature extraction module outputs fea-
 216 ture maps \mathcal{S}_c and \mathbf{Q}_q , where \mathcal{S}_c contains features from K support samples
 217 of class $c \in C$ and \mathbf{Q}_q ($q = 1, \dots, m$) is the feature of a single query sam-
 218 ple. Moreover, we pool all features from the same class, i.e., applying a
 219 reshaping function to \mathcal{S}_c to map all features of class c into a single matrix
 220 $\mathbf{S}_c: \mathbb{R}^{K \times hw \times l} \rightarrow \mathbb{R}^{Khw \times l}$, where h, w, l are the height, width and number of
 221 channels of the feature map respectively.

222 One core component of the SRM module is the reconstruction of query
 223 features, using both support features and their own query features, as elab-
 224 orated below.

225 The feature map \mathbf{Q}_q can be reconstructed from support features \mathbf{S}_c ac-
 226 cording to Eq. 4, generating the class-specific reconstructed query features
 227 $\bar{\mathbf{Q}}_c$:

$$\bar{\mathbf{Q}}_c = \rho \mathbf{A}(\lambda) \mathbf{S}_c = \rho \mathbf{Q}_q (\mathbf{S}_c^\top \mathbf{S}_c + \lambda \mathbf{I})^{-1} \mathbf{S}_c^\top \mathbf{S}_c, \quad (5)$$

228 where ρ is a learnable re-scaling parameter. Eq. 5 addresses the misalignment
 229 between query and support features.

230 Applying the same technique, we can reconstruct \mathbf{Q}_q from its own feature,
 231 so as to map the query feature to a reconstruction space. Note that the
 232 feature map $\mathbf{Q}_q \in \mathbb{R}^{hw \times l}$ comes from a single query image, rather than all
 233 query images. The self-reconstructed query feature $\bar{\mathbf{Q}}_q$ is obtained as follows:

$$\bar{\mathbf{Q}}_q = \rho \mathbf{A}(\lambda) \mathbf{Q}_q = \rho \mathbf{Q}_q (\mathbf{Q}_q^\top \mathbf{Q}_q + \lambda \mathbf{I})^{-1} \mathbf{Q}_q^\top \mathbf{Q}_q. \quad (6)$$

The parameters ρ and λ in Eqs. 5 and 6 are not shared and they are updated in the same way as parameters in the feature extraction module during training. In order to ensure ρ and λ are positive, they are converted to e^α and e^β respectively, with α and β initialized to zero:

$$\lambda = \frac{Khw}{l}e^\alpha, \quad (7)$$

$$\rho = e^\beta. \quad (8)$$

After obtaining the reconstructed features, we carry out distance calculation as in metric-based methods. In this step, the class-specific reconstructed query features $\bar{\mathbf{Q}}_c$ and the self-reconstructed query features $\bar{\mathbf{Q}}_q$ are used differently, where $\bar{\mathbf{Q}}_c$ serves as class-specific prototypes for spatially aligning support features with the original query features \mathbf{Q}_q and $\bar{\mathbf{Q}}_q$ serves as additional samples for expanding the representation capability of the learned query feature space. More specifically, we calculate the squared Euclidean distance $d_{c,1}$ between the support-reconstructed query features $\bar{\mathbf{Q}}_c$ and the original query feature \mathbf{Q}_q , the squared Euclidean distance $d_{c,2}$ between $\bar{\mathbf{Q}}_c$ and the self-reconstructed query feature $\bar{\mathbf{Q}}_q$, and finally sum up the two distances to get the distance $d_{c,q}$, which represents the distance between the query and the class c :

$$d_{c,1} = \|\bar{\mathbf{Q}}_c - \mathbf{Q}_q\|_F^2, \quad (9)$$

$$d_{c,2} = \|\bar{\mathbf{Q}}_c - \bar{\mathbf{Q}}_q\|_F^2, \quad (10)$$

$$d_{c,q} = d_{c,1} + d_{c,2}. \quad (11)$$

234 *3.5. Loss Functions*

235 By using the above distance, the final classification probability can be
236 obtained as:

$$P(y_q = c | x_q) = \frac{\exp(-\tau d_{c,q})}{\sum_{c' \in C} \exp(-\tau d_{c',q})}, \quad (12)$$

237 where τ is a learnable hyperparameter that controls the sharpness of the
238 metric distance.

239 In the training phase, we use a joint loss integrating two loss functions to
240 optimize the model. One is the widely-used cross-entropy (CE) loss:

$$L_{CE} = -\frac{1}{m} \sum_{q=1}^m (\mathbf{y}_q^\top \log(\mathbf{p}_q)), \quad (13)$$

241 where \mathbf{y}_q denotes the one-hot vector, \mathbf{p}_q denotes the vector of predicted
242 probability, and m is the number of query samples.

243 To further alleviate the overfitting problem, we propose to restrain the
244 cross-entropy loss on the training classes. Specifically, we design a new re-
245 strained cross-entropy (RCE) loss in Eq. 14, which has the opposite effect to
246 the CE loss and can prevent the learned model from being over-confident in
247 its predictions on the training set:

$$L_{RCE} = -\frac{1}{m} \sum_{q=1}^m \left((\mathbf{1} - \mathbf{y}_q)^\top \log(\mathbf{p}_q) \right). \quad (14)$$

248 The final loss is obtained by combining the CE loss and RCE loss:

$$\begin{aligned} Loss &= L_{CE} + kL_{RCE} \\ &= -\frac{1}{m} \sum_{q=1}^m \left(\left[\mathbf{y}_q^\top + k(\mathbf{1} - \mathbf{y}_q)^\top \right] \log(\mathbf{p}_q) \right), \end{aligned} \quad (15)$$

249 where $0 < k < 1$ adjusts the influence of the RCE loss. When combined
250 with the CE loss as in Eq. 15, the RCE loss can effectively restrict the

251 decrease of CE loss. Moreover, when $0 < k < 1$, the joint loss trades off
 252 between assigning the probability of 1 to the correct class and assigning equal
 253 probabilities to all classes. Therefore, the training process will encourage the
 254 query sample to be correctly classified as usual, but not necessarily with a
 255 high classification probability, thus preventing the model from overfitting to
 256 noises or non-generalizable features and improving generalization [34].

257 4. Experiments and Discussions

258 To evaluate the effectiveness of our proposed approach, we present in this
 259 section experiments designed for six purposes:

- 260 • To compare our proposed method with state-of-the-art methods for the
 261 task of few-shot fine-grained image classification (Sec. 4.2);
- 262 • To study the effectiveness of each branch of our network in image classi-
 263 fication and in mitigating the overfitting problem (Sec. 4.3, 4.5.1, 4.5.3);
- 264 • To investigate the stability of the classification accuracy (Sec. 4.4);
- 265 • To evaluate the discriminative power of the learned features (Sec. 4.5.4,
 266 4.5.2);
- 267 • To illustrate the effect of feature reconstruction (Sec. 4.5.5);
- 268 • To assess the computational complexity (Sec. 4.6).

269 4.1. Implementation Details

270 We conduct experiments on the following five popular benchmark datasets:
 271 CUB-200-2011 (CUB) [35], Stanford-Cars (Cars) [36], Stanford-Dogs (Dogs) [37],

272 Flowers [38], and FGVC-Aircraft (Aircraft) [39]. The CUB dataset is con-
273 sistent with that in [13], and the images are pre-cropped into the bounding
274 boxes provided. All datasets are divided into base, validation, and novel
275 datasets according to the ratio of 2:1:1.

276 For the feature extractor, we adopt two widely-used backbones: ResNet-
277 12 and ResNet-18 [1, 13]. The ResNet-12 structure has 4 residual blocks, and
278 each residual block contains 3 convolutional layers with 3×3 convolution
279 kernel. Each convolutional layer is followed by a batch normalization, and
280 the first convolutional layer is followed by a ReLU nonlinearity, while a 2×2
281 maximum pooling layer is added at the end of each residual block. The input
282 dimension of the network is $3 \times 84 \times 84$, and the output feature dimension
283 is $640 \times 5 \times 5$ after feature extraction.

284 Unlike the ResNet-18 network in [40], our ResNet-18 network is modified
285 on the ResNet-12 network. [40] uses a 7×7 convolution kernel in the first
286 convolutional layer, which is not conducive to fine-grained feature extraction
287 as the subtle discriminative features may locate in a tiny region in fine-grained
288 image classification. Our ResNet-18 network, like ResNet-12, has 4 residual
289 blocks, but the first two residual blocks are divided into two sub-residual
290 blocks; each sub-residual block contains 3 convolutional layers with 3×3
291 convolution kernel. The rest of the structure is similar to ResNet-12.

292 Throughout the experiments, we use the setting of 10-way 5-shot for
293 training, 5-way 5-shot for validating, and 5-way 1-shot and 5-way 5-shot for
294 testing. The query set contains 16 images in all three phases. In the training
295 phase, we use the SGD optimizer on all datasets with an initial learning
296 rate of 0.1 and momentum of 0.9, and train 1,200 epochs in total. In the

297 testing phase, the accuracy score is obtained by averaging over 10,000 trials.
 298 The coefficient k in Eq. 15 is chosen separately for each dataset according to
 299 the accuracy on the validation set. The code of our method is available at
 300 <https://github.com/liz-lut/SRM-main>.

301 4.2. Comparison with State-of-the-arts

302 To validate the efficacy of our method, we compare it with the following
 303 13 methods: 1) four classical few-shot learning methods, namely Match-
 304 ingNet [20], ProtoNet [8], RelationNet [9], Baseline++ [1]; 2) three state-
 305 of-the-art metric-based methods, namely DeepEMD [12], MCL-Katz [16],
 306 DeepBDC [25]; 3) five state-of-the-art fine-grained few-shot methods, namely
 307 VFD [41], FRN [13], TDM [30], AGPF-FSFG [17], LCCRN [29]; and 4)
 308 one method that targets the generalizability, namely Neg-margin [6]. Ta-
 309 ble 1 and Table 2 list the performance of these methods with ResNet-12
 310 and ResNet-18 as the backbone, respectively. From the experimental results,
 311 we observe that our approach achieves the highest classification accuracy on
 312 three datasets for both 5-way 1-shot and 5-way 5-shot classification and is the
 313 second-best in most other cases. In particular, compared with FRN which
 314 reconstructs query features only from support features, our approach, adding
 315 self-reconstruction and restrained cross-entropy loss, shows 2.44% gains in 5-
 316 way 1-shot classification and 2.75% gains in the 5-way 5-shot classification on
 317 the Flowers dataset with the ResNet-12 backbone; similar improvements can
 318 be observed with the ResNet-18 backbone. This demonstrates the effective-
 319 ness of our proposed method and encourages the use of self-reconstruction
 320 and restrained cross-entropy loss.

Table 1: Comparison of few-shot classification methods on CUB-200-2011, Flowers, FGVC Aircraft, Stanford-Cars, and Stanford-Dogs datasets. All experiments adopt ResNet-12 as the backbone network. Mean accuracy and 95% confidence interval are reported. The best-performing methods are shown in bold and the second best ones are underlined.

Method	<i>CUB</i>	<i>Flowers</i>	<i>Aircraft</i>	<i>Cars</i>	<i>Dogs</i>
	5-Way 1-Shot Accuracy (%)				
MatchingNet (NeurIPS 16')[20]	73.02 ± 0.88	75.70 ± 0.88	82.2 ± 0.80	75.70 ± 0.88	66.48 ± 0.88
ProtoNet (NeurIPS 17')[8]	79.64 ± 0.20	75.41 ± 0.22	86.57 ± 0.18	82.29 ± 0.20	72.85 ± 0.22
RelationNet (CVPR 18')[9]	63.94 ± 0.92	69.51 ± 1.01	74.2 ± 1.04	46.04 ± 0.91	47.35 ± 0.88
Baseline++ (CVPR 19')[1]	64.62 ± 0.98	69.03 ± 0.92	74.51 ± 0.90	67.92 ± 0.92	59.64 ± 0.89
DeepEMD (CVPR 20')[12]	71.11 ± 0.31	70.00 ± 0.35	69.86 ± 0.30	73.30 ± 0.29	67.59 ± 0.30
VFD (ICCV 21')[41]	79.12 ± 0.83	76.20 ± 0.92	76.88 ± 0.85	77.52 ± 0.85	76.24 ± 0.87
FRN (CVPR 21')[13]	83.16 ± 0.19	81.07 ± 0.20	88.04 ± 0.17	86.48 ± 0.18	76.49 ± 0.21
TDM (CVPR 22')[30]	82.41 ± 0.19	82.85 ± 0.19	88.35 ± 0.17	87.04 ± 0.17	76.20 ± 0.21
MCL-Katz (CVPR 22')[16]	85.63 ± 0.00	76.55 ± 0.00	87.69 ± 0.00	85.04 ± 0.00	71.49 ± 0.00
DeepBDC (CVPR 22')[25]	79.32 ± 0.43	80.57 ± 0.49	83.14 ± 0.41	83.24 ± 0.42	76.61 ± 0.46
AGPF-FSFG (PR 22')[17]	78.54 ± 0.83	77.92 ± 0.94	82.65 ± 0.89	83.94 ± 0.76	72.06 ± 0.91
LCCRN (TCSVT 23')[29]	82.97 ± 0.19	82.39 ± 0.19	88.48 ± 0.17	87.04 ± 0.17	75.87 ± 0.20
Ours	<u>83.82 ± 0.18</u>	83.51 ± 0.19	88.94 ± 0.16	88.02 ± 0.16	<u>76.54 ± 0.21</u>
5-Way 5-Shot Accuracy (%)					
MatchingNet (NeurIPS 16')[20]	85.17 ± 0.60	87.61 ± 0.55	88.99 ± 0.50	87.61 ± 0.55	79.57 ± 0.63
ProtoNet (NeurIPS 17')[8]	91.15 ± 0.11	89.46 ± 0.14	93.51 ± 0.09	93.11 ± 0.10	86.54 ± 0.13
RelationNet (CVPR 18')[9]	77.87 ± 0.64	86.84 ± 0.56	86.62 ± 0.55	68.52 ± 0.78	66.20 ± 0.74
Baseline++ (CVPR 19')[1]	81.15 ± 0.61	85.72 ± 0.63	88.06 ± 0.44	84.17 ± 0.58	77.36 ± 0.62
DeepEMD (CVPR 20')[12]	86.30 ± 0.19	83.63 ± 0.26	85.17 ± 0.28	88.37 ± 0.17	81.13 ± 0.20
VFD (ICCV 21')[41]	91.48 ± 0.39	89.90 ± 0.53	88.77 ± 0.46	90.76 ± 0.46	88.00 ± 0.47
FRN (CVPR 21')[13]	92.59 ± 0.10	92.52 ± 0.11	94.21 ± 0.08	94.78 ± 0.08	88.22 ± 0.12
TDM (CVPR 22')[30]	92.37 ± 0.10	93.60 ± 0.10	94.47 ± 0.08	96.11 ± 0.07	88.32 ± 0.12
MCL-Katz (CVPR 22')[16]	93.18 ± 0.00	90.31 ± 0.00	93.28 ± 0.00	93.92 ± 0.00	85.24 ± 0.00
DeepBDC (CVPR 22')[25]	92.10 ± 0.24	92.82 ± 0.24	93.25 ± 0.18	94.97 ± 0.19	90.22 ± 0.23
AGPF-FSFG (PR 22')[17]	89.85 ± 0.44	91.96 ± 0.45	89.25 ± 0.49	94.11 ± 0.36	84.83 ± 0.50
LCCRN (TCSVT 23')[29]	93.63 ± 0.10	94.24 ± 0.09	94.61 ± 0.08	96.19 ± 0.07	88.36 ± 0.11
Ours	<u>93.45 ± 0.10</u>	95.27 ± 0.08	94.88 ± 0.08	96.23 ± 0.07	<u>88.52 ± 0.12</u>

Table 2: Comparison of few-shot classification methods on CUB-200-2011, Flowers, FGVC Aircraft, Stanford-Cars, and Stanford-Dogs datasets. All experiments adopt ResNet-18 as the backbone network. Mean accuracy and 95% confidence interval are reported. The best-performing methods are shown in bold and the second best ones are underlined.

Method	<i>CUB</i>	<i>Flowers</i>	<i>Aircraft</i>	<i>Cars</i>	<i>Dogs</i>
	5-Way 1-Shot Accuracy (%)				
MatchingNet (NeurIPS 16')[20]	72.88 ± 0.89	76.07 ± 0.82	82.84 ± 0.81	75.03 ± 0.95	65.59 ± 0.95
RelationNet (CVPR 18')[9]	68.82 ± 1.04	69.04 ± 0.97	74.76 ± 0.97	64.08 ± 1.05	54.21 ± 1.00
Baseline++ (CVPR 19')[1]	65.67 ± 0.95	67.90 ± 0.96	75.92 ± 0.88	67.41 ± 0.99	62.54 ± 0.87
Neg-margin (CVPR 19')[6]	72.51 ± 0.82	76.34 ± 0.89	77.40 ± 0.86	76.04 ± 0.81	68.86 ± 0.83
FRN (CVPR 21')[13]	83.40 ± 0.19	81.22 ± 0.21	87.89 ± 0.18	87.63 ± 0.17	77.53 ± 0.21
TDM (CVPR 22')[30]	83.25 ± 0.19	82.31 ± 0.20	87.91 ± 0.17	87.69 ± 0.17	76.59 ± 0.21
MCL-Katz (CVPR 22')[16]	85.84 ± 0.00	76.57 ± 0.00	88.44 ± 0.00	86.12 ± 0.00	72.07 ± 0.00
DeepBDC (CVPR 22')[25]	81.85 ± 0.42	81.07 ± 0.50	85.22 ± 0.41	85.48 ± 0.40	78.81 ± 0.43
AGPF-FSFG (PR 22')[17]	79.02 ± 0.83	78.69 ± 0.84	85.02 ± 0.86	84.68 ± 0.78	73.61 ± 0.91
LCCRN (TCSVT 23')[29]	82.80 ± 0.19	82.86 ± 0.19	88.66 ± 0.18	86.24 ± 0.18	77.29 ± 0.20
Ours	<u>84.14 ± 0.18</u>	83.25 ± 0.19	89.14 ± 0.17	88.70 ± 0.16	<u>77.57 ± 0.20</u>
5-Way 5-Shot Accuracy (%)					
MatchingNet (NeurIPS 16')[20]	85.25 ± 0.57	87.46 ± 0.51	88.77 ± 0.54	87.02 ± 0.56	80.94 ± 0.60
RelationNet (CVPR 18')[9]	82.68 ± 0.58	85.46 ± 0.58	87.45 ± 0.55	91.45 ± 0.44	80.42 ± 0.62
Baseline++ (CVPR 19')[1]	81.53 ± 0.58	84.34 ± 0.62	88.13 ± 0.47	85.50 ± 0.58	79.04 ± 0.61
Neg-margin (CVPR 19')[6]	89.25 ± 0.43	90.83 ± 0.47	90.92 ± 0.39	93.06 ± 0.38	85.75 ± 0.52
FRN (CVPR 21')[13]	92.69 ± 0.10	92.33 ± 0.11	93.96 ± 0.09	95.35 ± 0.08	89.05 ± 0.11
TDM (CVPR 22')[30]	92.98 ± 0.10	93.46 ± 0.11	94.28 ± 0.08	96.06 ± 0.07	88.87 ± 0.11
MCL-Katz (CVPR 22')[16]	93.29 ± 0.00	90.06 ± 0.00	93.22 ± 0.00	93.35 ± 0.00	85.46 ± 0.00
DeepBDC (CVPR 22')[25]	93.00 ± 0.24	93.19 ± 0.24	94.26 ± 0.16	95.84 ± 0.16	91.33 ± 0.22
AGPF-FSFG (PR 22')[17]	89.92 ± 0.42	92.78 ± 0.40	91.94 ± 0.47	94.87 ± 0.33	85.68 ± 0.52
LCCRN (TCSVT 23')[29]	93.60 ± 0.10	93.87 ± 0.10	94.72 ± 0.07	96.34 ± 0.07	<u>89.54 ± 0.10</u>
Ours	<u>93.58 ± 0.09</u>	94.70 ± 0.09	94.79 ± 0.08	96.40 ± 0.07	89.20 ± 0.11

Table 3: Ablation studies on the Flowers and Cars datasets. All experiments adopt the ResNet-12 backbone. Mean accuracy and its 95% confidence interval are reported. ProtoNet with both SRM and L_{RCE} is our proposed network, noted in the table as Ours.

	Method	Flowers	Cars
5-Way 1-Shot	ProtoNet	75.41 \pm 0.22	82.29 \pm 0.20
	ProtoNet+ L_{RCE}	76.62 \pm 0.22	85.25 \pm 0.19
	ProtoNet+SRM	82.67 \pm 0.19	85.74 \pm 0.18
	Ours	83.51 \pm 0.18	88.02 \pm 0.16
5-Way 5-Shot	ProtoNet	89.46 \pm 0.14	93.11 \pm 0.10
	ProtoNet+ L_{RCE}	89.18 \pm 0.14	93.55 \pm 0.10
	ProtoNet+SRM	94.89 \pm 0.08	95.37 \pm 0.08
	Ours	95.27 \pm 0.08	96.23 \pm 0.07

321 4.3. Ablation Studies

322 To further demonstrate the effectiveness of our proposed network, we
 323 conduct 5-way 1-shot and 5-way 5-shot experiments on four model structures
 324 on Flowers and Cars datasets with ResNet-12 backbone. The four models are
 325 as follows: ProtoNet, ProtoNet with our proposed self-reconstruction metric
 326 module (SRM), ProtoNet with our proposed loss L_{RCE} , and ProtoNet with
 327 both SRM and L_{RCE} (i.e., Ours). As shown in Table 3, when using the
 328 SRM module on the ProtoNet, the accuracy is higher than ProtoNet in all
 329 cases. When using loss L_{RCE} on ProtoNet, the accuracy is slightly lower
 330 than ProtoNet on the 1-shot task on Flowers, while it improves in other
 331 cases. Finally, when the SRM module and L_{RCE} are used together (Ours), the
 332 accuracy increases dramatically. This highlights the importance of combining
 333 both techniques to get a better network.

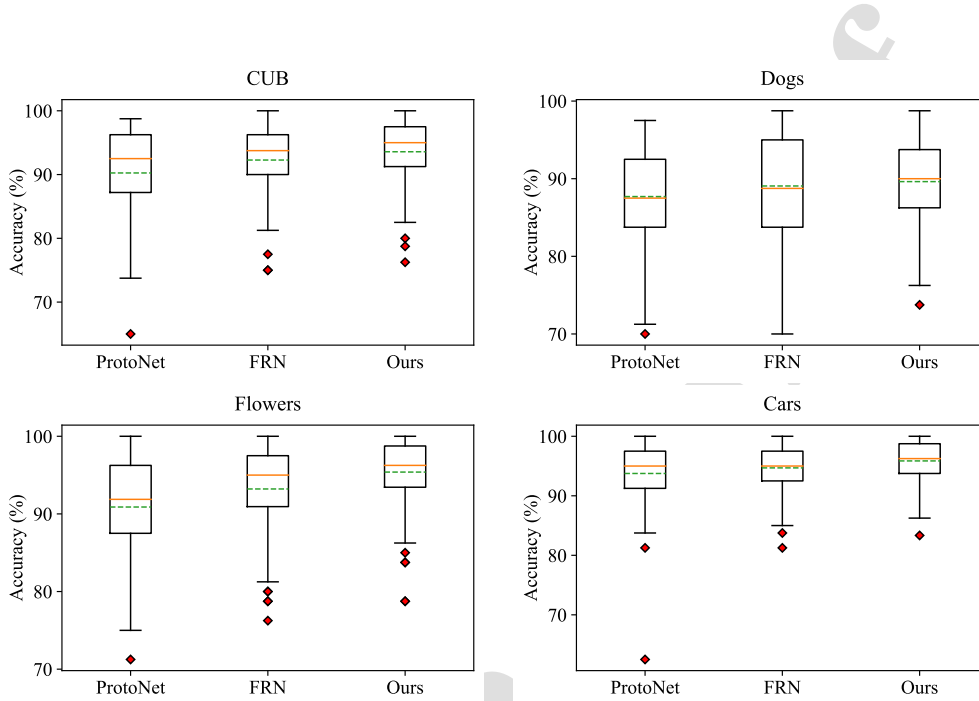


Figure 3: Boxplots of classification accuracy of ProtoNet, FRN, and our proposed method on CUB, Dogs, Flowers, and Cars datasets. All experiments are based on a 5-way 5-shot classification setup with the ResNet-12 backbone. Each method has been evaluated for 100 rounds, and the distributions of test accuracy are shown via boxplots. In each boxplot, the central orange line marks the median, and the green dashed line marks the mean; the edges of the box are the 25th and 75th percentiles, respectively; and the outliers are marked in red individually.

334 4.4. Stability Analysis

335 4.4.1. Boxplots of Classification Accuracy

336 Tables 1 and 2 list the mean value and 95% confidence interval of test
 337 accuracy. To provide a further insight into the classification performance, we
 338 present the boxplots of classification accuracy of ProtoNet, FRN, and our
 339 proposed method on CUB, Dogs, Flowers, and Cars datasets in Figure 3.

340 All experiments are based on a 5-way 5-shot classification setup with the
341 ResNet-12 backbone. On each dataset, 100 tasks were randomly selected
342 and used to evaluate the methods.

343 As can be seen from Figure 3, our method is obviously better than the
344 other two methods in terms of the median (orange lines) and mean (green
345 dashed lines). Moreover, the range of classification accuracy excluding out-
346 liers (i.e., the distance between whiskers) of our method is significantly nar-
347 rower than that of ProtoNet and FRN, indicating that our method is more
348 stable and has higher confidence. Furthermore, looking at the outliers (red
349 points), we can see that the classification accuracy of our method is also
350 better than the other two methods on the worst-performing tasks.

351 4.4.2. Classification Accuracy Under Different Shots

352 To further evaluate the stability of our method, we calculate the test
353 accuracy of ProtoNet, FRN, and our proposed method in different K -shot
354 settings on CUB, Dogs, Flowers, and Cars datasets. Figure 4 shows the
355 classification accuracy in different K -shot settings for the 5-way classification
356 task. While our method is only slightly better than ProtoNet and FRN on
357 the Dogs dataset, its advantage is more obvious on other datasets. The
358 advantage over ProtoNet is more pronounced for 1-shot classification, which
359 is reasonable since overfitting is more likely to occur when only one support
360 image is provided. The advantage over FRN holds across different number
361 of shots.

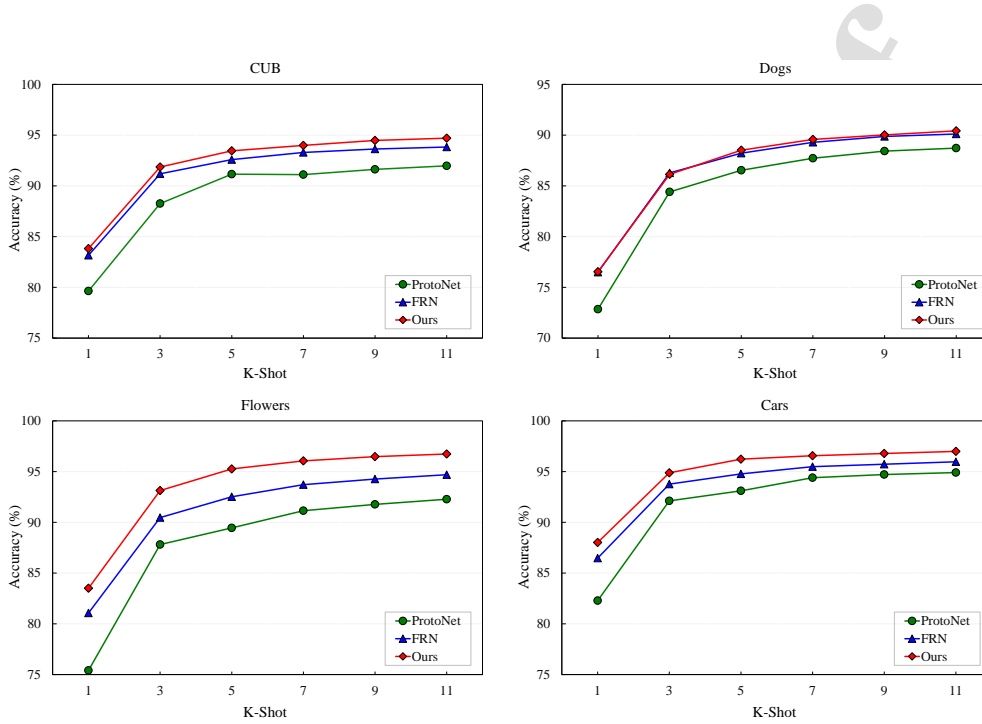


Figure 4: Test accuracy of ProtoNet, FRN, and our proposed method on CUB, Dogs, Flowers, and Cars datasets in different K -shot settings. All experiments are based on a 5-way classification setup with the ResNet-12 backbone.

362 4.5. Visualization Analysis

363 4.5.1. Visualization of Validation Loss

364 As discussed in the introduction, FRN suffers from overfitting and, to
 365 overcome this issue, we propose both self-reconstruction and restrained cross-
 366 entropy loss. Figure 5 shows the cross-entropy losses of four model struc-
 367 tures on Cars; the four models are FRN, ours without the proposed self-
 368 reconstruction metric module (SRM), ours without the proposed loss L_{RCE} ,
 369 and our proposed method. Each data point is calculated on the validation
 370 set after each epoch. As we can see, FRN reaches the minimum value of

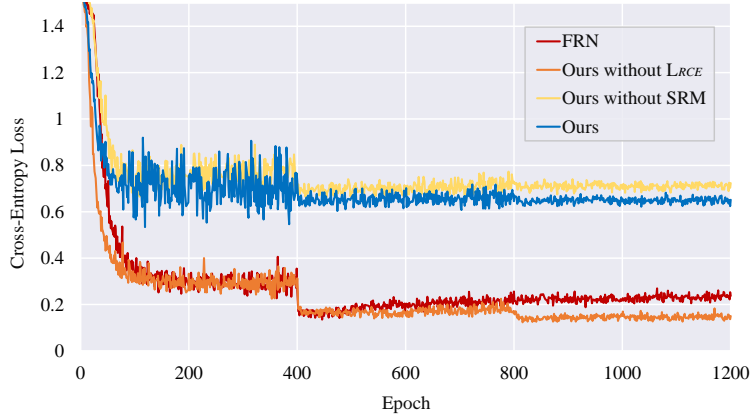


Figure 5: Validation loss of FRN and our method on the Cars dataset. All experiments are based on a 5-way 5-shot classification setup with the ResNet-12 backbone.

371 cross-entropy loss at around the 420th epoch and gradually increases after-
 372 ward, indicating the occurrence of overfitting. However, the cross-entropy
 373 loss does not tend to increase after the 400th epoch for our method and its
 374 two variants. Therefore, both the proposed SRM structure and restrained
 375 cross-entropy loss L_{RCE} can effectively mitigate the overfitting problem.

376 4.5.2. Visualization of Classification Probabilities

377 To further evaluate the effectiveness of the proposed RCE loss in prevent-
 378 ing over-confident predictions, we present a heat map of the classification
 379 probabilities of query samples predicted by ProtoNet, FRN, and our method
 380 (defined by Eq. 12) on the test set of Cars dataset. As shown in Figure 6, in
 381 each confusion matrix, the vertical axis shows 5 classes in a task, and the hor-
 382 izontal axis shows query samples in the 5 classes with each class containing
 383 16 query samples. The main diagonal line indicates the correct classification.
 384 Warmer color represents higher probability score.

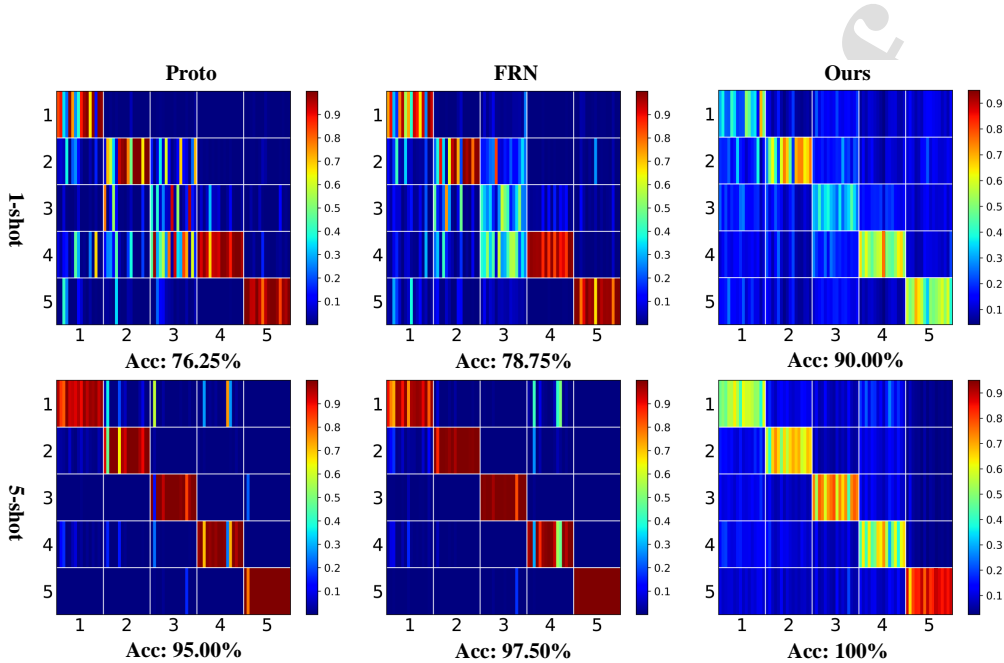


Figure 6: Visualization of classification probabilities predicted by ProtoNet, FRN and our proposed method on the test set of Cars dataset. In each confusion matrix, the vertical axis shows 5 classes in a task and the horizontal axis shows query samples in the 5 classes with each class containing 16 query samples. Warmer color means higher probability. Classification accuracy of each method is displayed below its respective figure.

385 Firstly, we notice that ProtoNet and FRN make a number of incorrect
 386 predictions for 1-shot classification. Taking class 3 as an example, the prob-
 387 abilities of assigning queries of class 3 to the correct class are often lower
 388 than the probabilities of assigning them to other classes, resulting in incor-
 389 rect predictions. These match the classification accuracy shown below their
 390 respective figures. Secondly, we see that our method makes correct predic-
 391 tions on more query samples as the correct class is often assigned with the
 392 highest classification probabilities. Moreover, the correct classification prob-
 393 abilities of our method are typically not very high, mostly ranging from 0.4

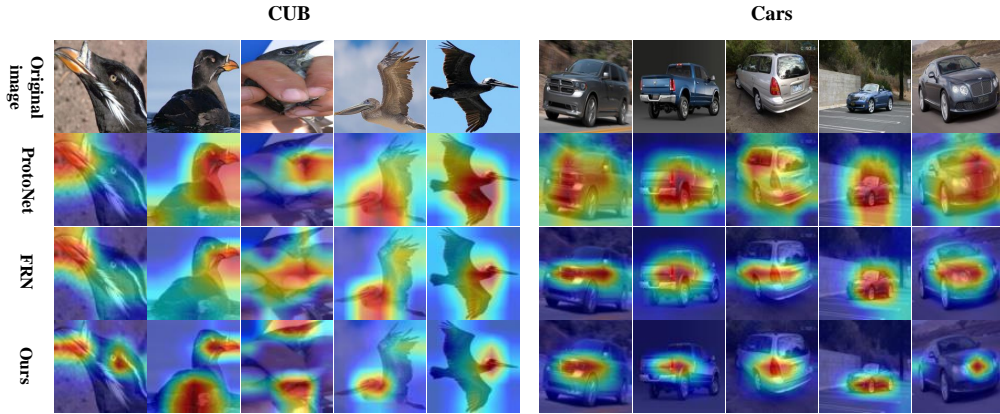


Figure 7: Feature visualization under ProtoNet, FRN, and our method on the CUB (left) and Cars (right) datasets. Red indicates the learned discriminative features. The redder the region, the more discriminative the learned features.

394 to 0.7. This is a consequence of including the restricted cross-entropy in the
 395 loss function to prevent over-confident predictions, aligning with the findings
 396 drawn from Figure 5.

397 4.5.3. Visualization of Discriminative Feature Regions

398 To demonstrate the effectiveness of our proposed method on feature ex-
 399 traction, we show the feature regions extracted from the original image by
 400 visualizing them using Grad-CAM [42].

401 Figure 7 shows that the discriminative feature regions extracted by our
 402 method are more concentrated compared with ProtoNet and FRN. It also ex-
 403 tracts discriminative features that ProtoNet and FRN do not capture well, so
 404 our method could generalize better and become more robust in the presence
 405 of noises in some spatial regions.

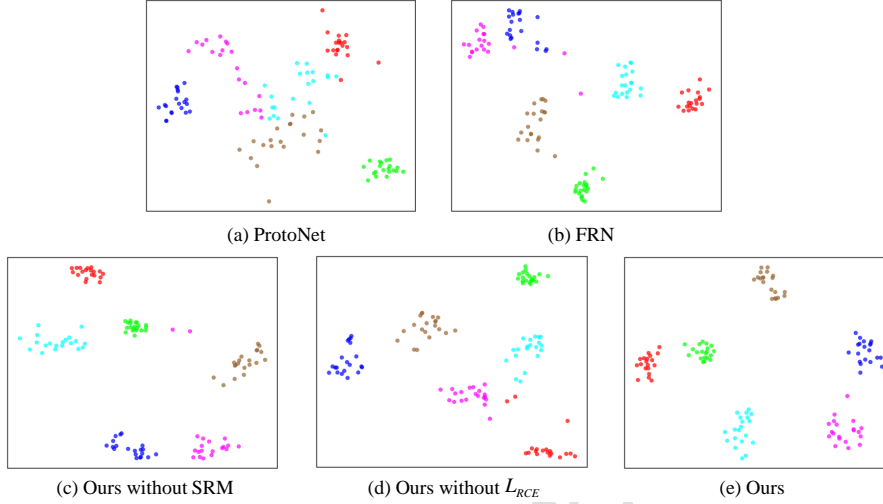


Figure 8: The t-SNE visualization of convolutional feature of samples by ProtoNet, FRN, and Ours with the ResNet-12 backbone on the Flowers dataset.

406 4.5.4. Visualization of Feature Embeddings

407 We also visualize the feature embeddings learned by ProtoNet, FRN, and
 408 our method by t -distributed stochastic neighbor embedding (t -SNE) [43] on
 409 the Flowers dataset and show the results in Figure 8. As can be seen from
 410 the figure, the distribution of feature embeddings of our method is relatively
 411 concentrated, and the inter-class margin is large. The two variants of our
 412 method, ours without self-reconstruction metric module (SRM) or loss L_{RCE} ,
 413 also show larger inter-class margin compared with FRN and ProtoNet. Thus,
 414 our approach improves separability between classes, which is more favorable
 415 for the task of few-shot classification.

416 *4.5.5. Visualization of Reconstructed Features*

417 To provide a deeper insight into the reconstruction module, we train an
 418 image generator to recover images from query features. Specifically, three
 419 types of query features are considered, namely original features obtained di-
 420 rectly after the feature extraction module, self-reconstructed query features,
 421 and support-reconstructed query features. To map features back to images,
 422 an inverted ResNet-12 decoder is trained on the 5-way 5-shot task to decode
 423 features of dimension $640 \times 5 \times 5$ into $3 \times 84 \times 84$. We use the Adam optimizer
 424 and L_1 loss to train the decoder. The initial learning rate of the optimizer
 425 is set as 0.01, and the batch size is set as 200. After training for 2,000
 426 epochs, we select the parameter with the minimum loss as the parameter of
 427 the decoder.

428 In Figure 9, panel (a) shows 5 query images of the same class, and panel
 429 (e) shows 5 support images from 5 classes, where images in the third column
 430 has the same class as the query. Panels (b)-(d) are images generated from
 431 query features after feature extraction, i.e., Q_q , from self-reconstructed query
 432 features, i.e., \bar{Q}_q , and from support-reconstructed query features, i.e., \bar{Q}_c .
 433 The (i, j) th image in panel (d) is the recovered image for i th query by using
 434 query features reconstructed from the j th support class.

435 As we can see from Figure 9, firstly, images recovered from self-reconstructed
 436 query features (panel (c)) are very similar to those recovered from the origi-
 437 nal query features (panel (b)), but there are some differences. For example,
 438 in the second image, feather is more blurred in the self-reconstructed case
 439 than in the original case; in the fifth image, beak is missing in the self-
 440 reconstructed case. This suggests that the self-reconstructed query features

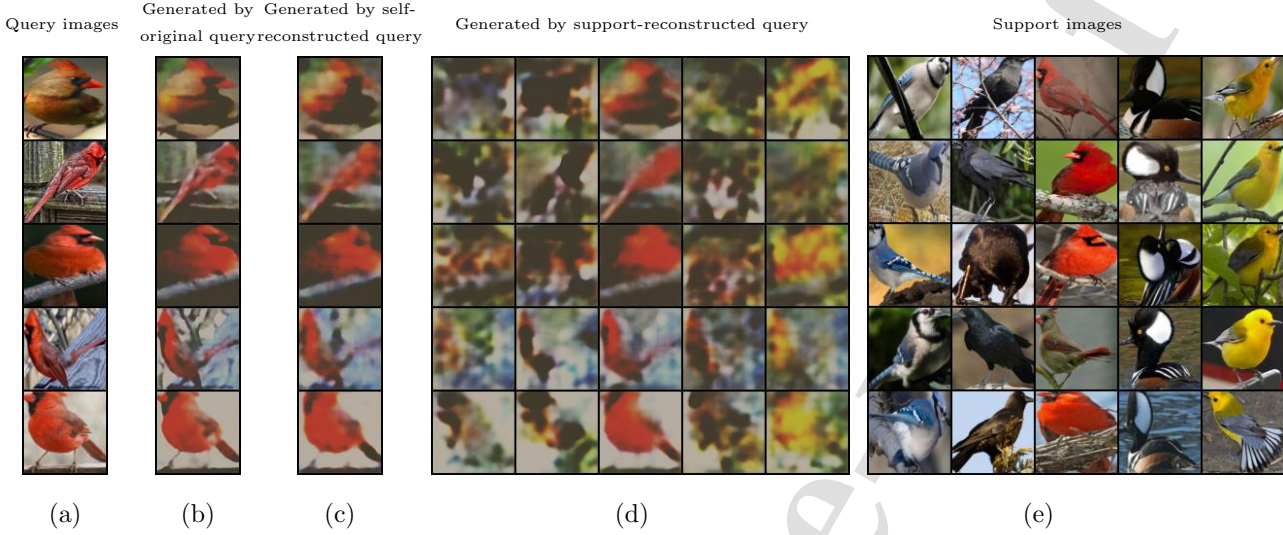


Figure 9: Visualization of image reconstruction for CUB. Panel (a) shows 5 query images. Panel (e) shows 25 support images (5 in each class on a column). Panels (b)-(d) are images generated by using original features after feature extraction, self-reconstructed query features, and support-reconstructed query features, respectively. Self-reconstruction increases the diversity, especially the hard samples, and reconstruction from the same class is more credible than that from different classes.

441 increase the feature diversity, especially generating some hard samples for
 442 classification, and therefore help alleviate overfitting. Secondly, as seen from
 443 panel (d), images generated from the query features that were reconstructed
 444 from same-class support features (i.e., the middle column in that panel) are
 445 much more similar to the ground-truth query images than those based on
 446 different-class reconstructions (i.e., other four columns in that panel). This
 447 shows that our proposed reconstruction module preserves the merit of FRN
 448 in reconstructing semantically faithful images from support images.

Table 4: Comparison of model complexity and computational cost.

Method	Number of parameters	Training time (per epoch)		
		CUB	Flowers	Cars
FRN	12,424,323	11.96s	12.34s	21.03s
Ours	12,424,325	12.45s	12.62s	22.23s

449 4.6. Computational Cost

450 We evaluate the computational complexity of the proposed method. Ta-
 451 ble 4 lists the model complexity in terms of the number of parameters and
 452 the computational cost in terms of the training time per epoch. All methods
 453 were implemented by using PyTorch on an NVIDIA RTX 3090 GPU. Com-
 454 pared with FRN, our method introduces an additional self-reconstruction
 455 step and distance calculation. The self-reconstruction step only adds two
 456 new parameters, namely ρ and λ in Eq. 6. Moreover, the increase in training
 457 time is marginal.

458 5. Conclusion

459 In this paper, we proposed a self-reconstruction network for few-shot fine-
 460 grained image classification. Our innovation includes enhancing feature di-
 461 versity by self-reconstructing query samples and introducing restrained cross-
 462 entropy loss to mitigate overfitting. Extensive experiments on five bench-
 463 mark fine-grained datasets demonstrate the efficacy of our method with the
 464 state-of-the-art performance achieved on both 5-way 1-shot and 5-way 5-shot
 465 classification tasks.

466 We can observe in Figure 4 a big performance increase from 1-shot clas-
467 sification to 3-shot classification. This is likely due to the lack of important
468 information in the support images in the 1-shot situation. If we can simu-
469 late support images by using some generative models or expand the support
470 feature set by dynamically utilizing query images in a semi-supervised way,
471 the accuracy of 1-shot classification may improve substantially.

472 **Acknowledgements**

473 This work was supported in part by the Beijing Natural Science Founda-
474 tion Project No. Z200002, in part by the Royal Society under International
475 Exchanges Award IEC\NSFC\201071, in part by the National Natural Sci-
476 ence Foundation of China (NSFC) No. 62111530146, 62176110, 61906080,
477 61922015, U19B2036, 62225601, in part by Young Doctoral Fund of Educa-
478 tion Department of Gansu Province under Grant 2021QB-038, Youth Innova-
479 tive Research Team of BUPT No. 2023QNTD02, and Hong-liu Distinguished
480 Young Talents Foundation of Lanzhou University of Technology.

481 **References**

- 482 [1] W.-Y. Chen, Y.-C. Liu, Z. Kira, Y.-C. F. Wang, J.-B. Huang, A closer
483 look at few-shot classification, in: International Conference on Learning
484 Representations, 2018.
- 485 [2] Y. Wang, Q. Yao, J. T. Kwok, L. M. Ni, Generalizing from a few ex-
486 amples: A survey on few-shot learning, ACM computing surveys (csur)
487 53 (3) (2020) 1–34.

- 488 [3] Y. Guo, R. Du, X. Li, J. Xie, Z. Ma, Y. Dong, Learning calibrated class
489 centers for few-shot classification by pair-wise similarity, *IEEE Transac-*
490 *tions on Image Processing* 31 (2022) 4543–4555.
- 491 [4] X. Li, X. Yang, Z. Ma, J.-H. Xue, Deep metric learning for few-shot im-
492 age classification: A review of recent developments, *Pattern Recognition*
493 (2023) 109381.
- 494 [5] C. Finn, P. Abbeel, S. Levine, Model-agnostic meta-learning for fast
495 adaptation of deep networks, in: *International Conference on Machine*
496 *Learning*, PMLR, 2017, pp. 1126–1135.
- 497 [6] B. Liu, Y. Cao, Y. Lin, Q. Li, Z. Zhang, M. Long, H. Hu, Negative
498 margin matters: Understanding margin in few-shot classification, in:
499 *European Conference on Computer Vision*, 2020, pp. 438–455.
- 500 [7] W. Chen, Z. Zhang, W. Wang, L. Wang, Z. Wang, T. Tan, Few-shot
501 learning with unsupervised part discovery and part-aligned similarity,
502 *Pattern Recognition* 133 (2023) 108986.
- 503 [8] J. Snell, K. Swersky, R. S. Zemel, Prototypical networks for few-shot
504 learning, *Advances in Neural Information Processing Systems* 30 (2017).
- 505 [9] F. Sung, Y. Yang, L. Zhang, T. Xiang, P. H. S. Torr, T. M.
506 Hospedales, Learning to compare: Relation network for few-shot learn-
507 ing, *IEEE/CVF Conference on Computer Vision and Pattern Recogni-*
508 *tion* (2018) 1199–1208.
- 509 [10] W. Li, L. Wang, J. Xu, J. Huo, Y. Gao, J. Luo, Revisiting local de-
510 scriptor based image-to-class measure for few-shot learning, *IEEE/CVF*

- 511 Conference on Computer Vision and Pattern Recognition (2019) 7253–
512 7260.
- 513 [11] Z. Wu, Y. Li, L. Guo, K. Jia, PARN: Position-aware relation networks
514 for few-shot learning, IEEE/CVF International Conference on Com-
515 puter Vision (2019) 6658–6666.
- 516 [12] C. Zhang, Y. Cai, G. Lin, C. Shen, DeepEMD: Few-shot image classifica-
517 tion with differentiable earth mover’s distance and structured classifiers,
518 IEEE/CVF Conference on Computer Vision and Pattern Recognition
519 (2020) 12200–12210.
- 520 [13] D. Wertheimer, L. Tang, B. Hariharan, Few-shot classification with fea-
521 ture map reconstruction networks, IEEE/CVF Conference on Computer
522 Vision and Pattern Recognition (2021) 8008–8017.
- 523 [14] Y. Liu, Y. Bai, X. Che, J. He, Few-shot fine-grained image classification:
524 A survey, in: International Conference on Natural Language Processing,
525 IEEE, 2022, pp. 201–211.
- 526 [15] Y. Yu, D. Zhang, S. Wang, Z. Ji, Z. Zhang, Local spatial alignment
527 network for few-shot learning, Neurocomputing 497 (2022) 182–190.
- 528 [16] Y. Liu, W. Zhang, C. Xiang, T. Zheng, D. Cai, X. He, Learning
529 to affiliate: Mutual centralized learning for few-shot classification, in:
530 IEEE/CVF Conference on Computer Vision and Pattern Recognition,
531 2022, pp. 14391–14400.
- 532 [17] H. Tang, C. Yuan, Z. Li, J. Tang, Learning attention-guided pyramidal

- 533 features for few-shot fine-grained recognition, *Pattern Recognition* 130
534 (2022) 108792.
- 535 [18] Y.-X. Wang, R. Girshick, M. Hebert, B. Hariharan, Low-shot learning
536 from imaginary data, in: *IEEE Conference on Computer Vision and*
537 *Pattern Recognition*, 2018, pp. 7278–7286.
- 538 [19] C. Wang, S. Song, Q. Yang, X. Li, G. Huang, Fine-grained few shot
539 learning with foreground object transformation, *Neurocomputing* 466
540 (2021) 16–26.
- 541 [20] O. Vinyals, C. Blundell, T. P. Lillicrap, K. Kavukcuoglu, D. Wierstra,
542 Matching networks for one shot learning, in: *Advances in Neural Infor-*
543 *mation Processing Systems*, 2016.
- 544 [21] H. Huang, Z. Wu, W. Li, J. Huo, Y. Gao, Local descriptor-based multi-
545 prototype network for few-shot learning, *Pattern Recognition* 116 (2021)
546 107935.
- 547 [22] K. Cao, M. Brbic, J. Leskovec, Concept learners for few-shot learning,
548 in: *International Conference on Learning Representation*, 2021.
- 549 [23] Y. Zhou, Y. Guo, S. Hao, R. Hong, Hierarchical prototype refinement
550 with progressive inter-categorical discrimination maximization for few-
551 shot learning, *IEEE Transactions on Image Processing* 31 (2022) 3414–
552 3429.
- 553 [24] X. Huang, S. H. Choi, SAPENet: self-attention based prototype en-
554 hancement network for few-shot learning, *Pattern Recognition* 135
555 (2023) 109170.

- 556 [25] J. Xie, F. Long, J. Lv, Q. Wang, P. Li, Joint distribution matters: Deep
557 brownian distance covariance for few-shot classification, in: IEEE/CVF
558 Conference on Computer Vision and Pattern Recognition, 2022, pp.
559 7972–7981.
- 560 [26] X. Li, J. Wu, Z. Sun, Z. Ma, J. Cao, J.-H. Xue, BSNet: Bi-similarity net-
561 work for few-shot fine-grained image classification, IEEE Transactions
562 on Image Processing 30 (2020) 1318–1331.
- 563 [27] W. Zhu, W. Li, H. Liao, J. Luo, Temperature network for few-shot
564 learning with distribution-aware large-margin metric, Pattern Recogni-
565 tion 112 (2021) 107797.
- 566 [28] H. Huang, J. Zhang, J. Zhang, J. Xu, Q. Wu, Low-rank pairwise align-
567 ment bilinear network for few-shot fine-grained image classification,
568 IEEE Transactions on Multimedia 23 (2021) 1666–1680.
- 569 [29] X. Li, Q. Song, J. Wu, R. Zhu, Z. Ma, J.-H. Xue, Locally-enriched
570 cross-reconstruction for few-shot fine-grained image classification, IEEE
571 Transactions on Circuits and Systems for Video Technology (2023) 1–1.
- 572 [30] S. Lee, W. Moon, J.-P. Heo, Task discrepancy maximization for fine-
573 grained few-shot classification, IEEE/CVF Conference on Computer Vi-
574 sion and Pattern Recognition (2022) 5321–5330.
- 575 [31] Z. Li, Z. Hu, W. Luo, X. Hu, Sabernet: Self-attention based effective
576 relation network for few-shot learning, Pattern Recognition 133 (2023)
577 109024.

- 578 [32] B. Munjal, A. Flaborea, S. Amin, F. Tombari, F. Galasso, Query-guided
579 networks for few-shot fine-grained classification and person search, *Pat-
580 tern Recognition* 133 (2023) 109049.
- 581 [33] K. B. Petersen, M. S. Pedersen, et al., The matrix cookbook, Technical
582 University of Denmark 7 (15) (2008) 510.
- 583 [34] G. Pereyra, G. Tucker, J. Chorowski, L. Kaiser, G. Hinton, Regular-
584 izing neural networks by penalizing confident output distributions, in:
585 International Conference on Learning Representations, 2017.
- 586 [35] C. Wah, S. Branson, P. Welinder, P. Perona, S. Belongie, The Caltech-
587 UCSD birds-200-2011 dataset (2011).
- 588 [36] J. Krause, M. Stark, J. Deng, L. Fei-Fei, 3d object representations
589 for fine-grained categorization, *IEEE International Conference on Com-
590 puter Vision Workshops* (2013) 554–561.
- 591 [37] A. Khosla, N. Jayadevaprakash, B. Yao, F.-F. Li, Novel dataset for fine-
592 grained image categorization: Stanford dogs, in: *CVPR workshop on
593 fine-grained visual categorization (FGVC)*, 2011.
- 594 [38] M.-E. Nilsback, A. Zisserman, Automated flower classification over a
595 large number of classes, *Indian Conference on Computer Vision, Graph-
596 ics & Image Processing* (2008) 722–729.
- 597 [39] S. Maji, E. Rahtu, J. Kannala, M. B. Blaschko, A. Vedaldi, Fine-grained
598 visual classification of aircraft, *ArXiv abs/1306.5151* (2013).

- 599 [40] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recog-
600 nition, in: IEEE Conference on Computer Vision and Pattern Recogni-
601 tion, 2016, pp. 770–778.
- 602 [41] J. Xu, H. Le, M. Huang, S. Athar, D. Samaras, Variational feature
603 disentangling for fine-grained few-shot classification, in: IEEE/CVF In-
604 ternational Conference on Computer Vision, 2021, pp. 8812–8821.
- 605 [42] R. R. Selvaraju, A. Das, R. Vedantam, M. Cogswell, D. Parikh, D. Batra,
606 Grad-CAM: Visual explanations from deep networks via gradient-based
607 localization, International Journal of Computer Vision 128 (2017) 336–
608 359.
- 609 [43] L. van der Maaten, G. E. Hinton, Visualizing data using t-SNE, Journal
610 of Machine Learning Research 9 (2008) 2579–2605.

Title: Self-Reconstruction Network for Fine-Grained Few-Shot Classification

Highlights:

- A novel fine-grained few-shot classification method is proposed to alleviate overfitting.
- The proposed method uses self-reconstruction to increase the diversity of query features.
- The proposed method introduces a restrained cross-entropy loss.
- The proposed method generalizes well across five benchmark datasets.

Xiaoxu Li received her Ph.D. degree from Beijing University of Posts and Telecommunications in 2012. She is a 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.

Zhen Li received his B.Eng. degree in Measurement and Control Technology and Instrument from Taiyuan University of Technology, China, in 2012. He is currently pursuing the master's degree at Lanzhou University of Technology. His research interests include computer vision and machine learning.

Jiyang Xie received his B.E. degree in information engineering from Beijing University of Posts and Telecommunications (BUPT), China, in 2017, where he received his Ph.D. degree in 2022. His research interests include pattern recognition and machine learning fundamentals with a focus on applications in image processing, data mining, and deep learning.

Xiaochen Yang received the Ph.D. degree in statistical science from University College London in 2020. She is a Lecturer with the School of Mathematics and Statistics, University of Glasgow. Her research interests include statistical classification, metric learning, and hyperspectral image analysis. She is an Associate Editor of Neurocomputing.

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 with the Department of Statistical Science at University College London. His research interests include multivariate and high-dimensional data analysis, statistical pattern recognition, machine learning and image processing. He is an Associate Editor of IEEE Transactions on Circuits and Systems for Video Technology, IEEE Transactions on Cybernetics, and IEEE Transactions on Neural Networks and Learning Systems.

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.

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-proof