

---

# Generalized and Discriminative Few-Shot Object Detection via SVD-Dictionary Enhancement

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Few-shot object detection (FSOD) aims to detect new objects based on few annotated  
2 samples. To alleviate the impact of few samples, enhancing the generalization  
3 and discrimination abilities of detectors on new objects plays an important role.  
4 In this paper, we explore employing Singular Value Decomposition (SVD) to  
5 boost both the generalization and discrimination abilities. In specific, we propose  
6 a novel method, namely, SVD-Dictionary enhancement, to build two separated  
7 spaces based on the sorted singular values. Concretely, the eigenvectors corre-  
8 sponding to larger singular values are used to build the generalization space in  
9 which localization is performed, as these eigenvectors generally suppress certain  
10 variations (e.g., the variation of styles) and contain intrinsic characteristics of  
11 objects. Meanwhile, since the eigenvectors corresponding to relatively smaller sin-  
12 gular values may contain richer category-related information, we can utilize them  
13 to build the discrimination space in which classification is performed. Dictionary  
14 learning is further leveraged to capture high-level discriminative information from  
15 the discrimination space, which is beneficial for improving detection accuracy. In  
16 the experiments, we separately verify the effectiveness of our method on PASCAL  
17 VOC and COCO benchmarks. Particularly, for the 2-shot case in VOC split1, our  
18 method significantly outperforms the baseline by 6.2%. Moreover, visualization  
19 analysis shows that our method is instrumental in doing FSOD.

## 20 1 Introduction

21 With the rejuvenation of deep neural networks, for object detection, many progresses [11, 12, 1, 26,  
22 22] have been achieved. Though these methods obtain outstanding detection performances, they  
23 usually require a large number of labeled samples for training, which are labored yet expensive to  
24 collect and annotate. On the contrary, human beings are born with the ability to learn a new visual  
25 concept with only few samples. To imitate such an ability of human beings, the task of few-shot  
26 object detection (FSOD) [2, 17] has been proposed, which aims to improve the detection performance  
27 for new objects that contain few annotated training samples.

28 The main challenge of FSOD lies in how to learn generalized and discriminative object features from  
29 both abundant samples in base object categories and few samples in new object categories, which  
30 can improve the representation ability of object features and alleviate overfitting on new objects.  
31 Following the popular methods for few-shot image classification, earlier attempts [37, 36, 33, 8] in  
32 FSOD utilize the meta-learning strategy [29, 31, 10], whose goal is to learn detectors across tasks and  
33 then transfer to the few-shot detection task. However, compared with traditional two-stage fine-tuning  
34 based approaches [34, 35, 30], the meta-learning strategy fails to effectively improve generalization  
35 and discrimination of object features and leads to weak performance. The reason may be that during

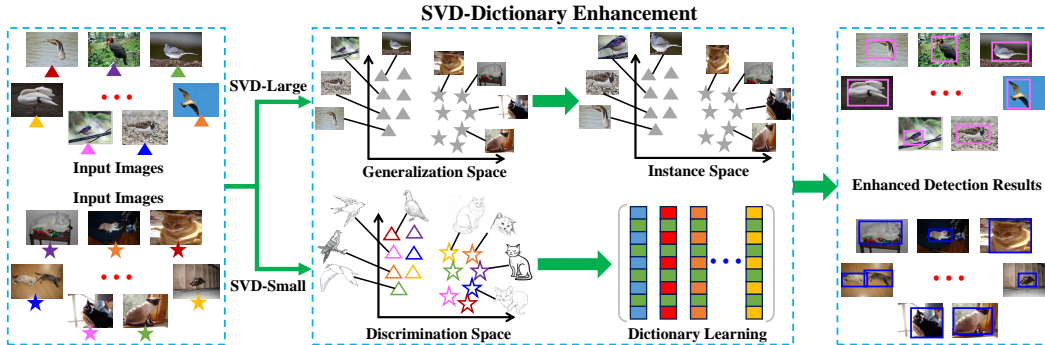


Figure 1: SVD-Dictionary enhancement for FSOD. ‘SVD-Large’ indicates that we use the eigenvectors corresponding to larger singular values to build the generalization space in which localization is performed. ‘SVD-Small’ indicates that we use the eigenvectors corresponding to smaller singular values to build the discrimination space in which classification is performed. Meanwhile, dictionary learning [40] is used to capture high-level discriminative information from the discrimination space, which is beneficial for improving the detection accuracy.

36 each training episode, meta-learning methods focus on transferability across different tasks and ignore  
 37 learning of generalized and discriminative feature representations.

38 For FSOD, the generalized representations may contain intrinsic characteristics of object features,  
 39 which is beneficial for adapting knowledge from base object categories to new object categories.  
 40 Meanwhile, the discriminative representations may contain certain category-related information,  
 41 which is helpful for boosting the detection accuracy. Furthermore, recent research [3] has shown that  
 42 from a spectral analysis perspective, the feature representations can be decomposed into eigenvectors  
 43 with importance quantified by the corresponding singular values. The eigenvectors corresponding to  
 44 larger singular values contribute to the generalization ability, as these eigenvectors could suppress  
 45 certain variations (e.g., the variations of style and texture). Meanwhile, since the eigenvectors  
 46 corresponding to relatively smaller singular values contain richer category-related information (e.g.,  
 47 the structures of objects), these eigenvectors are beneficial for discrimination. Therefore, in this  
 48 paper, we explore employing Singular Value Decomposition (SVD) (as shown in Fig. 1) to promote  
 49 detectors to learn generalized and discriminative object features.

50 Particularly, we propose a method named as SVD-Dictionary enhancement for FSOD. Given an  
 51 input image, a backbone network is first used to extract the corresponding feature map. Then, SVD  
 52 is performed on the feature map. Here, we select the eigenvectors corresponding to the first  $k$   
 53 largest singular values to compute a generalization map. And the generalization ability is enhanced  
 54 by a residual operation between the generalization map and the original feature map. Next, the  
 55 residual eigenvectors are used to calculate a discrimination map. Meanwhile, to further enhance  
 56 discrimination, we define a codebook containing multiple codewords and employ dictionary learning  
 57 [40] to capture high-level discriminative information from the discrimination map, which is good for  
 58 accurate detection. Compared with most methods [35, 33, 17] for FSOD, our method includes two  
 59 virtues. One is that during enhancing generalization, our method does not introduce extra parameters.  
 60 The other is that with the help of the discrimination map and dictionary learning, our method could  
 61 capture high-level discriminative information of different categories, which is conducive to reducing  
 62 the data-scarce impact on new object categories. During training, we first train the model on the  
 63 data-abundant base object categories. Then, the model is fine-tuned on a reconstructed training set  
 64 that contains a small number of balanced training samples from both base and new object categories.  
 65 Extensive experiments on two benchmarks demonstrate the superiorities of our method.

66 The contributions of our work are summarized as follows:

- 67 • To boost both the generalization and discrimination abilities, we propose to build the  
 68 generalization and discrimination spaces based on the sorted singular values.
- 69 • To further enhance the discrimination ability, we explore dictionary learning to capture  
 70 high-level discriminative information from the discrimination map.
- 71 • By plugging our method into two fine-tuning based two-stage methods, i.e., MPSR [35] and  
 72 FSCE [30], our method significantly improves their performances on PASCAL VOC [6, 7]  
 73 and COCO [20] benchmarks.

## 74 2 Related Work

75 **Few-shot image classification.** The goal of few-shot image classification [29, 24] is to recognize new  
76 categories with very few labeled samples. Recently, many progresses [5, 32, 41, 39, 13] have been  
77 achieved. Particularly, meta-learning [10] is a widely used method to solve few-shot classification,  
78 which aims to leverage task-level meta knowledge to help models adapt to new tasks with few labeled  
79 samples. Based on the meta-learning policy, Snell et al. [29] proposed a prototypical network to learn  
80 a metric space in which classification can be performed by computing distances to the prototype  
81 representation of each category. However, the performance of this method relies on the quality of the  
82 learned prototypes. When the training data is scarce, the learned prototypes could not represent the  
83 information of each category sufficiently, which affects the classification performance. Liu et al. [23]  
84 proposed a method of prototype rectification, which considers the intra-class bias and the cross-class  
85 bias and improves the performance significantly. Apart from these methods, more methods, e.g.,  
86 sample synthesis and augmentation, in few-shot learning can be seen in the work [24]. Whereas,  
87 these classification methods could not be directly applied to detection that requires localizing and  
88 recognizing objects simultaneously.

89 **Few-shot object detection.** Towards FSOD, most existing methods [18, 8, 25, 2, 38] employ a meta-  
90 learning or fine-tuning based mechanism. Particularly, Wang et al. [33] proposed a meta-learning  
91 framework to leverage meta-level knowledge from base object categories to facilitate the generation  
92 of a detector for new object categories. Based on this work [33], Kang et al. [17] further proposed  
93 a one-stage detection architecture that contains a meta feature learner and a reweighting module.  
94 In order to alleviate the impact of complex background and multiple objects on one image, Yan et  
95 al. [37] extended Faster R-CNN [27] and Mask R-CNN [16] by proposing meta-learning over RoI  
96 (Region-of-Interest) features. Recently, the two-stage fine-tuning based approach (TFA) [34] reveals  
97 a potential for addressing FSOD. By simply fine-tuning the box classifier and regressor, this method  
98 outperforms many meta-learning based methods. Wu et al. [35] considered the impact of the scale  
99 bias on the fine-tuning process, which further improves the detection performance.

100 Different from the above methods, in this paper, we explore enhancing both generalization and  
101 discrimination for FSOD. And we propose a method of SVD-Dictionary enhancement that combines  
102 SVD with dictionary learning. Experimental results and visualization analysis demonstrate the  
103 superiorities of the proposed method.

## 104 3 SVD-Dictionary Enhancement for FSOD

105 In this paper, we follow the same settings introduced in Kang et al. [17]. Concretely, there are a set  
106 of base object categories that contain abundant annotated samples and a set of new object categories  
107 that contain only few (usually less than 30) annotated samples per category. The main purpose is to  
108 improve the detection performance of new object categories.

### 109 3.1 SVD Enhancement

110 For FSOD, generalization and discrimination are two important criteria that characterize the goodness  
111 of feature representation. Particularly, enhancing generalization is beneficial for adapting the knowl-  
112 edge learned from base object categories to new object categories, which alleviates the data-scarce  
113 impact on new object categories. Meanwhile, discrimination refers to the ability to separate different  
114 categories based on the learned representations. And enhancing discrimination is helpful for reducing  
115 the overfitting risk on new object categories, which improves the detection accuracy. To this end, we  
116 explore SVD to enhance both the generalization and discrimination abilities of detectors.

117 Concretely, as shown in Fig. 2, we adopt a widely used two-stage object detector, i.e., Faster R-CNN  
118 [27], as the basic detection model. Given an input image, we first employ the feature extractor, e.g.,  
119 ResNet [15], to extract the corresponding feature map  $F \in \mathbb{R}^{m \times w \times h}$ , where  $m$ ,  $w$ , and  $h$  separately  
120 denote the number of channels, width, and height. Then,  $F$  is reshaped as  $\mathbf{F} \in \mathbb{R}^{m \times n}$ , where  
121  $n = w \times h$ . SVD is used to factorize the matrix  $\mathbf{F}$ , i.e.,  $\mathbf{F} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T \in \mathbb{R}^{m \times n}$ , into the product of  
122 three matrices, where  $\mathbf{U} \in \mathbb{R}^{m \times m}$  and  $\mathbf{V} \in \mathbb{R}^{n \times n}$  are orthogonal, and  $\mathbf{\Sigma}$  contains the sorted singular  
123 values along its main diagonal [4]. Since the eigenvectors corresponding to larger singular values  
124 contain more information of the original matrix  $\mathbf{F}$ , we select the eigenvectors corresponding to the  
125 first  $k$  largest singular values to compute the generalization map  $\mathbf{G} \in \mathbb{R}^{m \times n}$ . Next, by feat of the

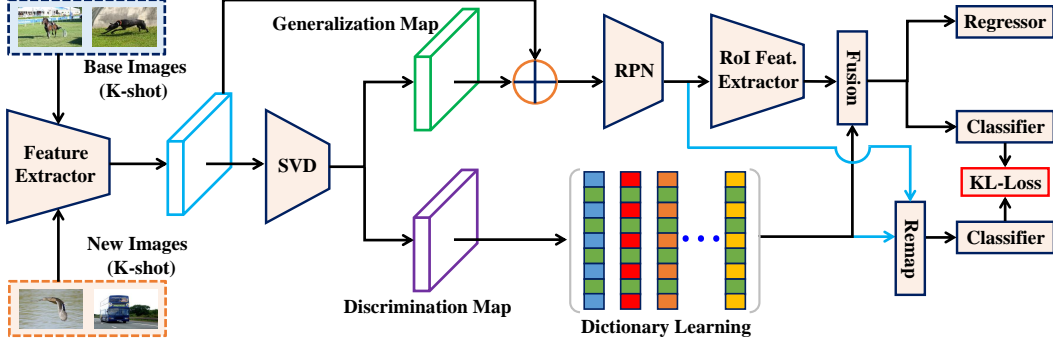


Figure 2: The architecture of generalized and discriminative FSOD via SVD-Dictionary enhancement. Here, ‘ $\oplus$ ’ indicates the residual operation. ‘RPN’ denotes Region-Proposal Network with RoI Pooling. After extracting the corresponding feature maps of input images, SVD is utilized to compute all singular values and eigenvectors. Then, the eigenvectors corresponding to larger singular values are used to compute the generalization map. And the eigenvectors corresponding to smaller singular values are used to calculate the discrimination map. Finally, dictionary learning is used to further capture high-level discriminative information, which helps improve the ability of accurate detection.

126 residual operation between  $G$  and  $F$ , the generalization ability of the extracted features is enhanced.  
 127 The processes are shown as follows:

$$G = U_{m \times k} \Sigma_{k \times k} V_{k \times n}^T, \quad E = G + F, \quad (1)$$

128 where  $U_{m \times k}$  and  $V_{k \times n}^T$  indicate that we select the first  $k$  columns and rows from the matrix  $U$  and  
 129  $V^T$ , respectively.  $\Sigma_{k \times k}$  is a diagonal matrix with the dimension  $k \times k$ .  $E \in \mathbb{R}^{m \times n}$  is the enhanced  
 130 matrix. Finally,  $E$  is reshaped as  $E \in \mathbb{R}^{m \times w \times h}$  that is used to perform the following RPN operation.  
 131 It is worth noting that in the process of enhancing generalization, we only perform the SVD operation  
 132 and do not introduce extra parameters. Besides, we utilize the residual operation to obtain the output  
 133  $E$ , which strengthens the generalization ability and retains the discriminative information in the output.  
 134 In the experiment, we observe that utilizing the operation of enhancing generalization improves the  
 135 detection performance effectively.

136 Next, the remaining eigenvectors and corresponding singular values are used to calculate the discrim-  
 137 ination map  $D \in \mathbb{R}^{m \times n}$ . The processes are the same as computing  $G$ . Since the map  $D$  contains  
 138 more category-related information [3], e.g., the structures of objects, it is helpful for enhancing the  
 139 discrimination ability. Similarly, for this process, we do not introduce extra parameters, either.

### 140 3.2 SVD-based Dictionary Learning

141 **Dictionary Learning.** Based on the map  $D$ , we explore employing dictionary learning [40, 14] to  
 142 capture high-level discriminative information, which is beneficial for strengthening the discrimination  
 143 ability of detectors. Concretely, we define a learned codebook  $C = \{c_j \in \mathbb{R}^m, j = 1, \dots, Q\}$   
 144 that contains  $Q$  codewords. Each element  $d_i \in \mathbb{R}^m$  of the map  $D$  can be assigned with a weight  
 145  $a_{ij}$  to each codeword  $c_j$  and the corresponding residual vector is denoted by  $r_{ij} = d_i - c_j$ , where  
 146  $i = 1, 2, \dots, n$ . Thus, dictionary learning can be calculated as follows:

$$x_j = \sum_{i=1}^n a_{ij} r_{ij}, \quad a_{ij} = \frac{\exp(-s_j \|r_{ij}\|^2)}{\sum_{j=1}^Q \exp(-s_j \|r_{ij}\|^2)}, \quad (2)$$

147 where  $s_j$  indicates the learnable smoothing factor for the corresponding codeword  $c_j$ . Finally, the  
 148 output of dictionary learning is a fixed length representation  $X = \{x_j \in \mathbb{R}^m, j = 1, \dots, Q\}$ . Next,  
 149 we take  $E$  as the input of the RPN module to obtain a set of object proposals  $P \in \mathbb{R}^{z \times m \times o \times o}$ , where  
 150  $z$  and  $o$  separately denote the number of proposals and their spatial size. And the fusion result of  $P$   
 151 and  $X$  is taken as the input of the classifier.

$$P = \text{RPN}(E), \quad y = \text{cls}([\phi(P), w_c X + b_c]), \quad (3)$$

152 where  $\phi$  consists of two fully-connected layers.  $w_c$  and  $b_c$  are learnable parameters. ‘[.]’ indicates  
 153 the fusion operation. Here, we use the concatenation operation. ‘cls’ denotes the classifier. By the

154 constraint of the classification loss, we can promote the learned representation  $X$  and codebook  $C$  to  
 155 absorb category-related information, which is good for enhancing detection accuracy.

156 **Dictionary-based Remap.** To further facilitate the learned codebook  $C$  to retain more category-  
 157 related characteristics, we try to remap  $P$  to the dictionary space and perform classification. Con-  
 158 cretely, each element  $p \in \mathbb{R}^m$  of  $P$  is remapped as a combination of codewords in the codebook  $C$ .  
 159 The processes are shown as follows:

$$rep = \sum_{j=1}^Q \frac{\exp(\psi(p)c_j^T)}{\sum_{j=1}^Q \exp(\psi(p)c_j^T)} c_j, \quad (4)$$

160 where  $\psi$  is a fully-connected layer that maps  $p$  to the dictionary space.  $rep \in \mathbb{R}^m$  indicates one  
 161 element of the remapping output  $Rep \in \mathbb{R}^{z \times m \times o \times o}$ . Next,  $Rep$  is taken as the input of the classifier  
 162 to output the probability:

$$y_{rep} = \text{cls}([\phi(P), \phi(Rep)]), \quad (5)$$

163 where  $y_{rep}$  indicates the output probability. Eq. (3) and Eq. (5) share the same classifier. Finally, the  
 164 KL-Divergence loss  $\mathcal{L}_{kl}$  is leveraged to enforce the prediction consistency between  $y_{rep}$  and  $y$ . By  
 165 performing classification in the dictionary space, the codebook  $C$  could be directly facilitated to learn  
 166 category-related characteristics, which is conducive to the improvement of the discrimination ability.

### 167 3.3 Two-Stage Fine-Tuning Mechanism

168 In this paper, we employ the commonly used detection loss [27] to optimize the model. Concretely,  
 169 the joint training loss is defined as follows:

$$\mathcal{L} = \mathcal{L}_{cls} + \mathcal{L}_{loc} + \mathcal{L}_{rpn} + \lambda \mathcal{L}_{kl}, \quad (6)$$

170 where  $\mathcal{L}_{cls}$  and  $\mathcal{L}_{loc}$  separately indicate the classification and bounding-box regression losses.  $\mathcal{L}_{rpn}$   
 171 is the RPN loss that is used to distinguish foreground from background and refine bounding-box  
 172 anchors. The hyper-parameter  $\lambda$  is set to 1.0 in the experiment.

173 During training, we employ the two-stage fine-tuning mechanism to optimize the proposed method.  
 174 Currently, there exist two fine-tuning training strategies. One is that during the base training and  
 175 fine-tuning stage, all the parameters of the detector are optimized simultaneously [35]. The other is  
 176 that during the fine-tuning stage, some important parameters of the detector are optimized. And the  
 177 remaining parameters are fixed [34, 30]. To demonstrate the effectiveness of the proposed method,  
 178 we separately utilize these two strategies to optimize the detector. Specifically, in the base training  
 179 stage, we employ the joint loss  $\mathcal{L}$  to optimize the entire model based on the data-abundant base object  
 180 categories. During fine-tuning, the last fully-connected layer (for classification) of the detection head  
 181 is replaced. The new classifier is randomly initialized. For the first strategy, we follow MPSR [35] to  
 182 optimize all the parameters of the model based on a balanced training set consisting of both the few  
 183 base and new object categories. For the second strategy, we follow FSCE [30] to jointly fine-tune the  
 184 FPN [21] pathway and RPN while fixing the backbone.

### 185 3.4 Further Discussion

186 In this section, we further discuss SVD and dictionary learning for few-shot object detection.

187 For FSOD, the two-stage fine-tuning mechanism can be regarded as a method that adapts the  
 188 knowledge from base object categories to new object categories, which is effective to alleviate the  
 189 data-scarce impact. Most existing methods [34, 30] focus on designing an effective optimizing  
 190 strategy and pay little attention to improving both the generalization and discrimination during the  
 191 fine-tuning stage. Recently, FSCE [30] brings contrastive learning [19] into FSOD, which is beneficial  
 192 for enhancing discrimination. However, the contrastive loss is calculated based on object proposals,  
 193 which neglects the impact of generalization on object localization.

194 For FSOD, we propose an SVD-Dictionary method to enhance both generalization and discrimination.  
 195 Particularly, the eigenvectors corresponding to larger singular values are directly used to enhance  
 196 generalization without introducing extra parameters. Meanwhile, we employ dictionary learning  
 197 to capture high-level discriminative information, which leads to accurate detection. Experimental  
 198 results and visualization analysis demonstrate the superiorities of our method.

## 199 4 Experiments

200 In the experiments, the proposed method is evaluated on PASCAL VOC [6, 7] and COCO [20]  
201 benchmarks. We strictly follow the consistent few-shot detection data construction and evaluation  
202 protocol [17, 35, 36, 34] to ensure fair and direct comparison. Meanwhile, since our method is trained  
203 based on the two-stage fine-tuning mechanism, we take two-stage methods, i.e., TFA [34], MPSR  
204 [35], and FSCE [30], as the compared baselines.

### 205 4.1 Implementation Details and Few-Shot Detection Benchmarks

206 **Implementation Details.** For the detection model, we use Faster R-CNN [27] with the RoI Align  
207 [16] layer. The backbone is ResNet-101 [15]. The parameters are pre-trained on ImageNet [28] for  
208 initialization. In Eq. (1), we select the first  $k$  largest singular values to compute the generalization  
209 map. Here,  $k$  is set to half of the total number of singular values. For dictionary learning, the  
210 number of codewords is set to 24. All newly introduced parameters are initialized randomly. All the  
211 experiments are trained using the standard SGD optimizer with a momentum of 0.9 and a weight  
212 decay of 0.0001. During inference, we take the output  $y$  of Eq. (3) as the classification result.

213 **FSOD Benchmarks.** For PASCAL VOC, the overall 20 categories are divided into 15 base object  
214 categories and 5 new object categories. All base object category data from PASCAL VOC 07+12  
215 trainval sets is available. For each new object category, there exist  $\mathbf{K}$  instances available and  $\mathbf{K}$  is  
216 set to 1, 2, 3, 5, and 10. Following existing methods [17, 34, 35], we utilize the same three random  
217 partitions of base and new object categories, referred to as New Split 1, 2, and 3. And for the  
218 predictions on PASCAL VOC 2007 test set, we separately report the results of nAP50 and nAP75.

219 For the 80 categories in COCO, 20 categories overlapped with PASCAL VOC are taken as new object  
220 categories. The remaining 60 categories are used as base object categories. The  $\mathbf{K} = 10$  and 30 shots  
221 detection performance is evaluated on 5,000 images from COCO 2014 validation set.

### 222 4.2 Performance Analysis of Few-Shot Detection

223 **PASCAL VOC Results.** Table 1 shows the results on three PASCAL VOC New Splits. We can see  
224 that as the number of object instances increases, the performance continually improves significantly.  
225 This shows that few samples affect the performance of object detection. Besides, compared with the  
226 two-stage fine-tuning training mechanism, the training process of the meta-learning mechanism is  
227 more complex. However, for FSOD, meta-learning based methods [37, 33, 36] fail to obtain superior  
228 performance. The reason may be that these methods focus on learning task-level transferability and  
229 ignore the learning of feature generalization and discrimination. Next, we can see that plugging  
230 our method into MPSR [35] and FSCE [30] improves their performances significantly. Particularly,  
231 based on nAP50 and nAP75, the performance of FSCE is significantly improved. These analyses  
232 demonstrate that the proposed method is helpful for enhancing the generalization and discrimination  
233 abilities of detectors, which is beneficial for FSOD.

234 In Fig. 3, we show some detection examples. We can see that compared with MPSR and FSCE,  
235 our method localizes and recognizes the objects in these images accurately. Particularly, there exist  
236 three types of error detections, i.e., missing detection that misses the detection of certain objects (e.g.,  
237 the fifth example in the first row), uncertain detection that classifies objects into multiple different  
238 categories (e.g., the second example in the first row), and mis-classifications of objects (e.g., the first  
239 example in the first row). For these examples, our method reduces the appearance of these errors,  
240 which shows improving generalization and discrimination is beneficial for accurate detection.

241 **COCO Results.** Table 2 shows the COCO results. We can also see that plugging our method into  
242 MPSR and FSCE leads to performance improvement. Particularly, for MPSR, based on the 30-shot  
243 case, plugging our method separately improves its performance by 2.4 % (AP), 2.4 % (AP75), and  
244 3.7 % (AP<sub>L</sub>). For FSCE, in terms of the five metrics, plugging our method is beneficial for boosting  
245 the detection performance. This further demonstrates the effectiveness of our method. Besides,  
246 FSOD-VE [36] is a recently proposed meta-learning method, which explores leveraging viewpoint  
247 estimation to solve FSOD. Though FSOD-VE’s performance outperforms fine-tuning based methods  
248 [35, 30], the training process of meta-learning is much more complex. And the performance on  
249 small objects is weaker. This shows that improving the generalization and discrimination during the  
250 fine-tuning process is an effective solution for FSOD.

Table 1: Few-shot detection performance (%) on PASCAL VOC New Split sets. ‘MPSR + Ours’ and ‘FSCE + Ours’ separately indicate that we plug our method into MPSR [35] and FSCE [30]. ‘ft’ denotes fine-tuning. ‘†’ represents meta-learning based methods. ‘\*’ indicates that we directly run the released code to obtain the results. The evaluation of the last two rows is based on nAP75. The evaluation of the other rows is based on nAP50.

Method (nAP50) / Shot	New Split 1					New Split 2					New Split 3				
	1	2	3	5	10	1	2	3	5	10	1	2	3	5	10
FRCN-ft [33]	13.8	19.6	32.8	41.5	45.6	7.9	15.3	26.2	31.6	39.1	9.8	11.3	19.1	35.0	45.1
FRCN+FPN-ft [34]	8.2	20.3	29.0	40.1	45.5	13.4	20.6	28.6	32.4	38.8	19.6	20.8	28.7	42.2	42.1
†Meta R-CNN [37]	19.9	25.5	35.0	45.7	51.5	10.4	19.4	29.6	34.8	45.4	14.3	18.2	27.5	41.2	48.1
†MetaDet [33]	18.9	20.6	30.2	36.8	49.6	21.8	23.1	27.8	31.7	43.0	20.6	23.9	29.4	43.9	44.1
†FSOD-VE [36]	24.2	35.3	42.2	49.1	57.4	21.6	24.6	31.9	37.0	45.7	21.2	30.0	37.2	43.8	49.6
TFA w/fc [34]	36.8	29.1	43.6	55.7	57.0	18.2	29.0	33.4	35.5	39.0	27.7	33.6	42.5	48.7	50.2
TFA w/cos [34]	39.8	36.1	44.7	55.7	56.0	23.5	26.9	34.1	35.1	39.1	30.8	34.8	42.8	49.5	49.8
Retentive R-CNN [9]	42.4	45.8	45.9	53.7	56.1	21.7	27.8	35.2	37.0	40.3	30.2	37.6	43.0	49.7	50.1
MPSR* [35]	40.7	41.2	48.9	53.6	60.3	24.4	29.3	39.2	39.9	47.8	32.9	34.4	42.3	48.0	49.2
MPSR + Ours	41.5	<b>47.4</b>	<b>51.5</b>	57.7	61.2	<b>29.4</b>	29.6	39.8	41.2	<b>51.5</b>	36.0	39.4	45.4	50.4	51.3
FSCE* [30]	44.2	43.2	45.7	58.3	61.0	25.4	29.5	42.1	43.6	48.7	37.2	43.5	45.8	53.3	55.8
FSCE + Ours	<b>46.1</b>	43.5	48.9	<b>60.0</b>	<b>61.7</b>	25.6	<b>29.9</b>	<b>44.8</b>	<b>47.5</b>	48.2	<b>39.5</b>	<b>45.4</b>	<b>48.9</b>	<b>53.9</b>	<b>56.9</b>
FSCE* (nAP75) [30]	21.9	21.2	20.1	32.7	38.8	6.9	8.4	14.7	20.3	<b>25.9</b>	16.3	18.3	18.9	25.4	29.6
FSCE + Ours (nAP75)	<b>25.1</b>	<b>21.4</b>	<b>25.1</b>	<b>36.5</b>	<b>39.8</b>	<b>9.4</b>	<b>11.3</b>	<b>18.5</b>	<b>24.1</b>	25.6	<b>18.4</b>	<b>20.5</b>	<b>24.2</b>	<b>26.8</b>	<b>30.5</b>

Table 2: Few-shot detection evaluation results (%) on COCO. Here, AP<sub>S</sub>, AP<sub>M</sub>, and AP<sub>L</sub> separately indicate the detection performances of the small, medium, and large objects.

Shots	Method	AP	AP75	AP <sub>S</sub>	AP <sub>M</sub>	AP <sub>L</sub>
10	†Meta R-CNN [37]	8.7	6.6	2.3	7.7	14.0
	†MetaDet [33]	7.1	6.1	1.0	4.1	12.2
	†FSOD-VE [36]	<b>12.5</b>	9.8	2.5	<b>13.8</b>	<b>19.9</b>
	TFA w/fc [34]	10.0	9.2	–	–	–
	TFA w/cos [34]	10.0	9.3	–	–	–
	MPSR* [35]	9.5	9.5	3.3	8.2	15.9
	MPSR + Ours	11.0	<b>10.6</b>	<b>4.4</b>	11.5	17.1
	FSCE* [30]	11.3	9.6	3.7	10.7	18.6
	FSCE + Ours	12.0	10.4	4.2	12.1	18.9
	†Meta R-CNN [37]	12.4	10.8	2.8	11.6	19.0
30	†MetaDet [33]	11.3	8.1	1.1	6.2	17.3
	†FSOD-VE [36]	14.7	12.2	3.2	15.2	23.8
	TFA w/fc [34]	13.4	13.2	–	–	–
	TFA w/cos [34]	13.7	13.4	–	–	–
	MPSR* [35]	13.8	13.5	4.0	12.9	22.9
	MPSR + Ours	<b>16.2</b>	<b>15.9</b>	4.6	14.6	<b>26.6</b>
	FSCE* [30]	15.4	14.2	5.5	14.9	24.4
	FSCE + Ours	16.0	15.3	<b>6.0</b>	<b>16.8</b>	24.9

### 251 4.3 Ablation Analysis

252 In this section, ablation analysis is performed based on  
 253 the New Split 1 of PASCAL VOC. And we plug our  
 254 method into MPSR to make the ablation analysis.

255 **Analysis of Hyper-parameter  $k$ .** In Eq. (1), we select  
 256 the first  $k$  columns and rows from  $U$  and  $V^T$  that cor-  
 257 respond to the first  $k$  largest singular values to build  
 258 the generalization map. To enhance the generalization  
 259 ability, the generalization map is expected to contain  
 260 much information reflecting intrinsic characteristics  
 261 of objects. In Table 3, we analyze the impact of  $k$ . Here,  
 262 we only change the setting of  $k$ . The other modules are kept unchanged. We can see that different

Table 3: The performance (%) of using a different number of singular values. Here, ‘proportion’ indicates the percentage of the total number of singular values.

proportion/shot	1	2	3	5	10
10%	40.6	43.9	49.1	55.6	<b>62.1</b>
25%	38.3	44.7	49.4	56.2	61.7
50%	<b>41.5</b>	<b>47.4</b>	<b>51.5</b>	<b>57.7</b>	61.2
75%	36.3	42.9	48.7	55.1	60.9
90%	37.9	41.8	48.1	55.8	60.2



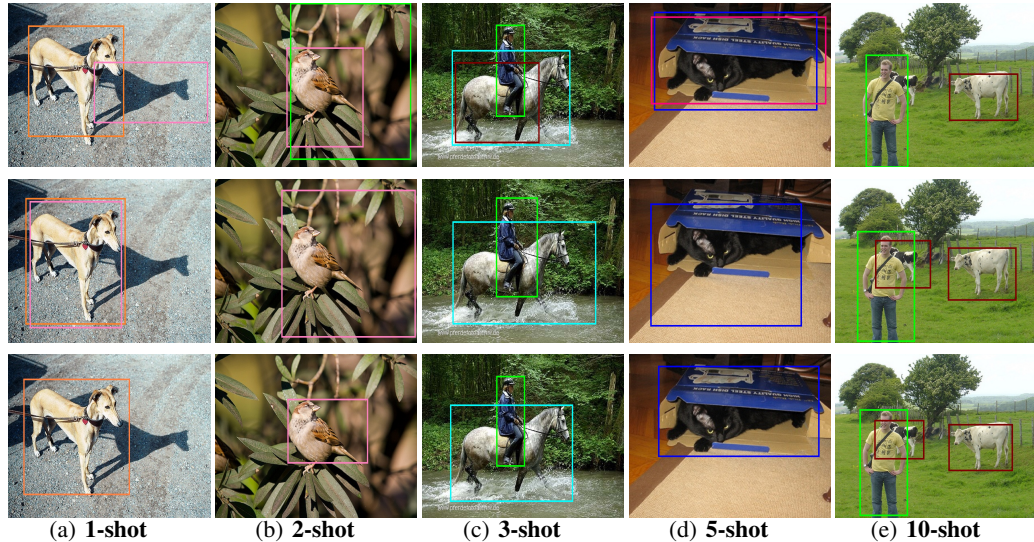


Figure 3: Detection examples based on different shots. The first, second, and third rows separately indicate detections based on MPSR [35], FSCE [30], and our method. We can see our method accurately detect ‘dog’, ‘bird’, ‘person’, ‘horse’, ‘cat’, and ‘cow’.

263 settings of  $k$  affect the performance of FSOD. Particularly, while  $k$  is set to a large value or a small  
 264 value, the performance decreases. The reason may be that using a large value of proportion introduces  
 265 much information that is not related to intrinsic characteristics, which weakens the generalization  
 266 ability. Meanwhile, using a small value of proportion may lead to the loss of certain object-related  
 267 information, which weakens the feature representation. We observe that the performance of using  
 268 50% proportion is the best.

269 **Analysis of SVD-based Generalization and Dictionary Learning.** To demonstrate the effectiveness of  
 270 the proposed method, we remove the module of dictionary learning and only keep the SVD-based generaliza-  
 271 tion. From the 1-shot to 10-shot case, the performance  
 272 is 41.2%, 44.3%, 49.7%, 54.8%, and 60.9%. We can  
 273 see that employing dictionary learning is helpful for  
 274 improving detection performance. Particularly, taking  
 275 the 2-shot case as an example, the performance is im-  
 276 proved by 3.1%. This indicates based on the output of  
 277 SVD operation, dictionary learning is able to leverage  
 278 multiple codewords to capture high-level discrimina-  
 279 tive information that is helpful for accurate detection. Meanwhile, this also shows that the learned  
 280 codewords contain category-related information, which enhances the discrimination ability of the  
 281 detector. Besides, we can see the current performance of only using SVD-based generalization still  
 282 outperforms MPSR. Taking the 2-shot and 5-shot cases as examples, the performance is separately  
 283 improved by 3.1% and 1.2%. This indicates that utilizing eigenvectors corresponding to larger  
 284 singular values to build the generalization map is beneficial for extracting generalized information  
 285 without introducing extra parameters, thereby boosting the performance of FSOD.

288 **Analysis of the Number of Codewords in Dictionary Learning.** In this paper, we define a codebook contain-  
 289 ing multiple codewords to sufficiently capture category-  
 290 related discriminative information from the discrimina-  
 291 tion map corresponding to relatively smaller singular  
 292 values. In Table 4, we analyze the impact of using a dif-  
 293 ferent number of codewords. We can see that when the  
 294 number is small, e.g., 16, the performance decreases.  
 295 The reason may be that a small number of codewords  
 296 are not sufficient to capture much discriminative information, which affects the detection performance.

Table 4: The performance (%) of using a different number of codewords in the codebook of dictionary learning.

number/shot	1	2	3	5	10
16	39.2	<b>48.3</b>	51.3	55.8	60.5
20	40.1	48.1	49.6	56.1	60.7
24	41.5	47.4	<b>51.5</b>	<b>57.7</b>	<b>61.2</b>
28	41.9	47.6	50.9	56.5	60.4
32	<b>42.1</b>	46.2	51.3	56.9	60.3

Table 5: The performance (%) of base and new object categories.

Method	Base AP50			New AP50		
	1	3	5	1	3	5
MPSR [35]	59.9	68.5	69.4	40.7	48.9	53.6
MPSR + Ours	61.3	69.4	69.8	41.5	<b>51.5</b>	57.7
FSCE [30]	78.3	74.2	76.6	44.2	45.7	58.3
FSCE + Ours	<b>78.6</b>	<b>74.8</b>	<b>77.8</b>	<b>46.1</b>	48.9	<b>60.0</b>



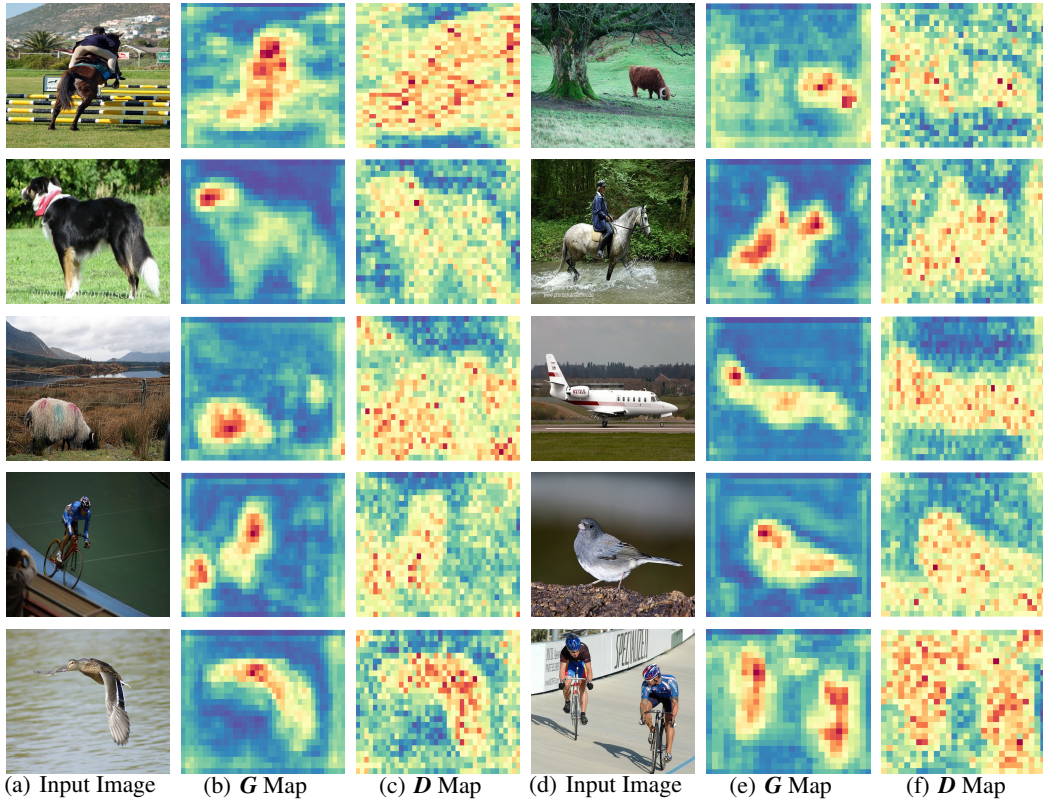


Figure 4: Visualization of the generalization ( $G$ ) map and discrimination ( $D$ ) map. The first, second, third, fourth, and fifth rows separately denote 1-shot, 2-shot, 3-shot, 5-shot, and 10-shot cases. For each feature map, the channels corresponding to the maximum value are selected for visualization.

298 Besides, when the number is large, e.g., 32, the performance also decreases. The reason may be  
 299 that employing more codewords increases the parameters, which leads to overfitting on new object  
 300 categories. For our method, the performance of using 24 codewords is the best.

301 **The Performance of Base Object Categories.** In Table 5, we can see that plugging our method into  
 302 MPSR [35] and FSCE [30] improves not only the performance of new object categories but also the  
 303 performance of base object categories. This further shows our method is beneficial for enhancing  
 304 generalization and discrimination, which is conducive to the improvement of detection performance.

#### 305 4.4 Visualization Analysis

306 In Fig. 4, we give visualization examples of the generalization ( $G$ ) map and discrimination ( $D$ ) map.  
 307 We can see that the generalization map focuses on the representative object characteristics, e.g., the  
 308 head of the dog and bird, which are helpful for improving generalization and accuracy of localization.  
 309 Meanwhile, the discrimination map contains rich information of background and object, which  
 310 enables the following dictionary learning to sufficiently capture high-level discriminative information.  
 311 This further shows that our method is effective to enhance both the generalization and discrimination  
 312 abilities for FSOD.

## 313 5 Conclusion

314 In this paper, for FSOD, we focus on improving generalization and discrimination via SVD-Dictionary  
 315 enhancement. Specifically, the eigenvectors corresponding to larger singular values are used to  
 316 calculate a generalization map. And the eigenvectors corresponding to relatively smaller singular  
 317 values are garnered to compute a discrimination map. Meanwhile, dictionary learning is employed to  
 318 capture high-level discriminative information from the discrimination map. The experimental results  
 319 and visualization analysis demonstrate the superiorities of our proposed method.

## References

- 320
- 321 [1] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey  
322 Zagoruyko. End-to-end object detection with transformers. *European Conference on Computer Vision*, 2020.
- 323 [2] Hao Chen, Yali Wang, Guoyou Wang, and Yu Qiao. Lstd: A low-shot transfer detector for object detection.  
324 In *AAAI-18 AAAI Conference on Artificial Intelligence*, pages 2836–2843, 2018.
- 325 [3] Xinyang Chen, Sinan Wang, Mingsheng Long, and Jianmin Wang. Transferability vs. discriminability:  
326 Batch spectral penalization for adversarial domain adaptation. In *International conference on machine  
327 learning*, pages 1081–1090. PMLR, 2019.
- 328 [4] Marc Peter Deisenroth, A Aldo Faisal, and Cheng Soon Ong. *Mathematics for machine learning*. Cambridge  
329 University Press, 2020.
- 330 [5] Guneet Singh Dhillon, Pratik Chaudhari, Avinash Ravichandran, and Stefano Soatto. A baseline for few-shot  
331 image classification. In *ICLR 2020 : Eighth International Conference on Learning Representations*, 2020.
- 332 [6] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal  
333 visual object classes (voc) challenge. *International journal of computer vision*, 88(2):303–338, 2010.
- 334 [7] Mark Everingham, SM Ali Eslami, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew  
335 Zisserman. The pascal visual object classes challenge: A retrospective. *International journal of computer  
336 vision*, 111(1):98–136, 2015.
- 337 [8] Qi Fan, Wei Zhuo, Chi-Keung Tang, and Yu-Wing Tai. Few-shot object detection with attention-rpn  
338 and multi-relation detector. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern  
339 Recognition*, pages 4013–4022, 2020.
- 340 [9] Zhibo Fan, Yuchen Ma, Zeming Li, and Sun Jian. Generalized few-shot object detection without forgetting.  
341 In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2021.
- 342 [10] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep  
343 networks. In *ICML'17 Proceedings of the 34th International Conference on Machine Learning - Volume 70*,  
344 pages 1126–1135, 2017.
- 345 [11] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate  
346 object detection and semantic segmentation. In *Proceedings of the IEEE conference on computer vision and  
347 pattern recognition*, pages 580–587, 2014.
- 348 [12] Ross Girshick. Fast r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages  
349 1440–1448, 2015.
- 350 [13] Micah Goldblum, Liam Fowl, and Tom Goldstein. Adversarially robust few-shot learning: A meta-learning  
351 approach. In *Advances in Neural Information Processing Systems*, volume 33, pages 17886–17895, 2020.
- 352 [14] Nima Hatami, Mohsin Bilal, and Nasir Rajpoot. Deep multi-resolution dictionary learning for histopathol-  
353 ogy image analysis. *arXiv preprint arXiv:2104.00669*, 2021.
- 354 [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition.  
355 In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016.
- 356 [16] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE  
357 international conference on computer vision*, pages 2961–2969, 2017.
- 358 [17] Bingyi Kang, Zhuang Liu, Xin Wang, Fisher Yu, Jiashi Feng, and Trevor Darrell. Few-shot object detection  
359 via feature reweighting. In *Proceedings of the IEEE International Conference on Computer Vision*, pages  
360 8420–8429, 2019.
- 361 [18] Leonid Karlinsky, Joseph Shtok, Sivan Harary, Eli Schwartz, Amit Aides, Rogerio Feris, Raja Giryes, and  
362 Alex M. Bronstein. Repmet: Representative-based metric learning for classification and few-shot object  
363 detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*,  
364 June 2019.
- 365 [19] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot,  
366 Ce Liu, and Dilip Krishnan. Supervised contrastive learning. In *Advances in Neural Information Processing  
367 Systems*, volume 33, pages 18661–18673, 2020.

- 368 [20] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár,  
369 and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer  
370 vision*, pages 740–755. Springer, 2014.
- 371 [21] Tsung-Yi Lin, Piotr Dollar, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature  
372 pyramid networks for object detection. In *2017 IEEE Conference on Computer Vision and Pattern Recognition  
373 (CVPR)*, pages 936–944, 2017.
- 374 [22] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and  
375 Alexander C Berg. Ssd: Single shot multibox detector. In *European conference on computer vision*, pages  
376 21–37. Springer, 2016.
- 377 [23] Jinlu Liu, Liang Song, and Yongqiang Qin. Prototype rectification for few-shot learning. *European  
378 Conference on Computer Vision*, 2020.
- 379 [24] Jiang Lu, Pinghua Gong, Jieping Ye, and Changshui Zhang. Learning from very few samples: A survey.  
380 *arXiv preprint arXiv:2009.02653*, 2020.
- 381 [25] Juan-Manuel Perez-Rua, Xiatian Zhu, Timothy M Hospedales, and Tao Xiang. Incremental few-shot object  
382 detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages  
383 13846–13855, 2020.
- 384 [26] Joseph Redmon and Ali Farhadi. Yolo9000: better, faster, stronger. In *Proceedings of the IEEE conference  
385 on computer vision and pattern recognition*, pages 7263–7271, 2017.
- 386 [27] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection  
387 with region proposal networks. In *Advances in neural information processing systems*, pages 91–99, 2015.
- 388 [28] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang,  
389 Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge.  
390 *International journal of computer vision*, 115(3):211–252, 2015.
- 391 [29] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. In *Advances  
392 in Neural Information Processing Systems*, pages 4077–4087, 2017.
- 393 [30] Bo Sun, Banghuai Li, Shengcai Cai, Ye Yuan, and Chi Zhang. Fsce: Few-shot object detection via  
394 contrastive proposal encoding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern  
395 Recognition*, 2021.
- 396 [31] Yonglong Tian, Yue Wang, Dilip Krishnan, Joshua B. Tenenbaum, and Phillip Isola. Rethinking few-shot  
397 image classification: a good embedding is all you need? *arXiv preprint arXiv:2003.11539*, 2020.
- 398 [32] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. Matching networks for one shot  
399 learning. In *Advances in neural information processing systems*, pages 3630–3638, 2016.
- 400 [33] Yu-Xiong Wang, Deva Ramanan, and Martial Hebert. Meta-learning to detect rare objects. In *Proceedings  
401 of the IEEE International Conference on Computer Vision*, pages 9925–9934, 2019.
- 402 [34] Xin Wang, Thomas E Huang, Trevor Darrell, Joseph E Gonzalez, and Fisher Yu. Frustratingly simple  
403 few-shot object detection. *International Conference on Machine Learning*, 2020.
- 404 [35] Jiayi Wu, Songtao Liu, Di Huang, and Yunhong Wang. Multi-scale positive sample refinement for few-shot  
405 object detection. *European Conference on Computer Vision*, 2020.
- 406 [36] Yang Xiao and Renaud Marlet. Few-shot object detection and viewpoint estimation for objects in the wild.  
407 *European Conference on Computer Vision*, 2020.
- 408 [37] Xiaopeng Yan, Ziliang Chen, Anni Xu, Xiaoxi Wang, Xiaodan Liang, and Liang Lin. Meta r-cnn: Towards  
409 general solver for instance-level low-shot learning. In *Proceedings of the IEEE International Conference on  
410 Computer Vision*, pages 9577–9586, 2019.
- 411 [38] Yukuan Yang, Fangyu Wei, Miaoqing Shi, and Guoqi Li. Restoring negative information in few-shot object  
412 detection. *Advances in Neural Information Processing Systems*, 2020.
- 413 [39] Zhongqi Yue, Hanwang Zhang, Qianru Sun, and Xian-Sheng Hua. Interventional few-shot learning. In  
414 *Advances in Neural Information Processing Systems*, volume 33, pages 2734–2746, 2020.
- 415 [40] Hang Zhang, Jia Xue, and Kristin Dana. Deep ten: Texture encoding network. In *Proceedings of the IEEE  
416 conference on computer vision and pattern recognition*, pages 708–717, 2017.
- 417 [41] Ruixiang Zhang, Tong Che, Zoubin Ghahramani, Yoshua Bengio, and Yangqiu Song. Metagan: An  
418 adversarial approach to few-shot learning. *NeurIPS*, 2:8, 2018.

419 **Checklist**

- 420 1. For all authors...
- 421 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's  
422 contributions and scope? [Yes]
- 423 (b) Did you describe the limitations of your work? [Yes] See Sec. 3.4
- 424 (c) Did you discuss any potential negative societal impacts of your work? [No] Our work  
425 is to improve the generalization and discrimination of detectors, which is helpful for  
426 promoting the development of detectors. Thus, our work do not have potential negative  
427 societal impacts.
- 428 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
429 them? [Yes]
- 430 2. If you are including theoretical results...
- 431 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 432 (b) Did you include complete proofs of all theoretical results? [N/A]
- 433 3. If you ran experiments...
- 434 (a) Did you include the code, data, and instructions needed to reproduce the main exper-  
435 imental results (either in the supplemental material or as a URL)? [No] We did not  
436 include the code, but we include the details in the paper.
- 437 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
438 were chosen)? [Yes] See Sec 4.1.
- 439 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
440 ments multiple times)? [N/A]
- 441 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
442 of GPUs, internal cluster, or cloud provider)? [N/A]
- 443 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 444 (a) If your work uses existing assets, did you cite the creators? [Yes] We cite many related  
445 works.
- 446 (b) Did you mention the license of the assets? [N/A]
- 447 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 448
- 449 (d) Did you discuss whether and how consent was obtained from people whose data you're  
450 using/curating? [N/A]
- 451 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
452 information or offensive content? [N/A]
- 453 5. If you used crowdsourcing or conducted research with human subjects...
- 454 (a) Did you include the full text of instructions given to participants and screenshots, if  
455 applicable? [N/A]
- 456 (b) Did you describe any potential participant risks, with links to Institutional Review  
457 Board (IRB) approvals, if applicable? [N/A]
- 458 (c) Did you include the estimated hourly wage paid to participants and the total amount  
459 spent on participant compensation? [N/A]