



## Few-shot Camouflaged Animal Detection and Segmentation

Thanh-Danh **Nguyen**<sup>a,d,\*\*</sup>, Anh-Khoa **Nguyen Vu**<sup>a,d,\*\*</sup>, Nhat-Duy **Nguyen**<sup>a,d,\*\*</sup>, Vinh-Tiep **Nguyen**<sup>a,d</sup>, Thanh Duc **Ngo**<sup>a,d</sup>,  
 Thanh-Toan **Do**<sup>e</sup>, Minh-Triet **Tran**<sup>b,c,d</sup>, Tam V. **Nguyen**<sup>f,\*\*</sup>

<sup>a</sup>University of Information Technology, Ho Chi Minh City, Vietnam

<sup>b</sup>University of Science, Ho Chi Minh City, Vietnam

<sup>c</sup>John von Neumann Institute, VNU-HCM, Vietnam

<sup>d</sup>Vietnam National University, Ho Chi Minh City, Vietnam

<sup>e</sup>Monash University, Clayton, VIC 3800, Australia

<sup>f</sup>University of Dayton, Dayton, OH 45469, United States

### ABSTRACT

Camouflaged object detection and segmentation is a new and challenging research topic in computer vision. There is a serious issue of lacking data of camouflaged objects such as camouflaged animals in natural scenes. In this paper, we address the problem of few-shot learning for camouflaged object detection and segmentation. To this end, we first collect a new dataset, CAMO-FS, for the benchmark. We then propose a novel method to efficiently detect and segment the camouflaged objects in the images. In particular, we introduce the instance triplet loss and the instance memory storage. The extensive experiments demonstrated that our proposed method achieves state-of-the-art performance on the newly collected dataset.

© 2023 Elsevier Ltd. All rights reserved.

### 1. Introduction

Camouflage is a defense mechanism that animals use to conceal their appearance by blending in with their environment (Singh et al., 2013). Autonomously detecting camouflaged animals is helpful in various applications, e.g., search-and-rescue missions (Le et al., 2019); wild species discovery and preservation activities (Le et al., 2019); and media forensics (manipulated image/video detection and segmentation (Le et al., 2021a)). Although image segmentation methods have been proposed for a long time, general detectors cannot deal with camouflaged animals (Kervrann and Heitz, 1995; Boykov and Funka-Lea, 2006; Li and Sahbi, 2011; Sulimowicz et al.,

2018). The detectors initially developed for camouflage detection (Galun et al., 2003a; Song and Geng, 2010a; Xue et al., 2016a; Pan et al., 2011; Liu et al., 2012; P. Sengottuvelan and Shanmugam, 2008; Yin et al., 2011; Gallego and Bertolino, 2014), which use handcrafted low-level features, are effective only for images with a simple and uniform background. More recently developed deep learning-based detectors (Le et al., 2019; Fan et al., 2020a) for camouflaged object segmentation rely on large-scale data.

With the approach of few-shot learning, we can perform machine learning tasks with given limited data. A few-shot method requires two stages of processing: (1) one base phase training for the model to gain concept knowledge of general domains with abundant data, and then (2) perform novel phase which can do the specific task on few-shot data. In case of the camouflaged object detection and segmentation task, we leverage

<sup>\*\*</sup>Thanh-Danh Nguyen, Anh-Khoa Nguyen Vu, and Nhat-Duy Nguyen contributed equally to this work. Tam V. Nguyen (tamnguyen@udayton.edu) is the corresponding author.

few-shot learning in the concept of camouflaged animals which is rare and hard to find in the wild. Thus, with limited data on camouflaged objects, the models can still well handle the given tasks. However, none of the aforementioned publications targets few-shot camouflaged object detection and segmentation despite its practical applications. The task of segmentation supports better to identify camouflaged objects in terms of specifying which pixels in the images concealing the objects in comparison to classification and detection. In fact, the research on camouflaged animals suffers due to the lack of data. There are not many object classes with rare instances captured in photos. Therefore, in this paper, we would like to address few-shot learning with camouflaged animals.

Our contributions in this work are two-fold:

- First, we build a new benchmark dataset, CAMO-FS, which is among the first datasets to support few-shot detection and instance segmentation on camouflaged instances in nature.
- Second, we propose a new framework to efficiently detect and segment camouflaged instances given a small number of training data for novel classes.

The remainder of this paper is organized as follows. **Section 2** summarizes related work. Next, **Section 3** introduces the newly constructed CAMO-FS dataset and presents our proposed framework for few-shot camouflaged object detection and segmentation. **Section 4** presents the results of our evaluation of baselines on the newly constructed dataset. Finally, **Section 5** summarizes the key points and mentions future work.

## 2. Related Work

### 2.1. Camouflage Research

Given any region (i.e. bounding boxes or polygon masks) presented for an object of interest (i.e. animals or artificial objects) in an image and then they tend to be classified as background, contents in that region can be qualified as camouflaged objects. Thus, a camouflaged object is defined as a set of bounding boxes or camouflaged pixels in an image without any further detailed information such as the number of objects or the

semantic meaning (Le et al., 2019). Although tasks related to camouflaged animals are performed in a wide range of applications, this research field has not been well explored in the literature, especially few-shot learning which is practically suitable to the context of scarce data as camouflaged animals.

**Binary camouflage segmentation.** Prior to the advancement of deep neural networks, most of the work exploit identical regions between camouflaged regions and the background by handcrafted or low-level features, specifically based on external characteristics (e.g., color, shape, orientation, and brightness). Particularly, early camouflage detection works had attention on the foreground region even when some of its texture was similar to the background (Galun et al., 2003a; Song and Geng, 2010b,a; Xue et al., 2016a). The foreground was distinguishable from the background via simple features, such as color, intensity, shape, orientation, and edge (Siricharoen et al., 2010; Galun et al., 2003b; Kavitha et al., 2011; Song and Geng, 2010b; Xue et al., 2016b). A few methods (Pan et al., 2011; Hou and Li, 2011; Liu et al., 2012; P. Sengottuvelan and Shanmugam, 2008; Yin et al., 2011; Gallego and Bertolino, 2014) based on handcrafted low-level features have been proposed for tackling the problem of camouflage detection. However, they are effective only for images with a simple and uniform background. Thus, their performances are unsatisfactory in camouflaged object segmentation due to the substantial similarity between the foreground and the background.

Until now, the convention of binary prefers binary ground truth camouflaged object datasets (Fan et al., 2020a; Le et al., 2019; Skurowski et al., 2018). Existing methods for camouflaged objects (Le et al., 2019; Zhai et al., 2021; Li et al., 2021a; Fan et al., 2020a; Mei et al., 2021; Yan et al., 2021; Zhu et al., 2021b; Lv et al., 2021) based on binary ground truth are considered as the binary camouflage segmentation. For example, Le *et al.* (Le et al., 2019) proposed an end-to-end Anabranched Network, dubbed ANet which includes two streams of classification and segmentation. The outputs of both streams are fused to improve the segmentation performance of camouflaged objects. This proposed network was also flexibly applied to any

fully convolutional networks. Similarly, motivated by the way of hunting strategies of predators, Fan *et al.* (Fan et al., 2020a) designed Search Identification Network (SINet) with two main modules to simulate this hunting behavior, namely a search module searching for targets and an identification module identifying the existence of targets then catching them. Yan *et al.* (Yan et al., 2021) recently introduced MirrorNet, a dual-stream network comprising a mainstream and a mirror stream. This mirror stream aimed to capture instinct information by horizontally flipping camouflaged objects to break their camouflaged nature and make them more distinguishable. Zhu *et al.* (Zhu et al., 2021b) presented the TINet, which interactively refines multi-level texture and segmentation features and thereby gradually enhances the segmentation of camouflaged objects. Lv *et al.* (Lv et al., 2021) simultaneously worked on ranking and localization to well-present camouflaged objects. As a result, they formed a triplet task with localizing, segmenting, and ranking the camouflaged objects. Besides, the authors also introduced the NC4K dataset for camouflaged segmentation. Such methods reveal the presence of the camouflaged objects with the high level of bounding boxes and contain corresponding pixel-wise ground truth belonging to camouflage. Further understanding of the camouflage level may help us to give comparative analyses, finding evidence for links between camouflage and other defensive strategies with aspects of habitat and life-history (Price et al., 2019).

**Camouflage instance segmentation.** Although several works have been proposed, there is still a difficulty in efficiently exploring the information of camouflage animals, especially at the instance level with more challenging detailed masks. Therefore, for ease of training methods with the challenging task of camouflaged instance segmentation, Le *et al.* (Le et al., 2021b) introduced a framework with several state-of-the-art methods and proposed a tool with user interactive cues to tune the segmentation mask on a website. Realizing that the semantic level is not detailed enough, Le *et al.* (Le et al., 2022) introduced a camouflage fusion learning (CFL) to utilize the strength of different instance segmentation methods by fusing various models

via learning image contexts.

**Camouflage datasets.** CamouflagedAnimals (Pia Bideau, 2016) and CHAMELEON (Skurowski et al., 2018) were the first two camouflage datasets with mask annotations. The two datasets do not contain enough images to train deep learning methods. Le *et al.* (Le et al., 2019) created the CAMO dataset, the first camouflage dataset with more than 1,000 annotated images. It contains 1,250 annotated images, which is a limited number of samples to train and evaluate deep learning methods. Then, Fan *et al.* (Fan et al., 2020a) collected the COD dataset, which comprises 10,000 images (both camouflage and non-camouflage) divided into 5 meta-categories. However, they annotated only 5,066 camouflage images. Lamdouar *et al.* (Lamdouar et al., 2020) recently developed the MoCA dataset for the camouflage object detection task; it contains only bounding box ground truths. Hence, these datasets limit their annotations at binary ground truth datasets which have a shortage of intensive annotations for multi-task camouflage problems. CAMO++ (Le et al., 2022) is different from the above-mentioned dataset, CAMO++ provides a benchmark for camouflaged instance segmentation with more comprehensive annotations and diverse meta-categories of 10. The dataset comprises 5,500 images with superiority over other datasets on instances including 32,756 instances for both camo and non-camo objects.

## 2.2. Few-shot Learning

**Few-shot object detection (FSOD).** When having some available samples of given classes with their corresponding bounding boxes, FSOD aims to learn from these limited data in order to help models adapt to the new classes. To date, several works (Fan et al., 2020b; Kang et al., 2019; Wang et al., 2020; Yan et al., 2019) have been proposed to deal with FSOD. Early works (Kang et al., 2019; Yan et al., 2019) mainly prefer to overcome the difficulties of the data scarcity of FSOD via meta-learning approaches by combining supportive information from meta-based streams with their main streams. Particularly, Bingyi (Kang et al., 2019) proposed a Feature Reweighting framework that leverages the free-proposal approach of a

well-known one-stage framework such as YOLO (Redmon and Farhadi, 2017) to boost FSOD performance. The network integrated a meta-model that aims to generate reweighting vectors from support samples for highlighting the attention to features from the YOLO network. Conversely, Meta RCNN (Yan et al., 2019) based on the two-stage proposal approach as Mask RCNN (He et al., 2017) and fed available annotations such as bounding boxes and segmented masks to train a meta-network called Predictor-head Remodeling Network for inferring attention features. Fan *et al.* (Fan et al., 2020b) recently proposed to take advantage of support images from a massive FSOD dataset to generate significant results combined with their proposed network called Attention-RPN, Multi-Relation Detectors. The Attention-RPN directed the trained model where to look on the image for the task of object detection. Differently, Wang *et al.* (Wang et al., 2020) simply adopted Faster RCNN with two-stage finetuning to transfer massive knowledge from abundant data in the base model to fine-tune the novel one by freezing the whole network except for the fully connected layer for object classification. Through this simple straightforward mechanism, this model significantly improved few-shot performance without a complex pipeline of training the model. Further, such works (Li et al., 2021b; Hu et al., 2021; Xiao and Marlet, 2020; Han et al., 2021, 2022a) presented advanced methods by applying class max-margin, multiple scale proposals, or feature alignment in FSOD. Other ones were based on transformed inputs (Li and Li, 2021; Wu et al., 2020), transformer approaches (Zhang et al., 2021; Han et al., 2022b), contrastive method (Sun et al., 2021), or kernels design (Zhang et al., 2022). Other methods (Wang et al., 2020; Qiao et al., 2021; Zhang and Wang, 2021; Zhu et al., 2021a) relied only on query images to deal with FSOD via extra text data (Zhu et al., 2021a), unlabeled image (Khandelwal et al., 2021), generated samples (Zhang and Wang, 2021), gradient scaling (Qiao et al., 2021).

**Few-shot object segmentation (FSOS).** Recently, the field of few-shot segmentation gained attention from the community. As mentioned above, the first work Meta RCNN originated from Mask RCNN, therefore, Meta RCNN simultaneously per-

formed detection and segmentation. Liu *et al.* (Liu et al., 2020) utilized a cross-reference network for generic image segmentation. The authors proposed a cross-reference mechanism and a mask refinement module to specifically support the task of segmentation. Before, Dong *et al.* (Dong and Xing, 2018) proposed a prototype learning component in a framework of semantic segmentation that learned to take discriminative information from features to help segment objects better. Also, Wang *et al.* (Wang et al., 2019) introduced a prototype align method that learns class-specific prototype representations from a few image samples to perform segmentation over the query images. Lately, Liu *et al.* proposed a dynamic prototype convolution network to address few-shot semantic segmentation. The work of (Saha et al., 2022) proposed context-aware prototype learning. (Tian et al., 2022) introduced generative models approach for this task. These methods focused on generic objects.

### 3. Proposed Method

#### 3.1. CAMO-FS Benchmark Dataset

Camouflaged data tends to be more difficult to collect in the real world rather than non-camouflaged ones. Generating intensive annotations with multi-task or hierarchical labels for camouflaged objects is also costly and complicated, especially with the pixel level as polygon masks. Particularly, the visual characteristics of a camouflaged object are extremely identical to the background. The external appearances (i.g. the intensity, color, and textures) are close to their surrounding environment, the boundary between camouflaged objects and the background or other identical-type camouflaged objects in case of being nearly or partly overlapped. Thus, it is really tough to provide the concurrence between annotators due to ambiguity in verifying camouflaged regions blended in surroundings. For ease of data preparation such as collections and annotations, one of the most common way is that inherits existing camouflaged datasets and CAMO++ (Le et al., 2022) is our selected dataset since it is a high-diversity dataset with a variety of camouflaged object categories. Furthermore, the key to few-shot learning lies in

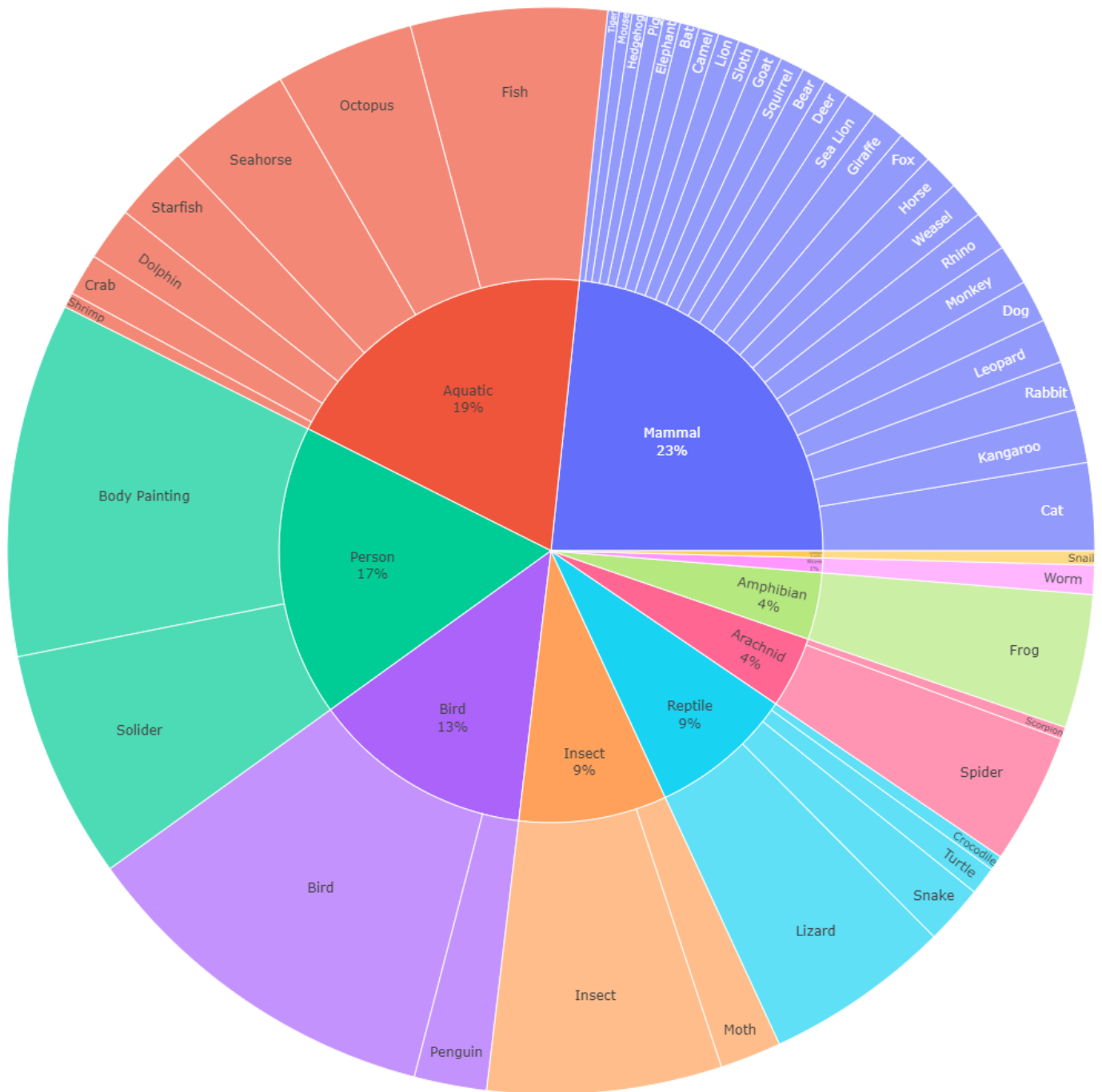


Fig. 1: Hierarchical taxonomic structure of our CAMO-FS dataset and the corresponding distribution of camouflaged coarse-grained classes. Best viewed online in color and zoomed in.

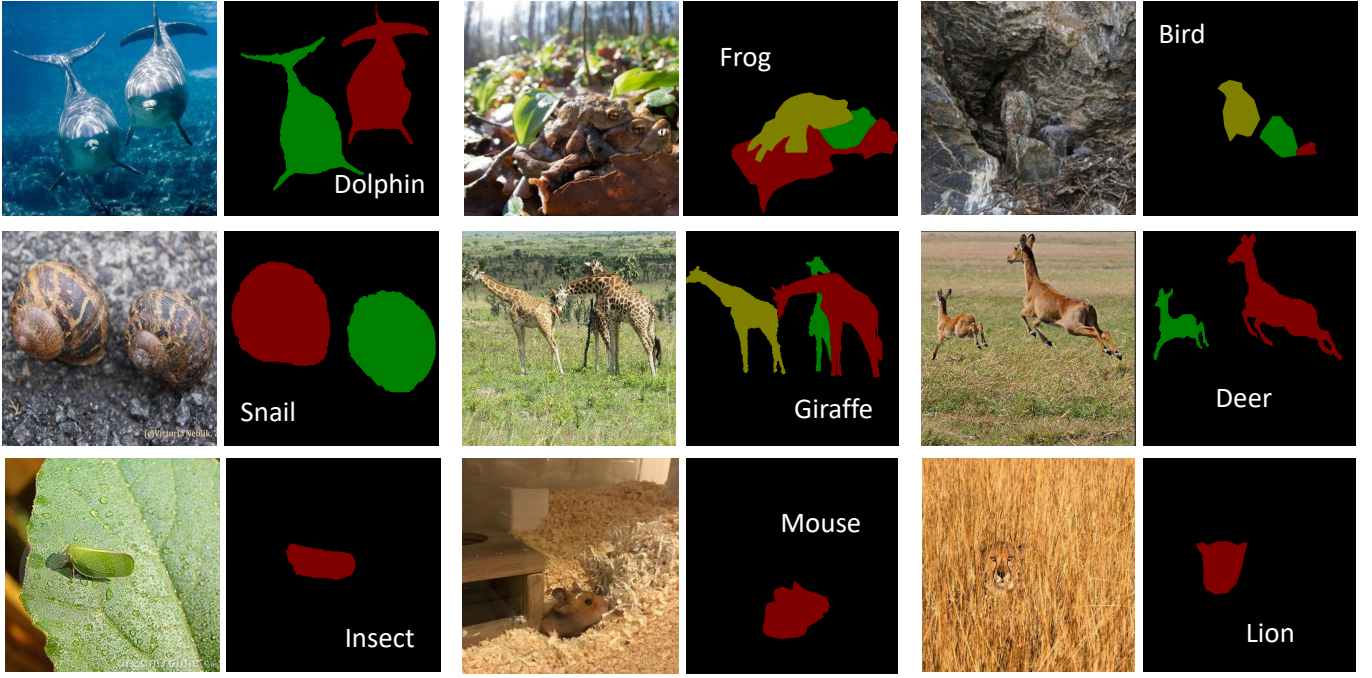


Fig. 2: Exemplary images with instance-level mask annotations from our proposed CAMO-FS dataset.

Table 1: Statistics of camouflage datasets (without non-camo images).

Dataset	Year	Publication	Type	#Annot. Camo. Img.	#Meta- Cat.	#Obj. Cat.	#Ins. or #Obj. per Img.	Bbox. GT	Obj. Mask GT	Ins. Mask GT	Few-shot
CamouflagedAnimals (Pia Bideau, 2016)	2016	ECCV	Video	181	-	6	1.238	×	✓	✓	×
MoCA (Lamdouar et al., 2020)	2020	ACCV	Video	7,617	-	67	1.000	✓	×	×	×
CHAMELEON (Skurowski et al., 2018)	2018	-	Image	76	-	-	1.000	×	✓	×	×
CAMO (Le et al., 2019)	2019	CVIU	Image	1,250	2	8	1.000	×	✓	×	×
COD (Fan et al., 2020a)	2020	CVPR	Image	5,066	5	69	1.171	✓	✓	✓	×
CAMO++ (Le et al., 2022)	2022	TIP	Image	2,695	10	47	1.171	✓	✓	✓	×
CAMO-FS (Ours)	2023	-	Image	2,858	10	47	1.172	✓	✓	✓	✓

the generalization ability of the pertinent model when presented with a few available samples. The context of camouflaged objects inherently match this understanding because the number of camouflaged images is often scarce in practice.

**CAMO++ Dataset.** CAMO++ generally contains camouflaged and non-camouflaged images with a total of 5,500 images corresponding to 32,756 instances (Le et al., 2022). The dataset contains 93 fine-grained classes assigned to 13 coarse-grained classes. However, in the case of camouflaged objects, there are 47 fine-grained classes designed with a hierarchical structure and assigned into 10 coarse-grained classes. In detail, CAMO++ contributes 2,695 camouflage images includ-

ing 1,250 existing camouflage images in the previous CAMO dataset with 1,450 newly collected camouflage images for CAMO++. In this scope of our paper, 2,800 remaining non-camouflage images are ignored. CAMO++ especially provides common ground truths such as bounding boxes, object masks, and instance masks which are suitable for many tasks of camouflage research.

**CAMO-FS Dataset.** We leverage the available CAMO++ to build our CAMO-FS dataset. In this way, we inherit the biology taxonomic and vision taxonomic structure of CAMO++ which helps us to reduce the burden of data collection. **Table 1** provides an overview of previous works done on camou-



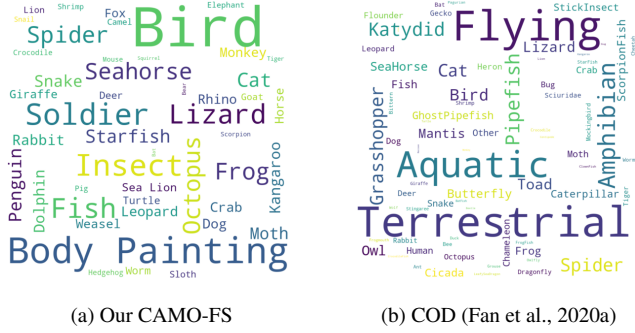


Fig. 3: Word cloud of class-distribution of camouflaged instances.

flage, which is mentioned in the related work, and our proposed CAMO-FS in terms of main characteristics. We exploit the diversity of CAMO++ by its 10 meta-categories to build up the few-shot concept for instance segmentation. To this end, our CAMO-FS not only keeps a good ratio of instances per image of 1.172 but also contributes as the very first dataset specific for few-shot research on camouflaged animals. Note that the large amount of images in some datasets does not mean they are all camouflaged images. Figure 3 illustrates the class-distribution of our CAMO-FS dataset and COD (Fan et al., 2020a). Our CAMO-FS shares similar categorical diversity with the original CAMO++ (Le et al., 2022).

However, imbalanced data and a shortage of the number of images of some classes inherently exist in CAMO++ posing problems of evaluation for few-shot tasks. Particularly, there are 11 classes (e.g. *Camel*, *Dolphin*, *Elephant*, *Horse*, *Kangaroo*, *Monkey*, *Penguin*, *Bat*, *Bear*, *Squirrel* and *Rhino*) even having a shortage of images that are needed to train a few-shot model. Hence, we hardly perform training or testing on these classes. As a result, we collect more data for these classes with 163 total images corresponding to 181 instances (an average of 15-16 instances per class). We also remove images with mistakes in the original dataset. The statistics of collected data are shown in Table 2. By gathering more camouflaged animals and combining them with the CAMO++ dataset, we conduct our CAMO-FS dataset for few-shot camouflaged animal detection and segmentation with 2,858 total images corresponding to 3,342 instances. Figure 1 shows the vision taxonomic structure of coarse-grained and corresponding fine-grained classes and il-

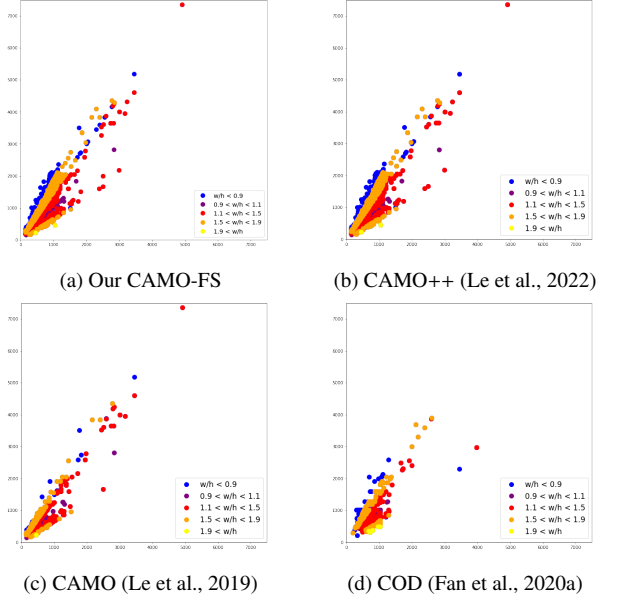


Fig. 4: Distribution of camouflage image resolution. Best viewed online in color and zoomed in.

lustrates the ratios of 10 coarse-grained classes in our proposed CAMO-FS dataset. Figure 2 visualizes the exemplary images with mask annotations from our proposed CAMO-FS.

In Table 3, we report the aggregated number of instances per image. The number of instances per image ranges from 1 to 25 and commonly falls into 1, then 2 and 3 while the remaining is beyond 3 instances. As can be seen, the number of images that contain 1 to 3 instances takes up a large proportion of the entire dataset. This also illustrates the problem of data imbalance between the number of instances and the ratio of images in the dataset, which reflects a problem that the presence of camouflaged animals captured in photos is often limited, i.e. mostly one animal per image. Additionally, although being claimed in (Le et al., 2022) that camouflaged objects in CAMO++ were localized over the entire image, after removing non-camouflage objects and adding new camouflaged images, we have the distributions of object centers in normalized image coordinates over all images in the CAMO-FS dataset as in Figure 5-a. This means camouflaged animals tend to be located in the center of images. Indeed, to capture images of camouflaged animals in the wild, photographers need to carefully focus on the animals, which leads to the central layout of collected images. Also in Figure 5, we illustrate the center bias of camouflaged images

Table 2: Extra collected number of images and instances in CAMO-FS dataset.

Classes	Bat	Bear	Camel	Dolphin	Elephant	Horse	Kangaroo	Monkey	Penguin	Rhino	Squirrel	Total
#images	12	14	14	13	14	16	22	16	11	14	17	163
#instances	12	14	15	19	14	17	25	20	14	14	17	181

Table 3: CAMO-FS dataset instances per image distribution.

No. of instances	Ratio (%)	#Images
1	90.5	2581
2	1.05	190
3	1.79	51
3+	6.66	30

in other CAMO (Le et al., 2019) and COD (Fan et al., 2020a) datasets for better visual comparison. In Figure 4, we present the image resolution among camouflage datasets. As we only consider camouflaged images of CAMO++ (Le et al., 2022) and COD (Fan et al., 2020a), the density of our CAMO-FS is slightly higher than CAMO++ as a result of our extra collection of images presented in Table 2. In comparison with the previous COD (Fan et al., 2020a) and CAMO (Le et al., 2019), our CAMO-FS image resolution distribution is more satisfying in diversity.

To effectively create the data for the few-shot problem, we get  $M$  instances from the CAMO-FS dataset to create training sets (in our setup,  $M = 5$ ) and use the remaining instances for testing. We only remove some objects of the higher-level training set if it exists to create the other few-shot settings. For example, we get all elements to generate 5-shot training data and discard 2 in 5 objects to make a 3-shot one. In this way, the 5-shot benchmark contains objects of the 3-shot dataset and the 3-shot setting contains the objects of the 2-shot one.

To the best of our knowledge, this is among the first works to address few-shot camouflaged instance segmentation and detection. Given the lack of a large-scale dataset for training and testing purposes on camouflaged animal issues, we build a

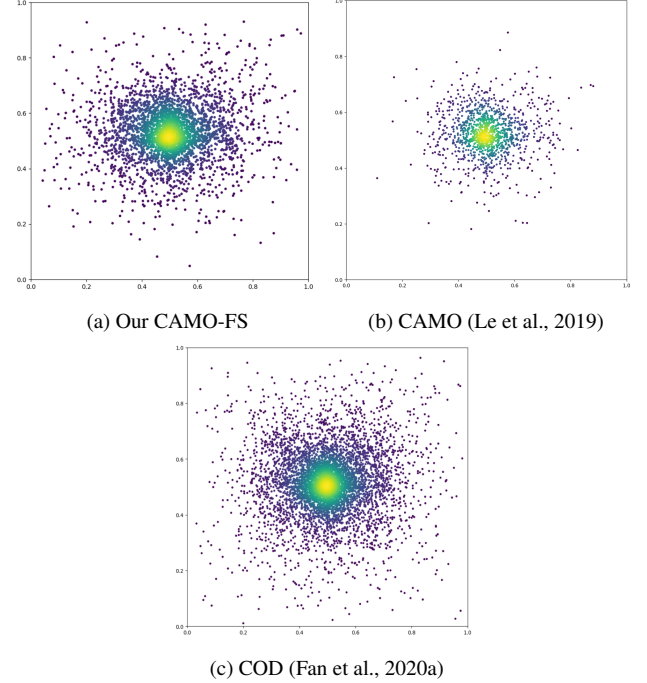


Fig. 5: Instance center bias camouflaged datasets. Best viewed online with color and zoomed in.

benchmark for the task of few-shot camouflaged instance segmentation and detection.

### 3.2. General Framework

**Few-shot instance segmentation formulation.** In few-shot learning, we have one set of base classes denoted  $C_{base}$  with a large amount of available training data, and one disjoint set of novel classes denoted  $C_{novel}$  containing a small amount of training data. This amount is small to a few samples. The ultimate goal is to train a model to predict well on the novel classes  $C_{test} = C_{novel}$  (Snell et al., 2017; Vinyals et al., 2016) or on both base and novel data  $C_{test} = C_{base} \cup C_{novel}$  (Gidaris and Komodakis, 2018). In few-shot classification, this work (Vinyals et al., 2016) introduces the method of episodic training. The method sets up a series of episodes  $E_i = (I_q, S_i)$  where  $S_i$  is a



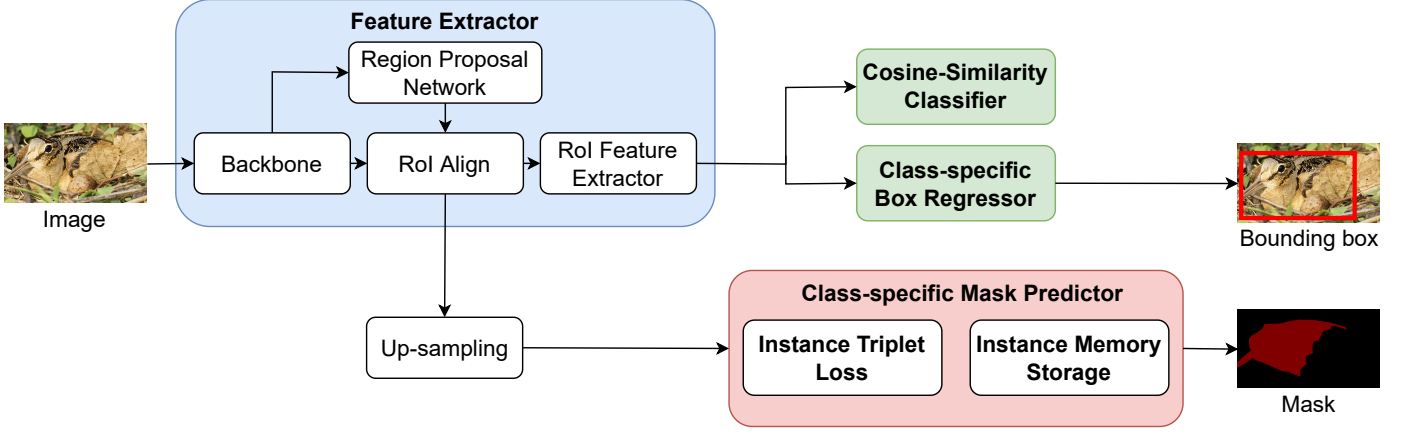


Fig. 6: General framework for few-shot camouflaged instance detection and segmentation.

support set that contains  $N$  classes from  $C_{train} = C_{novel} \cup C_{base}$  along with  $K$  examples per class (so-called  $N$ -way  $K$ -shot). A network is then trained to classify an input image  $I_q$ , termed query image, out of the classes in  $S_i$ . The key idea is that solving a different classification task for each episode leads to better generalization and results on  $C_{novel}$ . The extended versions of this method are FSOD (Kang et al., 2019) and FSIS (Fan et al., 2020c; Yan et al., 2019). Those proposals consider all objects in an image as queries and they have a single support set per image instead of per query. However, there exist challenges in FSIS which are not only classification task on the query objects, but also how to determine their localization and segmentation. Use an image  $I_q$  to query, FSIS returns labels  $y_i$ , bounding boxes  $b_i$ , and segmentation masks  $M_i$  for all objects in  $I_q$  that belong to the set of  $C_{test}$ .

**General framework.** Originated from TFA (Wang et al., 2020) which uses Faster R-CNN (Ren et al., 2015), MTFA (Ganea et al., 2021) employs a mask prediction branch to return the pixel-wise mask for the segmentation task. In this work, we leverage the architecture of MTFA model (Ganea et al., 2021) based on Mask R-CNN (He et al., 2017) which is a two-stage training and fine-tuning mechanism. We train the first stage of the framework on 80 classes from the COCO dataset. This stage results in the base model weights for the second stage of novel fine-tuning. In the fine-tuning stage, we apply the few-shot technique to learn the novel concepts of camouflaged instances in our proposed CAMO-FS dataset.

Similar to Mask R-CNN, the input images are fed into a feature extractor  $F$  consisting of backbone  $B$ , RoI Align, RoI feature extractor modules, and a region proposal network. There are three heads specifying three tasks that this scheme supports: a classification head  $C$ , a box regression head  $R$ , and a new attached mask prediction head  $M$ . In the first stage, the network is trained on the base classes  $C_{base}$ . Then in the second stage, we froze the backbone network  $B$  of the feature extractor  $F$  and only perform training on the prediction heads. Thus, only RoI classifier  $C$ , box regressor  $R$ , and mask predictor  $M$  are fine-tuned in the second stage. In Figure 6, there exists a branch called mask predictor  $M$ . We apply similarly to Ganea et al. (Ganea et al., 2021) by using this two-stage fine-tuning approach. Firstly, the network is trained on base classes with lots of abundant data and then fine-tuning all predictor head  $C$ ,  $R$ , and  $M$  on novel data of  $K$  shots for each class.

Not a simple mask predictor  $M$  that we use, we enhance the performance of the instance segmentation task by employing the two concepts of instance triplet loss and instance memory storage which are clearly described in the next section. The two improvements are inspired not only by the instance segmentation task in general but also by the camouflaged instance segmentation specifications.

### 3.3. Framework Improvement

One of the characteristics of camouflage instances is the camouflage texture similar to the background. This makes the pre-

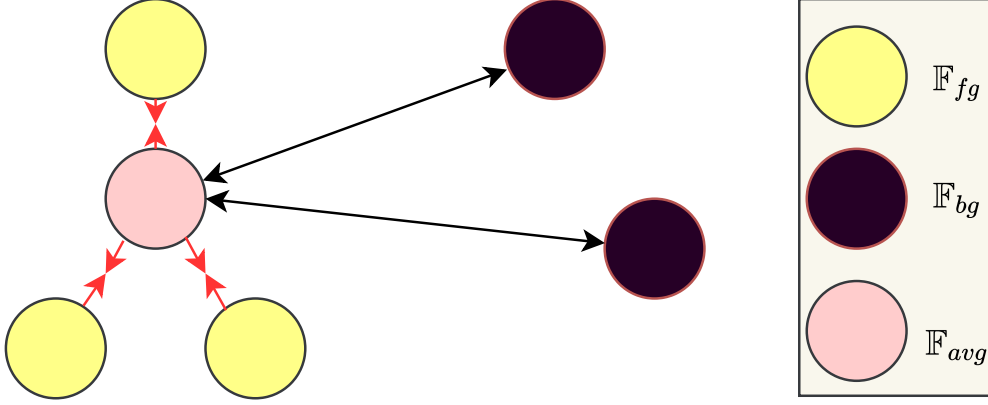


Fig. 7: Visualization of triplet loss for region proposal.

cise identification of the boundary areas difficult. It is more critical in the context of few-shot learning where the concepts of a class are represented by only a few samples.

In this work, we thus propose improvements to enhance distinguishable features between background and foreground areas. In particular, we explore two approaches that focus on loss functions. The first one is the triplet loss function which was known as a strong metric to support the network in creating discrimination features between anchor and negative. The second approach is the idea of memory bank, which is used to enhance the distance between foreground and background not only for individual instances but also for each novel class.

To calculate the loss function, we employ the mask annotation for RoI features to collect the  $\mathbb{F}_{bg}$  background and  $\mathbb{F}_{fg}$  foreground features by location on each RoI. Both  $\mathbb{F}_{bg}$  and  $\mathbb{F}_{fg}$  for each proposal are used to calculate the respective loss functions which are presented in the following sections.

### 3.3.1. Instance triplet loss

With the idea of enhancing the discrimination between camouflaged instances and their backgrounds, we leverage the power of the triplet loss function (Balntas et al., 2016). Specifically, we treat the pixels of an object as positive points and the background as negative ones. Accordingly, we force the model to learn the distinguished features among the foreground and background representatives. The more distinguished among features, the better a model can do to detect or segment camouflaged instances. In this way, we “highlight” the camouflaged

instances so that the model is able to “see” them.

For each RoI, we consider the average foreground features  $\mathbb{F}_{avg} = \frac{1}{|\mathbb{F}_{fg}|} \sum \mathbb{F}_{fg}$  as anchors with the foreground feature  $\mathbb{F}_{fg}$  as positive and the background feature  $\mathbb{F}_{bg}$  as negative to apply the triplet loss function (Balntas et al., 2016). In this way, the model tries to learn to minimize the distance between foreground representatives and maximize the distance between background representatives as shown in Figure 7. We use cosine similarity to calculate the distance instead of Euclidean distance. The loss function is defined as:

$$\mathcal{L}_{triplet} = \max\{d(\mathbb{F}_{avg}, \mathbb{F}_{fg}) - d(\mathbb{F}_{avg}, \mathbb{F}_{bg}) + margin, 0\} \quad (1)$$

$$d(x, y) = 1 - \frac{x \cdot y}{\|x\| \cdot \|y\|}$$

, where *margin* controls the discrimination between foreground and background features. In our experiments, we set *margin* of 0.5. The final loss is shown below:

$$\mathcal{L}_{final} = \mathcal{L}_{mrcnn} + \alpha \mathcal{L}_{triplet} \quad (2)$$

, where  $\mathcal{L}_{mrcnn}$  is the loss of Mask R-CNN (He et al., 2017) and  $\alpha$  is the balance weight of  $\mathcal{L}_{triplet}$  to train the model.

### 3.3.2. Instance memory storage

The memory bank is designed to store information within a class and the class information is updated during the training. Still, the model can learn information at a global level and high consistency for each class. On the other hand, storing

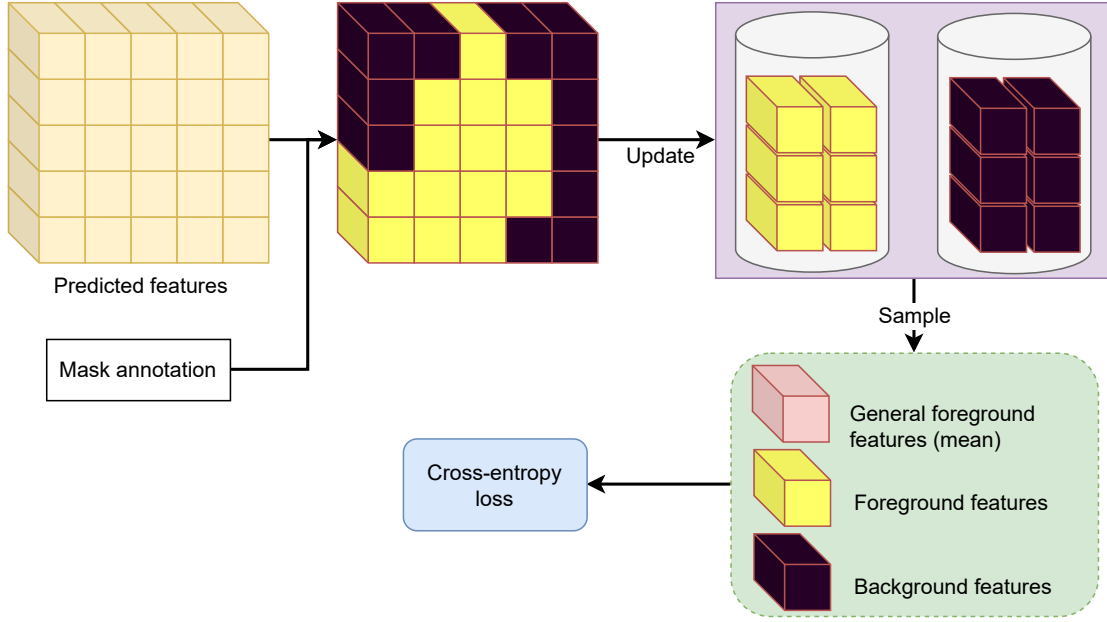


Fig. 8: Visualization of memory loss for region proposal.

and updating the features in the memory bank for each iteration during training also create more variants. By leveraging these advantages, we propose the memory bank for few-shot camouflage instance segmentation. To be specific, we use the memory bank to contain the background and foreground features per each class, and make use of features to calculate the discrimination between areas of object and no object in region proposals (shown in Figure 8).

**Storing and updating:** The memory bank for each class contains  $2N$  features including  $N$  of foreground features and  $N$  of background features. While the memory bank receives new features, the module concatenates them with existing old features. In case the number of features is greater than the given  $N$  features, the memory bank releases the oldest features to maintain the number of features to  $N$ . This process updates the features in the memory bank and keeps the quantity of the stored features appropriate to the memory size (also known as the memory capacity).

**Sampling:** To calculate the loss value, the memory bank has to provide three elements  $\mathbb{F}_{fg}$ ,  $\mathbb{F}_{bg}$ , and  $\mathbb{F}_{general}$ .  $\mathbb{F}_{fg}$  and  $\mathbb{F}_{bg}$  are all foreground and background features that module storing. The  $\mathbb{F}_{general}$  is the general foreground feature, and it is created for each class by averaging the  $\mathbb{F}_{fg}$ .

Let  $\mathbb{F}_{fg}^i$  be the  $i$ -th foreground feature and  $\tau$  be a temperature hyper-parameter in (Wu et al., 2018). The memory loss function for camouflaged instances is introduced as follows:

$$\mathcal{L}_{memory} = -\log \frac{\exp(\mathbb{F}_{general} \cdot \mathbb{F}_{fg}^i / \tau)}{\sum_{j=0}^{|\mathbb{F}_{bg}|} \exp(\mathbb{F}_{general} \cdot \mathbb{F}_{bg}^j / \tau) + \exp(\mathbb{F}_{general} \cdot \mathbb{F}_{fg}^i / \tau)} \quad (3)$$

In our experiments, we set  $\tau$  as 1. The final loss for memory approach is defined as:

$$\mathcal{L}_{final} = \mathcal{L}_{mrcnn} + \beta \mathcal{L}_{memory}. \quad (4)$$

Here, the parameter  $\beta$  of  $\mathcal{L}_{memory}$  is used during the training process to keep the balance between two loss functions.

## 4. Experiments

We first overview the metrics and the experiment settings and the implementation details in Section 4.1 and then we evaluate and discuss our improvement on the general framework, as well as ablation study for our core proposed methods in Section 4.2.

### 4.1. Overview

As specified in this work, we utilize the proposed CAMO-FS dataset containing images of camouflaged animals in the wild to establish the evaluation of our baseline and proposed improvement. We follow the concept procedure published in

Table 4: State-of-the-art comparison among the baseline model of MTFA (Ganea et al., 2021), Mask RCNN<sup>†</sup>, and IFS-RCNN (Nguyen and Todorovic, 2022) and our proposed methods of instance triplet loss and instance memory storage. The best, and second best performances are marked in **red**, and **blue**, respectively.

Model			Novel AP									
Method	Base model	Our proposal	Instance Segmentation					Object Detection				
			1	2	3	5	Mean	1	2	3	5	Mean
MTFA	COCO-80 R-50	-	3.88	6.95	7.30	8.66	6.70	3.43	6.99	7.21	7.82	6.36
M-RCNN <sup>†</sup>			4.08	6.79	6.90	8.29	6.52	2.82	5.09	5.46	6.18	4.89
IFS-RCNN			4.17	6.26	5.73	6.38	5.64	3.92	6.06	5.47	6.60	5.51
MTFA	COCO-80 R-101	-	3.66	6.21	6.16	5.95	5.50	2.93	5.9	5.84	5.84	5.13
M-RCNN <sup>†</sup>			4.39	7.69	7.94	<b>10.09</b>	7.53	3.03	5.80	6.20	7.79	5.71
IFS-RCNN			4.27	6.55	6.07	7.80	6.17	3.79	6.28	6.01	8.08	6.04
MTFA	COCO-80 R-101	Instance Triplet Loss	4.46	5.57	6.41	8.48	6.23	4.04	7.28	7.49	<b>9.76</b>	7.14
M-RCNN <sup>†</sup>			<b>5.73</b>	<b>7.97</b>	<b>8.52</b>	<b>9.92</b>	<b>8.04</b>	<b>5.08</b>	<b>7.56</b>	<b>7.85</b>	9.67	<b>7.54</b>
IFS-RCNN			5.35	6.01	7.80	6.23	6.35	4.71	5.66	7.10	6.06	5.88
MTFA	COCO-80 R-101	Instance Memory Storage	5.46	6.95	7.36	9.61	7.35	4.50	6.95	7.55	<b>10.36</b>	7.34
M-RCNN <sup>†</sup>			<b>5.52</b>	<b>7.84</b>	<b>8.65</b>	9.82	<b>7.96</b>	<b>4.92</b>	<b>7.39</b>	<b>7.96</b>	9.52	<b>7.45</b>
IFS-RCNN			2.99	6.83	6.14	9.03	6.25	2.74	6.39	5.94	8.44	5.88

M-RCNN<sup>†</sup> is Mask R-CNN (He et al., 2017) with sigmoid classifier.

FSOD (Wang et al., 2020; Kang et al., 2019; Yan et al., 2019). In the first stage of the base phase, we train our model with abundant data from 80 classes of the COCO dataset as proposed in (Kang et al., 2019). In the second stage of the novel phase, we evaluate the performance of having  $K = 1, 2, 3, 5$  shots per each novel class.

To report our results on detection and instance segmentation tasks, we use average precision (AP) and average recall (AR). To be detailed, we report AP@50 and AP@75, along with AR@10. Besides, we also report AP and AR at small, medium, and large scales of the instances to further understand the model performance. For more details, readers can visit the homepage of the COCO dataset for detection and segmentation evaluation metrics <sup>1</sup>.

Our MTFA (Ganea et al., 2021) baseline is implemented using Detectron2 framework (Wu et al., 2019). Our backbone is ResNet-101 (He et al., 2016) with Feature Pyramid Network

(Lin et al., 2017). Each experiment is set up with a single GPU GeForce RTX 2080Ti with a batch size of 2 images. The novel phase has a learning rate of 0.00125 inferred from the MTFA configuration. We set the balance parameters  $\alpha = 1e^{-1}$  and  $\beta = 1e^{-2}$  when we train the model with instance triplet and instance memory loss function, respectively. Please visit the publication (Wang et al., 2020) or (Wu et al., 2019) for more details on other parameters.

#### 4.2. Results and Discussion

To prove the effectiveness of our proposed methods, we conducted experiments on our proposed CAMO-FS dataset. We tested with  $K = 1, 2, 3, 5$  shots, respectively. Table 4 presents the evaluation of the performance of our method of triplet loss and memory storage over our baseline MTFA (Ganea et al., 2021), the model of Mask R-CNN with sigmoid classifier, and the state-of-the-art method IFS-RCNN (Nguyen and Todorovic, 2022) in the approach of few-shot instance segmentation. We reported experiments on those three models and chose COCO-

<sup>1</sup><https://cocodataset.org/#detection-eval>

Table 5: Our improvement of instance triplet loss and instance memory storage on MTFA (Wang et al., 2020). The best performance is marked in **boldface**. # is Number of shots, “Memory” is Instance Memory Storage and “Triplet” is Instance Triplet Loss.

#	Method	AP	AP50	AP75	APs	APm	API	AR1	AR10	ARs	ARm	ARl
<b>Instance Segmentation</b>												
<b>1</b>	<b>Baseline MTFA</b>	3.66	5.37	4.09	22.42	4.35	2.01	11.30	13.58	25.97	<b>12.96</b>	12.53
	<b>MTFA + Triplet</b>	4.46	8.21	4.60	21.33	4.13	<b>4.01</b>	12.36	15.04	23.17	9.49	16.67
	<b>MTFA + Memory</b>	<b>5.46</b>	<b>9.20</b>	<b>6.17</b>	<b>27.79</b>	<b>6.20</b>	<b>4.01</b>	<b>17.08</b>	<b>19.99</b>	<b>29.41</b>	11.45	<b>20.89</b>
<b>2</b>	<b>Baseline MTFA</b>	6.21	8.92	7.28	32.64	<b>7.75</b>	3.50	18.88	21.12	<b>35.82</b>	<b>15.49</b>	20.14
	<b>MTFA + Triplet</b>	5.57	9.45	6.04	25.83	3.01	5.37	15.67	17.33	26.13	7.37	17.50
	<b>MTFA + Memory</b>	<b>6.95</b>	<b>10.72</b>	<b>7.60</b>	<b>33.62</b>	5.73	<b>6.44</b>	<b>20.00</b>	<b>22.15</b>	34.25	13.86	<b>20.92</b>
<b>3</b>	<b>Baseline MTFA</b>	6.16	8.95	6.68	33.74	6.19	5.08	20.25	22.95	36.83	16.31	21.63
	<b>MTFA + Triplet</b>	6.41	10.67	6.72	30.39	5.17	5.30	20.69	22.98	31.90	15.69	22.53
	<b>MTFA + Memory</b>	<b>7.36</b>	<b>11.23</b>	<b>8.49</b>	<b>37.03</b>	<b>6.24</b>	<b>5.64</b>	<b>24.40</b>	<b>27.69</b>	<b>38.44</b>	<b>17.02</b>	<b>26.71</b>
<b>5</b>	<b>Baseline MTFA</b>	5.95	8.67	6.94	34.71	<b>6.25</b>	4.85	21.29	24.42	36.86	<b>14.51</b>	24.83
	<b>MTFA + Triplet</b>	8.48	13.43	9.80	36.66	5.75	8.04	23.83	26.66	37.03	11.62	25.91
	<b>MTFA + Memory</b>	<b>9.61</b>	<b>14.61</b>	<b>11.73</b>	<b>38.60</b>	5.79	<b>10.40</b>	<b>26.65</b>	<b>30.37</b>	<b>39.21</b>	12.26	<b>30.02</b>
<b>Object Detection</b>												
<b>1</b>	<b>Baseline MTFA</b>	2.93	5.86	2.20	20.95	4.18	2.03	9.25	10.84	21.74	<b>11.49</b>	8.77
	<b>MTFA + Triplet</b>	4.04	8.65	2.98	20.50	4.90	<b>4.22</b>	12.89	<b>15.53</b>	20.73	11.45	17.46
	<b>MTFA + Memory</b>	<b>4.50</b>	<b>9.14</b>	<b>3.45</b>	<b>22.88</b>	<b>5.61</b>	3.54	<b>13.14</b>	15.22	<b>23.14</b>	8.78	<b>16.33</b>
<b>2</b>	<b>Baseline MTFA</b>	5.90	8.87	6.83	33.04	9.74	3.10	17.26	19.25	34.04	<b>15.74</b>	19.61
	<b>MTFA + Triplet</b>	<b>7.28</b>	<b>11.22</b>	<b>8.25</b>	32.31	<b>10.72</b>	<b>6.83</b>	<b>20.52</b>	<b>22.69</b>	32.34	14.88	<b>23.52</b>
	<b>MTFA + Memory</b>	6.95	10.88	7.75	<b>33.93</b>	7.49	6.81	19.84	22.01	<b>34.10</b>	15.04	21.47
<b>3</b>	<b>Baseline MTFA</b>	5.84	8.98	6.29	34.56	7.78	4.31	19.13	21.83	35.80	15.93	21.09
	<b>MTFA + Triplet</b>	7.49	<b>11.51</b>	8.23	<b>38.45</b>	8.61	<b>6.38</b>	<b>24.88</b>	<b>27.52</b>	<b>38.55</b>	17.66	27.44
	<b>MTFA + Memory</b>	<b>7.55</b>	11.45	<b>8.50</b>	38.07	<b>9.21</b>	5.70	24.20	27.29	38.50	<b>18.10</b>	<b>27.56</b>
<b>5</b>	<b>Baseline MTFA</b>	5.84	9.13	6.04	35.44	8.17	4.22	19.67	22.96	35.94	<b>14.16</b>	22.58
	<b>MTFA + Triplet</b>	9.76	14.37	11.12	<b>40.05</b>	<b>8.82</b>	9.89	25.93	29.28	<b>40.05</b>	12.53	30.32
	<b>MTFA + Memory</b>	<b>10.36</b>	<b>16.27</b>	<b>11.79</b>	39.32	8.08	<b>11.36</b>	<b>26.34</b>	<b>30.30</b>	39.35	12.37	<b>30.91</b>

80 R-101 as their base model to apply our proposed methods. The details of this decision are declared in the ablation section. As can be inferred from Table 4, with the improvement of our methods, such models increased the performance in AP. In terms of instance segmentation, we improved MTFA (Ganea et al., 2021), Mask RCNN<sup>†</sup>, IFS-RCNN (Nguyen and Todorovic, 2022) on average percentages of 0.73, 0.51, and 0.18, respectively thanks to instance triplet loss and 1.85, 0.43, and

0.08, respectively thanks to instance memory storage. Regarding object detection, we got average amounts of 2.01, 1.83, and −0.16, respectively with instance triplet loss and 2.21, 1.74, and −0.16, respectively with instance memory storage. Meanwhile, the performance of IFS-RCNN (Nguyen and Todorovic, 2022) decreases in some cases by increasing the training data. A possible reason is that the probit classifier of IFS-RCNN (Nguyen and Todorovic, 2022) is based on the distribu-

Table 6: Ablation study on the base model with 1-shot results. The best, and second best performances are marked in **red**, and **blue**, respectively. “Memory” is Instance Memory Storage and “Triplet” is Instance Triplet Loss.

Method	Base Model	Segmentation			Detection		
		AP	AP50	AP75	AP	AP50	AP75
<b>Triplet</b>	COCO-80 R-101	<b>4.46</b>	<b>8.21</b>	<b>4.60</b>	<b>4.04</b>	<b>8.65</b>	<b>2.98</b>
	COCO-80 R-50	3.68	<b>6.79</b>	3.81	2.85	<b>6.67</b>	1.65
	COCO-60 R-101	<b>3.87</b>	6.26	<b>3.90</b>	<b>3.37</b>	6.51	<b>2.69</b>
	COCO-60 R-50	2.56	4.25	2.79	2.28	4.13	2.26
<b>Memory</b>	COCO-80 R-101	<b>5.46</b>	<b>9.20</b>	<b>6.17</b>	<b>4.50</b>	<b>9.14</b>	<b>3.45</b>
	COCO-80 R-50	<b>3.87</b>	<b>6.81</b>	<b>3.91</b>	<b>3.40</b>	<b>6.94</b>	2.76
	COCO-60 R-101	2.89	4.50	3.26	2.76	4.66	<b>2.81</b>
	COCO-60 R-50	2.63	4.50	3.02	2.25	4.50	1.65

Table 7: Ablation study on the margin and the  $\alpha$  ratio of the instance triplet loss in 1-shot settings. The best, and second best performances are marked in **red**, and **blue**, respectively.

AP $\alpha$	Segmentation					Detection				
	0	0.25	0.50	0.75	1.00	0	0.25	0.50	0.75	1.00
1	3.89	4.50	3.92	<b>5.16</b>	4.43	3.34	3.65	3.22	<b>4.22</b>	<b>3.68</b>
$1e^{-1}$	<b>4.82</b>	<b>4.74</b>	4.46	<b>4.58</b>	<b>4.57</b>	<b>4.36</b>	<b>4.27</b>	<b>4.04</b>	<b>4.16</b>	<b>3.79</b>
$1e^{-2}$	<b>4.29</b>	<b>4.74</b>	<b>4.69</b>	4.46	<b>4.39</b>	<b>4.02</b>	<b>3.97</b>	<b>4.24</b>	4.06	3.71

tion of base classes. Meanwhile, the novel classes distribution of our CAMO-FS has large gaps with base classes in the COCO dataset. Therefore, the performance of IFS-RCNN (Nguyen and Todorovic, 2022) is unstable in the context of scarce training data on camouflaged instances.

In **Table 5**, we also present the results of the baseline MTFA (Ganea et al., 2021) with its original default configuration along with our proposed improvements. On top of the baseline MTFA (Ganea et al., 2021), we establish fine-tuning configuration on this model by training all heads of classification, box regression, and mask prediction on few-shot novel data. The reported results prove the performance of the proposed instance triplet loss and instance memory storage.

In general, our approaches achieve outstanding results in comparison with the baseline. Our improvements surpass MTFA by a remarkable margin. Regarding the instance seg-

Table 8: Ablation study on the capacity of the instance memory storage in 1-shot settings. The best, and second best performances are marked in **red**, and **blue**, respectively.

Capacity	Segmentation			Detection		
	AP	AP50	AP75	AP	AP50	AP75
32	4.56	7.30	5.02	3.85	7.72	2.91
64	4.51	<b>7.67</b>	4.49	3.94	<b>8.37</b>	2.84
128	4.53	7.55	4.87	4.13	7.98	3.62
256	4.56	7.50	5.02	4.01	8.22	3.39
512	<b>4.76</b>	7.57	<b>5.37</b>	<b>4.48</b>	8.25	<b>4.44</b>
1024	<b>4.72</b>	<b>8.06</b>	<b>5.20</b>	<b>4.14</b>	<b>8.45</b>	<b>3.79</b>

mentation, our method improves 1.9%, 3.5%, and 2.3% in terms of AP, AP@50, and AP@75, respectively. These results manifest the efficiency of our methods in the context of few-shot camouflaged instances. Both loss functions enhance the



discrimination between foreground and background features which strongly supports the model to segment pixels that belong to the camouflaged animals. Regarding the memory bank and the triplet loss function, the results of the memory loss function are higher than those of the triplet loss function by about 1%. We realize that storing representatives for each class is a crucial element in few-shot learning. This technique not only expands the variants during training but also increases the consistency per class, thereby model can segment difficult objects better. In these ways, we also improve the corresponding results in camouflage object detection.

In Table 5, our improvements help the model segment animals in various sizes. Specifically, all three metrics including APs, APm, and API improve in comparison with the baseline model, which demonstrates that our model well segments small, medium, and large animals. This situation also happens in the detection problem. When data is very scarce as in a 1-shot or 2-shot setting, the triplet loss function has comparative results with the memory loss function. However, in the context of 3-shot or 5-shot settings, the memory loss function demonstrates outstanding efficiency thanks to storing and updating the memory via iterations to create discriminative features on a global level. Figure 9 illustrates the qualitative comparison among the results of 5-shot settings of the baseline MTFA (Ganea et al., 2021) and our proposed methods of Instance Triplet Loss and Instance Memory Storage. We chose to visualize the images with the confidence threshold of 0.5, which released a huge number of predictions with low confidence from the models. The two final rows are exemplary cases that either triplet loss or memory storage cannot well handle these camouflaged instances.

**Base model ablation study.** We also conduct ablation experiments on different backbone base models of the COCO settings including general and few-shot concepts. To be detailed, we report the performance of our proposed method of instance triplet loss and instance memory storage over four different backbones. The considered backbones are ResNet-50 and ResNet-101 (He et al., 2016). The two base datasets are MS-

Table 9: Ablation study on the  $\beta$  ratio of the instance memory loss (Eq. 4) in 1-shot settings. The best, and second best performances are marked in **red**, and **blue**, respectively.

$\beta$	Segmentation			Detection		
	AP	AP50	AP75	AP	AP50	AP75
$1e^{-1}$	3.36	6.58	2.91	3.69	8.02	2.90
$1e^{-2}$	<b>4.57</b>	<b>8.02</b>	<b>4.74</b>	3.73	7.78	2.76
$1e^{-3}$	4.51	7.15	4.67	3.87	7.16	3.51
$1e^{-4}$	<b>5.12</b>	<b>8.71</b>	<b>5.54</b>	<b>4.58</b>	<b>9.23</b>	<b>3.69</b>
$1e^{-5}$	4.44	7.58	3.89	<b>4.06</b>	<b>7.99</b>	<b>3.63</b>

COCO with 80 classes and 60 classes, respectively. Thus, it led to the combination of four different base models (i.e. COCO-80 R-101, COCO-80 R-50, COCO-60 R-101, and COCO-60 R-50). As can be seen from Table 6, the performance of applying COCO-80 R-101 base weight yields better results among others evaluated on AP, AP@50, and AP@75 in both segmentation and detection tasks. In both cases of our two proposed improvements, the ablation results demonstrate our selection of COCO-80 R-101 is the best among the tested backbones of the base phase. For the segmentation task, we achieve 4.46 and 5.46 of AP reported for triplet loss and memory storage, respectively. For the detection task, we reach 4.04 and 4.50 also of AP for the two proposals, respectively. In summary, the chosen backbone of the base weight presents a higher performance of around 1% to 2% of evaluated on common metrics as in the table. To be explained, the base from COCO-80 contains more semantic concepts in comparison with the COCO-60 base, which leads to the higher performance reported. Note that all the results in this ablation section are reported for the 1-shot setting.

**Ablation on instance triplet loss component.** In terms of the instance triplet loss described in Eq. 1, we establish ablation experiments to evaluate the performance of the model with different configurations of margin and  $\alpha$  value. To this end, we set up the margin varying from 0 to 1, with a step of 0.25. For the  $\alpha$  ratio of the loss function (Eq. 2), we check out  $\alpha = \{1, 1e^{-1}, 1e^{-2}\}$ . To be enhanced, the margin value indicates how



Fig. 9: Qualitative comparison among the baseline MTFA (Ganea et al., 2021) and our proposed methods. The results are from 5-shot settings. “Memory” is Instance Memory Storage and “Triplet” is Instance Triplet Loss.



distinguished foreground and background features are. Meanwhile, the  $\alpha$  controls the effect of the instance triplet loss on the total loss of the framework. Table 7 presents the evaluation of both detection and segmentation issues in 1-shot manner. As can be inferred from the table, the effect of  $\alpha$  decides which margin should be selected for the triplet loss. With  $\alpha = 1$  meaning we keep the original ratio of the loss, the segmentation result in 1-shot setting yields the highest performance of 5.16% mAP with a 0.75 margin value. Meanwhile, the detection result gets the highest performance of 4.36% with  $\alpha = 1e^{-1}$  and zero margin. This table offers a better understanding of the impact of  $\alpha$  and the margin over the total performance.

**Ablation on instance memory storage component.** As for the instance memory storage, as introduced in Eq. 3, and Eq. 4, there are several parameters that need analyzing, listed as the amount of capacity in the memory storage and the  $\beta$  ratio controlling the effect of the memory storage loss in the total loss. Table 8, and Table 9 present the ablation experimental results of those issues, respectively. In terms of the capacity of the memory storage, we establish experiments on a range of memory capacity  $= 2^i$  where  $i = \{5, 6, 7, 8, 9, 10\}$ . The reported results figure out that the performance on both segmentation and detection tasks increases with a larger capacity of memory storage. To be detailed, with a capacity of 512, the mAP metric achieves the highest value among configurations, i.e. 4.76% and 4.48% for segmentation and detection, respectively. Empirically, we select 512 to be the suitable capacity of the memory storage, not the largest. To this end, the larger capacity can confuse the model in the process of learning when retrieving information in such a large memory bank. Besides, Table 9 expresses the effectiveness of the memory loss to the total loss function. As can be inferred,  $\beta = 1e^{-4}$  gives the best performance evaluated on mAP, AP50, and AP75 among all configurations.

## 5. Conclusion

In this work, we investigated the interesting yet challenging problem of few-shot learning for camouflaged animal detection and segmentation. We first collect a new dataset, CAMO-FS,

for benchmarking purpose. We then propose a novel method to efficiently detect and segment the camouflaged animals in the images. In particular, we introduce the instance triplet loss and the instance memory storage. The extensive experiments demonstrated that our proposed method achieves state-of-the-art performance on the newly constructed dataset. We expect our work will encourage more research work in this field. In the future, we would like to extend our work with more shots for new classes. In addition, we aim to improve the computational model by taking the context into consideration.

## Acknowledgments

This research was supported by the VNUHCM-University of Information Technology’s Scientific Research Support Fund.

## References

- Balntas, V., Riba, E., Ponsa, D., Mikolajczyk, K., 2016. Learning local feature descriptors with triplets and shallow convolutional neural networks., in: Bmvc, p. 3.
- Boykov, Y., Funka-Lea, G., 2006. Graph cuts and efficient nd image segmentation. IJCV 70, 109–131.
- Dong, N., Xing, E.P., 2018. Few-shot semantic segmentation with prototype learning., in: BMVC.
- Fan, D.P., Ji, G.P., Sun, G., Cheng, M.M., Shen, J., Shao, L., 2020a. Camouflaged object detection, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 2777–2787.
- Fan, Q., Zhuo, W., Tang, C.K., Tai, Y.W., 2020b. Few-shot object detection with attention-rpn and multi-relation detector, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR).
- Fan, Z., Yu, J.G., Liang, Z., Ou, J., Gao, C., Xia, G.S., Li, Y., 2020c. Fgn: Fully guided network for few-shot instance segmentation, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 9172–9181.
- Gallego, J., Bertolino, P., 2014. Foreground object segmentation for moving camera sequences based on foreground-background probabilistic models and prior probability maps, in: ICIP, pp. 3312–3316.
- Galun, M., Sharon, E., Basri, R., Brandt, A., 2003a. Texture segmentation by multiscale aggregation of filter responses and shape elements, in: ICCV, pp. 716–723.
- Galun, M., Sharon, E., Basri, R., Brandt, A., 2003b. Texture segmentation by multiscale aggregation of filter responses and shape elements., in: ICCV, p. 716.

- Ganea, D.A., Boom, B., Poppe, R., 2021. Incremental few-shot instance segmentation, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1185–1194.
- Gidaris, S., Komodakis, N., 2018. Dynamic few-shot visual learning without forgetting, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4367–4375.
- Han, G., He, Y., Huang, S., Ma, J., Chang, S.F., 2021. Query adaptive few-shot object detection with heterogeneous graph convolutional networks, in: *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 3263–3272.
- Han, G., Huang, S., Ma, J., He, Y., Chang, S.F., 2022a. Meta faster r-cnn: Towards accurate few-shot object detection with attentive feature alignment, in: *Proceedings of the AAAI Conference on Artificial Intelligence*, pp. 780–789.
- Han, G., Ma, J., Huang, S., Chen, L., Chang, S.F., 2022b. Few-shot object detection with fully cross-transformer, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5321–5330.
- He, K., Gkioxari, G., Dollár, P., Girshick, R., 2017. Mask r-cnn, in: *ICCV*, pp. 2980–2988.
- He, K., Zhang, X., Ren, S., Sun, J., 2016. Deep residual learning for image recognition, in: *CVPR*.
- Hou, J.Y.Y.H.W., Li, J., 2011. Detection of the mobile object with camouflage color under dynamic background based on optical flow. *Procedia Engineering* 15, 2201–2205.
- Hu, H., Bai, S., Li, A., Cui, J., Wang, L., 2021. Dense relation distillation with context-aware aggregation for few-shot object detection, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Kang, B., Liu, Z., Wang, X., Yu, F., Feng, J., Darrell, T., 2019. Few-shot object detection via feature reweighting, in: *ICCV*.
- Kavitha, C., Rao, B.P., Govardhan, A., 2011. An efficient content based image retrieval using color and texture of image sub-blocks. *International Journal of Engineering Science and Technology (IJEST)* 3, 1060–1068.
- Kervrann, C., Heitz, F., 1995. A markov random field model-based approach to unsupervised texture segmentation using local and global spatial statistics. *IEEE TIP* 4, 856–862.
- Khandelwal, S., Goyal, R., Sigal, L., 2021. Unit: Unified knowledge transfer for any-shot object detection and segmentation, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Lamdouar, H., Yang, C., Xie, W., Zisserman, A., 2020. Betrayed by motion: Camouflaged object discovery via motion segmentation, in: *ACCV*.
- Le, T.N., Cao, Y., Nguyen, T.C., Le, M.Q., Nguyen, K.D., Do, T.T., Tran, M.T., Nguyen, T.V., 2022. Camouflaged instance segmentation in-the-wild: Dataset, method, and benchmark suite. *IEEE Transactions on Image Processing* 31, 287–300.
- Le, T.N., Nguyen, H.H., Yamagishi, J., Echizen, I., 2021a. Openforensics: Large-scale challenging dataset for multi-face forgery detection and segmentation in-the-wild, in: *ICCV*.
- Le, T.N., Nguyen, T.V., Nie, Z., Tran, M.T., Sugimoto, A., 2019. Anabran network for camouflaged object segmentation. *CVIU* 184, 45–56.
- Le, T.N., Nguyen, V., Le, C., Nguyen, T.C., Tran, M.T., Nguyen, T.V., 2021b. Camouflander: Finding camouflaged instances in images, in: *Proceedings of the AAAI Conference on Artificial Intelligence*, pp. 16071–16074.
- Li, A., Li, Z., 2021. Transformation invariant few-shot object detection, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Li, A., Zhang, J., Lv, Y., Liu, B., Zhang, T., Dai, Y., 2021a. Uncertainty-aware joint salient object and camouflaged object detection, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10071–10081.
- Li, B., Yang, B., Liu, C., Liu, F., Ji, R., Ye, Q., 2021b. Beyond max-margin: Class margin equilibrium for few-shot object detection, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Li, X., Sahbi, H., 2011. Superpixel-based object class segmentation using conditional random fields, in: *ICASSP*, pp. 1101–1104.
- Lin, T.Y., Dollár, P., Girshick, R., He, K., Hariharan, B., Belongie, S., 2017. Feature pyramid networks for object detection, in: *CVPR*.
- Liu, W., Zhang, C., Lin, G., Liu, F., 2020. Crnet: Cross-reference networks for few-shot segmentation, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Liu, Z., Huang, K., Tan, T., 2012. Foreground object detection using top-down information based on em framework. *IEEE TIP* 21, 4204–4217.
- Lv, Y., Zhang, J., Dai, Y., Li, A., Liu, B., Barnes, N., Fan, D.P., 2021. Simultaneously localize, segment and rank the camouflaged objects, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11591–11601.
- Mei, H., Ji, G.P., Wei, Z., Yang, X., Wei, X., Fan, D.P., 2021. Camouflaged object segmentation with distraction mining, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8772–8781.
- Nguyen, K., Todorovic, S., 2022. ifs-rcnn: An incremental few-shot instance segmenter, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7010–7019.
- P. Sengottuvelan, A.W., Shanmugam, A., 2008. Performance of decamouflaging through exploratory image analysis, in: *ICETET*, pp. 6–10.
- Pan, Y., Chen, Y., Fu, Q., Zhang, P., Xu, X., 2011. Study on the camouflaged target detection method based on 3d convexity. *Modern Applied Science* 5, 152.
- Pia Bideau, E.L.M., 2016. It's moving! a probabilistic model for causal motion segmentation in moving camera videos, in: *ECCV*.
- Price, N., Green, S., Troscianko, J., Tregenza, T., Stevens, M., 2019. Background matching and disruptive coloration as habitat-specific strategies for camouflage. *Scientific reports* 9, 1–10.
- Qiao, L., Zhao, Y., Li, Z., Qiu, X., Wu, J., Zhang, C., 2021. Defrcn: Decoupled faster r-cnn for few-shot object detection, in: *Proceedings of the IEEE/CVF*

- International Conference on Computer Vision, pp. 8681–8690.
- Redmon, J., Farhadi, A., 2017. Yolo9000: better, faster, stronger, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 7263–7271.
- Ren, S., He, K., Girshick, R., Sun, J., 2015. Faster r-cnn: Towards real-time object detection with region proposal networks, in: *NeurIPS*, pp. 91–99.
- Saha, O., Cheng, Z., Maji, S., 2022. Ganorcon: Are generative models useful for few-shot segmentation?, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 9991–10000.
- Singh, S., Dhawale, C., Misra, S., 2013. Survey of object detection methods in camouflaged image. *IERI Procedia* 4, 351 – 357.
- Siricharoen, P., Aramvith, S., Chalidabhongse, T., Siddhichai, S., 2010. Robust outdoor human segmentation based on color-based statistical approach and edge combination, in: *The 2010 International Conference on Green Circuits and Systems, IEEE*. pp. 463–468.
- Skurowski, P., Abdulameer, H., Baszczyk, J., Depta, T., Kornacki, A., Kozie, P., 2018. Animal camouflage analysis: Chameleon database. Unpublished Manuscript .
- Snell, J., Swersky, K., Zemel, R., 2017. Prototypical networks for few-shot learning. *Advances in neural information processing systems* 30.
- Song, L., Geng, W., 2010a. A new camouflage texture evaluation method based on wssim and nature image features, in: *International Conference on Multimedia Technology*, pp. 1–4.
- Song, L., Geng, W., 2010b. A new camouflage texture evaluation method based on wssim and nature image features, in: *2010 International Conference on Multimedia Technology, IEEE*. pp. 1–4.
- Sulimowicz, L., Ahmad, I., Aved, A., 2018. Superpixel-enhanced pairwise conditional random field for semantic segmentation, in: *ICIP*, pp. 271–275.
- Sun, B., Li, B., Cai, S., Yuan, Y., Zhang, C., 2021. Fsce: Few-shot object detection via contrastive proposal encoding, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7352–7362.
- Tian, Z., Lai, X., Jiang, L., Liu, S., Shu, M., Zhao, H., Jia, J., 2022. Generalized few-shot semantic segmentation, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 11563–11572.
- Vinyals, O., Blundell, C., Lillicrap, T., Wierstra, D., et al., 2016. Matching networks for one shot learning. *Advances in neural information processing systems* 29.
- Wang, K., Liew, J.H., Zou, Y., Zhou, D., Feng, J., 2019. Panet: Few-shot image semantic segmentation with prototype alignment, in: *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*.
- Wang, X., Huang, T.E., Darrell, T., Gonzalez, J.E., Yu, F., 2020. Frustratingly simple few-shot object detection, in: *ICML*.
- Wu, J., Liu, S., Huang, D., Wang, Y., 2020. Multi-scale positive sample refinement for few-shot object detection, in: *ECCV*.
- Wu, Y., Kirillov, A., Massa, F., Lo, W.Y., Girshick, R., 2019. Detectron2. <https://github.com/facebookresearch/detectron2>.
- Wu, Z., Xiong, Y., Yu, S.X., Lin, D., 2018. Unsupervised feature learning via non-parametric instance discrimination, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3733–3742.
- Xiao, Y., Marlet, R., 2020. Few-shot object detection and viewpoint estimation for objects in the wild, in: *European conference on computer vision*, Springer. pp. 192–210.
- Xue, F., Yong, C., Xu, S., Dong, H., Luo, Y., Jia, W., 2016a. Camouflage performance analysis and evaluation framework based on features fusion. *Multimedia Tools and Applications* 75, 4065–4082.
- Xue, F., Yong, C., Xu, S., Dong, H., Luo, Y., Jia, W., 2016b. Camouflage performance analysis and evaluation framework based on features fusion. *Multimedia Tools and Applications* 75, 4065–4082.
- Yan, J., Le, T.N., Nguyen, K.D., Tran, M.T., Do, T.T., Nguyen, T.V., 2021. Mirrornet: Bio-inspired camouflaged object segmentation. *IEEE Access* 9, 43290–43300.
- Yan, X., Chen, Z., Xu, A., Wang, X., Liang, X., Lin, L., 2019. Meta r-cnn: Towards general solver for instance-level low-shot learning, in: *ICCV*.
- Yin, J., Han, Y., Hou, W., Li, J., 2011. Detection of the mobile object with camouflage color under dynamic background based on optical flow. *Procedia Engineering* 15, 2201 – 2205.
- Zhai, Q., Li, X., Yang, F., Chen, C., Cheng, H., Fan, D.P., 2021. Mutual graph learning for camouflaged object detection, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12997–13007.
- Zhang, G., Luo, Z., Cui, K., Lu, S., 2021. Meta-detr: Image-level few-shot object detection with inter-class correlation exploitation. *arXiv preprint arXiv:2103.11731* .
- Zhang, S., Wang, L., Murray, N., Koniusz, P., 2022. Kernelized few-shot object detection with efficient integral aggregation, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 19207–19216.
- Zhang, W., Wang, Y.X., 2021. Hallucination improves few-shot object detection, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 13008–13017.
- Zhu, C., Chen, F., Ahmed, U., Shen, Z., Savvides, M., 2021a. Semantic relation reasoning for shot-stable few-shot object detection, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Zhu, J., Zhang, X., Zhang, S., Liu, J., 2021b. Inferring camouflage objects by texture-aware interactive guidance network, in: *AAAI*.