

Rethinking Few-Shot Object Detection on a Multi-Domain Benchmark

Kibok Lee^{1,2*}, Hao Yang^{1†}, Satyaki Chakraborty¹, Zhaowei Cai¹,
Gurumurthy Swaminathan¹, Avinash Ravichandran¹, and Onkar Dabeer¹

¹AWS AI Labs ²Yonsei University

{kibok,haoyng,satyaki,zhaoweic,gurumurs,ravinash,onkardab}@amazon.com
kibok@yonsei.ac.kr

Abstract. Most existing works on few-shot object detection (FSOD) focus on a setting where both pre-training and few-shot learning datasets are from a similar domain. However, few-shot algorithms are important in multiple domains; hence evaluation needs to reflect the broad applications. We propose a Multi-domain Few-Shot Object Detection (MoFSOD) benchmark consisting of 10 datasets from a wide range of domains to evaluate FSOD algorithms. We comprehensively analyze the impacts of freezing layers, different architectures, and different pre-training datasets on FSOD performance. Our empirical results show several key factors that have not been explored in previous works: 1) contrary to previous belief, on a multi-domain benchmark, fine-tuning (FT) is a strong baseline for FSOD, performing on par or better than the state-of-the-art (SOTA) algorithms; 2) utilizing FT as the baseline allows us to explore multiple architectures, and we found them to have a significant impact on down-stream few-shot tasks, even with similar pre-training performances; 3) by decoupling pre-training and few-shot learning, MoFSOD allows us to explore the impact of different pre-training datasets, and the right choice can boost the performance of the down-stream tasks significantly. Based on these findings, we list possible avenues of investigation for improving FSOD performance and propose two simple modifications to existing algorithms that lead to SOTA performance on the MoFSOD benchmark. The code is available [here](#).

Keywords: Few-Shot Learning, Object Detection

1 Introduction

Convolutional neural networks have led to significant progress in object detection by learning with a large number of training images with annotations [30,28,3,4]. However, humans can easily localize and recognize new objects with only a

*Work done at AWS. †Corresponding author.

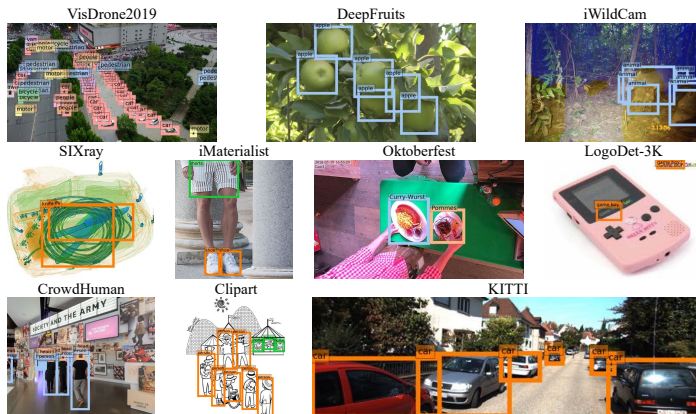


Fig. 1: Sample images in the proposed FSOD benchmark.

few examples. Few-shot object detection (FSOD) is a task to address this setting [19,41,8,35,26]. FSOD is desirable for many real-world applications in diverse domains due to lack of training data, difficulties in annotating them, or both, *e.g.*, identifying new logos, detecting anomalies in the manufacturing process or rare animals in the wild, *etc.* These diverse tasks naturally have vast differences in class distribution and style of images. Moreover, large-scale pre-training datasets in the same domain are not available for many of these tasks. In such cases, we can only rely on existing natural image datasets, such as COCO [24] and OpenImages [21] for pre-training.

Despite the diverse nature of FSOD tasks, FSOD benchmarks used in prior works are limited to a homogeneous setting [19,41,8,35,26], such that the pre-training and few-shot test sets in these benchmarks are from the same domain, or even the same dataset, *e.g.*, VOC [7] 15 + 5 and COCO [24] 60 + 20 splits. The class distributions of such few-shot test sets are also fixed to be balanced. While they provide an artificially balanced environment for evaluating different algorithms, it might lead to skewed conclusions for applying them in more realistic scenarios. Note that few-shot classification suffered from the same problem in the past few years [39,29]; Meta-dataset [37] addressed the problem with 10 different domains and a sophisticated scheme to sample imbalanced few-shot episodes.

Inspired by Meta-dataset [37], we propose a Multi-domain FSOD (MoFSOD) benchmark consisting of 10 datasets from 10 different domains, as shown in Figure 1. The diversity of MoFSOD datasets can be seen in Figure 2 where the domain distance of each dataset to COCO [24] is depicted. Our benchmark enables us to estimate the performance of FSOD algorithms across domains and settings and helps in better understanding of various factors, such as pre-training, architectures, *etc.*, that influence the algorithm performance. In addition, we propose a simple natural K -shot sampling algorithm that encourages more diversity in class distributions than balanced sampling.

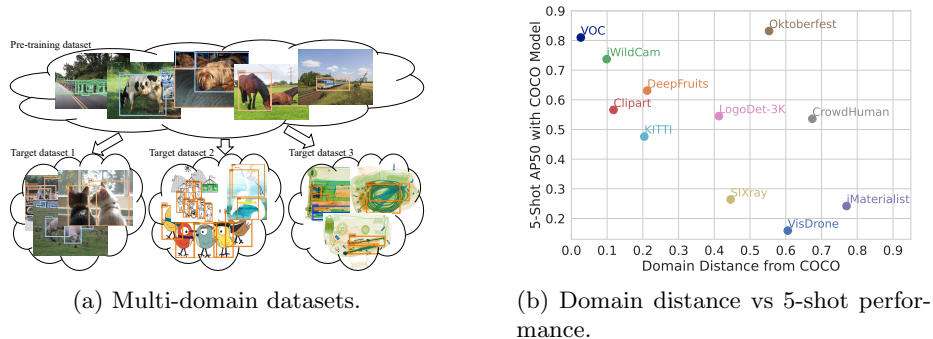


Fig. 2: (a) Real-world applications of FSOD are not limited to the natural image domain; we propose to pre-train models on large-scale natural image datasets and transfer to target domains for few-shot learning. (b) We measure the domain distance between datasets (see Section 3.2 for details) in the benchmark and COCO, and plot against 5-shot AP50 of these datasets fine-tuned from a model pre-trained on COCO. VOC is added for reference. The figure shows the benchmark covers a wide range of domains.

Building on our benchmark, we extensively study the impact of freezing parameters, detection architectures, pre-training datasets, and the effectiveness of several state-of-the-art (SOTA) FSOD algorithms [41,35,26]. Our empirical study leads to rethinking conventions in the field and interesting findings that can guide future research on FSOD.

Conventionally, in FSOD or general few-shot learning, it is believed that freezing parameters to avoid overfitting is helpful or even crucial for good performance [41,35,34,43]. If we choose to tune more parameters, specific components or designs must be added, such as weight imprinting [44] or decoupled gradient [26], to prevent overfitting. Our experiments in the MoFSOD show that these design choices might be helpful when pre-training and few-shot learning are in similar domains, as in previous benchmarks. However, if we consider a broader spectrum of domains, unfreezing the whole network results in better overall performance, as the network has more freedom to adapt. We further demonstrate a correlation between the performance gain of tuning more parameters and domain distance (see Figure 3a). Overall, *fine-tuning (FT) is a strong baseline* for FSOD on MoFSOD without any bells and whistles.

Using FT as a baseline allows us to explore the impact of different architectures on FSOD tasks. Previous FSOD methods [41,8,35,26] need to make architecture-specific design choices; hence focus on a single architecture – mostly Faster R-CNN [30], while we conduct extensive study on the impact of different architectures, *e.g.*, recent development of anchor-free [51,53] and transformer-based architectures [4,56], on few-shot performance. Surprisingly, we find that even with similar performance on COCO, different architectures have very different downstream few-shot performances. This finding suggests the potential benefits of specifically designed few-shot architectures for improved performance.

Moreover, unlike previous benchmarks, which split the pre-training and few-shot test sets from the same datasets (VOC or COCO), MoFSOD allows us to freely choose different pre-training datasets and explore the potential benefits of large-scale pre-training. To this end, we systematically study the effect of pre-training datasets with ImageNet [5], COCO [24], LVIS [13], FSOD-Dataset [8], Unified [52], and the integration of large-scale language-vision models. Similar to observations in recent works in image classification [20] and NLP [2,27], we find that large-scale pre-training can play a crucial role for downstream few-shot tasks.

Finally, motivated by the effectiveness of the unfreezing parameters and language-vision pre-training, we propose two extensions: FSCE+ and LVIS+. FSCE+ extends FSCE [35] to fine-tune more parameters with a simplified fully-connected (FC) detection head. LVIS+ follows the idea of using CLIP embedding of class names as the classifier, but instead of using it in zero-shot/open-vocabulary setting as in [50,11], we extend it to few-shot fine-tuning. Both methods achieve SOTA results with/without extra pre-training data.

We summarize our contributions as follows:

- We propose a Multi-domain Few-Shot Object Detection (MoFSOD) benchmark to simulate real-world applications in diverse domains.
- We conduct extensive studies on the effect of architectures, pre-training datasets, and hyperparameters with fine-tuning and several SOTA methods on the proposed benchmark. We summarize the observations below:
 - **Unfreezing more layers** do not lead to detrimental overfitting and improve the FSOD performance across different domains.
 - **Object detection model architectures** have a significant impact on the FSOD performance even when the architectures have a similar performance on the pre-training dataset.
 - **Pre-training datasets** play an important role in the downstream FSOD performance. Effective utilization of the pre-training dataset can significantly boost performance.
- Based on these findings, we propose two extensions that outperform SOTA methods by a significant margin on our benchmark.

2 Related Work

Meta-learning-based methods for FSOD are inspired by few-shot classification. Kang et al. [19] proposed a meta feature extractor with feature reweight module, which maps support images to mean features and reweight query features with the mean features, inspired by prototypical networks [34]. Meta R-CNN [46] extended the idea with an extra predictor-head remodeling network to extract class-attentive vectors. MetaDet [43] proposed meta-knowledge transfer with weight prediction module.

Two-stream methods take one query image and support images as inputs, and use the correlations between query and support features as the final features to the detection head and the Region Proposal Network (RPN). Several works in

this direction [8,48,14] have shown competitive results. These methods require all classes to have at least one support image to be fed to the model, which makes the overall process slow.

Fine-tuning-based methods update only the linear classification and regression layers [41], the whole detection head and RPN with an additional contrastive loss [35], or decoupling the gradient of RPN and the detection head while updating the whole network [26]. These methods are simple yet have shown competitive results. We focus on benchmarking them due to their simplicity, efficiency, and higher performance than other types.

Multi-domain few-shot classification benchmarks. In few-shot classification, miniImageNet [39] and tieredImageNet [29] have been used as standard benchmarks. Similar to benchmarks in FSOD, they are divided into two splits and used for pre-training and few-shot learning, respectively, such that they are in the same natural image domain. Recent works have proposed new benchmarks to address this issue: Tseng *et al.* [38] proposed a cross-domain few-shot classification benchmark with five datasets from different domains. Triantafyllou *et al.* [37] proposed Meta-Dataset, which is a large-scale few-shot classification benchmark with ten datasets and a sophisticatedly designed sampling algorithm to sample realistically imbalanced few-shot training datasets. Although not specifically catering to few-shot applications, Wang *et al.* [42] proposed universal object detection, which aims to cover multi-domain datasets with a single model for high-shot object detection.

3 MoFSOD: A Multi-Domain FSOD Benchmark

In this section, we first describe existing FSOD benchmarks and their limitations. Then, we propose a Multi-domain FSOD (MoFSOD) benchmark.

3.1 Existing Benchmarks and Limitations

Recent FSOD works have evaluated their methods in PASCAL VOC 15 + 5 and MS COCO 60 + 20 benchmarks proposed by Kang *et al.* [19]. From the original VOC [7] with 20 classes and COCO [24] with 80 classes, they took 25% of classes as novel classes for few-shot learning, and the rest of them as base classes for pre-training. For VOC 15 + 5, three splits were made, where each of them consists of 15 base classes for pre-training, and the other 5 novel classes for few-shot learning. For each novel class, $K = \{1, 3, 5, 10\}$ object instances are sampled, which are referred to as shot numbers. For COCO, 20 classes overlapped with VOC are considered novel classes, and $K = \{10, 30\}$ -shot settings are used. Different from classification, as an image usually contains multiple annotations in object detection, sampling exactly K annotations per class is difficult. Kang *et al.* [19] proposed pre-defined support sets for few-shot training, which would cause overfitting [17]. Wang *et al.* [41] proposed to sample few-shot training datasets with different random seeds to mitigate this issue, but the resulting

sampled datasets often contain more than K instances. While these benchmarks contributed to the research progress in FSOD, they have several limitations.

First, these benchmarks do not capture the breadth of few-shot tasks and domains as they sample few-shot task instances from a single dataset, as we discussed in Section 1. Second, these benchmarks contain only a fixed number of classes, 5 or 20. However, real-world applications might have a varying number of classes, ranging from one class, *e.g.*, face/pedestrian detection [47,49], to thousands of classes, *e.g.*, logo detection [40]. Last but not least, these benchmarks are constructed with the balanced K -shot sampling. For example, in the 5-shot setting, a set of images containing exactly 5 objects [19] is pre-defined. Such a setting is unlikely in real-world few-shot tasks. We also demonstrate that such a sampling strategy can lead to high variances in the performance of multiple episodes (see Table 2b). Moreover, different from classification, object detection datasets tend to be imbalanced due to the multi-label nature of the datasets. For example, COCO [24] and OpenImages [21] have a more dominant number of person instances than any other objects. The benchmark datasets should also explore these imbalanced scenarios.

3.2 Multi-Domain Benchmark Datasets

While FSOD applications span a wide range of domains, gathering enough pre-training data from these domains might be difficult. Hence, it becomes important to test the few-shot algorithm performance in settings where the pre-training and few-shot domains are different. Similar to Meta-dataset [37] in few-shot classification, we propose to extend the benchmark with datasets from a wide range of domains rather than a subset of natural image datasets. Our proposed benchmark consists of 10 datasets from 10 domains: VisDrone [54] in aerial images, DeepFruits [31] in agriculture, iWildCam [1] in animals in the wild, Clipart in cartoon, iMaterialist [12] in fashion, Oktoberfest [57] in food, LogoDet-3K [40] in logo, CrowdHuman [33] in person, SIXray [25] in security, and KITTI [10] in traffic/autonomous driving. We provide statistics of these datasets in Table 1a. The number of classes varies from 1 to 352, and that of boxes per image varies from 1.2 to 54.4, covering a wide range of scenarios.

In Figure 2b, we illustrate the diversity of domains in our benchmark by computing the domain distances between these datasets and COCO [24] and plotting against the 5-shot performance of fine-tuning (FT) on each dataset from a model pre-trained on COCO. Specifically, we measure the domain similarity by calculating the recall of a pre-trained COCO model on each dataset in a class-agnostic fashion, similar to the measurement of unsupervised object proposals [16]. Intuitively, if a dataset is in a domain similar to COCO, then objects in the dataset are likely to be localized well by the model pre-trained on COCO. As a reference, VOC has a recall of 97%. For presentation purpose, we define $(1 - \text{recall})$ as the domain distance. We can see diverse domain distances in the benchmark, ranging from 0.1 to 0.8. Interestingly, the domain distance also correlates with the FSOD performance. Although this is not the only deciding factor, as the intrinsic properties (such as the similarity between training

Table 1: Statistics of pre-training and MoFSOD datasets (Table 1a) and performance of different architectures on the pre-training datasets (Table 1b).

| Domain | Dataset | # classes | # training images |
|---------|---------|-----------|-------------------|
| Natural | COCO | 80 | 117k |
| | FSODD | 800 | 56k |
| Image | LVIS | 1203 | 100k |
| | Unified | 723 | 2M |

| Domain | Dataset | # classes | # bboxes per image |
|-------------|--------------|-----------|--------------------|
| Aerial | VisDrone | 10 | 54.4 |
| Agriculture | DeepFruits | 7 | 5.6 |
| Animal | iWildCam | 1 | 1.5 |
| Cartoon | Clipart | 20 | 3.3 |
| Fashion | iMaterialist | 46 | 7.3 |
| Food | Oktoberfest | 15 | 2.4 |
| Logo | LogoDet-3K | 352 | 1.2 |
| Person | CrowdHuman | 2 | 47.1 |
| Security | SIXray | 5 | 2.1 |
| Traffic | KITTI | 4 | 7.0 |

| Dataset | Architecture | AP |
|---------|-------------------|------|
| COCO | Faster R-CNN | 42.7 |
| | Cascade R-CNN | 45.1 |
| | CenterNet2 | 45.3 |
| | RetinaNet | 39.3 |
| | Deformable DETR | 46.3 |
| | Cascade R-CNN-P67 | 45.9 |
| LVIS | Faster R-CNN | 24.2 |
| | CenterNet2 | 28.3 |
| | Cascade R-CNN-P67 | 26.2 |

(a) **Top:** statistics of pre-training datasets.
Bottom: statistics of MoFSOD datasets.

(b) Performance of benchmark architectures pre-trained on COCO and LVIS. FSODD and Unified do not have a pre-defined validation/test set, so we do not measure their pre-training performances.

and test datasets) of a dataset also play an important role, we can still see the linear correlation between the domain distance and 5-shot performance with the Pearson correlation coefficient -0.43. Oktoberfest and CrowdHuman are outliers in our analysis possibly as they are relatively easy.

Natural K -Shot Sampling. We use a natural K -shot sampling algorithm to maintain the original class distribution for this benchmark. Specifically, we sample $C \times K$ images from the original dataset without worrying about class labels, where C is the number of classes of the original dataset. Then, we check missing images to ensure we have at least one image for each class of all classes. We provide the details in the supplementary material. The comparison between the balanced K -shot and natural K -shot sampling shows that our conclusions do not change based on the sampling algorithm, but the performance of the natural K -shot sampling is more consistent on different episodes (see Table 2b) and covers imbalanced class distributions existing in the real-world applications.

Evaluation Protocol. To evaluate the scalability of methods, we experiment with four different average shot numbers, $K = \{1, 3, 5, 10\}$. We first sample a few-shot training dataset from the original training dataset with the natural K -shot sampling algorithm for each episode. Then, we initialize the object detection model with pre-trained model parameters and train the model with FSOD methods. For evaluation, we randomly sample 1k images from the original test set if the test set is larger than 1k. We repeat this episode 10 times with different random seeds for all multi-domain datasets and report the average of the mean and standard deviation of the performance.

Metrics. As evaluation metrics, we use AP50 and the average rank among compared methods [37]. AP50 stands for the average precision of predictions where the IoU threshold is 0.5, and the rank is an integer ranging from 1 to the

number of compared methods, where the method with the highest AP50 gets rank 1. We first take the best AP50 among different hyperparameters at the end of training for each episode, compute the mean and standard deviation of AP50 and the rank over 10 episodes, and then average them over different datasets and/or shots, depending on the experiments.

4 Experiments

In this section, we conduct extensive experiments on MoFSOD and discuss the results. For better presentations, we highlight compared methods and architectures in *italics* and pre-training/few-shot datasets in **bold**.

4.1 Experimental Setup

Model architecture. We conduct experiments on six different architectures. For simplicity, we use ResNet-50 [15] as the backbone of all architectures. We also employ deformable convolution v2 [55] in the last three stages of the backbone. Specifically, we benchmark two-stage object detection architectures: 1) *Faster R-CNN* [30] 2) *Cascade R-CNN* [3], and the newly proposed 3) *CenterNet2* [51], and one-stage architectures: 4) *RetinaNet* [23], as well as transformer-based 5) *Deformable-DETR* [56]. Note that all architectures utilize Feature Pyramid Networks (FPN) [22] or similar multi-scale feature aggregation techniques. In addition, we also experiment the combination of the FPN-P67 design from *RetinaNet* and *Cascade R-CNN*, dubbed 6) *Cascade R-CNN-P67* [52]. We conduct our architecture analysis pre-trained on **COCO** [24] and **LVIS** [13]. Table 1b summarizes the architectures and their pre-training performance.

Freezing parameters. Based on the design of detectors, we can think of three different levels of fine-tuning the network: 1) only the last classification and regression fully-connected (FC) layer [41], 2) the detection head consisting of several FC and/or convolutional layers [8], and 3) the whole network, *i.e.*, standard fine-tuning.¹ We study the effects of these three ways of tuning on different domains in MoFSOD with *Faster R-CNN* [30].

Pre-training datasets. To explore the effect of pre-training dataset, we conduct experiments on five pre-training datasets: **ImageNet**² [5], **COCO** [24], **FSOD**³ [8], **LVIS** [13], and **Unified**, which is a union of OpenImages v5 [21], Object365 v1 [32], Mapillary [6], and **COCO**, combined as in [52]. To reduce the combinations of different architectures and pre-training datasets, we conduct most of the studies on the best performing architecture *Cascade R-CNN-P67*.

¹ When training an object detection model, the batch normalization layers [18] and the first two macro blocks of the backbone (**stem** and **res2**) are usually frozen, even for large-scale datasets. We follow this convention in our paper.

² This is ImageNet-1K for classification, which is commonly used for pre-training standard object detection methods, *i.e.*, we omit pre-training on an object detection task.

³ The name of the dataset is also FSOD, so we introduce an additional D to distinguish the dataset from the task.

Also, inspired by [50,11], we experiment with the effect of CLIP [27] embeddings to initialize the final classification layer of detector when pre-training on **LVIS** dubbed **LVIS+**. In addition, **LVIS++** uses the backbone pre-trained on the ImageNet-21K classification task instead of ImageNet-1K before pre-training on **LVIS**. All experiments are done with Detectron2 [45].

Hyperparameters. For pre-training, we mostly follow standard hyperparameters of the corresponding method, with the addition of deformable convolution v2 [55]. On **COCO** and **LVIS**, for *Faster R-CNN*, *Cascade R-CNN*, and *Cascade R-CNN-P67*, we use the $3\times$ scheduler with 270k iterations, the batch size of 16, the SGD optimizer with initial learning rate of 0.02 decaying by the factor of 0.1 at 210k and 250k. For *RetinaNet*, the initial learning rate is 0.01 [23]. For *CenterNet2*, following [51], we use the *CenterNet* [53] as the first stage and the Cascade R-CNN head as the second stage, where the other hyperparameters are the same as above. For *Deformable-DETR* [56], we follow the two-stage training of 50 epochs, the AdamW optimizer, and the initial learning rate of 0.0002 decaying by the factor of 0.1 at 40 epochs. On **FSODD**, we train for 60 epochs with the learning rate of 0.02 decaying by the factor of 0.1 at 40 and 54 epochs and the batch size of 32. On **Unified**, following [52], the label space of four datasets are unified, the dataset-aware sampling and equalization loss [36] are applied to handle long-tailed distributions, and training is done for 600k iterations with the learning rate of 0.02 decaying by the factor of 0.1 at 400k and 540k iterations, and the batch size of 32.

For few-shot training, we train models for 2k iterations with the batch size of 4 on a single V100 GPU.⁴ For *Faster R-CNN*, *Cascade R-CNN*, *Cascade R-CNN-P67*, and *CenterNet2*, we train models with the SGD optimizer and different initial learning rates in $\{0.0025, 0.005, 0.01\}$ and choose the best, where the learning rate is decayed by the factor of 0.1 after 80% of training. For *RetinaNet*, we halve the learning rates to $\{0.001, 0.0025, 0.005\}$, as we often observe training diverges with the learning rate of 0.01. For *Deformable-DETR* and *CenterNet2 with CLIP*, we use the AdamW optimizer and initial learning rates in $\{0.0001, 0.0002, 0.0004\}$.⁵

Compared methods. *TFA* [41] or *Two-stage Fine-tuning Approach* has shown to be a simple yet effective method for FSOD. TFA fine-tunes the box regressor and classifier on the few-shot dataset while freezing the other parameters of the model. For this method, we use the FC head, such that TFA is essentially the same as tuning the final FC layer.⁶

FSCE [35] or *Few-Shot object detection via Contrastive proposals Encoding* improves TFA by 1) additionally unfreezing detection head in the setting of TFA, 2) doubling the number of proposals kept after NMS and halving the number of sampled proposals in the RoI head, and 3) optimizing the contrastive proposal

⁴ The batch size could be less than 4 if the sampled dataset size is less than 4, e.g., when the number of classes and shot number K is 1, and the batch size has to be 1.

⁵ The learning rates are chosen from our initial experiments on three datasets. Note that training *CenterNet2* with AdamW results in worse performance than SGD.

⁶ Replacing the FC head with the cosine-similarity results in a similar performance.

Table 2: The K -shot performance on MoFSOD with $K = \{1, 3, 5, 10\}$.

| Method | Pre-training | 1-shot | 3-shot | 5-shot | 10-shot |
|----------------|--------------|-------------------|-------------------|-------------------|-------------------|
| TFA [41] | | 23.4 ± 4.6 | 29.2 ± 1.8 | 32.0 ± 1.3 | 35.2 ± 1.2 |
| FSCE-base [35] | | 30.5 ± 5.3 | 39.4 ± 1.6 | 43.9 ± 1.1 | 50.2 ± 1.1 |
| FSCE-con [35] | | 29.4 ± 2.5 | 38.8 ± 1.5 | 43.6 ± 1.2 | 50.4 ± 1.0 |
| DeFRCN [26] | COCO | 29.3 ± 4.2 | 37.8 ± 3.1 | 41.6 ± 1.9 | 48.2 ± 1.9 |
| Ours-FT | | 31.5 ± 2.0 | 41.1 ± 1.8 | 46.1 ± 1.4 | 52.6 ± 1.8 |
| Ours-FSCE+ | | 31.2 ± 2.4 | 41.3 ± 1.5 | 46.4 ± 1.1 | 53.2 ± 1.5 |
| Ours-FT+ | COCO | 35.4 ± 1.8 | 44.7 ± 1.6 | 49.9 ± 1.2 | 56.4 ± 1.1 |
| Ours-FT++ | LVIS++ | 35.8 ± 3.4 | 47.2 ± 2.1 | 52.6 ± 1.4 | 59.4 ± 1.1 |

(a) Performance with the natural K -shot.

| Architecture | Pre-training | 1-shot | 3-shot | 5-shot | 10-shot |
|-------------------|--------------|-------------------|-------------------|-------------------|-------------------|
| Faster R-CNN | | 31.5 ± 2.0 | 41.1 ± 1.8 | 46.1 ± 1.4 | 52.6 ± 1.8 |
| Cascade R-CNN | | 31.5 ± 2.1 | 41.2 ± 1.5 | 45.6 ± 1.4 | 52.7 ± 1.2 |
| CenterNet2 | COCO | 29.1 ± 5.2 | 40.2 ± 1.9 | 45.3 ± 1.8 | 52.5 ± 1.9 |
| RetinaNet | | 25.4 ± 4.0 | 34.8 ± 2.2 | 40.6 ± 1.7 | 48.8 ± 1.2 |
| Deformable-DETR | | 32.0 ± 2.9 | 42.3 ± 1.6 | 47.4 ± 1.4 | 54.7 ± 1.0 |
| Cascade R-CNN-P67 | | 35.4 ± 1.8 | 44.7 ± 1.6 | 49.9 ± 1.2 | 56.4 ± 1.1 |
| Faster R-CNN | | 31.7 ± 1.7 | 41.6 ± 1.5 | 46.4 ± 1.2 | 53.6 ± 1.0 |
| CenterNet2 | LVIS | 28.1 ± 3.3 | 39.0 ± 1.7 | 44.3 ± 1.2 | 51.6 ± 1.0 |
| Cascade R-CNN-P67 | | 34.4 ± 1.2 | 44.0 ± 1.6 | 48.7 ± 1.4 | 55.6 ± 1.0 |

(c) The effect of different architectures.

| Method | Pre-training | 1-shot | 3-shot | 5-shot | 10-shot |
|---------------|--------------|-------------------|-------------------|-------------------|-------------------|
| TFA [41] | | 22.9 ± 5.4 | 27.6 ± 2.1 | 28.3 ± 5.9 | 31.6 ± 2.7 |
| FSCE-con [35] | | 26.8 ± 3.9 | 33.6 ± 4.9 | 36.3 ± 5.0 | 40.2 ± 5.4 |
| DeFRCN [26] | COCO | 27.5 ± 5.1 | 34.8 ± 2.2 | 37.4 ± 1.7 | 40.1 ± 6.5 |
| Ours-FT | | 28.8 ± 2.3 | 34.8 ± 3.1 | 37.4 ± 4.4 | 42.1 ± 5.4 |
| Ours-FSCE+ | | 28.7 ± 3.8 | 36.5 ± 5.3 | 38.2 ± 5.0 | 42.7 ± 5.7 |

(b) Performance with the balanced K -shot.

| Architecture | Pre-training | 1-shot | 3-shot | 5-shot | 10-shot |
|--------------|--------------|-------------------|-------------------|-------------------|-------------------|
| ImageNet | | 13.5 ± 1.6 | 23.2 ± 1.5 | 29.5 ± 1.2 | 37.7 ± 1.2 |
| COCO | | 35.4 ± 1.8 | 44.7 ± 1.6 | 49.9 ± 1.2 | 56.4 ± 1.1 |
| FSODD | | 26.7 ± 2.9 | 36.9 ± 1.5 | 42.3 ± 1.2 | 49.1 ± 1.0 |
| LVIS | | 34.4 ± 2.0 | 44.0 ± 1.6 | 48.7 ± 1.4 | 55.6 ± 1.0 |
| Unified | | 33.3 ± 2.2 | 44.3 ± 1.4 | 49.6 ± 1.2 | 56.6 ± 1.1 |
| LVIS+ | | 34.7 ± 4.2 | 46.6 ± 1.6 | 52.0 ± 1.0 | 59.0 ± 1.0 |
| COCO | | 29.1 ± 5.2 | 40.2 ± 1.9 | 45.3 ± 1.8 | 52.5 ± 1.9 |
| LVIS | | 28.1 ± 3.3 | 39.0 ± 1.7 | 44.3 ± 1.2 | 51.6 ± 1.0 |
| LVIS+ | | 34.9 ± 3.2 | 46.5 ± 1.9 | 51.8 ± 1.3 | 58.8 ± 1.0 |
| LVIS++ | | 35.8 ± 3.4 | 47.2 ± 2.1 | 52.6 ± 1.4 | 59.4 ± 1.1 |

(d) The effect of different pre-training.

encoding loss. While the original work did not apply the contrastive loss for extremely few-shot settings (less than 3), we explicitly compare two versions in all shots: without (*FSCE-base*, the same as tuning the detection head) and with the contrastive loss (*FSCE-con*).

DeFRCN [26] or *Decoupled Faster R-CNN* can be distinguished with other methods by 1) freezing only the R-CNN head, 2) decoupling gradients to suppress gradients from RPN while scaling those from the R-CNN head, and 3) calibrating the classification score from an offline prototypical calibration block (PCB), which is a CNN-based prototype classifier pre-trained on ImageNet [5]. We note that PCB does not re-scale input images in their original implementation, unlike the object detector, so we manually scaled images to avoid GPU memory overflow if they are too large.

FT or *Fine-tuning* does not freeze model parameters as done for other methods. Though it is undervalued in prior works, we found that this simple baseline outperforms state-of-the-art methods in our proposed benchmark. All experiments are done with this method unless otherwise specified.

4.2 Experimental Results and Discussions

Effect of tuning more parameters. We first analyze the effect of tuning more or fewer parameters on MoFSOD. In Table 2a and Table 3a, we examine three methods freezing different number of parameters when fine-tuning: *TFA* as tuning the last FC layers only, *FSCE-base* as tuning the detection head, and *FT* as tuning the whole network. We observe that freezing fewer parameters improves the average performance: tuning the whole network (*FT*) shows better performance than others, while tuning the last FC layers only (*TFA*) shows lower performance than others. Also, the performance gap becomes larger as the size of few-shot training datasets increases. For example, *FT* outperforms

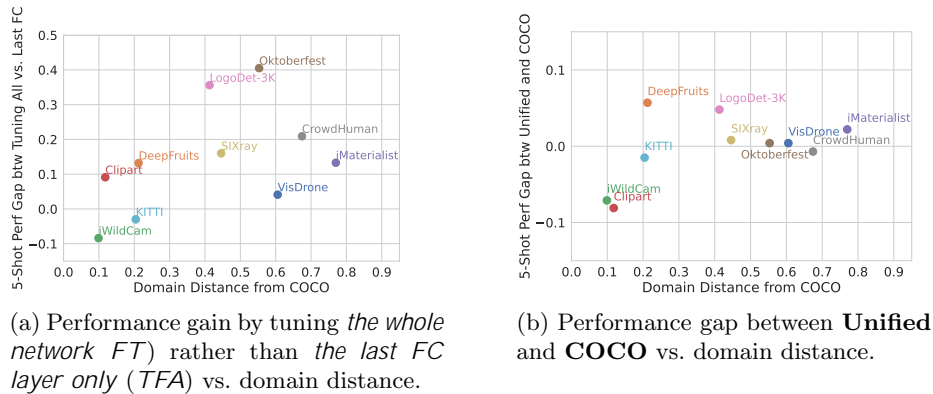


Fig. 3: We demonstrate the correlation between tuning more parameters and domain distance in Figure 3a and the correlation between pre-training datasets and domain distance in Figure 3b.

FSCE-base and *TFA* by 1.0% and 8.1% in 1-shot, and 2.4%, 17.4% in 10-shot, respectively. This contrasts with the conventional belief that freezing most of the parameters generally improves the performance of few-shot learning, as it prevents overfitting [34,9,41,35]. However, this is not necessarily true for FSOD. For example, in the standard two-stage object detector training, RPN is class-agnostic, such that its initialization for training downstream few-shot tasks can be the one pre-trained on large-scale datasets, preserving the pre-trained knowledge on objectness. Also, the detection head utilizes thousands of examples even in few-shot scenarios, because RPN could generate 1–2k proposals per image. Hence, the risk of overfitting is relatively low.

However, fine-tuning more parameters does not always improve performance. Figure 3a illustrates the performance gain by tuning more parameters with respect to the domain distance. There is a linear correlation, *i.e.*, the performance gain by fine-tuning more parameters increases when the domain distance increases. This implies that fine-tuning fewer parameters to preserve the pre-trained knowledge is better when the few-shot dataset is close to the pre-training dataset. Hence, for datasets close to **COCO**, such as **KITTI** and **iWildCam**, *tuning Last FC layers (TFA)* is the best performing method. From these observations, an interesting research direction might be exploring a sophisticated tuning of layers based on the few-shot problem definition and the domain gap between pre-training and few-shot tasks.

One way to design such sophisticated tuning is to develop a better measure for domain distance. In fact, the proposed measure with class-agnostic object recall has limitations. If we decouple the object detection task into localization (background vs. foreground) and classification (among foreground classes), then the proposed domain distance is biased towards measuring localization gaps. Therefore, it ignores the potential classification gaps that would also require tuning more layers. For example, although we can get a good coverage with the propos-

Table 3: Per-dataset 5-shot performance of the effects of tuning different parameters, different architectures and pre-training datasets.

| 5-shot | Aerial | Agriculture | Animal | Cartoon | Fashion | Food | Logo | Person | Security | Traffic | Mean | Rank |
|---------------------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|------------------|
| Unfrozen | VisDrone | DeepFruits | iWildCam | Clipart | iMaterialist | Oktoberfest | LogoDet-3K | CrowdHuman | SIXray | KITTI | | |
| Last FC layers (TFA [41]) | 10.1 ± 0.6 | 47.5 ± 1.8 | 71.7 ± 2.4 | 40.2 ± 2.6 | 8.0 ± 0.9 | 41.1 ± 4.8 | 14.2 ± 3.7 | 30.6 ± 0.8 | 7.9 ± 2.0 | 48.3 ± 3.4 | 32.0 ± 1.3 | 2.7 ± 0.2 |
| Detection head (FSCE-base [35]) | 13.0 ± 0.7 | 59.2 ± 3.2 | 70.6 ± 1.7 | 43.5 ± 2.5 | 20.5 ± 0.9 | 70.7 ± 3.5 | 47.2 ± 3.5 | 51.6 ± 2.0 | 15.9 ± 2.2 | 47.1 ± 4.3 | 43.9 ± 1.1 | 1.9 ± 0.2 |
| Whole network (Ours-FT) | 14.2 ± 0.8 | 60.7 ± 3.8 | 63.3 ± 3.7 | 49.3 ± 3.5 | 21.3 ± 0.8 | 81.6 ± 4.0 | 49.8 ± 3.5 | 51.5 ± 2.2 | 23.9 ± 3.9 | 45.3 ± 3.7 | 46.1 ± 1.4 | 1.5 ± 0.2 |

(a) Fine-tuning different number of parameters with Faster R-CNN pre-trained on COCO.

| 5-shot | Aerial | Agriculture | Animal | Cartoon | Fashion | Food | Logo | Person | Security | Traffic | Mean | Rank | |
|-------------------|--------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|------------------|
| Architecture | Pre-training | VisDrone | DeepFruits | iWildCam | Clipart | iMaterialist | Oktoberfest | LogoDet-3K | CrowdHuman | SIXray | KITTI | | |
| Faster R-CNN | | 14.2 ± 0.8 | 60.7 ± 3.8 | 63.3 ± 3.7 | 49.3 ± 3.5 | 21.3 ± 0.8 | 81.6 ± 4.0 | 49.8 ± 3.5 | 51.5 ± 2.2 | 23.9 ± 3.9 | 45.3 ± 3.7 | 46.1 ± 1.4 | 3.4 ± 0.3 |
| Cascade R-CNN | | 13.0 ± 0.8 | 58.9 ± 3.3 | 66.8 ± 3.7 | 50.8 ± 2.8 | 20.5 ± 0.6 | 80.5 ± 2.2 | 48.5 ± 4.7 | 51.4 ± 2.0 | 21.6 ± 4.4 | 44.4 ± 4.1 | 45.6 ± 1.4 | 4.1 ± 0.4 |
| CenterNet2 | | 13.5 ± 0.7 | 59.0 ± 4.7 | 61.6 ± 4.7 | 49.2 ± 7.6 | 22.2 ± 1.9 | 79.7 ± 3.5 | 51.0 ± 4.2 | 51.4 ± 3.7 | 22.2 ± 4.8 | 43.7 ± 4.9 | 45.3 ± 1.8 | 3.9 ± 0.3 |
| RetinaNet | COCO | 9.9 ± 0.6 | 55.8 ± 2.7 | 59.2 ± 6.8 | 26.8 ± 2.2 | 17.6 ± 0.4 | 79.5 ± 4.0 | 49.2 ± 3.5 | 47.6 ± 1.9 | 20.6 ± 3.3 | 39.7 ± 3.3 | 40.6 ± 1.7 | 5.4 ± 0.5 |
| Deformable-DETR | | 15.0 ± 0.6 | 68.8 ± 4.3 | 66.0 ± 4.3 | 43.7 ± 1.6 | 22.4 ± 1.1 | 77.2 ± 4.6 | 50.3 ± 3.5 | 56.7 ± 2.4 | 26.3 ± 3.8 | 47.0 ± 3.2 | 47.4 ± 1.4 | 2.6 ± 0.5 |
| Cascade R-CNN-P67 | | 15.9 ± 0.8 | 63.1 ± 2.3 | 73.7 ± 2.8 | 56.6 ± 2.3 | 24.2 ± 0.8 | 83.2 ± 2.4 | 54.5 ± 4.4 | 53.6 ± 1.8 | 26.4 ± 3.9 | 47.6 ± 3.6 | 49.9 ± 1.2 | 1.5 ± 0.3 |
| Faster R-CNN | | 14.2 ± 0.7 | 66.2 ± 3.4 | 69.3 ± 3.7 | 39.8 ± 1.8 | 28.0 ± 0.6 | 80.7 ± 2.7 | 51.3 ± 4.2 | 48.3 ± 2.2 | 24.7 ± 3.6 | 41.3 ± 3.5 | 46.4 ± 1.2 | 2.1 ± 0.3 |
| CenterNet2 | LVIS | 13.3 ± 0.7 | 64.6 ± 3.5 | 50.2 ± 25.3 | 34.1 ± 2.5 | 25.6 ± 0.5 | 81.6 ± 2.8 | 53.7 ± 3.8 | 46.6 ± 2.0 | 22.8 ± 3.3 | 38.9 ± 3.9 | 43.1 ± 6.9 | 2.7 ± 0.3 |
| Cascade R-CNN-P67 | | 15.2 ± 0.7 | 65.4 ± 2.0 | 71.7 ± 2.0 | 46.0 ± 2.5 | 30.2 ± 0.6 | 81.9 ± 3.0 | 56.6 ± 5.0 | 49.6 ± 1.7 | 25.9 ± 4.3 | 44.7 ± 3.7 | 48.7 ± 1.4 | 1.2 ± 0.3 |

(b) Performance of different architectures pre-trained on COCO and LVIS.

| 5-shot | Aerial | Agriculture | Animal | Cartoon | Fashion | Food | Logo | Person | Security | Traffic | Mean | Rank | |
|-------------------|--------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|------------------|
| Architecture | Pre-training | VisDrone | DeepFruits | iWildCam | Clipart | iMaterialist | Oktoberfest | LogoDet-3K | CrowdHuman | SIXray | KITTI | | |
| ImageNet | | 9.8 ± 0.5 | 53.0 ± 3.5 | 7.4 ± 3.2 | 13.1 ± 2.5 | 19.2 ± 0.6 | 77.4 ± 3.4 | 46.4 ± 4.4 | 32.0 ± 3.1 | 13.1 ± 3.2 | 23.8 ± 3.4 | 29.5 ± 1.2 | 6.0 ± 0.0 |
| COCO | | 15.9 ± 0.8 | 63.1 ± 2.3 | 73.7 ± 2.8 | 56.6 ± 2.3 | 24.2 ± 0.8 | 83.2 ± 2.4 | 54.5 ± 4.4 | 53.6 ± 1.8 | 26.4 ± 3.9 | 47.6 ± 3.6 | 49.9 ± 1.2 | 2.8 ± 0.4 |
| FSODD | | 10.6 ± 0.6 | 67.5 ± 3.7 | 64.5 ± 3.4 | 26.5 ± 2.0 | 21.2 ± 0.5 | 82.9 ± 2.9 | 55.7 ± 3.8 | 38.4 ± 2.2 | 23.2 ± 3.7 | 32.3 ± 3.6 | 42.3 ± 1.2 | 4.3 ± 0.4 |
| Cascade R-CNN-P67 | LVIS | 15.2 ± 0.7 | 65.4 ± 2.0 | 71.7 ± 2.0 | 46.0 ± 2.5 | 30.2 ± 0.6 | 81.9 ± 3.0 | 56.6 ± 5.0 | 49.6 ± 1.7 | 25.9 ± 4.3 | 44.7 ± 3.7 | 47.4 ± 1.4 | 3.4 ± 0.4 |
| Unified | | 16.3 ± 0.9 | 68.8 ± 3.0 | 66.6 ± 4.0 | 48.5 ± 2.9 | 26.4 ± 0.7 | 83.6 ± 2.9 | 59.3 ± 3.7 | 52.9 ± 1.6 | 27.2 ± 3.6 | 46.1 ± 3.9 | 49.6 ± 1.2 | 2.5 ± 0.4 |
| LVIS+ | | 18.5 ± 1.0 | 76.8 ± 3.3 | 59.5 ± 4.1 | 52.6 ± 2.0 | 31.5 ± 0.7 | 82.5 ± 3.3 | 61.2 ± 2.7 | 52.3 ± 1.9 | 36.0 ± 2.5 | 49.6 ± 3.9 | 52.0 ± 1.1 | 1.9 ± 0.5 |
| COCO | | 13.5 ± 0.7 | 59.0 ± 4.7 | 61.6 ± 4.7 | 49.2 ± 7.6 | 22.2 ± 1.9 | 79.7 ± 3.5 | 51.0 ± 4.2 | 51.4 ± 3.7 | 22.2 ± 4.8 | 43.7 ± 4.9 | 45.3 ± 1.8 | 3.2 ± 0.3 |
| LVIS | | 13.3 ± 0.7 | 64.6 ± 3.5 | 61.5 ± 4.3 | 34.1 ± 2.5 | 25.6 ± 0.5 | 81.6 ± 2.8 | 53.7 ± 3.8 | 46.6 ± 2.0 | 22.8 ± 3.3 | 38.9 ± 3.9 | 44.3 ± 1.2 | 3.2 ± 0.2 |
| LVIS+ | | 18.2 ± 0.9 | 74.0 ± 3.3 | 63.7 ± 3.8 | 50.5 ± 1.8 | 31.5 ± 0.8 | 80.4 ± 4.0 | 62.3 ± 3.2 | 54.0 ± 2.1 | 37.3 ± 4.1 | 46.5 ± 4.7 | 51.8 ± 1.3 | 2.0 ± 0.4 |
| LVIS++ | | 18.1 ± 0.9 | 77.5 ± 2.4 | 64.4 ± 4.7 | 52.1 ± 1.8 | 33.1 ± 0.7 | 80.4 ± 3.5 | 61.0 ± 4.5 | 54.7 ± 2.0 | 37.8 ± 3.2 | 46.9 ± 4.1 | 52.6 ± 1.4 | 1.5 ± 0.3 |

(c) Performance of Cascade R-CNN-P67 and CenterNet2 pre-trained on different datasets.

als of **COCO** pre-training for **Clipart** (classic domain adaption dataset) and **DeepFruits** (infrared images), resulting in relatively small domain distances, there exist significant gaps of feature discrimination for fine-grained classification. In this case, we need to tune more parameters for better performance.

Effect of model architectures. A benefit of using *FT* as the baseline is that we can systematically study the effects of model architectures without the constraints of specifically designed components. Many different architectures have been proposed to solve object detection problems; each has its own merits and drawbacks. One-stage methods, such as YOLO [28] and RetinaNet [23], are known for their fast inference speed. However, different from two-stage methods, they do not have the benefit of inheriting the pre-trained class-agnostic RPN. Specifically, the classification layers for discriminating background/foreground and foreground classes need to be reinitialized, as they are often tied together in one-stage methods. We validate this hypothesis by comparing with two-stage methods in Table 2c. Compared to *Faster R-CNN*, *RetinaNet* has 4–6% low performance on MoFSOD. Per-dataset performance in Table 3b shows that *RetinaNet* is worse in all cases.

The same principle of preserving as much information as possible from pre-training also applies for two-stage methods, *i.e.*, reducing the number of randomly reinitialized parameters is better. Specifically, we look into the performance of *Cascade R-CNN* vs *Faster R-CNN*. For *Cascade R-CNN*, we need to

reinitialize and learn three FC layers as there are three stages in the cascade detection head, while we only need to reinitialize the last FC layer for *Faster R-CNN*. However, *Cascade R-CNN* is known to have better performance, as demonstrated in Table 1b. In FSOD, these two factors cancelled out, such that their performance is on a par with each other as shown in Table 2c.

Based on these insights, we extend *Cascade R-CNN* by applying the FPN-P67 architecture [23], similar to [52,51]. Specifically, assuming ResNet-like architecture [15], we use the last three stages of the backbone, namely [res3, res4, res5], instead of the last four in standard FPN. Then, we add P6 and P7 of the FPN features from P5 with two different FC layers, such that RPN takes in features from [P3, P4, P5, P6, P7] to improve the class-agnostic coverage, which can be inherited for down-stream tasks. While the detection head still uses [P3, P4, P5] only. As shown in Table 2c, the resulting architecture, *Cascade R-CNN-P67* improves *Cascade R-CNN* by 3–4% on the downstream few-shot tasks.

Moreover, recent works proposed new directions of improvement, such as utilizing point-based predictions [53,51] or transformer-based set predictions [4,56]. These methods are unknown quantities in FSOD, as no previous FSOD work has studied them. In our experiments, while *CenterNet2* [51] outperforms *Faster R-CNN* on **COCO** by 2.6% as shown in Table 1b, its FSOD performance on MoFSOD is lower, *e.g.*, 2.4% in 1-shot. In the case of *Deformable-DETR* [56], it outperforms *Faster R-CNN* in both pre-training and few-shot learning, by 3.6% on **COCO** and 2.1% on MoFSOD in 10-shot. These results show that the upstream performance might not necessarily translate to the downstream FSOD performance. We note that we could not observe a significant correlation between the performance gap of different architectures and domain distances.

Effect of pre-training datasets. MoFSOD consists of datasets from a wide range of domains, allowing us to freely explore different pre-training datasets while ensuring domain shifts between pre-training and few-shot learning. We examine the impact of the pre-training datasets with the best performing *Cascade R-CNN-P67* architecture.

In Table 2d, we first observe that pre-training on **ImageNet** for classification results in low performance, as it does not provide a good initialization for downstream FSOD, especially for RPN. On the other hand, compared to **COCO**, **FSODD**, **LVIS**, and **Unified** have more classes and/or more annotations, while they have a fewer, similar, and more number of images, respectively. Pre-training on these larger object detection datasets does not improve the FSOD performance significantly, as shown in Table 2d. For example, pre-training on **Unified** improves the performance over **COCO** when the domain distance from **COCO** is large, such as **Deepfruits** and **LogoDet-3K** as shown in Figure 3b. However, pre-training on **Unified** results in lower performance for few-shot datasets close to **COCO**, such that the overall performance is similar. We hypothesize that this could be due to the non-optimal pre-training of **LVIS** and **Unified**, as these two datasets are highly imbalanced and difficult to train. It could also be the case that even **LVIS** and **Unified** do not have better coverage for these datasets.

On the other end of the spectrum, we can combine the idea of preserving knowledge and large-scale pre-training by utilizing a large-scale language-vision model. Following [11,50], we use CLIP to extract text features from each class name and build a classifier initialized with the text features. In this way, we initialize the classifier with the CLIP text embeddings for downstream few-shot tasks, such that it has strong built-in knowledge of text-image alignment, better than random initialization. As demonstrated in Table 2d, For **LVIS+**, we can see this improves performance significantly by 7.5% for *CenterNet2*, and 3.3% for *Cascade R-CNN-P67* in 5-shot. **LVIS++** pre-trains the backbone on ImageNet-21K instead of ImageNet-1K (before pre-training on **LVIS**), and it further improves over **LVIS+** by 0.8% in 5-shot. However, the benefit of CLIP initialization is valid only when class names are matched with texts presented in CLIP; an exceptional case is **Oktoberfest**, which has German class names, such that **LVIS+** does not help.

Comparison with SOTA methods. Table 2a and Table 2b compare SOTA methods and our proposed methods. For balanced K -shot sampling, we follow Wang *et al.* [41] to sample K instances for each class whenever possible greedily. Here *FT* and *FSCE+* employ a similar backbone/architecture and pre-trained data as all SOTA methods for a fair comparison. We confirm that *FT* is indeed a strong baseline, such that it performs better than other SOTA methods in both natural K -shot and balanced K -shot settings. In addition to *FT*, based on the insights above, we propose several extensions: 1) *FSCE+* is an extension of *FSCE* by tuning the whole network parameters, similar to *FT*. We keep the contrastive proposal encoding loss, but we simplify the classification head from the cosine similarity head to the FC head. We can see the improvement by 2–3% compared to *FSCE* for both natural K -shot and balanced K -shot scenarios and performs slightly better than *FT*. 2) *FT+* replaces *Faster R-CNN* with *Cascade R-CNN-P67*, and it improves over *FT* by 3–4% without sacrificing inference speed or memory consumption much. 3) *FT++* replaces *Faster R-CNN* with *CenterNet2* and uses **LVIS++** for initialization, and it further improves the performance by around 3% in 3-, 5-, and 10-shot. We also observe that while the overall trend of performance is similar for both natural and balanced K -shot sampling, the standard deviation of the natural K -shot performance is less than that of the balanced K -shot.

5 Conclusion

We present the Multi-domain Few-Shot Object Detection (MoFSOD) benchmark consisting of 10 datasets from different domains to evaluate FSOD methods. Under the proposed benchmark, we conducted extensive experiments on the impact of freezing parameters, different architectures, and different pre-training datasets. Based on our findings, we proposed simple extensions of the existing methods and achieved state-of-the-art results on the benchmark. In the future, we would like to go beyond empirical studies and modifications, to designing architectures and smart-tuning methods for a wide range of FSOD tasks.

References

1. Beery, S., Agarwal, A., Cole, E., Birodkar, V.: The iwildcam 2021 competition dataset. arXiv preprint arXiv:2105.03494 (2021) [6](#)
2. Brown, T.B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D.M., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., Amodei, D.: Language models are few-shot learners. In: NeurIPS (2020) [4](#)
3. Cai, Z., Vasconcelos, N.: Cascade R-CNN: Delving into high quality object detection. In: CVPR (2018) [1](#), [8](#)
4. Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov, A., Zagoruyko, S.: End-to-end object detection with transformers. In: ECCV (2020) [1](#), [3](#), [13](#)
5. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: ImageNet: A large-scale hierarchical image database. In: CVPR. pp. 248–255. Ieee (2009) [4](#), [8](#), [10](#)
6. Ertler, C., Mislej, J., Ollmann, T., Porzi, L., Kuang, Y.: Traffic sign detection and classification around the world. arXiv preprint arXiv:1909.04422 (2019) [8](#)
7. Everingham, M., Van Gool, L., Williams, C.K., Winn, J., Zisserman, A.: The pascal visual object classes (VOC) challenge. IJCV **88**(2), 303–338 (2010) [2](#), [5](#)
8. Fan, Q., Zhuo, W., Tang, C.K., Tai, Y.W.: Few-shot object detection with attention-rpn and multi-relation detector. In: CVPR (2020) [2](#), [3](#), [4](#), [5](#), [8](#)
9. Finn, C., Abbeel, P., Levine, S.: Model-agnostic meta-learning for fast adaptation of deep networks. In: ICML (2017) [11](#)
10. Geiger, A., Lenz, P., Urtasun, R.: Are we ready for autonomous driving? the kitti vision benchmark suite. In: CVPR (2012) [6](#)
11. Gu, X., Lin, T., Kuo, W., Cui, Y.: Open-vocabulary object detection via vision and language knowledge distillation. In: ICLR (2022) [4](#), [9](#), [14](#)
12. Guo, S., Huang, W., Zhang, X., Srikhanta, P., Cui, Y., Li, Y., R.Scott, M., Adam, H., Belongie, S.: The imaterialist fashion attribute dataset. arXiv preprint arXiv:1906.05750 (2019) [6](#)
13. Gupta, A., Dollár, P., Girshick, R.B.: LVIS: A dataset for large vocabulary instance segmentation. In: CVPR (2019) [4](#), [8](#)
14. Han, G., He, Y., Huang, S., Ma, J., Chang, S.F.: Query adaptive few-shot object detection with heterogeneous graph convolutional networks. In: ICCV (2021) [5](#)
15. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: CVPR (2016) [8](#), [13](#)
16. Hosang, J.H., Benenson, R., Schiele, B.: How good are detection proposals, really? In: Valstar, M.F., French, A.P., Pridmore, T.P. (eds.) BMVC (2014) [6](#)
17. Huang, G., Laradji, I., Vazquez, D., Lacoste-Julien, S., Rodriguez, P.: A survey of self-supervised and few-shot object detection. arXiv preprint arXiv:2110.14711 (2021) [5](#)
18. Ioffe, S., Szegedy, C.: Batch normalization: Accelerating deep network training by reducing internal covariate shift. In: ICML (2015) [8](#)
19. Kang, B., Liu, Z., Wang, X., Yu, F., Feng, J., Darrell, T.: Few-shot object detection via feature reweighting. In: ICCV (2019) [2](#), [4](#), [5](#), [6](#)
20. Kolesnikov, A., Beyer, L., Zhai, X., Puigcerver, J., Yung, J., Gelly, S., Houlsby, N.: Big transfer (bit): General visual representation learning. In: Vedaldi, A., Bischof, H., Brox, T., Frahm, J. (eds.) ECCV (2020) [4](#)

21. Krasin, I., Duerig, T., Alldrin, N., Ferrari, V., Abu-El-Haija, S., Kuznetsova, A., Rom, H., Uijlings, J., Popov, S., Kamali, S., Mallocci, M., Pont-Tuset, J., Veit, A., Belongie, S., Gomes, V., Gupta, A., Sun, C., Chechik, G., Cai, D., Feng, Z., Narayanan, D., Murphy, K.: OpenImages: A public dataset for large-scale multi-label and multi-class image classification. Dataset available from <https://storage.googleapis.com/openimages/web/index.html> (2017) 2, 6, 8
22. Lin, T.Y., Dollár, P., Girshick, R., He, K., Hariharan, B., Belongie, S.: Feature pyramid networks for object detection. In: CVPR (2017) 8
23. Lin, T., Goyal, P., Girshick, R.B., He, K., Dollár, P.: Focal loss for dense object detection. In: ICCV (2017) 8, 9, 12, 13
24. Lin, T.Y., Maire, M., Belongie, S., Bourdev, L., Girshick, R., Hays, J., Perona, P., Ramanan, D., Zitnick, C.L., Dollár, P.: Microsoft COCO: Common objects in context. arXiv:1405.0312 (2014) 2, 4, 5, 6, 8
25. Miao, C., Xie, L., Wan, F., Su, C., Liu, H., Jiao, J., Ye, Q.: Sixray: A large-scale security inspection x-ray benchmark for prohibited item discovery in overlapping images. In: CVPR (2019) 6
26. Qiao, L., Zhao, Y., Li, Z., Qiu, X., Wu, J., Zhang, C.: DeFRCN: Decoupled faster r-cnn for few-shot object detection. In: ICCV (2021) 2, 3, 5, 10
27. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., Krueger, G., Sutskever, I.: Learning transferable visual models from natural language supervision. In: ICML (2021) 4, 9
28. Redmon, J., Farhadi, A.: Yolo9000: better, faster, stronger. In: CVPR. pp. 7263–7271 (2017) 1, 12
29. Ren, M., Triantafillou, E., Ravi, S., Snell, J., Swersky, K., Tenenbaum, J.B., Larochelle, H., Zemel, R.S.: Meta-learning for semi-supervised few-shot classification. In: ICLR (2018) 2, 5
30. Ren, S., He, K., Girshick, R., Sun, J.: Faster R-CNN: Towards real-time object detection with region proposal networks. NeurIPS (2015) 1, 3, 8
31. Sa, I., Ge, Z., Dayoub, F., Upcroft, B., Perez, T., McCool, C.: DeepFruits: A fruit detection system using deep neural networks. *sensors* 16(8), 1222 (2016) 6
32. Shao, S., Li, Z., Zhang, T., Peng, C., Yu, G., Zhang, X., Li, J., Sun, J.: Objects365: A large-scale, high-quality dataset for object detection. In: ICCV (2019) 8
33. Shao, S., Zhao, Z., Li, B., Xiao, T., Yu, G., Zhang, X., Sun, J.: CrowdHuman: A benchmark for detecting human in a crowd. arXiv preprint arXiv:1805.00123 (2018) 6
34. Snell, J., Swersky, K., Zemel, R.S.: Prototypical networks for few-shot learning. In: NeurIPS (2017) 3, 4, 11
35. Sun, B., Li, B., Cai, S., Yuan, Y., Zhang, C.: Fsce: Few-shot object detection via contrastive proposal encoding. In: CVPR (2021) 2, 3, 4, 5, 9, 10, 11, 12
36. Tan, J., Wang, C., Li, B., Li, Q., Ouyang, W., Yin, C., Yan, J.: Equalization loss for long-tailed object recognition. In: CVPR (2020) 9
37. Triantafillou, E., Zhu, T., Dumoulin, V., Lamblin, P., Evci, U., Xu, K., Goroshin, R., Gelada, C., Swersky, K., Manzagol, P.A., et al.: Meta-dataset: A dataset of datasets for learning to learn from few examples. In: ICLR (2020) 2, 5, 6, 7
38. Tseng, H.Y., Lee, H.Y., Huang, J.B., Yang, M.H.: Cross-domain few-shot classification via learned feature-wise transformation. In: ICLR (2020) 5
39. Vinyals, O., Blundell, C., Lillicrap, T., Kavukcuoglu, K., Wierstra, D.: Matching networks for one shot learning. In: NeurIPS (2016) 2, 5
40. Wang, J., Min, W., Hou, S., Ma, S., Zheng, Y., Jiang, S.: LogoDet-3K: A large-scale image dataset for logo detection. arXiv preprint arXiv:2008.05359 (2020) 6

41. Wang, X., Huang, T.E., Darrell, T., Gonzalez, J.E., Yu, F.: Frustratingly simple few-shot object detection. In: ICML (2020) [2](#), [3](#), [5](#), [8](#), [9](#), [10](#), [11](#), [12](#), [14](#)
42. Wang, X., Cai, Z., Gao, D., Vasconcelos, N.: Towards universal object detection by domain attention. In: CVPR (2019) [5](#)
43. Wang, Y.X., Ramanan, D., Hebert, M.: Meta-learning to detect rare objects. In: ICCV (2019) [3](#), [4](#)
44. Wu, X., Sahoo, D., Hoi, S.: Meta-rcnn: Meta learning for few-shot object detection. In: Proceedings of the 28th ACM International Conference on Multimedia (2020) [3](#)
45. Wu, Y., Kirillov, A., Massa, F., Lo, W.Y., Girshick, R.: Detectron2. <https://github.com/facebookresearch/detectron2> (2019) [9](#)
46. Yan, X., Chen, Z., Xu, A., Wang, X., Liang, X., Lin, L.: Meta R-CNN: Towards general solver for instance-level low-shot learning. In: ICCV (2019) [4](#)
47. Yang, S., Luo, P., Loy, C.C., Tang, X.: Wider face: A face detection benchmark. In: CVPR (2016) [6](#)
48. Zhang, L., Zhou, S., Guan, J., Zhang, J.: Accurate few-shot object detection with support-query mutual guidance and hybrid loss. In: CVPR (2021) [5](#)
49. Zhang, S., Xie, Y., Wan, J., Xia, H., Li, S.Z., Guo, G.: WiderPerson: A diverse dataset for dense pedestrian detection in the wild. IEEE Transactions on Multimedia **22**(2), 380–393 (2019) [6](#)
50. Zhou, X., Girdhar, R., Joulin, A., Krähenbühl, P., Misra, I.: Detecting twenty-thousand classes using image-level supervision. arXiv preprint arXiv:2201.02605 (2021) [4](#), [9](#), [14](#)
51. Zhou, X., Koltun, V., Krähenbühl, P.: Probabilistic two-stage detection. arXiv preprint arXiv:2103.07461 (2021) [3](#), [8](#), [9](#), [13](#)
52. Zhou, X., Koltun, V., Krähenbühl, P.: Simple multi-dataset detection. In: CVPR (2022) [4](#), [8](#), [9](#), [13](#)
53. Zhou, X., Wang, D., Krähenbühl, P.: Objects as points. arXiv preprint arXiv:1904.07850 (2019) [3](#), [9](#), [13](#)
54. Zhu, P., Wen, L., Du, D., Bian, X., Fan, H., Hu, Q., Ling, H.: Detection and tracking meet drones challenge. IEEE TPAMI (2021). <https://doi.org/10.1109/TPAMI.2021.3119563> [6](#)
55. Zhu, X., Hu, H., Lin, S., Dai, J.: Deformable convnets V2: more deformable, better results. In: CVPR (2019) [8](#), [9](#)
56. Zhu, X., Su, W., Lu, L., Li, B., Wang, X., Dai, J.: Deformable DETR: Deformable transformers for end-to-end object detection. In: ICLR (2021) [3](#), [8](#), [9](#), [13](#)
57. Ziller, A., Hansjakob, J., Rusinov, V., Zügner, D., Vogel, P., Günnemann, S.: Oktoberfest food dataset. arXiv preprint arXiv:1912.05007 (2019) [6](#)