# AFD-Net: Adaptive Fully-Dual Network for Few-Shot Object Detection

Longyao Liu, Bo Ma, Yulin Zhang, Xin Yi and Haozhi Li
Beijing Institute of Technology
roel_liu@bit.edu.cn

## ABSTRACT

Few-shot object detection (FSOD) aims at learning a detector that can fast adapt to previously unseen objects with scarce annotated examples, which is challenging and demanding. Existing methods solve this problem by performing subtasks of classification and localization utilizing a shared component (*e.g.*, RoI head) in the detector, yet few of them take the distinct preferences of two subtasks towards feature embedding into consideration. In this paper, we carefully analyze the characteristics of FSOD, and present that a general few-shot detector should consider the explicit decomposition of two subtasks, as well as leveraging information from both of them to enhance feature representations. To the end, we propose a simple yet effective Adaptive Fully-Dual Network (AFD-Net). Specifically, we extend Faster R-CNN by introducing Dual Query Encoder and Dual Attention Generator for separate feature extraction, and Dual Aggregator for separate model reweighting. Spontaneously, separate state estimation is achieved by the R-CNN detector. Besides, for the acquisition of enhanced feature representations, we further introduce Adaptive Fusion Mechanism to adaptively perform feature fusion in different subtasks. Extensive experiments on PASCAL VOC and MS COCO in various settings show that, our method achieves new state-of-the-art performance by a large margin, demonstrating its effectiveness and generalization ability.

## CCS CONCEPTS

• **Computing methodologies → Object detection**.

## KEYWORDS

Few-Shot Object Detection, Meta-Learning, Task Decomposition, Adaptive Feature Fusion

## 1 INTRODUCTION

Recent years have witnessed impressive advances of convolutional neural networks (CNNs) [9, 18, 19] in object detection [4, 8, 14, 15, 24, 31, 32, 34], due to the availability of large-scale benchmarks with accurate annotations [10, 25, 35]. However, training general object detection models from scratch typically requires rich labeled data, which is extremely expensive to obtain or even hard to collect, such as endangered animals or certain medical data. Thus, detectors significantly suffer a performance drop when training examples are inadequate [47]. On the contrary, humans exhibit a strong ability to address this issue: even a child can easily learn to recognize novel characteristics from only a few instances [36].

This triggers recent researches on few-shot learning (FSL). It is considered promising to enhance the generalization ability of deep networks from limited training examples [1, 12, 23, 27, 39]. Concretely, FSL aims at recognizing instances from novel classes
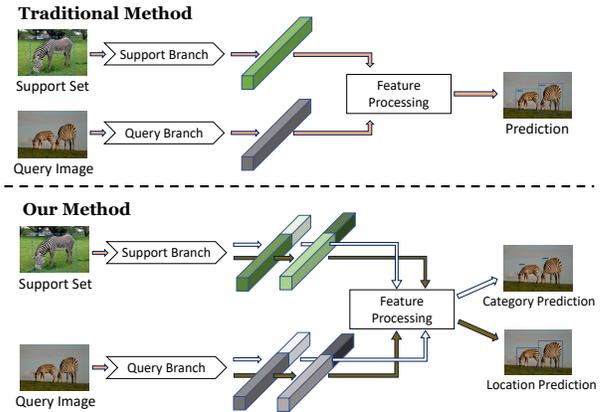


**Figure 1: Comparison of our framework and traditional methods. Traditional methods use shared encoded features to perform the estimation of object categories and locations, while our Adaptive Fully-Dual Network decomposes the tasks of classification (white arrow) and localization (brown arrow). Besides, feature fusion mechanism (concatenating encoded features with different colors) is introduced for enhancing feature representations.**

given only a few annotated data per category in the inference stage, with the availability of abundant labeled training samples from base classes. Most researches in few-shot learning community focus on image classification [13, 22, 29, 37–39], while far less progress has been made in the field of object detection [20, 44, 46], which is generally considered much more challenging due to the existence of an additional subtask, *i.e.*, few-shot localization, besides the subtask of object recognition.

A trend to solve few-shot object detection is to conjoin a reweighting module with a base object detector, *e.g.*, Faster R-CNN [34]. Concretely, these few-shot detectors [44, 46] introduce an additional branch to extract discriminative features of support set from novel classes, and then use these features to reweight RoI (Region-of-Interest) head, which is shared for subsequent estimation of object categories and locations. However, classification and localization are two fundamentally different subtasks in few-shot object detection. The former subtask focuses on providing a coarse location of the object via classification, while the latter aims at estimating an accurate object state by a refined bounding box [43, 45]. This leads to distinct preferences towards feature representation of these two subtasks. In other words, generated features for classification are probably not suitable for bounding box regression, and using the same RoI representations in two subtasks is suboptimal. Owing to the explicit task decomposition, in the process of a specific subtask, the obtained features for state estimation can be enriched by the

information from the unfocused subtask. Therefore, it is crucial to exploit an efficient feature fusion mechanism between tasks of classification and location regression.

Motivated by the abovementioned analysis, we present that a general few-shot detector should: 1) decouple the process of classification from that of localization, involving feature representations, model reweighting and state estimation; 2) introduce efficient feature fusion mechanism between two subtasks into the encoder of few-shot detector, to enhance feature representations for both query image and support set, as shown in Fig. 1.

In this paper, we propose a novel and intuitive framework for few-shot object detection, namely *Adaptive Fully-Dual Network (AFD-Net)*, as illustrated in Fig.2. Specifically, we extend Faster R-CNN by introducing three modules, *i.e.*, Dual Query Encoder (DQE), Dual Attention Generator (DAG), and Dual Aggregator (DA), for separate process of two subtasks in various key stages of FSOD. Besides, Adaptive Fusion Mechanism (AFM), guiding the design of DQE and DAG, is introduced for the acquisition of enhanced feature representations. Query features generated by backbone network are encoded by DQE into two groups of RoI vectors. These vectors contain meta information generalizable to detect novel objects within RoIs, and each group is utilized for a specific subtask. In parallel with DQE, DAG takes support images as input and encodes them into class-attentive vectors in the same scheme. DA, consisting of two aggregators for two subtasks, performs separate model reweighting. Each RoI vector from query image and each class-attentive vector from support set assigned the same subtask are aggregated in the corresponding aggregator. In this way, some meta features of query image informative for detecting novel objects would be activated. Finally, the aggregated features are fed into the R-CNN detector to estimate the object location and category, respectively, achieving separate state estimation, spontaneously.

Extensive experiments on two public datasets, *i.e.*, PASCAL VOC and MS COCO, in various settings, show that despite its simplicity, our method outperforms state-of-the-art approaches by a large margin, demonstrating its effectiveness and generalization ability.

In summary, the main contributions of this paper are three-fold:

- We propose a simple yet effective framework for few-shot object detection. To the best of our knowledge, we are among the first to solve FSOD by decomposing the process of classification and localization.
- We further introduce three modules, *i.e.*, Dual Query Encoder, Dual Attention Generator and Dual Aggregator, to perform decomposition in multiple components of our network, along with Adaptive Fusion Mechanism for enhancing feature representations.
- The experimental results demerstrate that our proposed method achieves new state-of-the-art performance on multiple benchmarks.

## 2 RELATED WORK

**General Object Detection.** Recent object detectors based on deep CNNs can be mainly divided into two steams, *i.e.*, one-stage detectors and two-stage ones. R-CNN series [8, 14, 15, 17, 24, 34] belong to the first category, which generate region proposals [34] of potential objects in the first stage, and then perform category and bounding box estimation at the proposal-level. On the contrary, YOLO [31] and the variants [3, 26, 32, 33] dominate the one-stage steam. These methods use a single CNN to predict categories and locations of the objects directly, without explicitly generating proposals. These two branches of general object detectors unanimously depend on a huge amount of data with elaborate bounding box annotations. When training samples are limited, they struggle heavily.

**Few-shot Learning.** Few-shot learning refers to learning to learn general knowledge that can be easily transferred to new tasks with only a handful of annotated examples [1, 23, 39]. Few-shot classification has recently been widely investigated as a representative task of few-shot learning. Generally, solutions to this problem involve two groups: meta-learning based and fine-tuning based methods. The former can be further categorized into: 1) metric-learning based methods [22, 37, 38] that focus on the similarity of input images in the embedding space; 2) optimization-based ones [12, 27, 30], where a meta-learner is designed to simulate the optimization process for the fast adaptation to novel classes with limited samples; 3) model-based methods [2, 13, 41] that aim to estimate the network parameters for novel tasks with the help of a learned predictor. Fine-tuning based approaches [6, 7] demonstrate that simple fine-tuning techniques are crucial and effective towards few-shot learning.

**Few-shot Object Detection.** While substantial progress has been made in few-shot classification, the problem of few-shot object detection is relatively unsolved, since object localization should be additionally processed in the scenario where distracting irrelevant objects may exist or even there is no object within an image. LSTD [5] applies transfer learning technique in object detection with limited examples and demonstrates its effectiveness. TFA [40] proposes that only fine-tuning the last layer of existing detectors is crucial and effective despite its simplicity. Metric learning widely explored in few-shot classification can be extended to few-shot object detection. For example, RepMet [21] scores the pair-wise similarity between embedded features of input query and support images. [11] incorporates attention mechanism into RPN (Region Proposal Network) and explicitly defines a multi-relation detector to suppress detection of the background. Recently, methods based on meta-learning have been a popular trend. FSRW [20] attaches a reweighting module with YOLOv2 [32] to adjust the full input image features using class-attentive vectors of support images, while Meta R-CNN [46] applies reweighting on RoI features based on Faster R-CNN [34]. MetaDet [42] presents a meta-model to predict the parameters of category-agnostic and category-specific components in a detector, separately. FSDetView [44] applys a slightly more complicated feature aggregation scheme to further improve the detection performance and evaluation reliability.

## 3 METHODOLOGY

### 3.1 Problem Setup

As in previous work [20, 44, 46], we adopt the following few-shot object detection settings. In the training phase, we are provided with two training sets, *i.e.*, a base set $D_{base}$ from base classes $C_{base}$ with abundant instances and a novel set $D_{novel}$ from novel classes $C_{novel}$ with only a few samples per category, where $C_{base}$ and $C_{novel}$ are non-overlapping, *i.e.*, $C_{base} \cap C_{novel} = \varnothing$. For each sample $(x, y) \in$
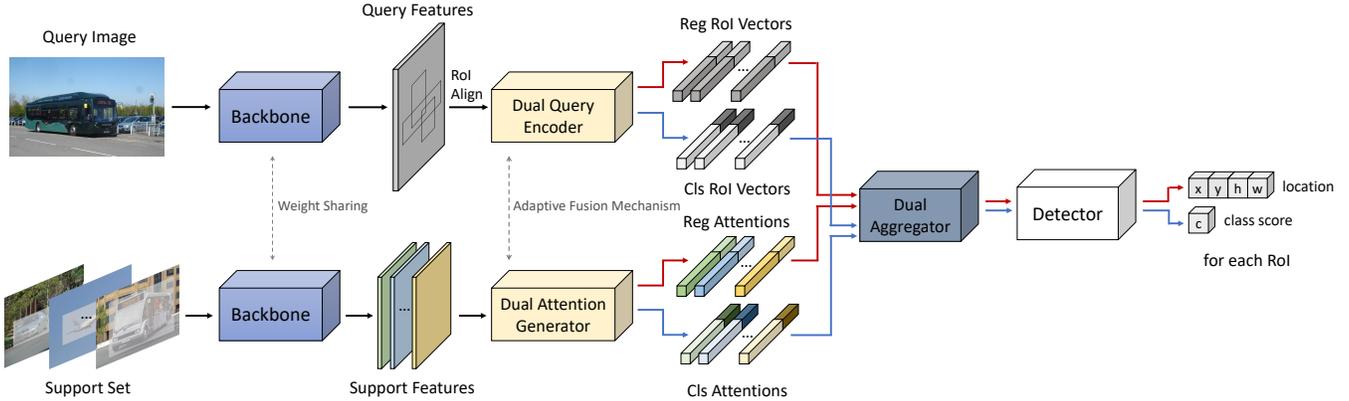
**Figure 2: The pipeline of Adaptive Fully-Dual Network for few-shot object detection. It consists of three key components: 1) the top query branch receives a query image and generates its RoI vectors; 2) the bottom support branch encodes support set into its class-attentive vectors; 3) two groups of features from two branches are aggregated for subsequent estimation of object categories and locations with a R-CNN detector. We further introduce three modules, *i.e.*, Dual Query Encoder, Dual Attention Generator and Dual Aggregator in three components, along with Adaptive Fusion Mechanism for enhancing feature representations. All these three modules are in a dual architecture, where the specific subtask is performed in the assigned path. Black line represents shared path for two subtasks, red and blue lines represent subtasks of bounding box regression and category classification, respectively.**

$D_{base} \cup D_{novel}$, $x = \{obj_i\}_{i=1}^{N}$ is an image containing $N$ objects, and $y = \{(cls_i, box_i)\}_{i=1}^{N}$ denotes $N$ categories $\{cls_i\}_{i=1}^{N}$ with each $cls_i \in C_{base} \cup C_{novel}$ along with $N$ structured annotations $\{box_i\}_{i=1}^{N}$ of the $N$ objects in the image $x$. The aim of the few-shot object detector is to classify and locate objects from both novel and base classes in an image with only $K$ (usually less than 10) available instances per class in the phase of inference, with the help of transferable knowledge learned from abundant examples from base classes.

## 3.2 Adaptive Fully-Dual Network

**Pipeline.** Our proposed Adaptive Fully-Dual Network decomposes the process of few-shot classification and localization, as well as enhancing feature representations. We adopt the widely used two-stage detector Faster R-CNN [34] as the base model, which first generates RoIs of potential objects, and then performs state estimation using a single RoI head. The pipeline of our proposed network is demonstrated in Fig. 2. Concretely, we extend Faster R-CNN to a dual Siamese architecture, where a query image is encoded by the top Faster R-CNN branch, and the bottom branch is designed for support set. Each branch further contains two paths, and feature extraction for the specific subtask is performed in the corresponding path. Generated features from two branches are then aggregated, achieving the reweighting of Faster R-CNN for the subsequent detection of novel instances. We introduce Dual Query Encoder (DQE) and Dual Attention Generator (DAG) for separate feature representations, along with Dual Aggregator (DA) for separate model reweighting, in our proposed architecture. Furthermore, we present Adaptive Fusion Mechanism (AFM) to efficiently encode both the query image and support set.

**Dual Feature Representations.** Adaptive Fully-Dual Network consists of two branches, *i.e.*, the top Faster R-CNN branch for

processing the query image and the bottom support branch for support set, as shown in Fig. 2. The Faster R-CNN branch aims to learn meta features generalizable to detect novel objects within the input query image. The support branch encodes the input support set into discriminative vectors, which will be subsequently utilized to adjust the contribution of meta features generated from Faster R-CNN branch. In this way, meta features informative for detecting novel objects would be activated. Therefore, these support vectors, namely class-attentive vectors, can be seen as the attention coefficients of meta features.

In two-stage object detectors, for an input image $Q$, RPN is applied to generate $n$ class-agnostic RoIs. These RoIs are then embedded by the RoI head into RoI vectors for subsequent state estimation. Instead of generating a single group of RoI vectors $\{r_i\}_{i=1}^{n}$ for the input query image $Q$ by a single RoI head in previous work [44, 46], we propose Dual Query Encoder $\mathcal{E}$ in Faster R-CNN branch, to obtain two groups of $m$-dimensional RoI vectors $\left\{r_i^{cls}, r_i^{reg}\right\}_{i=1}^{n} \in \mathbb{R}^{2n \times m}$ containing meta information for subtasks of classification and bounding box regression respectively, guided by our observation that these two subtasks in few-shot detection should be treated separately. This procedure is formulated as below:

$$\left\{r_i^{cls}, r_i^{reg}\right\}_{i=1}^{n} = \mathcal{E}(\mathcal{R}(\mathcal{B}(Q))) \tag{1}$$

where $\mathcal{B}$ denotes the backbone network of Faster R-CNN, and $\mathcal{R}$ denotes the RoIAlign operation.

As for the support branch, it takes support images as input and generates their class-attentive vectors for subsequently reweighting meta features from Faster R-CNN branch. To effectively capture the support information, the input support image $S$ is concatenated with a structured binary mask $M$ indicating the bounding box annotation of the target object to detect [20]. Supposing there are
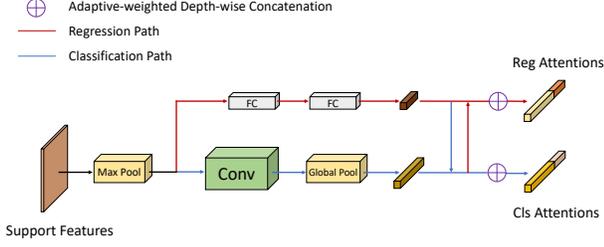
Figure 3: Illustration of Adaptive Fusion Mechanism.



Figure 4: Illustration of feature aggregation scheme.

$s$ categories in the support set. Support features of input support image concatenated with its associated mask $[S_j, M_j]$ from class $j, j = 1, 2, \ldots, s$, are generated from backbone network $\mathcal{B}$ sharing weights with Faster R-CNN branch. Then our proposed Dual Attention Generator $\mathcal{G}$ embeds them into two $m$-dimensional class-attentive vectors $\left\{a_j^{cls}, a_j^{reg}\right\} \in \mathbb{R}^{2 \times m}$ for category classification and bounding box regression, respectively:

$$\left\{a_j^{cls}, a_j^{reg}\right\} = \mathcal{G}(\mathcal{B}[S_j, M_j]) \tag{2}$$

Each RoI vector of query image and each class-attentive vector of support set assigned the same subtask will be aggregated by our proposed Dual Aggregator $\mathcal{D}$ before being fed into R-CNN detector for state estimation. We will discuss the details next.

In this paper, the input query image $Q$ and support image $S$ are resized into $224 \times 224$. The backbone network $\mathcal{B}$ is the first four ResNet blocks, and the input dimension of support branch is $224 \times 224 \times 4$. The size of each query RoI is $7 \times 7 \times 1024$, and that of each support feature is $14 \times 14 \times 1024$. All of the RoI vectors and class-attentive vectors are 3072-dimensional.

**Adaptive Fusion Mechanism.** To enhance feature representations for both the query and support set, we propose Adaptive Fusion Mechanism. Its key philosophy is performing feature extraction for two subtasks in distinct manners, as well as using the information from the unfocused subtask for further enhancing feature representations. Adaptive Fusion Mechanism guides the design of Dual Query Encoder and Dual Attention Generator. Actually these two modules share weights except for the additional Max Pool layer in Dual Attention Generator, and we take this module as an example to elaborate this mechanism, as shown in Fig. 3.

Dual Attention Generator firstly applies a Max Pool layer to ensure the consistency in the size of input feature with Dual Query Encoder. Then it involves two parallel branches and each branch encodes input support images into class-attentive vectors for the specific subtask $t$, *i.e.*, $cls$ for category classification and $reg$ for bounding box regression. We adopt a two-layer fully-connected (2-fc) encoder $\mathcal{G}_{fc}$ in regression branch and a convolution ($conv$) encoder $\mathcal{G}_{conv}$ in classification branch. Output task-specific class-attentive vectors from two branches are the concatenation of weighted $fc$ features and $conv$ features. The weights can be adaptively adjusted according to the specific subtask, and their values indicate the contributions to the integrated task-specific outputs. Therefore, Eq. (2) can be rewritten as:

$$\left\{a_j^t\right\} = \left[\lambda_{conv}^t \mathcal{G}_{conv}\mathcal{B}([S_j, M_j]), \lambda_{fc}^t \mathcal{G}_{fc}\mathcal{B}([S_j, M_j])\right], t \in \{cls, reg\} \tag{3}$$
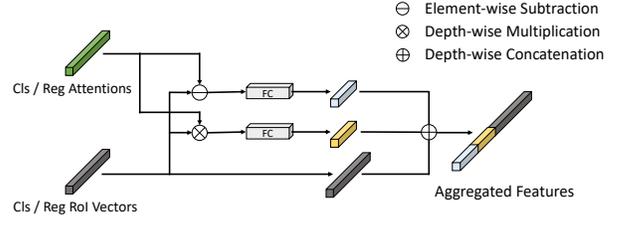
where $\lambda_{conv}^t$ and $\lambda_{fc}^t$ denote the learnable weights of corresponding task-specific class-attentive vectors and are initialized to 1, $[\cdot, \cdot]$ denotes the depth-wise concatenation operation.

In this paper, the kernel size of Max Pool layer is $2 \times 2$. The $conv$ encoder $\mathcal{G}_{conv}$ is the last ResNet block, and the Global Pool layer encodes the output $conv$ features into a 2048-dimensional vector. Two layers of $fc$ encoder $\mathcal{G}_{fc}$ are both 1024-dimensional, thus the output vector of $\mathcal{G}_{fc}$ is 1024-dimensional. Output task-specific class-attentive vectors are 3072-dimensional.

**Dual Feature Aggregation.** The class-attentive vectors are responsible for reweighting RoI vectors and activating those that encode potential novel objects. The aggregated RoI vectors are then fed into the R-CNN detector for state estimation. In this paper, besides separate feature representations for two subtasks, we decompose the process of feature aggregation to realize separate model reweighting, leading to our Dual Aggregator $\mathcal{A}$. It consists of a classification aggregator and a bounding box regression aggregator in parallel. In each aggregator, each RoI vector of query image is aggregated with each class-attentive vector of support features following the scheme introduced in [44], as shown in Fig. 4, which can be formulated as:

$$
\begin{aligned}
r_{i,j}^t &= \mathcal{A}\left(r_i^t, a_j^t\right) \\
&= \left[f_m(r_i^t \otimes a_j^t), f_s(r_i^t - a_j^t), r_i^t\right], t \in \{cls, reg\} \\
&\quad \text{for each } i \in \{1, 2, \cdots, n\}, j \in \{1, 2, \cdots, s\}
\end{aligned}
\tag{4}
$$

where $r_{i,j}^t$ denotes the $i$-th RoI vector reweighted by the $j$-th class-attentive vector from the aggregator assigned the specific subtask $t$, and $\otimes$ denotes the depth-wise multiplication implemented through $1 \times 1$ depth-wise convolution. $f_m$ and $f_s$ are FC layers in Fig. 4, and they are both 1536-dimensional. Therefore, the output aggregated features are 6144-dimensional.

**Dual State Estimation.** As the input features of R-CNN detector are divided into two groups $r_{i,j}^{cls}$ and $r_{i,j}^{reg}$ where $i = 1, 2, \ldots, n$, $j = 1, 2, \ldots, s$, the state estimation in our framework is decomposed as well. Specifically, the aggregated RoI vector $r_{i,j}^{cls}$ is utilized for predicting the probability of the $i$-th RoI containing an object from novel class $j$. Therefore, the classifier will produce $s$ outcomes for each RoI, and the category $h$ with the highest confidence score will be assigned. Then the associated RoI vector $r_{i,h}^{reg}$ will be used for location estimation of the $i$-th RoI by the regressor. If the highest confidence score is lower than the threshold, this RoI will be regarded as background and discarded, consistent with the process in [44, 46].

## 3.3 Training Procedure

**Training Phase.** Following the common practice in [20, 40, 44, 46], we train our network in two phases. Base data $D_{base}$ from base classes $C_{base}$ with abundant samples per class are used to train the model in the first base training phase. Then the balanced base data $D_{base}$ and novel data $D_{novel}$ with only $K$ samples per class are fed into the network in the second fine-tuning phase for the fast adaptation to novel classes.

**Training Data Organization.** In the training phase, distinct from general object detectors that applying an image $x_i$ as the training mini-batch, our few-shot detector applies a task $T_i$ as an input training data in the meta-learning paradigm [20, 44, 46]. Each input task $T_i = S_i \cup Q_i$ is the union of a query set $Q_i$ and a support set $S_i$, where $Q_i$ provided to query branch in Fig.2 is a query image $q_i$ containing objects from $m$ classes $C_i^{meta} \subseteq C_{base} \cup C_{novel}$, and $S_i$ for support branch contains $m$ $K$-shot ($K = 200$ in the base training phase and $K = \{1, 2, 3, 5, 10\}$ in the fine-tuning phase) clusters $\{g_i^p\}_{p=1}^m$, where each cluster $g_i^p$ includes $K$ images $\{s_{i,j}^p\}_{j=1}^K$ in the category $p$. Each support image $s_{i,j}^p$ is depth-wise concatenated with a structured binary $mask_{i,j}^p$ (see the input support set in Fig. 2, only one object is considered when multiple objects are present within an image).

**Loss Function.** We optimize our network in both training phases using the loss function introduced in [46]:

$$L = L_{Faster\ R\text{-}CNN} + L_{meta} \tag{5}$$

where $L_{Faster\ R\text{-}CNN}$ denotes the loss function of base detector Faster R-CNN, involving the RPN loss, and classification loss together with bounding box regression loss for the state estimation, $L_{meta}$ denotes a meta loss aiming at encouraging class-attentive vectors of support set to distinguish with each other.

Since our proposed Dual Attention Generator and Dual Query Encoder are in a dual architecture, $L_{meta}$ in this paper is a combination of two components for classification and regression, respectively. Thus the loss function in Eq. (5) can be precisely written as:

$$L = L_{Faster\ R\text{-}CNN} + L_{meta\text{-}cls} + L_{meta\text{-}reg} \tag{6}$$

where $L_{meta\text{-}cls}$ and $L_{meta\text{-}reg}$ denote the meta losses for classification and boundinng box regression, respectively, implemented by the cross-entropy loss.

## 4 EXPERIMENTS

### 4.1 Experimental Settings

**Benchmarks.** We evaluate our method on general object detection benchmarks, *i.e.*, PASCAL VOC 2007, 2012 and MS COCO, in few-shot detection settings. As for PASCAL VOC, we follow the common practice in previous work [20, 40, 44, 46], and use train/val sets of VOC 07 and 12 for training while the test set of VOC 07 for testing. This benchmark covers 20 object categories, where 15 of them are regarded as base classes for the base training phase and the remaining 5 categories with only $K$ ($K = \{1, 2, 3, 5, 10\}$) samples per class are considered as novel classes for few-shot fine-tuning phase. For the fair quantitative comparison, we adopt the same three class splits provided in [20]. For the evaluation protocol, we use mean Average Precision (*mAP*) of novel objects and the

Intersection of Union (*IoU*) is set as 0.5 ($AP_{50}$). Another benchmark we evaluate on is MS COCO, which has 80k train images and 40k validation images, covering 80 classes. Among them, we denote the 20 classes overlapped with PASCAL VOC as novel classes with $K$ ($K = \{10, 30\}$) samples per category and the rest 60 classes as base classes. We use 5k images from validation set for evaluation with the standard COCO-style evaluation metrics [26, 33] and the rest for training. To compare empirically, on both datasets, we evaluate on novel classes over $n$ ($n = 30$ for PASCAL VOC and $n = 10$ for MS COCO) repeated runs unless otherwise specified, and report the average performance.

**Implementation Details.** Our proposed approach is implemented in the *PyTorch* [28] library. It employs the backbone network of ResNet [18] pre-trained on ImageNet [35] along with a RoIAlign [16] layer. Specifically, we use ResNet-101 on PASCAL VOC and ResNet-50 on MS COCO. Our model is trained end-to-end using a batch size of 4 on a single Nvidia GeForce GTX 1080Ti with 13 GB memory. For optimization, we keep the setups in [44, 46] on both benchmarks. Concretely, we adopt the SGD optimizer with an initial learning rate of 0.001, weight decay of 0.0005 and momentum of 0.9. The model is trained for 20 epochs in the base training phase and 9 epochs in the $K$-shot fine-tuning phase. The learning rate is reduced by 0.1 every 5 epochs in two training phases. During inference, the support branch will be removed for directly performing detection without requiring support images as input. It is because that class-attentive vectors for model reweighting can be obtained in the $K$-shot fine-tuning phase by averaging the outputs from support branch over the input $K$ support images.

### 4.2 Comparison with State-of-the-art Methods

**PASCAL VOC.** Our evaluation results are presented in Table 1. The comparison experiments cover few-shot object detection scenarios in five setups ($K = \{1, 2, 3, 5, 10\}$) across three base/novel set splits. As can be observed, our method outperforms recent state-of-the-art methods by a large margin and obtains the best performance in all 15 cases. We notice that in the majority of cases, the improvements are much larger than the gap among previous approaches, which indicates the strong generalization ability of our model. Surprisingly, although the existence of high variance of support data in the extreme few-shot setup ($K = 1$), we obtain large improvements (+6.4% in the first split and +6.2% in the third split), showing the robustness of our model in tough scenarios. Besides, in the robust 10-shot setup, the improvements (+2.9% in the first split, +1.2% in the second split and +3.9% in the third split) are lower than those in other setups. This can be explained by that, as the number of instances per novel class increases, the sample variance decreases and the detection performance stabilizes.

Taking evaluation performance on base classes into consideration, we provide detailed results on the first base/novel split in Table 2. Note that our proposed method achieves a new state-of-the-art performance on both base and novel classes. Additionally, the performance improvements are much larger in novel classes compared with that in base classes. As for the performance on the single category, our approach has the best detection performance for most of categories except several base classes in 3-shot scenario, such as "table" and "train". We will give detailed analysis in the ablation study, from the perspective of sample variance.

**Table 1: Few-shot detection performance ($AP_{50}$) for novel categories on PASCAL VOC dataset. We evaluate baselines on three different novel sets. Our approach consistently outperforms other methods by a large margin. [*n]Reported results are averaged over n repeated runs.**

| Method / Shot | Novel Set 1 | | | | | Novel Set 2 | | | | | Novel Set 3 | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 5 | 10 | 1 | 2 | 3 | 5 | 10 | 1 | 2 | 3 | 5 | 10 |
| LSTD [5] | 8.2 | 11.0 | 12.4 | 29.1 | 38.5 | 11.4 | 3.8 | 5.0 | 15.8 | 31.0 | 12.6 | 8.5 | 15.0 | 27.3 | 36.3 |
| FSRW [20] | 14.8 | 15.5 | 26.7 | 33.9 | 47.2 | 15.7 | 15.2 | 22.7 | 30.1 | 40.5 | 21.3 | 25.6 | 28.4 | 42.8 | 45.9 |
| MetaDet[*5] [42] | 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 |
| Meta R-CNN[*5] [46] | 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 |
| TFA[*30] w/ fc [40] | 22.9 | 34.5 | 40.4 | 46.7 | 52.0 | 16.9 | 26.4 | 30.5 | 34.6 | 39.7 | 15.7 | 27.2 | 34.7 | 40.8 | 44.6 |
| TFA[*30] w/ cos [40] | 25.3 | 36.4 | 42.1 | 47.9 | 52.8 | 18.3 | 27.5 | 30.9 | 34.1 | 39.5 | 17.9 | 27.2 | 34.3 | 40.8 | 45.6 |
| FSDetView[*10] [44] | 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 |
| AFD-Net[*30] (Ours) | **31.7** | **41.4** | **49.5** | **54.6** | **60.3** | **23.2** | **31.3** | **38.4** | **41.9** | **46.9** | **27.4** | **35.3** | **41.7** | **46.7** | **53.5** |

**Table 2: Few-shot detection performance ($AP_{50}$) for base and novel categories on Novel Set 1 of PASCAL VOC dataset. Our approach outperforms other methods on both base and novel classes. [*n]Reported results are averaged over n repeated runs.**

| Shot | Method | Novel classes | | | | | | Base classes | | | | | | | | | | | | | | mAP |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | bird | bus | cow | mbike | sofa | mean | aero | bike | boat | bottle | car | cat | chair | table | dog | horse | person | plant | sheep | train | tv | mean | |
| 3 | LSTD [5] | 23.1 | 22.6 | 15.9 | 0.4 | 0.0 | 12.4 | 74.8 | 68.7 | 57.1 | 44.1 | 78.0 | 83.4 | 46.9 | 64.0 | 78.7 | **79.1** | 70.1 | 39.2 | 58.1 | **79.8** | **71.9** | 66.3 | 52.8 |
| | FSRW [20] | 26.1 | 19.1 | 40.7 | 20.4 | 27.1 | 26.7 | 73.6 | 73.1 | 56.7 | 41.6 | 76.1 | 78.7 | 42.6 | **66.8** | 72.0 | 77.7 | 68.5 | **42.0** | 57.1 | 74.7 | 70.7 | 64.8 | 55.2 |
| | Meta R-CNN[*5] [46] | 30.1 | 44.6 | **50.8** | 38.8 | 10.7 | 35.0 | 67.6 | 70.5 | **59.8** | 50.0 | 75.7 | 81.4 | 44.9 | 57.7 | 76.3 | 74.9 | 76.9 | 34.7 | 58.7 | 74.7 | 67.8 | 64.8 | 57.3 |
| | AFD-Net[*30] (Ours) | **51.8** | **60.3** | 43.8 | **60.5** | **31.3** | **49.5** | 69.9 | **75.2** | 56.9 | **57.9** | **79.5** | **84.2** | 47.9 | 60.5 | **82.4** | 76.7 | **77.4** | 40.8 | **68.6** | 75.2 | 71.0 | **68.3** | **63.6** |
| 10 | LSTD [5] | 22.8 | 52.5 | 31.3 | 45.6 | 40.3 | 38.5 | 70.9 | 71.3 | 59.8 | 41.1 | 77.1 | 81.9 | 45.1 | **67.2** | 78.0 | 78.9 | 70.7 | 41.6 | 63.8 | **79.7** | 66.8 | 66.3 | 59.4 |
| | FSRW [20] | 30.0 | 62.7 | 43.2 | 60.6 | 39.6 | 47.2 | 65.3 | 73.5 | 54.7 | 39.5 | 75.7 | 81.1 | 35.3 | 62.5 | 72.8 | 78.8 | 68.6 | 41.5 | 59.2 | 76.2 | 69.2 | 63.6 | 59.5 |
| | Meta R-CNN[*5] [46] | 52.5 | 55.9 | 52.7 | 54.6 | 41.6 | 51.5 | 68.1 | 73.9 | 59.8 | 54.2 | 80.1 | 82.9 | 48.8 | 62.8 | 80.1 | **81.4** | 77.2 | 37.2 | 65.7 | 75.8 | 70.6 | 67.9 | 63.8 |
| | AFD-Net[*30] (Ours) | **62.7** | **67.9** | **60.2** | **68.1** | **42.7** | **60.3** | **73.0** | **77.9** | **60.4** | **59.7** | **81.8** | **85.4** | **51.0** | 66.5 | **84.7** | 80.3 | **78.1** | **44.6** | **70.1** | 77.6 | **72.7** | **70.9** | **68.3** |

**Table 3: Few-shot detection performance for novel categories on MS COCO dataset. Our approach achieves an significant improvement over other methods with notably smaller standard deviations. [*n]Reported results are averaged over n repeated runs.**

| Shots | Method | Average Precision | | | | | | Average Recall | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 0.5:0.95 | 0.5 | 0.75 | S | M | L | 1 | 10 | 100 | S | M | L |
| 10 | LSTD [5] | 3.2 | 8.1 | 2.1 | 0.9 | 2.0 | 6.5 | 7.8 | 10.4 | 10.4 | 1.1 | 5.6 | 19.6 |
| | FSRW [20] | 5.6 | 12.3 | 4.6 | 0.9 | 3.5 | 10.5 | 10.1 | 14.3 | 14.4 | 1.5 | 8.4 | 28.2 |
| | MetaDet[*5] [42] | 7.1 | 14.6 | 6.1 | 1.0 | 4.1 | 12.2 | 11.9 | 15.1 | 15.5 | 1.7 | 9.7 | 30.1 |
| | Meta R-CNN[*5] [46] | 8.7 | 19.1 | 6.6 | 2.3 | 7.7 | 14.0 | 12.6 | 17.8 | 17.9 | 7.8 | 15.6 | 27.2 |
| | TFA[*10] w/ fc [40] | 9.1 ± 0.5 | 17.3±1.0 | 8.5±0.5 | - | - | - | - | - | - | - | - | - |
| | TFA[*10] w/ cos [40] | 9.1 ± 0.5 | 17.1±1.1 | 8.8±0.5 | - | - | - | - | - | - | - | - | - |
| | FSOD [11] | 11.1 | 20.4 | 10.6 | - | - | - | - | - | - | - | - | - |
| | FSDetView[*10] [44] | 12.5 | 27.3 | 9.8 | 2.5 | 13.8 | 19.9 | 20.0 | 25.5 | 25.7 | 7.5 | 27.6 | 38.9 |
| | AFD-Net[*10] (Ours) | **17.3 ± 0.1** | **33.4 ± 0.2** | **16.6 ±0.1** | **5.7±0.2** | **18.6±0.2** | **27.4±0.2** | **24.5±0.1** | **31.5±0.1** | **31.8±0.1** | **12.7±0.2** | **33.8±0.1** | **45.8±0.2** |
| 30 | LSTD [5] | 6.7 | 15.8 | 5.1 | 0.4 | 2.9 | 12.3 | 10.9 | 14.3 | 14.3 | 0.9 | 7.1 | 27.0 |
| | FSRW [20] | 9.1 | 19.0 | 7.6 | 0.8 | 4.9 | 16.8 | 13.2 | 17.7 | 17.8 | 1.5 | 10.4 | 33.5 |
| | MetaDet[*5] [42] | 11.3 | 21.7 | 8.1 | 1.1 | 6.2 | 17.3 | 14.5 | 18.9 | 19.2 | 1.8 | 11.1 | 34.4 |
| | Meta R-CNN[*5] [46] | 12.4 | 25.3 | 10.8 | 2.8 | 11.6 | 19.0 | 15.0 | 21.4 | 21.7 | 8.6 | 20.0 | 32.1 |
| | TFA[*10] w/fc [40] | 12.0±0.4 | 22.2±0.6 | 11.8±0.4 | - | - | - | - | - | - | - | - | - |
| | TFA[*10] w/cos [40] | 12.1±0.4 | 22.0±0.7 | 12.0±0.5 | - | - | - | - | - | - | - | - | - |
| | FSDetView[*10] [44] | 14.7 | 30.6 | 12.2 | 3.2 | 15.2 | 23.8 | 22.0 | 28.2 | 28.4 | 8.3 | 30.3 | 42.1 |
| | AFD-Net[*10] (Ours) | **19.1±0.04** | **35.8±0.1** | **18.7±0.2** | **5.7±0.1** | **20.6±0.1** | **29.2±0.2** | **26.0±0.2** | **33.4±0.2** | **33.7±0.2** | **14.3±0.1** | **36.1±0.2** | **47.6±0.4** |

**MS COCO.** We show the evaluation of 20 novel classes on MS COCO in 10/30-shot setups, and report the standard COCO-style metrics in Table 3. Obviously, our approach significantly outperforms recent state-of-the-art methods with much smaller confidence intervals, despite the complexity and a huge amount of data in MS COCO, verifying its effectiveness and robustness. Note that our method performs much better in detection of small objects in comparison with other state-of-the-arts, indicating that our proposed network indeed enhances feature representations and obtains high-quality image information.

In Fig. 5, we provide the comparison of qualitative 30-shot detection results between our approach and baseline method FS-DetView [44]. We can find that our method performs better on detecting small and occluded objects, e.g., the first two columns. This is consistent with the observations from Table 3. Besides, our model provides accurate locations of target objects, e.g., the third column. The last two columns show the ability of our approach in classifying objects from both base and novel categories, including correctly classifying target objects and avoiding duplicate estimations of object categories.

**Figure 5: Comparison of qualitative 30-shot detection results on base and novel classes from MS COCO between our approach and baseline method FSDetView [44].**
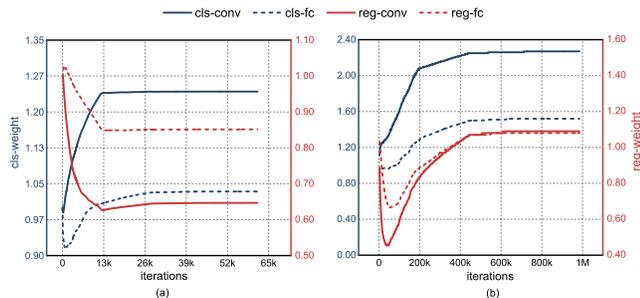


**Figure 6: Performance visualization of our proposed Adaptive Fusion Mechanism. We plot values of four learnable weights, *i.e.*, $\lambda_{conv}^{cls}$, $\lambda_{fc}^{cls}$, $\lambda_{conv}^{reg}$, and $\lambda_{fc}^{reg}$ against number of training iterations on: (a) PASCAL VOC in the first split, (b) MS COCO, in the base training phase. *cls-conv, cls-fc, reg-conv, reg-fc* indicate $\lambda_{conv}^{cls}$, $\lambda_{fc}^{cls}$, $\lambda_{conv}^{reg}$, and $\lambda_{fc}^{reg}$, respectively.**

## 4.3 Feature Fusion Analysis

In this subsection, we discuss the performance of Adaptive Fusion Mechanism. We introduce four learnable weights, *i.e.*, $\lambda_{conv}^{cls}$, $\lambda_{fc}^{cls}$, $\lambda_{conv}^{reg}$, and $\lambda_{fc}^{reg}$ in Eq. (3), to control the contributions of encoded *conv* and *fc* feature components from two paths to the fused output features in Dual Query Encoder and Dual Attention Generator shown in Fig.3. Obviously, larger $\lambda_i^j$ with $i \in \{conv, fc\}$, $j \in \{cls, reg\}$ indicates that extractor $i$ dominates the subtask $j$ in feature extraction.

We visualize the variation of these four weights in base training phase ($K = 200$) on PASCAL VOC and MS COCO benchmarks, as shown in Fig. 6. As for PASCAL VOC, it is illustrated that the stable

value of $\lambda_{conv}^{cls}$ is larger than $\lambda_{fc}^{cls}$ while $\lambda_{conv}^{reg}$ is lower compared with $\lambda_{fc}^{reg}$, representing that *conv* extractor dominates the classification subtask but is less important than *fc* extractor in bounding box regression in this experimental setup and implementation. Besides, we note that $\lambda_i^{cls}$ with $i \in \{conv, fc\}$ in classification subtask is larger than $\lambda_i^{reg}$ with $i \in \{conv, fc\}$ in regression subtask respectively, indicating that performing classification generally needs more information than location estimation. When it comes to MS COCO, *conv* extractor consistently leads the classification subtask while these two kinds of extractors almost have equal contributions to bounding box regression. Besides, after stabilizing, all these four weights are larger than those in the experiments on PASCAL VOC respectively, showing that training on complex MS COCO dataset generally requires richer information. Analysis above validates the efficacy of Adaptive Fusion Mechanism when facing distinct subtasks and datasets.

## 4.4 Ablation Study

In this subsection, we conduct relative ablations in 10-shot scenario on the first base/novel split of PASCAL VOC. Results in all experiments are obtained by averaging over 10 random runs, as 10 is sufficient to provide statistically stable results in this scenario.

**Effect of Feature Fusion.** To examine the effectiveness of Adaptive Fusion Mechanism in the design of Dual Query Encoder and Dual Attention Generator, we compare the performance of various feature fusion combinations and report $AP_{50}$ of novel classes in 10-shot setup, as shown in Tabel 4, where *j-i*, $j \in \{cls, reg\}$, $i \in \{conv, fc\}$ in the header represents applying extractor $i$ in subtask *j*. The existence of two extractors within a branch indicates that the information fusion is applied (*e.g.*, the regression branch in
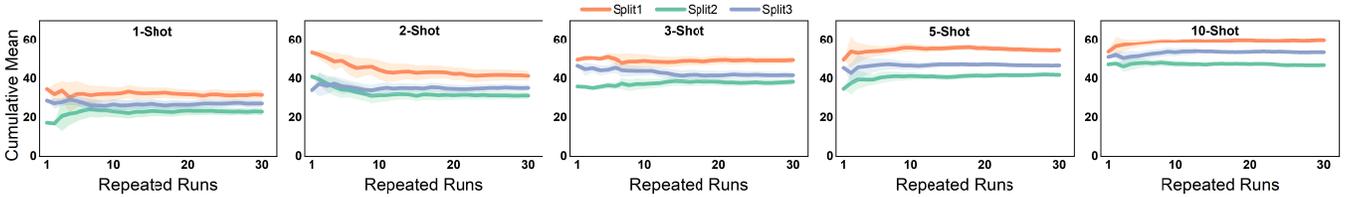
**Figure 7: Ablation study on effect of multiple runs on PASCAL VOC.** We plot cumulative means of $AP_{50}$ with 95% confidence intervals across 30 repeated runs on novel classes of three novel sets. The means and variances have been stable before 30 runs.

**Table 4: Ablation study on effect of feature fusion.** Comparison results among various fusion combinations in terms of $AP_{50}$ for novel categories on Novel Set 1 of PASCAL VOC dataset are reported.

| Ablation | cls-conv | cls-fc | reg-conv | reg-fc | Novel $AP_{50}$ |
|---|---|---|---|---|---|
| 1 | | ✓ | ✓ | | 49.9 |
| 2 | ✓ | | | ✓ | 56.2 |
| 3 | ✓ | | ✓ | | 55.5 |
| 4 | | ✓ | | ✓ | 55.1 |
| 5 | ✓ | ✓ | | ✓ | 59.1 |
| 6 | ✓ | ✓ | ✓ | | 57.6 |
| 7 | ✓ | | ✓ | ✓ | 59.0 |
| 8 | ✓ | ✓ | ✓ | ✓ | **60.7** |

**Table 5: Ablation study on effect of meta losses.** Comparison results among various meta loss combinations in terms of $AP_{50}$ for both base and novel categories on Novel Set 1 of PASCAL VOC dataset are reported.

| Ablation | meta-cls | meta-reg | Base $AP_{50}$ | Novel $AP_{50}$ |
|---|---|---|---|---|
| 1 | | | **70.8** | 57.3 |
| 2 | ✓ | | 70.3 | 58.9 |
| 3 | | ✓ | **70.8** | 58.4 |
| 4 | ✓ | ✓ | 70.6 | **60.7** |

Fig.3 can be represented by "*reg-fc & reg-conv*"). In *Ablation 1-4*, we adopt only one extractor in each branch without feature fusion. Among these four experiments, "*cls-conv & reg-fc*" achieves the best performance, while "*cls-fc & reg-conv*" performs the worst, indicating that *conv* extractor prefers classification and *fc* extractor is more suitable for bounding box regression in the first split on PASCAL VOC dataset. This is consistent with the observations in Section 4.3. In *Ablation 5-7*, feature fusion is adopted in only one branch and boost the detection performance, suggesting that leveraging information from unfocused tasks is essential in each subtask. Applying feature fusion in both branches, as shown in *Ablation 8*, obtains the best performance, further verifying the effectiveness of Adaptive Fusion Mechanism.

**Effect of Meta Losses.** In our approach, we further divide meta loss $L_{meta}$ introduced in [46] into task-specific meta losses, *i.e.*, $L_{meta-cls}$ for classification and $L_{meta-reg}$ for bounding box regression. Similarly, we evaluate the effect of meta losses in Table 5. Obviously, both the introduced $L_{meta-cls}$ and $L_{meta-reg}$ can indeed improve the performance on novel classes. Besides, they provide similar contributions individually. Interestingly, the existence of meta losses almost have no positive influence on the detection of base classes. We believe the reason is that, the few-shot detector has already been able to detect objects from base classes after base training phase with abundant annotations, despite the lack of meta losses.
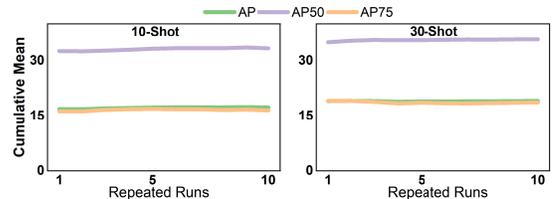


**Figure 8: Ablation study on effect of multiple runs on MS COCO.** We plot cumulative means with 95% confidence intervals across 10 repeated runs on novel classes. The means are stable and the variances are notably small, consistently.

**Effect of Multiple Runs.** In this ablation study, we focus on the effect of repeated runs in obtaining reported results on PASCAL VOC and MS COCO, as shown in Fig. 7 and Fig. 8, respectively. For PASCAL VOC, the confidence intervals are large across the first few runs (even up to 10 runs, in 3-shot setting), indicating the existence of large sample variance in low-shot scenarios. Therefore, overestimating the actual performance would occur when using only the first few random samples of training shots, as can be observed in the plot of 2-shot case in Fig. 7, leading to the unreliability of performance comparisons. In this paper, we adopt 30 repeated runs to obtain stable reported results with smaller variance. For MS COCO, we can observe from Fig. 8 that, the cumulative means are stable and the confidence intervals are quite small consistently, indicating the robustness of our model. In this paper, we adopt 10 repeated runs for fair comparisons with other methods.

## 5 CONCLUSION

This work targets the problem of few-shot object detection (FSOD). We carefully analyzed the characteristics of FSOD and presented that a general few-shot detector should: 1) explicitly decompose the process of category classification and localization; 2) utilize information from both subtasks for enhancing feature representations. Based on our observations, we proposed Adaptive Fully-Dual Network (AFD-Net) that performs two subtasks separately, involving feature representations, model reweighting, and state estimation. For the acquisition of enhanced features, we proposed a novel Adaptive Fusion Mechanism to guide the design of the feature extractor. Despite its simplicity, our approach achieved significant performance on multiple benchmarks, demonstrating its effectiveness and generalization ability. It is worth noting that, our proposed framework and feature fusion mechanism are general and simple, thus can be potentially applied into other two-stage models, and the performance could be further improved with carefully designed architectures.

# REFERENCES

[1] Marcin Andrychowicz, Misha Denil, Sergio Gomez, Matthew W Hoffman, David Pfau, Tom Schaul, Brendan Shillingford, and Nando De Freitas. 2016. Learning to learn by gradient descent by gradient descent. In *Advances in neural information processing systems*. 3981–3989.

[2] Luca Bertinetto, João F Henriques, Jack Valmadre, Philip Torr, and Andrea Vedaldi. 2016. Learning feed-forward one-shot learners. In *Advances in neural information processing systems*. 523–531.

[3] Alexey Bochkovskiy, Chien-Yao Wang, and Hong-Yuan Mark Liao. 2020. YOLOv4: Optimal Speed and Accuracy of Object Detection. *arXiv preprint arXiv:2004.10934* (2020).

[4] Zhaowei Cai and Nuno Vasconcelos. 2018. Cascade r-cnn: Delving into high quality object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 6154–6162.

[5] Hao Chen, Yali Wang, Guoyou Wang, and Yu Qiao. 2018. Lstd: A low-shot transfer detector for object detection. *arXiv preprint arXiv:1803.01529* (2018).

[6] Wei-Yu Chen, Yen-Cheng Liu, Zsolt Kira, Yu-Chiang Frank Wang, and Jia-Bin Huang. 2019. A closer look at few-shot classification. *arXiv preprint arXiv:1904.04232* (2019).

[7] Yinbo Chen, Xiaolong Wang, Zhuang Liu, Huijuan Xu, and Trevor Darrell. 2020. A new meta-baseline for few-shot learning. *arXiv preprint arXiv:2003.04390* (2020).

[8] Jifeng Dai, Yi Li, Kaiming He, and Jian Sun. 2016. R-fcn: Object detection via region-based fully convolutional networks. In *Advances in neural information processing systems*. 379–387.

[9] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. 2017. Deformable convolutional networks. In *Proceedings of the IEEE international conference on computer vision*. 764–773.

[10] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. 2010. The pascal visual object classes (voc) challenge. *International journal of computer vision* 88, 2 (2010), 303–338.

[11] Qi Fan, Wei Zhuo, Chi-Keung Tang, and Yu-Wing Tai. 2020. Few-shot object detection with attention-RPN and multi-relation detector. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4013–4022.

[12] Chelsea Finn, Pieter Abbeel, and Sergey Levine. 2017. Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks. In *ICML*.

[13] Spyros Gidaris and Nikos Komodakis. 2018. Dynamic few-shot visual learning without forgetting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 4367–4375.

[14] Ross Girshick. 2015. Fast r-cnn. In *Proceedings of the IEEE international conference on computer vision*. 1440–1448.

[15] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. 2014. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 580–587.

[16] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. 2017. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*. 2961–2969.

[17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2015. Spatial pyramid pooling in deep convolutional networks for visual recognition. *IEEE transactions on pattern analysis and machine intelligence* 37, 9 (2015), 1904–1916.

[18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.

[19] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. 2017. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 4700–4708.

[20] Bingyi Kang, Zhuang Liu, Xin Wang, Fisher Yu, Jiashi Feng, and Trevor Darrell. 2019. Few-shot object detection via feature reweighting. In *Proceedings of the IEEE International Conference on Computer Vision*. 8420–8429.

[21] Leonid Karlinsky, Joseph Shtok, Sivan Harary, Eli Schwartz, Amit Aides, Rogerio Feris, Raja Giryes, and Alex M Bronstein. 2019. Repmet: Representative-based metric learning for classification and few-shot object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 5197–5206.

[22] Gregory Koch, Richard Zemel, and Ruslan Salakhutdinov. 2015. Siamese neural networks for one-shot image recognition. In *ICML deep learning workshop*, Vol. 2. Lille.

[23] Brenden M Lake, Ruslan Salakhutdinov, and Joshua B Tenenbaum. 2015. Human-level concept learning through probabilistic program induction. *Science* 350, 6266 (2015), 1332–1338.

[24] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. 2017. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2117–2125.

[25] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *European conference on computer vision*. Springer, 740–755.

[26] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. 2016. Ssd: Single shot multibox detector. In *European conference on computer vision*. Springer, 21–37.

[27] Alex Nichol, Joshua Achiam, and John Schulman. 2018. On first-order meta-learning algorithms. *arXiv preprint arXiv:1803.02999* (2018).

[28] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. 2017. Automatic differentiation in pytorch. (2017).

[29] Siyuan Qiao, Chenxi Liu, Wei Shen, and Alan L Yuille. 2018. Few-shot image recognition by predicting parameters from activations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 7229–7238.

[30] Sachin Ravi and Hugo Larochelle. 2016. Optimization as a model for few-shot learning. (2016).

[31] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. 2016. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 779–788.

[32] Joseph Redmon and Ali Farhadi. 2017. YOLO9000: better, faster, stronger. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 7263–7271.

[33] Joseph Redmon and Ali Farhadi. 2018. Yolov3: An incremental improvement. *arXiv preprint arXiv:1804.02767* (2018).

[34] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. 2015. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*. 91–99.

[35] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. 2015. Imagenet large scale visual recognition challenge. *International journal of computer vision* 115, 3 (2015), 211–252.

[36] Larissa K Samuelson and Linda B Smith. 2005. They call it like they see it: Spontaneous naming and attention to shape. *Developmental Science* 8, 2 (2005), 182–198.

[37] Jake Snell, Kevin Swersky, and Richard Zemel. 2017. Prototypical networks for few-shot learning. In *Advances in neural information processing systems*. 4077–4087.

[38] Flood Sung, Yongxin Yang, Li Zhang, Tao Xiang, Philip HS Torr, and Timothy M Hospedales. 2018. Learning to compare: Relation network for few-shot learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 1199–1208.

[39] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. 2016. Matching networks for one shot learning. In *Advances in neural information processing systems*. 3630–3638.

[40] Xin Wang, Thomas E Huang, Trevor Darrell, Joseph E Gonzalez, and Fisher Yu. 2020. Frustratingly Simple Few-Shot Object Detection. *arXiv preprint arXiv:2003.06957* (2020).

[41] Xin Wang, Fisher Yu, Ruth Wang, Trevor Darrell, and Joseph E Gonzalez. 2019. Tafe-net: Task-aware feature embeddings for low shot learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 1831–1840.

[42] Yu-Xiong Wang, Deva Ramanan, and Martial Hebert. 2019. Meta-learning to detect rare objects. In *Proceedings of the IEEE International Conference on Computer Vision*. 9925–9934.

[43] Yue Wu, Yinpeng Chen, Lu Yuan, Zicheng Liu, Lijuan Wang, Hongzhi Li, and Yun Fu. 2020. Rethinking Classification and Localization for Object Detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10186–10195.

[44] Yang Xiao and Renaud Marlet. 2020. Few-Shot Object Detection and Viewpoint Estimation for Objects in the Wild. In *European Conference on Computer Vision (ECCV)*.

[45] Yinda Xu, Zeyu Wang, Zuoxin Li, Ye Yuan, and Gang Yu. 2020. SiamFC++: Towards Robust and Accurate Visual Tracking with Target Estimation Guidelines.. In *AAAI*. 12549–12556.

[46] Xiaopeng Yan, Ziliang Chen, Anni Xu, Xiaoxi Wang, Xiaodan Liang, and Liang Lin. 2019. Meta r-cnn: Towards general solver for instance-level low-shot learning. In *Proceedings of the IEEE International Conference on Computer Vision*. 9577–9586.

[47] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. 2016. Understanding deep learning requires rethinking generalization. *arXiv preprint arXiv:1611.03530* (2016).