# Advancing Incremental Few-shot Semantic Segmentation via Semantic-guided Relation Alignment and Adaptation

Yuan Zhou
Hefei University of Technology
Hefei, China
2018110971@mail.hfut.edu.cn

Xin Chen
Huawei Inc.
Shenzhen, China
chenxin061@gmail.com

Yanrong Guo
Hefei University of Technology
Hefei, China
yrguo@hfut.edu.cn

Shijie Hao
Hefei University of Technology
Hefei, China
hfut.hsj@gmail.com

Richang Hong
Hefei University of Technology
Hefei, China
hongrc.hfut@gmail.com

Qi Tian
Huawei Inc.
Shenzheng, China
tian.qi1@huawei.com

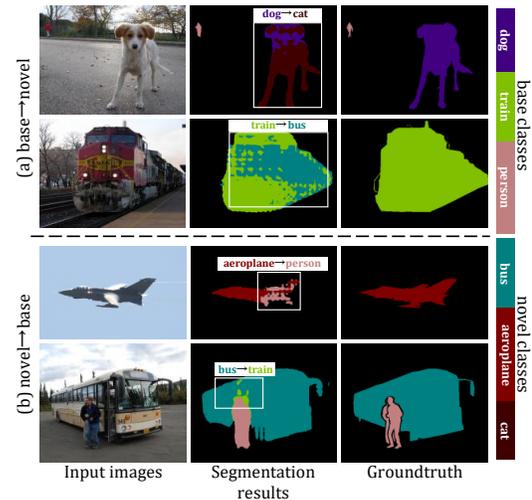arXiv:2305.10868v1 [cs.CV] 18 May 2023

## ABSTRACT

Incremental few-shot semantic segmentation (IFSS) aims to incrementally extend a semantic segmentation model to novel classes according to only a few pixel-level annotated data, while preserving its segmentation capability on previously learned base categories. This task faces a severe semantic-aliasing issue between base and novel classes due to data imbalance, which makes segmentation results unsatisfactory. To alleviate this issue, we propose the Semantic-guided Relation Alignment and Adaptation (SRAA) method that fully considers the guidance of prior semantic information. Specifically, we first conduct Semantic Relation Alignment (SRA) in the base step, so as to semantically align base class representations to their semantics. As a result, the embeddings of base classes are constrained to have relatively low semantic correlations to categories that are different from them. Afterwards, based on the semantically aligned base categories, Semantic-Guided Adaptation (SGA) is employed during the incremental learning stage. It aims to ensure affinities between visual and semantic embeddings of encountered novel categories, thereby making the feature representations be consistent with their semantic information. In this way, the semantic-aliasing issue can be suppressed. We evaluate our model on the PASCAL VOC 2012 and the COCO dataset. The experimental results on both these two datasets exhibit its competitive performance, which demonstrates the superiority of our method.

Figure 1: The typical examples for the semantic-aliasing issue in IFSS, which are obtained from the PASCAL VOC 2012 dataset on the 1-shot task. "A→B" indicates that regions belonging to A are incorrectly segmented as B.

## 1 INTRODUCTION

In recent years, semantic segmentation has achieved impressive performance by using deep neural networks. However, conventional semantic segmentation models generally have a fixed output space. Therefore, when encountering new categories, they need to be re-trained from scratch to update their segmentation capability. Moreover, these models require large-scale pixel-level labeled data, which are expensive and laborious to obtain. These issues limit their applicability in open-ended real-world scenarios. In this context, Cermelli et al. [4] proposed the Incremental Few-shot Semantic Segmentation (IFSS) task. It aims to effectively extend a semantic segmentation model to new categories using a few labeled samples, while maintaining its segmentation capability on previously learned old ones. In this way, the extendibility and flexibility of the model can be improved, which is critical for many real-world applications, such as autonomous driving and human-machine interaction.

More specifically, in the IFSS task, a base set with relatively more training samples is first provided to initialize the learnable parameters of a semantic segmentation model. Then, a few pixel-level annotated training samples of novel categories are given, helping incrementally expand the segmentation capability of the model to the encountered novel ones. However, the IFSS model is prone to fall into the semantic-aliasing issue due to data imbalance between base and novel classes. As shown in Figure 1, the semantic confusion between the base class "dog" and the encountered novel category "cat" misleads the model to draw the incorrect segmentation results, making the model performance unsatisfactory. Recently, semantic information has been successfully introduced in the few-shot classification task [15, 37–39], aiming to make feature embeddings more representative, e.g., GloVe [23] or word2vec [21] was employed in [15, 39] to provide prior semantic information while [37, 38] additionally considered the semantic guidance of CLIP [24].

Inspired by these methods, we propose to suppress the semantic-aliasing issue in IFSS by fully considering the guidance of visual semantics. Therefore, we propose the Semantic-guided Relation Alignment and Adaptation (SRAA) method in this paper, which is shown in Figure 2. On one hand, we propose to conduct Semantic Relation Alignment (SRA) in the base step, aiming to semantically align base class representations in latent semantic space. Therefore, the embeddings of base classes are constrained to have relatively low semantic correlations to categories that are different from them. Moreover, the cross-entropy loss is employed during this process to measure discrepancy between segmentation results and groundtruth label maps. As a result, the model is trained to segment base classes, while being aware of their semantic information. Based on the aligned base classes, Semantic-Guided Adaptation (SGA) is employed to incrementally adapt the model to novel classes. It aims to ensure affinities between visual and semantic embeddings of novel categories, thereby making the feature representations be consistent with their semantic information. By considering the semantic information of both the base and the novel classes, the semantic-aliasing issue can be alleviated. We evaluate our method on the public semantic segmentation datasets PASCAL VOC 2012 and COCO, following the cross validation used in [4]. On both these datasets, our method presents competitive performance.

All-in-all, the contributions of this paper can be summarized below:

- In this paper, we propose to suppress the semantic-aliasing issue in IFSS by fully considering the guidance of semantic information, thereby making segmentation results more accurate. To realize this goal, we accordingly propose the Semantic-guided Relation Alignment and Adaptation (SRAA) method.
- We propose to conduct Semantic Relation Alignment (SRA) in the base step, aiming to semantically align the representations of base categories. Therefore, the base class embeddings are guided to have relatively low semantic correlations to categories that are different from them.
- Based on the aligned base classes, we propose to conduct Semantic-Guided Adaptation (SGA) during the incremental learning stage, guiding the embeddings of novel classes to be consistent with their semantic information. In this way,

the semantic aliasing between the base and the encountered novel categories can be alleviated.

## 2 RELATED WORK

In this section, we review methods that are relevant to our research. We first briefly introduce typical methods of semantic segmentation, few-shot learning, and incremental learning in section 2.1, section 2.2, and section 2.3. Then, we review related incremental few-shot semantic segmentation methods in section 2.4, and introduce their differences to our work.

### 2.1 Semantic Segmentation

Semantic segmentation, a pixel-level image recognition technique, has achieved remarkable progress in recent years with development of deep learning. [19] is a typical deep-learning-based semantic segmentation method that uses the fully convolutional layer to realize efficient end-to-end dense predications for input images. Inspired by [19], many semantic segmentation models have been proposed. Zhao et al. [41] further introduced the pyramid pooling module, aiming to fully aggregate global context information of visual scenes. Chen et al. [5, 6] proposed to aggregate multi-scale context information using the atrous convolution, so as to make segmentation results more accurate. Based on [5, 6, 19], Zhang et al. [40] learned an inherent dictionary to aggregate semantic context information of a whole dataset, which help the model understand visual scenes in a more global way. Huang et al. [13] enhanced a semantic segmentation model with the proposed criss-cross attention layer. Therefore, sufficient context information is aggregated for each pixel, while the model is maintained with high efficiency. Recently, the methods [30, 36, 42] have successfully built semantic segmentation models upon the transformer [7, 32], thereby further boosting visual representations of input images.

### 2.2 Few-shot Learning

Few-shot learning aims to quickly transfer models to novel unseen categories according to only one or a few training instances, thereby reducing expenses cost on data preparation. Currently, few-shot learning methods are mainly based on metric learning, aiming at learning an effective metric classifier from given few-shot training instances. For example, Vinyals et al. [33] proposed a matching network that classifies query samples by measuring instance-wise consine similarity between queries and supports. Snell et al. [29] advanced the matching network by further introducing prototypical representations, thereby constructing global category representations for support samples. Santoro et al. [26] proposed a memory-augmented neural network that utilizes stored memory to make query categorization more accurate. Li et al. [15] and Zhang et al. [39] proposed to additionally consider semantic attributes encoded by GloVe[23] or word2vec [21], so as to further improve visual representations of input images. Besides, Xu et al. [37] and Yang et al. [38] proposed to further exploit the semantic guidance from CLIP [24], as they found semantic information encoded by CLIP is more effective in learning representative feature embeddings of visual scenes.

## 2.3 Incremental Learning

Incremental learning aims to effectively transfer a model to new categories, while maintaining its previously learned old knowledge as much as possible. Knowledge distillation [12] has shown its advantages in overcoming a catastrophic forgetting problem [16]. Aiming to incorporate the knowledge distillation with the data-replay strategy, Rebuffi et al. [25] introduced the exemplar-based knowledge distillation at the cost of extra small storage expenses. Castro et al. [2] and Kang et al. [14] pointed out that it is necessary to achieve a good balance between old class knowledge maintenance and new class adaptation. Therefore, the cross-distillation loss and the balanced finetune strategy were utilized in [2], while [14] employed the adaptive feature consolidation strategy to restrict the representation drift of critical old class feature embeddings. Recently, Wang et al. [34] advanced the incremental learning model with the gradient boosting strategy, so as to guide the model to effectively learn its residuals to the target one. Liu et al. [18] further enhanced the data-replay strategy by designing the reinforced memory management mechanism. It dynamically adjusts the stored memory information in each incremental step, thereby helping to overcome the sub-optimal memory allocation problem.

## 2.4 Incremental Few-shot Semantic Segmentation

Incremental Few-shot Semantic Segmentation (IFSS), proposed by Cermelli et al. [4], aims at enduing semantic segmentation models with the capability of few-shot incremental learning, thereby making them more suitable to be deployed in open-ended real-world applications. Aiming to address this task, Cermelli et al. [4] proposed the prototype-based knowledge distillation. It relieves the catastrophic forgetting issue by constraining the invariance of old class segmentation scores. Moreover, the overfitting to novel categories is suppressed by boosting the consistency between old and updated models. Shi et al. [27] proposed to build hyper-class feature representations, thereby helping to relieve the representation drift during the incremental learning. Furthermore, they adopted a different evaluation protocol than the one employed in [4]. Despite the success achieved by these methods, the guidance of visual semantics is ignored in them, which has been proven to play an important role in low-data tasks. Therefore, different from these works, in this paper, we study on how to exploit prior semantic information in IFSS to make segmentation results more accurate.

## 3 METHODOLOGY

We elaborate on our proposed method in this section. We first give the preliminaries in Section 3.1. Then, the details about our Semantic-guided Relation Alignment and Adaptation (SRAA) are provided in Section 3.2.

## 3.1 Preliminaries

The semantic space of the IFSS model is expanded over time. We define $C^t$ as the categories encountered at the step $t$. Therefore, after learning in this step, the semantic space of the model is expanded to $S^t = S^{t-1} \bigcup C^t$, where $S^{t-1} = \bigcup_{i=0}^{t-1} C^i$ denotes the semantic space learned after the step $t-1$. In each step, the dataset $D^t =$

$\{X_n^t, Y_n^t\}_{n=1}^{N_t}$ is provided to update learnable parameters, in which $X_n^t$ denotes the $n$-$th$ training image and $Y_n^t$ is the label map of $X_n^t$. In the IFSS task, the base dataset $D^0$ is provided in the base step (*i.e.*, $t = 0$) to initialize the parameters of the model, which contains relatively more training samples. After the base step, the dataset $D^t$ is only in the few-shot setting, i.e., each category contains one or a few labeled training instances, which satisfies the condition $N_t \ll N_0$ for $\forall t >= 1$. For brevity, in this paper, the categories given in the base step are called base categories, while the categories encountered in the incremental learning stage are termed novel categories. In the step $t$, the model only has access to the dataset $D^t$, and the datasets in the previous steps are unavailable.

## 3.2 Semantic-guided Relation Alignment and Adaptation

As described in Figure 2, our method consists of the two components that are Semantic Relation Alignment (SRA) and Semantic-Guided Adaptation (SGA). These two components are incorporated together to help the model be aware of semantic information of base and novel categories, thereby alleviating the semantic-aliasing issue between them. We elaborate on these two components in the following.

*3.2.1 Semantic Relation Alignment.* The goal of our SRA is to semantically align base classes in latent semantic space and guide the model to generate semantic-consistent visual representations. We first extract the visual embeddings $\{F_n^0\}_{n=1}^{N_b}$ from the input images $\{X_n^0\}_{n=1}^{N_b} \subset D^0$ using the visual encoder $f_v(\cdot|\Theta_v^0)$,

$$\{F_n^0\}_{n=1}^{N_b} = f_v(\{X_n^0\}_{n=1}^{N_b}|\Theta_v^0) \tag{1}$$

where $\Theta_v^0$ indicates the learnable parameters of the visual encoder in the base step, and $N_b$ denotes the number of images in a mini-batch. Meanwhile, the semantic encoder $f_s(\cdot|\Theta_s^0)$ encodes the semantic vectors $\{s_c^0\}_{c=1}^{|C_b|}$ of the categories $C_b$ involved in the inputs,

$$\{s_c^0\}_{c=1}^{|C_b|} = f_s(C_b|\Theta_s^0) \tag{2}$$
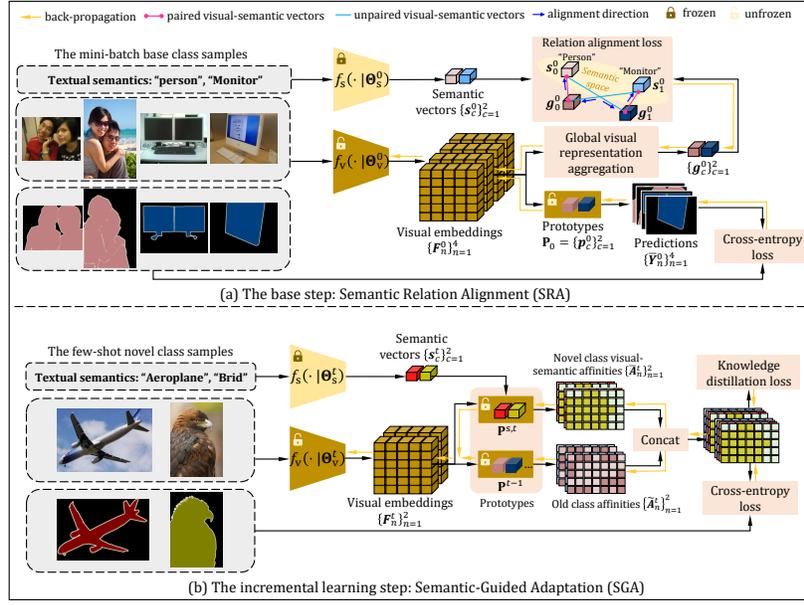
in which $\Theta_s^0$ represents the parameters of the semantic encoder, and $s_c^0$ denotes the encoded semantic vector about the class $c \in C_b$. After that, the global visual representations $\{g_c^0\}_{c=1}^{|C_b|}$ are aggregated for each of the categories $C_b$,

$$g_c^0 = \frac{\sum_{n=1}^{N_b} \sum_{i=1}^H \sum_{j=1}^W (F_{n,[i,j]}^0 * I_{n,[i,j]}^{0,c})}{\sum_{n=1}^{N_b} \sum_{i=1}^H \sum_{j=1}^W I_{n,[i,j]}^{0,c}} \tag{3}$$

$$s.t. \quad I_n^{0,c} = Y_n^0 == c$$

where $F_n^0 \in \mathbb{R}^{H \times W \times D}$ denotes the feature map encoded by the visual encoder $f_v(\cdot|\Theta_v^0)$, and $I_n^{0,c} \in \mathbb{R}^{H \times W}$ indicates the binary mask about the category $c$. Of note, if the pixel at the position $[i, j]$ of $X_n^0$ belongs to the category $c$, $I_{n,[i,j]}^{0,c} = 1$; otherwise, $I_{n,[i,j]}^{0,c} = 0$.

Aiming to align base class features with their semantics, the relation alignment loss $\mathcal{L}_{align}$ is employed. It jointly considers the correlations between visual and semantic embeddings that are paired

Figure 2: The illustration for our Semantic-guided Relation Alignment and Adaptation (SRAA) method. SRA aims to semantically align the representations of base classes in latent semantic space, and SGA aims at ensuring the visual-semantic affinities of encountered novel categories. Also, "$f_s(\cdot|\cdot)$" and "$f_v(\cdot|\cdot)$" represent a semantic encoder and a visual encoder respectively.

and unpaired,

$$\mathcal{L}_{\text{align}} = \underbrace{\sum_{c_1=1}^{|\text{C}_\text{b}|} \sum_{c_2=1, c_2 \neq c_1}^{|\text{C}_\text{b}|} \frac{\boldsymbol{g}_{c_1}^0 * \boldsymbol{s}_{c_2}^0}{|\text{C}_\text{b}| \times (|\text{C}_\text{b}| - 1)}}_{\text{Unpaired}} - \underbrace{\sum_{c=1}^{|\text{C}_\text{b}|} \frac{\boldsymbol{g}_c^0 * \boldsymbol{s}_c^0}{|\text{C}_\text{b}|}}_{\text{Paired}}. \quad (4)$$

The paired visual-semantic embeddings indicate that the visual vector $\boldsymbol{g}_c^0$ and the semantic vector $\boldsymbol{s}_c^0$ belong to the same class, and thus $\boldsymbol{g}_c^0$ should be aligned to match $\boldsymbol{s}_c^0$ in latent space. On the contrary, if the visual embeddings $\boldsymbol{g}_{c_1}^0$ and the semantic embeddings $\boldsymbol{s}_{c_2}^0$ ($c_2 \neq c_1$) are unpaired, the correlations between them should be suppressed to ensure representation discrimination between categories. We minimize the relation alignment loss $\mathcal{L}_{\text{align}}$ w.r.t. the learnable parameters of the visual encoder. Thereby, the model is guided to generate semantic-consistent visual representations, and the base classes embeddings are aligned with their semantic information.

Moreover, the segmentation results $\left\{ \bar{Y}_n^0 \in \mathbb{R}^{H \times W \times |\text{C}^0|} \right\}_{n=1}^{N_\text{b}}$ with respect to the semantic space $\text{C}^0$ are drawn by using the base class prototypical classifier $\text{P}^0 = \left\{ \boldsymbol{p}_c^0 \right\}_{c=1}^{|\text{C}^0|}$, which is shown below:

$$\bar{Y}_{n,[i,j,c]}^0 = P(c|X_{n,[i,j]}^0, \text{P}^0, \Theta_\text{v}^0) \quad (5)$$

$$= \frac{\exp(Sim(F_{n,[i,j]}^0, \boldsymbol{p}_c^0))}{\sum_{c' \in \text{C}^0} \exp(Sim(F_{n,[i,j]}^0, \boldsymbol{p}_{c'}^0))}.$$

In the above equation, $P(c|X_{n,[i,j]}^0, \text{P}^0, \Theta_\text{v}^0)$ indicates the probability that the pixel $X_{n,[i,j]}^0$ is inferred as the category $c$ according to $\text{P}^0$ and $\Theta_\text{v}^0$. $Sim(\cdot, \cdot)$ is a similarity metric function, which aims to measures consine similarity between feature embeddings. The

cross-entropy loss $\mathcal{L}_{\text{ce}}$ is used to measure the discrepancy between the segmentation results and the groundtruth labels of the inputs,

$$\mathcal{L}_{\text{ce}} = \frac{1}{N_\text{b}} \sum_{n=1}^{N_\text{b}} CE(\bar{Y}_n^0, Y_n^0). \quad (6)$$

By jointly minimizing the relation alignment loss $\mathcal{L}_{\text{align}}$ and the cross-entropy loss $\mathcal{L}_{\text{ce}}$ during the training process, the model learns to segment base categories while being aware of their semantic information.

*3.2.2 Semantic-Guided Adaptation.* After the base step, the model is incrementally extended to novel classes. We hope the model can also be aware of the semantic information of encountered novel categories. Therefore, we propose to ensure affinities between visual and semantic embeddings of encountered novel ones. Taking the step $t$ as an example, we first extract the visual embeddings from the images given in the few-shot dataset $\text{D}^t$ using the visual encoder $f_v(\cdot|\Theta_\text{v}^t)$,

$$\left\{ F_n^t \right\}_{n=1}^{N_t} = f_v(\left\{ X_n^t \right\}_{n=1}^{N_t}|\Theta_\text{v}^t). \quad (7)$$

In the equation, $\Theta_\text{v}^t$ indicates the parameters of the visual encoder in the step $t$. Meanwhile, the semantic encoder $f_s(\cdot|\Theta_\text{s}^t)$ encodes the semantic embeddings of the encountered new categories $\text{C}^t$,

$$\left\{ \boldsymbol{s}_c^t \right\}_{c=1}^{|\text{C}^t|} = f_s(\text{C}^t|\Theta_\text{s}^t). \quad (8)$$

Afterwards, these semantic vectors are used to imprint the weights of the semantic prototypes $\text{P}^{s,t} = \left\{ \boldsymbol{p}_c^{s,t} \right\}_{c=1}^{|\text{C}^t|}$, which are used to guide the finetune process on the novel categories. Specifically, we first calculate the affinities $\left\{ \bar{A}_n^t \in \mathbb{R}^{H \times W \times |\text{C}^t|} \right\}_{n=1}^{N_t}$ between the visual embeddings of the given images and the semantics of the

novel classes according to Eq. 9,

$$\bar{A}_{n,[i,j,c]}^t = \frac{F_{n,[i,j]}^t * \boldsymbol{p}_c^{s,t}}{|F_{n,[i,j]}^t| * |\boldsymbol{p}_c^{s,t}|}, \ s.t., 0 < c <= |C^t| \quad (9)$$

where $\bar{A}_n^t$ denotes the dense visual-semantic affinities about the sample $X_n^t$, and $\bar{A}_{n,[i,j,c]}^t$ indicates the affinity between the visual features at the position $[i, j]$ and the semantic embeddings of the category $c \in C^t$. The dense visual-semantic affinities reflect the relation between the visual embeddings and the semantics of the encountered novel classes.

Moreover, the prototypes $\mathbf{P}^{t-1} = \{\boldsymbol{p}_i^{t-1}\}_{i=1}^{|\cup_{j=0}^{t-1} C^j|}$ learned in the previous steps are utilized to compute the affinities of the current feature maps to the old categories $\{\tilde{A}_n^t \in \mathbb{R}^{H \times W \times |\cup_{j=0}^{t-1} C^j|}\}_{n=1}^{N_t}$,

$$\tilde{A}_{n,[i,j,c]}^t = \frac{F_{n,[i,j]}^t * \boldsymbol{p}_c^{t-1}}{|F_{n,[i,j]}^t| * |\boldsymbol{p}_c^{t-1}|}, s.t., 0 < c <= |\bigcup_{j=0}^{t-1} C^j|. \quad (10)$$

The affinity maps $\tilde{A}_n^t$ and $\bar{A}_n^t$ are concatenated together for each sample

$$A_n^t = \tilde{A}_n^t \oplus \bar{A}_n^t, \quad (11)$$

thereby producing $\{A_n^t \in \mathbb{R}^{H \times W \times |\cup_{j=0}^t C^j|}\}_{n=1}^{N_t}$. We use the cross entropy to constrain the correctness of the affinity maps $\{A_n^t\}_{n=1}^{N_t}$

$$\mathcal{L}_{\text{aff}} = \frac{1}{N_t} \sum_{n=1}^{N_t} CE(A_n^t, Y_n^t), \quad (12)$$

so as to guide the visual embeddings of the novel class images to have high correlations to their visual semantics while suppressing the affinities to the old classes. As a result, the feature embeddings of the novel classes can be consistent with their semantic information. Moreover, knowledge distillation is adopted to suppress the overfitting to encountered novel categories

$$\mathcal{L}_{\text{kd}} = \frac{1}{N_t} \sum_{n=1}^{N_t} CE(A_n^t, A_n^{t-1}), \quad (13)$$

where $A_n^{t-1}$ denotes the affinities drawn by the model after being trained in the step $t-1$. The joint minimization of $\mathcal{L}_{\text{aff}}$ and $\mathcal{L}_{\text{kd}}$ w.r.t. $\Theta_s^t$, $\mathbf{P}^{t-1}$, and $\mathbf{P}^{s,t}$ guides the model to be aware of the visual semantics of the encountered novel categories. Meanwhile, the prototypes $\mathbf{P}^{t-1}$ and $\mathbf{P}^{s,t}$ are optimized to be consistent to reflect the relationships between images and categories, so as to help accurately segment out the objects that belong to the encountered categories from images. After learning in the step $t$, the prototypes are updates: $\mathbf{P}^t \leftarrow \hat{\mathbf{P}}^{t-1} \bigcup \hat{\mathbf{P}}^{s,t}$, where $\hat{\mathbf{P}}^{t-1}$ and $\hat{\mathbf{P}}^{s,t}$ indicate the prototypes $\mathbf{P}^{t-1}$ and $\mathbf{P}^{s,t}$ after being optimized in the current step. The updated prototypes and visual encoder are employed to draw segmentation results for all the encountered classes, as same as the process shown in Eq. 5.

## 4 EXPERIMENTS

Experiments are provided in this section to validate our proposed method. We first introduce the datasets in Section 4.1 and give our implementation details in Section 4.2. Then, we report the main experimental results in Section 4.3, and conduct the ablation study in Section 4.4.

**Table 1: The dataset split on the PASCAL VOC 2012 dataset.**

| Split | Categories |
|---|---|
| 5-0 | aeroplane, bicycle, bird, boat, bottle |
| 5-1 | bus, car, cat, chair, cow |
| 5-2 | table, dog, horse, motorbike, person |
| 5-3 | plant, sheep, sofa, train, tv-monitor |

**Table 2: The dataset split on the COCO dataset.**

| Split | Categories |
|---|---|
| 20-0 | person, airplane, boat, parking meter, dog, elephant, refrigerator, backpack, suitcase, sports ball, skateboard, wine glass, spoon, sandwich, hot dog, chair, dining table, mouse, microwave, scissors |
| 20-1 | bicycle, bus, traffic light, bench, horse, bear, umbrella, frisbee, kite, surfboard, cup, bowl, orange, pizza,couch, toilet, remote, oven, book, teddy bear |
| 20-2 | car, train, fire hydrant, bird, sheep, zebra, handbag, skis, baseball bat, tennis racket, fork, banana, broccoli, donut, potted plant, tv, keyboard,toaster, clock, hair drier |
| 20-3 | motorcycle, truck, stop sign, cat, cow, giraffe, tie, snowboard, baseball glove, bottle, knife, apple, carrot, cake, bed, laptop, cell phone, sink, vase, toothbrush |

### 4.1 Datasets

We evaluate the proposed method on the two public semantic segmentation datasets that are PASCAL VOC 2012 [8, 10] and COCO [1, 17]. The PASCAL VOC 2012 dataset consists of 10582 training images and 1449 test images, which are collected from 20 different categories. Following the previous work [4], we divide these 20 categories into four folds and each fold includes five categories, which is summarized in Table 1. In addition, on the COCO dataset, 80 categories are used to evaluate the performance of the model, including about 110k training samples and 5k test samples. As presented in Table 2, the 80 categories of this dataset are split into four folds as well, which is the same as [4] does. For the cross validation on both these datasets, we use the categories of three folds to build the base set, while the categories of the rest one fold are used for testing.

### 4.2 Implementation Details

Our codes are implemented using PyTorch and run on the tesla V100 GPU card. In the experiments, the SGD optimizer is adopted to optimize the learnable parameters of our model based on the poly learning rate. For the experiments on the PASCAL VOC 2012 and the COCO dataset, the training details are slightly different. Specifically, on the PASCAL VOC 2012 dataset, we set the number of the epochs as 30 on the base step and 400 during the incremental learning. Also, in each phase, the initial learning rate of the optimizer is set as 0.01. On the COCO dataset, we train the model on the base set for 50 epochs with the initial poly learning rate 0.02. Moreover, during the incremental learning stage, the epochs are set as 400, and the initial learning rate is set as 0.01. Following the previous work [4], we evaluate our method in both the single few-shot step setting and the multiple few-shot step setting based on the cross validation protocol.

**Table 3: The experimental results on the PASCAL VOC 2012 dataset. In the table, "FT" indicates directly finetuning the model on novel classes using the cross-entropy loss, and "HM" indicates the harmonic mean of the mIoU on base and novel classes. Also, the first-place and the second-place result in each column are highlighted in bold font and underscore respectively.**

(a) The results under the single few-shot step setting.

| Method | 1-shot mIoU (%) | | | 2-shot mIoU (%) | | | 5-shot mIoU (%) | | |
|---|---|---|---|---|---|---|---|---|---|
| | base | novel | HM | base | novel | HM | base | novel | HM |
| FT | 58.3 | 9.7 | 16.7 | 59.1 | 19.7 | 29.5 | 55.8 | 29.6 | 38.7 |
| WI [22] | 62.7 | 15.5 | 24.8 | 63.3 | 19.2 | 29.5 | 63.3 | 21.7 | 32.3 |
| DWI [9] | 64.3 | 15.4 | 24.8 | **64.8** | 19.8 | 30.4 | 64.9 | 23.5 | 34.5 |
| RT [31] | 59.1 | 12.1 | 20.1 | 60.9 | 21.6 | 31.9 | 60.4 | 27.5 | 37.8 |
| AMP [28] | 57.5 | 16.7 | 25.8 | 54.4 | 18.8 | 27.9 | 51.9 | 18.9 | 27.7 |
| SPN [35] | 59.8 | 16.3 | 25.6 | 60.8 | 26.3 | 36.7 | 58.4 | 33.4 | 42.5 |
| LWF [16] | 61.5 | 10.7 | 18.2 | 63.6 | 18.9 | 29.2 | 59.7 | 30.9 | 40.8 |
| ILT [20] | 64.3 | 13.6 | 22.5 | 64.2 | 23.1 | 34.0 | 61.4 | 32.0 | 42.1 |
| MIB [3] | 61.0 | 5.2 | 9.7 | 63.5 | 12.7 | 21.1 | **65.0** | 28.1 | 39.3 |
| PIFS [4] | 60.9 | 18.6 | 28.4 | 60.5 | 26.4 | 36.8 | 60.0 | 33.4 | 42.8 |
| Ours | **65.2** | **19.1** | **29.5** | 62.7 | **27.4** | **38.1** | 63.8 | **36.7** | **46.6** |

(b) The results under the multiple few-shot step setting.

| Method | 1-shot mIoU (%) | | | 2-shot mIoU (%) | | | 5-shot mIoU (%) | | |
|---|---|---|---|---|---|---|---|---|---|
| | base | novel | HM | base | novel | HM | base | novel | HM |
| FT | 47.2 | 3.9 | 7.2 | 53.5 | 4.4 | 8.1 | 58.7 | 7.7 | 13.6 |
| WI [22] | 66.6 | 16.1 | 25.9 | 66.6 | 19.8 | 30.5 | 66.6 | 21.9 | 33.0 |
| DWI [9] | **67.2** | 16.3 | 26.2 | **67.5** | 21.6 | 32.7 | **67.6** | 25.4 | 36.9 |
| RT [31] | 49.2 | 5.8 | 10.4 | 36.0 | 4.9 | 8.6 | 45.1 | 10.0 | 16.4 |
| AMP [28] | 58.6 | 14.5 | 23.2 | 58.4 | 16.3 | 25.5 | 57.1 | 17.2 | 26.4 |
| SPN [35] | 49.8 | 8.1 | 13.9 | 56.4 | 10.4 | 17.6 | 61.6 | 16.3 | 25.8 |
| LWF [16] | 42.1 | 3.3 | 6.2 | 51.6 | 3.9 | 7.3 | 59.8 | 7.5 | 13.4 |
| ILT [20] | 43.7 | 3.3 | 6.1 | 52.2 | 4.4 | 8.1 | 59.0 | 7.9 | 13.9 |
| MIB [3] | 43.9 | 2.6 | 4.9 | 51.9 | 2.1 | 4.0 | 60.9 | 5.8 | 10.5 |
| PIFS [4] | 64.1 | 16.9 | 26.7 | 65.2 | 23.7 | 34.8 | 64.5 | 27.5 | 38.6 |
| Ours | 66.4 | **18.8** | **29.3** | 65.1 | **26.4** | **37.6** | 64.3 | **28.7** | **39.7** |

**Table 4: The experimental results on the COCO dataset.**

(a) The results under the single few-shot step setting.

| Method | 1-shot mIoU (%) | | | 2-shot mIoU (%) | | | 5-shot mIoU (%) | | |
|---|---|---|---|---|---|---|---|---|---|
| | base | novel | HM | base | novel | HM | base | novel | HM |
| FT | 41.2 | 4.1 | 7.5 | 41.5 | 7.3 | 12.4 | 41.6 | 12.3 | 19.0 |
| WI [22] | 43.8 | 6.9 | 11.9 | 44.2 | 7.9 | 13.5 | 43.6 | 8.7 | 14.6 |
| DWI [9] | 44.5 | 7.5 | 12.8 | 45.0 | 9.4 | 15.6 | 44.9 | 12.1 | 19.1 |
| RT [31] | **46.2** | 5.8 | 10.2 | **46.7** | 8.8 | 14.8 | **46.9** | 13.7 | 21.2 |
| AMP [28] | 37.5 | 7.4 | 12.4 | 35.7 | 8.8 | 14.2 | 34.6 | 11.0 | 16.7 |
| SPN [35] | 43.5 | 6.7 | 11.7 | 43.7 | 10.2 | 16.5 | 43.7 | 15.6 | 22.9 |
| LWF [16] | 43.9 | 3.8 | 7.0 | 44.3 | 7.1 | 12.3 | 44.6 | 12.9 | 20.1 |
| ILT [20] | **46.2** | 4.4 | 8.0 | 46.3 | 6.5 | 11.5 | 47.0 | 11.0 | 17.8 |
| MIB [3] | 43.8 | 3.5 | 6.5 | 44.4 | 6.0 | 10.6 | 44.7 | 11.9 | 18.8 |
| PIFS [4] | 40.8 | 8.2 | 13.7 | 40.9 | 11.1 | 17.5 | 42.8 | 15.7 | 23.0 |
| Ours | 41.2 | **9.3** | **15.2** | 42.1 | **12.7** | **19.5** | 42.6 | **17.1** | **24.4** |

(b) The results under the multiple few-shot step setting.

| Method | 1-shot mIoU (%) | | | 2-shot mIoU (%) | | | 5-shot mIoU (%) | | |
|---|---|---|---|---|---|---|---|---|---|
| | base | novel | HM | base | novel | HM | base | novel | HM |
| FT | 38.5 | 4.8 | 8.5 | 40.3 | 6.8 | 11.6 | 39.5 | 11.5 | 17.8 |
| WI [22] | **46.3** | 8.3 | 14.1 | 46.5 | 9.3 | 15.5 | 46.3 | 10.3 | 16.9 |
| DWI [9] | 46.2 | 9.2 | 15.3 | 46.5 | 11.4 | 18.3 | **46.6** | 14.5 | 22.1 |
| RT [31] | 38.4 | 5.2 | 9.2 | 43.8 | 10.1 | 16.4 | 44.1 | 16.0 | 23.5 |
| AMP [28] | 36.6 | 7.9 | 13.0 | 36.0 | 9.2 | 14.7 | 33.2 | 11.0 | 16.5 |
| SPN [35] | 40.3 | 8.7 | 14.3 | 41.7 | 12.5 | 19.2 | 41.4 | 18.2 | 25.3 |
| LWF [16] | 41.0 | 4.1 | 7.5 | 42.7 | 6.5 | 11.3 | 42.3 | 12.6 | 19.4 |
| ILT [20] | 43.7 | 6.2 | 10.9 | **47.1** | 10.0 | 16.5 | 45.3 | 15.3 | 22.9 |
| MIB [3] | 40.4 | 3.1 | 5.8 | 42.7 | 5.2 | 9.3 | 43.8 | 11.5 | 18.2 |
| PIFS [4] | 40.4 | 10.4 | 16.5 | 40.1 | 13.1 | 19.8 | 41.1 | 18.3 | 25.3 |
| Ours | 40.7 | **11.3** | **17.7** | 40.5 | **15.2** | **22.1** | 41.0 | **19.7** | **26.6** |

The single few-shot step setting indicates that all novel categories are given at once in an incremental step, while the multiple few-shot step setting means novel categories are progressively given in multiple steps. We employ the mean Intersection-over-Union (mIoU) metric in our experiments to evaluate the performance of the method. Besides, we build our semantic encoder using CLIP, due to its powerful capability in encoding semantic information [37, 38]. Following the previous methods [37, 38], we freeze the parameters of the semantic encoder during the training process, i.e., $\Theta_s^t = \Theta_s^0$ for $\forall t >= 1$. Meanwhile, as same as [4], we build our visual encoder by using resnet101 [11].

### 4.3 Main Results

The results of our method on the PASCAL VOC 2012 and the COCO dataset are summarized in Table 3 and Table 4 respectively. According to these results, we have the following observations. On the PASCAL VOC 2012 dataset, our method achieves higher mIoU on both base and novel categories than that of FT, RT, AMP, SPN, and PIFS under the single few-shot step setting. Despite the performance of MIB, ILT, LWF, DWI, and WI on base categories is better than that of ours in some cases, our method obviously achieves

higher mIoU on novel categories. For example, on the 2-shot task, the novel class mIoU of our method is 14.7%, 4.3%, 8.5%, 7.6%, and 8.2% higher than that of MIB, ILT, LWF, DWI, and WI respectively. In the meantime, our method shows its superiority on the PASCAL VOC 2012 dataset under the multiple few-shot step setting as well. For example, on all the 1-shot, the 2-shot, and the 5-shot task, the novel class mIoU of PIFS is lower than that of our proposed method. Similar results can also be found on the experiments of the COCO dataset. For example, in the single few-shot step setting, the performance of our method is better than that of PIFS and AMP on both base and novel categories. Although the base class mIoU of our method is lower than that of the several compared methods, it consistently shows higher mIoU on encountered novel categories, e.g., our method's novel class mIoU is 5.8%, 4.9%, 5.5%, 2.6%, 3.5%, 1.8%, and 2.4% higher than that of MIB, ILT, LWF, SPN, RT, DWI, and WI on the 1-shot task. Moreover, under the multiple few-shot step setting, the novel class mIoU of our proposed method is higher than that of all the compared methods in the table. For example, our method's novel class mIoU is 0.9%, 2.1%, and 1.4% higher than that of the second-place method PIFS on the 1-shot, the 2-shot, and the 5-shot task. In Figure 3, we give our step-by-step
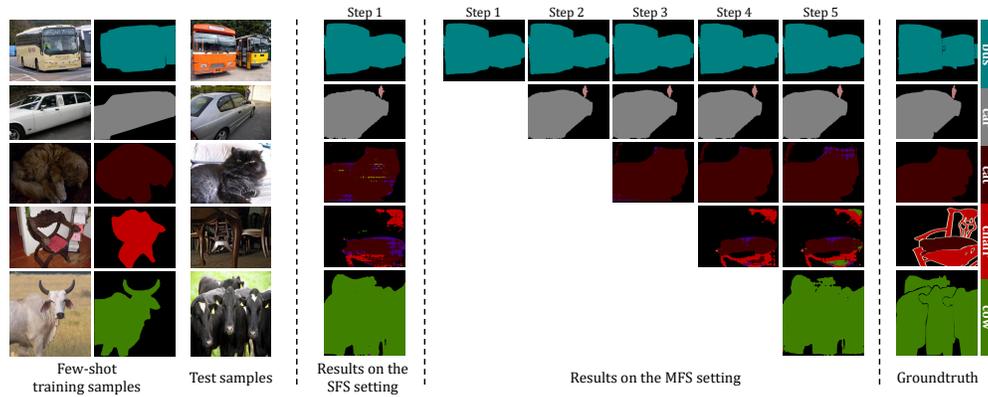
**Figure 3: The visualization for the step-by-step segmentation results of our method on both the Single Few-shot Step (SFS) setting and the Multiple Few-shot Step (MFS) setting according to a labeled sample per category. The model is progressively extended to the novel classes in the MFS setting. In the SFS setting, the novel classes are given at once in an incremental step.**
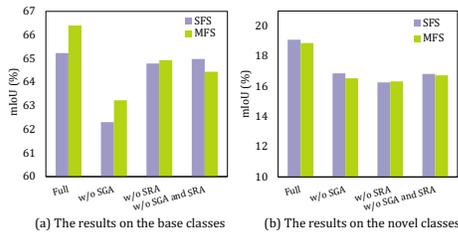


**Figure 4: The ablation study for our method on both the Single Few-shot Step (SFS) setting and the Multiple Few-shot Step (MFS) setting, which is conducted on the 1-shot task of the PASCAL VOC 2012 dataset.**
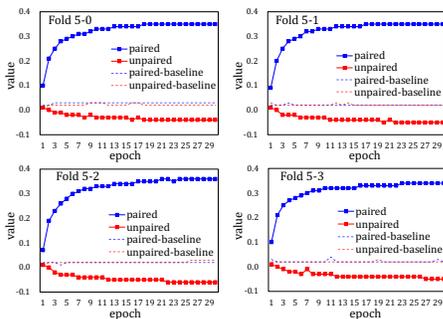


**Figure 5: The values of the paired and the unpaired term in our relation alignment loss $\mathcal{L}_{\text{align}}$ during the training process. To better reflect value change, the baseline curves are also drawn in the figure, which reflect the values of these two terms when $\mathcal{L}_{\text{align}}$ is not employed during the training. The above experiments are conducted on the PASCAL VOC 2012 dataset.**

segmentation results for the encountered novel categories under the two different settings. The results indicate that according to only one training instance per novel category, our method can still
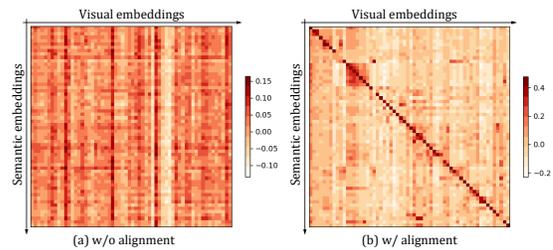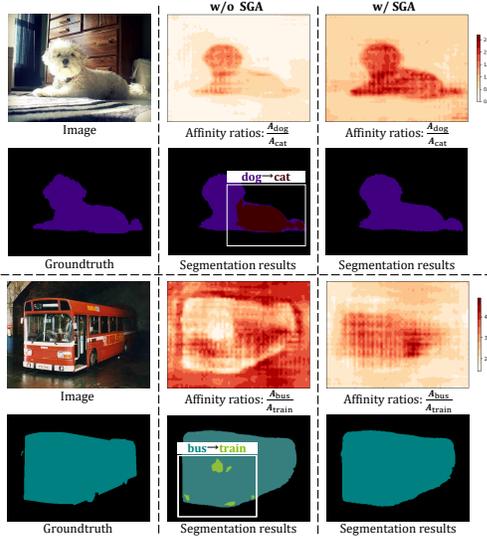


**Figure 6: The mean affinities between the visual and the semantic embeddings with or without being aligned by our SRA, which is conducted on the COCO dataset. The values in the diagonal are the affinities between the visual and the semantic embeddings that belong to the same class, while the others are the affinities between them that are unpaired.**

achieve the promising semantic segmentation results. Moreover, when encountering new classes, it can still maintain high effectiveness in segmenting the categories learned in the previous few-shot learning steps.

## 4.4 Ablation Study

In this subsection, we first study the influence of SGA and SRA on accuracy in Figure 4. "w/o SGA" indicates that the semantic guidance is not considered during the adaptation procedure. Thus, the prototypes about novel categories are imprinted by the mean visual embeddings of given samples. "w/o SRA" indicates that SRA is not employed, and base class embeddings are not aligned with their semantics. On one hand, the cooperation of our SRA and SGA (i.e., "Full") can achieve higher mIoU than that of the baseline model ("w/o SGA and SRA") on both base and novel classes under the two different settings, which demonstrates that the appropriate use of prior semantic information can make segmentation results more accurate. On the other hand, the removal of SGA or SRA (i.e., "w/o SGA" or "w/o SRA") leads to an obvious performance drop, thereby validating the importance of these two components. The results also indicate that semantic information should be considered in
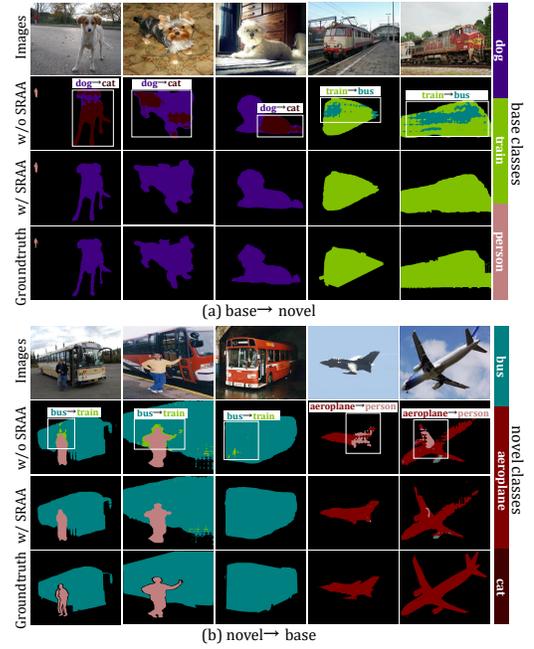
**Figure 7: The visualized analysis for the influence of our SGA. In the figure, $A_{\mathrm{dog}}$, $A_{\mathrm{cat}}$, $A_{\mathrm{bus}}$, or $A_{\mathrm{train}}$ denotes the affinity map of an image to the category "dog", "cat", "bus", or "train". $\frac{A_{\mathrm{dog}}}{A_{\mathrm{cat}}}$ or $\frac{A_{\mathrm{bus}}}{A_{\mathrm{train}}}$ indicates the ratio map between $A_{\mathrm{dog}}$ and $A_{\mathrm{cat}}$ or $A_{\mathrm{bus}}$ and $A_{\mathrm{train}}$**

both the base and the incremental learning stage. Otherwise, the training inconsistency between phases will reduce segmentation accuracy.

Then, in Figure 5, we visualize the values of the paired and the unpaired term in the relation alignment loss $\mathcal{L}_{\mathrm{align}}$ during the training process. The results indicate that $\mathcal{L}_{\mathrm{align}}$ can be optimized stably. On one hand, with the epoch increases, the paired term is maximized progressively, thereby constraining that visual and semantic embeddings belonging to the same category have relatively high correlations. On the other hand, the minimization of the unpaired term suppresses the similarity between visual and semantic embeddings that are unpaired. As a result, the visual embeddings belonging to the same class are guided to have high semantic correlations, while the semantic correlations of different categories are limited. We also visualize the mean affinities between the visual and the semantic embeddings that are aligned by our method. The results in Figure 6 validate the effectiveness of our SRA again, e.g., SRA can obviously rectify visual embeddings to better match their semantic information.

The analysis for the influence of our SGA is provided in Figure 7. As can be seen from the figure, without SGA, the sample about "dog" is incorrectly segmented as the class "cat" due to the incorrect affinity information, e.g., the affinity ratios $\frac{A_{\mathrm{dog}}}{A_{\mathrm{cat}}}$ have low values in the target regions. In contrast, the use of SGA can obviously boost the affinities to the target class, while suppressing the affinities to "cat". In this way, the segmentation results can be more accurate. For the instance "bus", the affinity ratios $\frac{A_{\mathrm{bus}}}{A_{\mathrm{train}}}$ show low values for a part of the target regions when not employing our SGA. Moreover, in the background areas, the affinity ratios have the incorrect high values. By leveraging the guidance of visual semantics, our method



**Figure 8: The visualized analysis for the influence of our method on alleviating the semantic-aliasing issue in both the "base→novel" and the "novel→base" scenario, which is conducted on the PASCAL VOC 2012 dataset under the 1-shot setting. Notice that "A→B" indicates regions belonging to A are incorrectly segmented as B.**

can rectify these incorrect affinities, thereby making segmentation results more accurate. Finally, in Figure 8, we also provide the qualitative analysis for the influence of our SRAA method on final segmentation results. The experimental results consistently validate the superiority of our method as well.

## 5 CONCLUSION

In this paper, we propose to alleviate the semantic-aliasing issue in IFSS by conducting Semantic-guided Relation Alignment and Adaptation (SRAA). On one hand, we introduce Semantic Relation Alignment (SRA) in the base step, aiming to semantically align the representations of base categories and guide the model to generate semantic-consistent feature representations. On the other hand, we employ Semantic-Guided Adaptation (SGA) to incrementally adapt the model to novel classes. It ensures the visual-semantic affinities of encountered novel categories, so as to make their feature embeddings be consistent with the corresponding semantic information. By considering the semantic information of both the base and the novel categories, the semantic-aliasing issue can be relieved. Currently, it is still very challenging to incrementally achieve accurate segmentation results for objects with complex and varied boundaries in IFSS. In the future, we plan to overcome this problem by fully considering the fine-grained information of local features.

# REFERENCES

[1] Holger Caesar, Jasper Uijlings, and Vittorio Ferrari. 2018. Coco-stuff: Thing and stuff classes in context. In *Proceedings of the Conference on Computer Vision and Pattern Recognition*. 1209–1218.

[2] Francisco M Castro, Manuel J Marín-Jiménez, Nicolás Guil, Cordelia Schmid, and Karteek Alahari. 2018. End-to-end incremental learning. In *Proceedings of the European Conference on Computer Vision*. 233–248.

[3] Fabio Cermelli, Massimiliano Mancini, Samuel Rota Bulo, Elisa Ricci, and Barbara Caputo. 2020. Modeling the background for incremental learning in semantic segmentation. In *Proceedings of the Conference on Computer Vision and Pattern Recognition*. 9233–9242.

[4] Fabio Cermelli, Massimiliano Mancini, Yongqin Xian, Zeynep Akata, and Barbara Caputo. 2021. Prototype-based Incremental Few-Shot Semantic Segmentation. In *Proceedings of the British Machine Vision Conference*. BMVC 2021, 484.

[5] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. 2017. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 40, 4 (2017), 834–848.

[6] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. 2017. Rethinking atrous convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587* (2017).

[7] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929* (2020).

[8] Mark Everingham, SM Ali Eslami, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. 2015. The pascal visual object classes challenge: A retrospective. *International Journal of Computer Vision* 111 (2015), 98–136.

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

[10] Bharath Hariharan, Pablo Arbeláez, Lubomir Bourdev, Subhransu Maji, and Jitendra Malik. 2011. Semantic contours from inverse detectors. In *Proceedings of the International Conference on Computer Vision*. IEEE, 991–998.

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

[12] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531* (2015).

[13] Zilong Huang, Xinggang Wang, Lichao Huang, Chang Huang, Yunchao Wei, and Wenyu Liu. 2019. Ccnet: Criss-cross attention for semantic segmentation. In *Proceedings of the International Conference on Computer Vision*. 603–612.

[14] Minsoo Kang, Jaeyoo Park, and Bohyung Han. 2022. Class-incremental learning by knowledge distillation with adaptive feature consolidation. In *Proceedings of the Conference on Computer Vision and Pattern Recognition*. 16071–16080.

[15] Aoxue Li, Weiran Huang, Xu Lan, Jiashi Feng, Zhenguo Li, and Liwei Wang. 2020. Boosting few-shot learning with adaptive margin loss. In *Proceedings of the Conference on Computer Vision and Pattern Recognition*. 12576–12584.

[16] Zhizhong Li and Derek Hoiem. 2017. Learning without forgetting. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 40, 12 (2017), 2935–2947.

[17] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *Proceedings of the European Conference on Computer Vision*. Springer, 740–755.

[18] Yaoyao Liu, Bernt Schiele, and Qianru Sun. 2021. RMM: Reinforced memory management for class-incremental learning. *Advances in Neural Information Processing Systems* 34 (2021), 3478–3490.

[19] Jonathan Long, Evan Shelhamer, and Trevor Darrell. 2015. Fully convolutional networks for semantic segmentation. In *Proceedings of the Conference on Computer Vision and Pattern Recognition*. 3431–3440.

[20] Umberto Michieli and Pietro Zanuttigh. 2019. Incremental learning techniques for semantic segmentation. In *Proceedings of the International Conference on Computer Vision Workshops*. 0–0.

[21] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781* (2013).

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

[23] Jeffrey Pennington, Richard Socher, and Christopher D Manning. 2014. Glove: Global vectors for word representation. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*. 1532–1543.

[24] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *Proceedings of the International Conference on Machine Learning*. PMLR, 8748–8763.

[25] Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. 2017. icarl: Incremental classifier and representation learning. In *Proceedings of the Conference on Computer Vision and Pattern Recognition*. 2001–2010.

[26] Adam Santoro, Sergey Bartunov, Matthew Botvinick, Daan Wierstra, and Timothy Lillicrap. 2016. Meta-learning with memory-augmented neural networks. In *Proceedings of the International Conference on Machine Learning*. PMLR, 1842–1850.

[27] Guangchen Shi, Yirui Wu, Jun Liu, Shaohua Wan, Wenhai Wang, and Tong Lu. 2022. Incremental few-shot semantic segmentation via embedding adaptive-update and hyper-class representation. In *Proceedings of the ACM International Conference on Multimedia*. 5547–5556.

[28] Mennatullah Siam, Boris Oreshkin, and Martin Jagersand. 2019. Adaptive masked proxies for few-shot segmentation. *arXiv preprint arXiv:1902.11123* (2019).

[29] Jake Snell, Kevin Swersky, and Richard Zemel. 2017. Prototypical networks for few-shot learning. *Advances in Neural Information Processing Systems* 30 (2017).

[30] Robin Strudel, Ricardo Garcia, Ivan Laptev, and Cordelia Schmid. 2021. Segmenter: Transformer for semantic segmentation. In *Proceedings of the International Conference on Computer Vision*. 7262–7272.

[31] Yonglong Tian, Yue Wang, Dilip Krishnan, Joshua B Tenenbaum, and Phillip Isola. 2020. Rethinking few-shot image classification: a good embedding is all you need?. In *Proceedings of the European Conference on Computer Vision*. Springer, 266–282.

[32] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in Neural Information Processing Systems* 30 (2017).

[33] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. 2016. Matching networks for one shot learning. *Advances in Neural Information Processing Systems* 29 (2016).

[34] Fu-Yun Wang, Da-Wei Zhou, Han-Jia Ye, and De-Chuan Zhan. 2022. Foster: Feature boosting and compression for class-incremental learning. In *Proceedings of the European Conference on Computer Vision*. Springer, 398–414.

[35] Yongqin Xian, Subhabrata Choudhury, Yang He, Bernt Schiele, and Zeynep Akata. 2019. Semantic projection network for zero-and few-label semantic segmentation. In *Proceedings of the Conference on Computer Vision and Pattern Recognition*. 8256–8265.

[36] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. 2021. SegFormer: Simple and efficient design for semantic segmentation with transformers. *Advances in Neural Information Processing Systems* 34 (2021), 12077–12090.

[37] Jingyi Xu and Hieu Le. 2022. Generating representative samples for few-shot classification. In *Proceedings of the Conference on Computer Vision and Pattern Recognition*. 9003–9013.

[38] Fengyuan Yang, Ruiping Wang, and Xilin Chen. 2023. Semantic Guided Latent Parts Embedding for Few-Shot Learning. In *Proceedings of the Winter Conference on Applications of Computer Vision*. 5447–5457.

[39] Baoquan Zhang, Xutao Li, Yunming Ye, Zhichao Huang, and Lisai Zhang. 2021. Prototype completion with primitive knowledge for few-shot learning. In *Proceedings of the Conference on Computer Vision and Pattern Recognition*. 3754–3762.

[40] Hang Zhang, Kristin Dana, Jianping Shi, Zhongyue Zhang, Xiaogang Wang, Ambrish Tyagi, and Amit Agrawal. 2018. Context encoding for semantic segmentation. In *Proceedings of the Conference on Computer Vision and Pattern Recognition*. 7151–7160.

[41] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. 2017. Pyramid scene parsing network. In *Proceedings of the Conference on Computer Vision and Pattern Recognition*. 2881–2890.

[42] Sixiao Zheng, Jiachen Lu, Hengshuang Zhao, Xiatian Zhu, Zekun Luo, Yabiao Wang, Yanwei Fu, Jianfeng Feng, Tao Xiang, Philip HS Torr, et al. 2021. Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers. In *Proceedings of the Conference on Computer Vision and Pattern Recognition*. 6881–6890.