# Self-Support Few-Shot Semantic Segmentation

Qi Fan[1], Wenjie Pei[2†], Yu-Wing Tai[1,3], and Chi-Keung Tang[1]

[1] HKUST, [2] Harbin Institute of Technology, Shenzhen, [3] Kuaishou Technology
`fanqics@gmail.com, wenjiecoder@outlook.com, yuwing@gmail.com,`
`cktang@cs.ust.hk`

**Abstract.** Existing few-shot segmentation methods have achieved great progress based on the support-query matching framework. But they still heavily suffer from the limited coverage of intra-class variations from the few-shot supports provided. Motivated by the simple Gestalt principle that pixels belonging to the same object are more similar than those to different objects of same class, we propose a novel self-support matching strategy to alleviate this problem, which uses query prototypes to match query features, where the query prototypes are collected from high-confidence query predictions. This strategy can effectively capture the consistent underlying characteristics of the query objects, and thus fittingly match query features. We also propose an adaptive self-support background prototype generation module and self-support loss to further facilitate the self-support matching procedure. Our self-support network substantially improves the prototype quality, benefits more improvement from stronger backbones and more supports, and achieves SOTA on multiple datasets. Codes are at https://github.com/fanq15/SSP.

**Keywords:** few-shot semantic segmentation, self-support prototype (SSP), self-support matching, adaptive background prototype generation.

## 1 Introduction

Semantic segmentation has achieved remarkable advances tapping into deep learning networks [29,42,34] and large-scale datasets such as [15,5,93]. However, current high-performing semantic segmentation methods rely heavily on laborious pixel-level annotations, which has expedited the recent development of few-shot semantic segmentation (FSS).

Few-shot semantic segmentation aims to segment arbitrary novel classes using only a few support samples. The dilemma is that the support images are limited and fixed (usually $\{1, 3, 5, 10\}$ supports per class), while the query images can be massive and arbitrary. Limited few-shot supports can easily fail to cover underlying appearance variations of the target class in query images, regardless of the support quality. This is clearly caused by the inherent data scarcity and diversity, two long standing issues in few-shot learning.
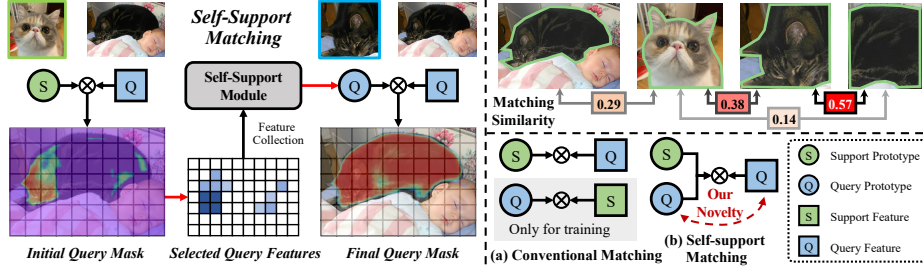
[†]Corresponding author.

**Fig. 1.** The *left* image illustrates the core idea of our self-support matching. We use the initial query mask prediction to collect query features in high-confidence regions and then use the generated **query prototype** to perform **self-matching** with **query features**. The *right top* image illustrates the motivation of our self-support matching: pixels/regions of the same objects are more similar than those from different objects. The numbers in boxes represent the cosine similarities between two objects. The *right bottom* image illustrates that our self-support matching is fundamentally distinct from conventional matching methods.

Existing methods try to solve the problem by making full use of the limited supports, such as proposing better matching mechanism [70,54,83,74,87,32,96] or generating representative prototypes [71,45,64,56,82,38,62,81,88,27]. Despite their success, they still cannot fundamentally solve the appearance discrepancy problem, bounded by the scarce few-shot supports.

We propose a novel self-support matching strategy to narrow the matching appearance discrepancy. This strategy uses query prototypes to match query features, or in other words, use the query feature to self-support itself. We thus call the query prototype as *self-support prototype* because of its self-matching property. This new idea is motivated by the classical Gestalt law [41] that pixels belonging to the same object are more similar than those to different objects.

Refer to Figure 1 for a high-level understanding of our novel self-support matching. First we generate the initial mask predictions by directly matching the support prototype and query features. Based on the initial query mask, we collect confident query features to generate the self-support prototype, which is used to perform matching with query features. Our *self-support module (SSM)* collects confident features of the cat head which are used to segment the entire black cat. Our model is optimized on base classes to retrieve other object parts supported by object fragments, *i.e.*, self-support prototype.

We apply our self-support module on both foreground and background prototypes for self-support matching. While SSM directly benefits foreground prototypes, note that the background is usually cluttered, which does not have the global semantic commonality shared among all background pixels. Thus, rather than generating a global background prototype by aggregating all the background pixels, we propose to adaptively generate self-support background prototypes for each query pixel, by dynamically aggregating similar background pixels in the query image. The *adaptive self-support background prototype (ASBP)*

is motivated by the fact that separate background regions have local semantic similarity. Finally, we propose a *self-support loss (SSL)* to further facilitate the self-support procedure.

Our self-support matching strategy is thus fundamentally different than conventional support-query matching. We use the flexible self-support prototypes to match query features, which can effectively capture the consistent underlying characteristics of the query objects, and thus fittingly match query features. As shown in Figure 1, the cats in the query and support images are very different in color, parts and scales, The garfield cat support has large appearance discrepancy to the black cat query, and undoubtedly conventional support-query matching produces inferior segmentation. In our self-support matching, our self-support prototype (the black cat head) is more consistent to the query (the entire black cat), and thus our method produces satisfactory results.

We are the first to perform self-support matching between query prototype and query features. As shown in Figure 1, our self-support matching fundamentally differs from conventional matching. Other methods learn better support prototypes for support-query matching from extra unlabeled images (PPNet [57] and MLC [82]) or builds various support prototype generation modules [71,45,81] or feature priors (PFENet [73]) based on support images. Although PANet [75] and CRNet [55] also explore query prototypes, they use *query prototypes* to match *support features* as a query-support matching only for auxiliary training, and cannot solve the appearance discrepancy.

Our self-support method significantly improves the prototype quality by alleviating the intra-class appearance discrepancy problem, evidenced by the performance boost on multiple datasets in our experimental validation. Despite the simple idea, our self-support method is very effective and has various advantages, such as benefiting more from stronger backbone and more supports, producing high-confidence predictions, more robustness to weak support labels, higher generalization to other methods and higher running efficiency. We will substantiate these advantages with thorough experiments. In summary, our contributions are:

- We propose novel self-support matching and build a novel self-support network to solve the appearance discrepancy problem in FSS.
- We propose self-support prototype, adaptive self-support background prototype and self-support loss to facilitate our self-support method.
- Our self-support method benefits more improvement from stronger backbones and more supports, and outperforms previous SOTAs on multiple datasets with many desirable advantages.

## 2 Related Works

**Semantic Segmentation.** Semantic segmentation is a fundamental computer vision task to produce pixel-wise dense semantic predictions. The state-of-the-art has recently been greatly advanced by the end-to-end fully convolutional network (FCN) [58]. Subsequent works have since followed this FCN paradigm

and contributed many effective modules to further promote the performance, such as encoder-decoder architectures [4,13,11,67], image and feature pyramid modules [39,91,10,49,50], context aggregation modules [25,26,33,92,95,36,85,89] and advance convolution layers [84,9,14,63]. Nevertheless, the above segmentation methods rely heavily on abundant pixel-level annotations. This paper aims to tackle the semantic segmentation problem in the few-shot scenario.

**Few-Shot Learning.** Few-shot learning targets at recognizing new concepts from very few samples. This low cost property has attracted a lot of research interests over the last years. There are three main approaches. The first is the transfer-learning approach [12,28,16,65] by adapting the prior knowledge learned from base classes to novel classes in a two-stage finetuning procedure. The second is the optimized-based approach [24,6,44,30,43,2,31,68], which rapidly updates models through meta-learning the optimization procedures from a few samples. The last is the metric-based approach [1,17,35,40,46,47], which applies a siamese network [40] on support-query pairs to learn a general metric for evaluating their relevance. Our work, including many few-shot works [23,37,90,80,22] on various high-level computer vision tasks, are inspired by the metric-based approach.

**Few-Shot Semantic Segmentation.** Few-shot semantic segmentation is pioneered by Shaban *et al.* [69]. Later works have mainly adopted the metric-based mainstream paradigm [18] with various improvements, *e.g.*, improving the matching procedure between support-query images with various attention mechanisms [70,54,83], better optimizations [94,53], memory modules [77,79], graph neural networks [78,74,87], learning-based classifiers [72,59], progressive matching [32,96], or other advanced techniques [86,52,61,48].

We are the first to perform self-support matching between query prototype and query features. Our self-support matching method is also related to the prototype generation methods. Some methods leverage extra unlabeled data [82,57] or feature priors [73] for further feature enhancement. Other methods generate representative support prototypes with various techniques, *e.g.*, attention mechanism [88,27], adaptive prototype learning [71,45,64], or various prototype generation approaches [62,81]. Although the query prototype has been explored in some methods [75,55], they only use query prototypes to match support features for prototype regularization. Finally, existing methods heavily suffer from the intra-class discrepancy problem in the support-query matching. On the other hand, we propose a novel self-support matching strategy to effectively address this matching problem.

## 3   Self-Support Few-Shot Semantic Segmentation

Given only a few support images, few-shot semantic segmentation aims to segment objects of novel classes using the model generalized from base classes. Existing mainstream few-shot semantic segmentation solution can be formulated as follows: The input support and query images $\{I_s, I_q\}$ are processed by a weight-shared backbone to extract image features $\{\mathcal{F}_s, \mathcal{F}_q\} \in \mathbb{R}^{C \times H \times W}$, where $C$ is the channel size and $H \times W$ is the feature spatial size. Then the support

**Table 1.** Cosine similarity for cross/intra object pixels.

| FG Pixels Similarity | | BG Pixels Similarity | |
|---|---|---|---|
| cross-object | intra-object | cross-image | intra-image |
| 0.308 | $0.416_{\uparrow 0.108}$ | 0.298 | $0.365_{\uparrow 0.067}$ |

feature $\mathcal{F}_s$ and its groundtruth mask $\mathcal{M}_s$ are fed into the masked average pooling layer to generate the support prototype vectors $\mathcal{P}_s = \{\mathcal{P}_{s,f}, \mathcal{P}_{s,b}\} \in \mathbb{R}^{C \times 1 \times 1}$ for foreground and background regions respectively. Finally, two distance maps $\mathcal{D} = \{\mathcal{D}_f, \mathcal{D}_b\}$ are generated by evaluating the cosine similarity between $\mathcal{P}_s$ and $\mathcal{F}_q$, which is then processed by a softmax operation as the final prediction $\mathcal{M}_1 = \text{softmax}(\mathcal{D})$.

### 3.1 Motivation

Current FSS methods rely heavily on the support prototype to segment query objects, by densely matching each query pixel with the support prototype. However, such cross-object matching severely suffers from intra-class appearance discrepancy, where objects in support and query can look very different even belonging to the same class. Such high intra-class variation cannot be reconciled by only a few supports, thus leading to poor matching results due to the large appearance gap between the query and supports.

To validate the relevance of Gestalt law [41] in narrowing such appearance discrepancy, we statistically analyze the feature cosine similarity of cross-object and intra-object pixels of Pascal VOC [19], where the pixel features are extracted from the ImageNet [15]-pretrianed ResNet-50 [34]. Table 1 shows that pixels belonging to the same object are much more similar than the cross-object pixels. Notably, background pixels share similar characteristics on their own, where intra-image background pixels are much more similar than cross-image pixels.

Thus, we propose to leverage the query feature to generate self-support prototypes to match the query feature itself. Notably, such prototype aligns the query along the homologous query features and thus can significantly narrow the feature gap between the support and query. In hindsight, the crucial reason the self-support matching works better than traditional support-query matching is that for a given visual object class, the intra-object similarities are much higher than the cross-object similarities.

### 3.2 Self-Support Prototype

Our core idea (Figure 2) is to aggregate query features to generate the query prototype and use it to self-support the query feature itself.

To recap, the regular support prototype generation procedure is:

$$\mathcal{P}_s = MAP(\mathcal{M}_s, \mathcal{F}_s), \tag{1}$$

where $MAP$ is the masked average pooling operation, which is used to generate the matching prediction with query feature $\mathcal{F}_q$:

$$\mathcal{M}_1 = \text{softmax}(\text{cosine}(\mathcal{P}_s, \mathcal{F}_q)), \tag{2}$$
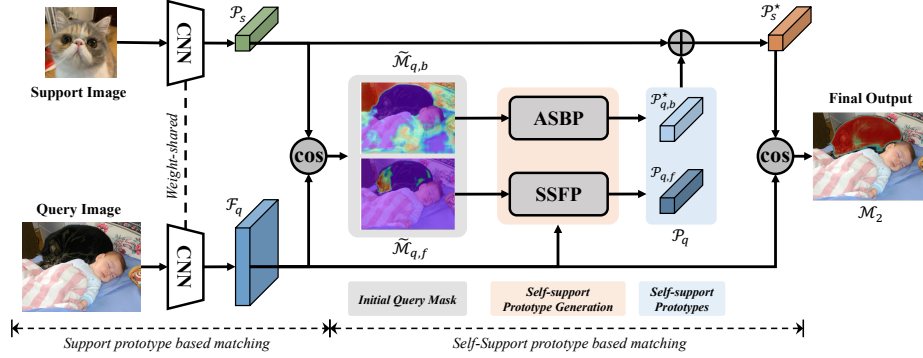
**Fig. 2.** Overall self-support network architecture. We first generate the initial mask predictions using the traditional support prototype based matching network. Then we leverage the initial query mask to aggregate query features to generate self-support prototypes, *i.e.*, the self-support foreground prototype (SSFP) and adaptive self-support background prototype (ASBP). Finally, we combine the support prototype and self-support prototypes to perform matching with query features.

where cosine is the cosine similarity metric.

Now, we can generate the query prototype $\mathcal{P}_q$ in the same manner, except the groundtruth masks of query images $\mathcal{M}_q$ are unavailable during inference. Thus, we need to use a predicted query mask $\widetilde{\mathcal{M}}_q$ to aggregate query features. The query prototype generation procedure can be formulated as:
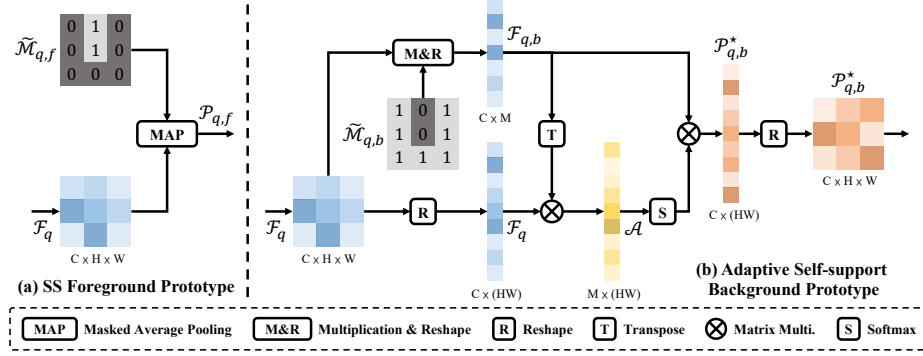
$$\mathcal{P}_q = MAP(\widetilde{\mathcal{M}}_q, \mathcal{F}_q), \tag{3}$$

where $\widetilde{\mathcal{M}}_q = \mathbb{1}(\mathcal{M}_1 > \tau)$, and $\mathcal{M}_1$ is the estimated query mask generated by Equation 2, $\mathbb{1}$ is the indicator function. The mask threshold $\tau$ is used to control the query feature sampling scope which is set as $\{\tau_{fg} = 0.7, \tau_{bg} = 0.6\}$ for foreground and background query masks respectively. The estimated self-support prototype $\mathcal{P}_q = \{\mathcal{P}_{q,f}, \mathcal{P}_{q,b}\}$ will be utilized to match query features.

We understand the reader's natural concern about the quality of self-support prototype, which is generated based on the estimated query mask, *i.e.*, whether the estimated mask is capable of effective self-support prototype generation. We found that even the estimated query mask is not perfect, as long as it covers some representative object fragments, it is sufficient to retrieve other regions of the same object. To validate partial object or object fragment is capable of supporting the entire object, we train and evaluate models with partial prototypes, which are aggregated from randomly selecting features based on the groundtruth mask labels. We conduct the 1-shot segmentation experiments on Pascal VOC dataset with the ResNet-50 backbone. As shown in Table 2, while reducing the aggregated object regions for prototype generation, our self-support prototype consistently achieves high segmentation performance. By contrast, the traditional support prototype consistently obtains much inferior performance, even using perfect support features from the entire object.

**Table 2.** The 1-shot matching results (mIoU) of support/self-support prototypes aggregated from full/partial objects.

| Object Ratio | full | 10% | 1% | 1%+noise |
|---|---|---|---|---|
| Support Prototype | 58.2 | 57.1 | 52.4 | 48.7 |
| Self-support Prototype | 83.0 | 82.5 | 79.2 | 74.6 |



**Fig. 3.** Prototype generations of (a) self-support (SS) foreground prototype and (b) adaptive self-support background prototype.

We further introduce noisy features (with 20% noise ratio) into partial prototypes to mimic realistic self-support generation during inference, by randomly selecting image features from non-target regions and aggregating these features into the above partial prototypes. To our pleasant surprise, our self-support prototype still works much better than the traditional support prototype in such noisy situation. Note that each image may contain multiple objects, thus the good performance indicates that our self-support prototype can also handle well the multiple objects scenarios. These results confirm the practicability and advantages of our self-support prototypes in the realistic applications.

### 3.3 Adaptive Self-Support Background Prototype

Foreground pixels share semantic commonalities [20,21], which constitutes the rationale behind our self-support prototype generation and matching procedure between query feature and support prototypes for foreground objects. Therefore, we can utilize a masked average pooling to generate the self-support foreground prototype (Figure 3 (a)):

$$\mathcal{P}_{q,f} = MAP(\widetilde{\mathcal{M}}_{q,f}, \mathcal{F}_q), \tag{4}$$

where $\widetilde{\mathcal{M}}_{q,f}$ is the aforementioned estimated query mask.

On the other hand, background can be cluttered, where commonalities can be reduced to local semantic similarities in disjoint regions, without a global semantic commonality shared among all background pixels. For example, for a query image with dog as the target class, other objects such as person and car

are both treated as background, but they are different in both appearance and semantic levels. This observation is also validated by the smaller background pixel similarity compared to foreground pixels as shown in Table 1, especially in the intra-object/image situation. This motivates us to generate multiple self-support background prototypes for different query semantic regions.

A straightforward solution is to directly group multiple background prototypes using a clustering algorithm, and then choose the most similar prototype at each query pixel for background matching. This explicit background grouping heavily relies on the clustering algorithm, which is unstable and time-consuming. Therefore, we propose a more flexible and efficient method to adaptively generate self-support background prototypes for each query pixel (Figure 3 (b)).

The idea is to dynamically aggregate similar background pixels for each query pixel to generate adaptive self-support background prototypes. Specifically, we first gather the background query features $\mathcal{F}_{q,b} \in \mathbb{R}^{C \times M}$ through the masked multiplication on the query feature $\mathcal{F}_q$ with the background mask $\widetilde{\mathcal{M}}_{q,b}$, where $M$ is the pixel number of the background region. Then we can generate the affinity matrix $\mathcal{A}$ between pixels of the reshaped background query feature $\mathcal{F}_{q,b}$ and full query feature $\mathcal{F}_q$ through a matrix multiplication operation $MatMul$:

$$\mathcal{A} = MatMul(\mathcal{F}_{q,b}{}^T, \mathcal{F}_q), \tag{5}$$

where $\mathcal{A}$ is in size of $\mathbb{R}^{M \times (H \times W)}$. The affinity matrix is normalized through a softmax operation along the first dimension, which is used to weighted aggregate background query features for each query pixel to generate the adaptive self-support background prototypes $\mathcal{P}_{q,b}^{\star} \in \mathbb{R}^{C \times H \times W}$:

$$\mathcal{P}_{q,b}^{\star} = MatMul(\mathcal{F}_{q,b}, \mathrm{softmax}(\mathcal{A})). \tag{6}$$

The self-support prototype is updated with the adaptive self-support background prototype: $\mathcal{P}_q = \{\mathcal{P}_{q,f}, \mathcal{P}_{q,b}^{\star}\}$.

### 3.4   Self-Support Matching

We weighted combine the support prototype $\mathcal{P}_s$ and self-support prototype $\mathcal{P}_q$:

$$\mathcal{P}_s^{\star} = \alpha_1 \mathcal{P}_s + \alpha_2 \mathcal{P}_q, \tag{7}$$

where $\alpha_1$ and $\alpha_2$ are the tuning weights and we set $\alpha_1 = \alpha_2 = 0.5$ in our experiments. Then we compute the cosine distance between the augmented support prototype $\mathcal{P}_s^{\star}$ and query feature $\mathcal{F}_q$ to generate the final matching prediction:

$$\mathcal{M}_2 = \mathrm{softmax}(\mathrm{cosine}(\mathcal{P}_s^{\star}, \mathcal{F}_q)). \tag{8}$$

Then we apply the training supervision on the generated distance maps:

$$\mathcal{L}_m = BCE(\mathrm{cosine}(\mathcal{P}_s^{\star}, \mathcal{F}_q), \mathcal{G}_q), \tag{9}$$

where $BCE$ is the binary cross entropy loss and $\mathcal{G}_q$ is the groundtruth mask of the query image.

To further facilitate the self-support matching procedure, we propose a novel query self-support loss. For the query feature $\mathcal{F}_q$ and its prototype $\mathcal{P}_q$, we apply the following training supervision:

$$\mathcal{L}_q = BCE(\mathrm{cosine}(\mathcal{P}_q, \mathcal{F}_q), \mathcal{G}_q). \tag{10}$$

We can apply the same procedure on the support feature to introduce the support self-matching loss $\mathcal{L}_s$.

Finally, we train the model in an end-to-end manner by jointly optimizing all the aforementioned losses:

$$\mathcal{L} = \lambda_1 \mathcal{L}_m + \lambda_2 \mathcal{L}_q + \lambda_3 \mathcal{L}_s, \tag{11}$$

where $\lambda_1 = 1.0, \lambda_2 = 1.0, \lambda_3 = 0.2$ are the loss weights.

## 4  Experiments

**Datasets.** We conduct experiments on two FSS benchmark datasets: PASCAL-$5^i$ [19] and COCO-$20^i$ [51]. We follow previous works [73,82] to split the data into four folds for cross validation, where three folds are used for training and the remaining one for evaluation. During inference, we randomly sample 1,000/4,000 support-query pairs to perform evaluation for PASCAL-$5^i$ and COCO-$20^i$, respectively. We use the popular mean Intersection-over-Union (mIoU, $\uparrow$[1]) as the default metric to evaluate our model under 1-shot and 5-shot settings. We also apply the Mean Absolute Error (MAE, $\downarrow$) to evaluate our prediction quality. By default, all analyses are conducted on PASCAL-$5^i$ dataset with ResNet-50 backbone in the 5-shot setting.

**Implementation details.** We adopt the popular ResNet-50/101 [34] pre-trained on ImageNet [15] as the backbone. Following previous work MLC [82], we discard the last backbone stage and the last ReLU for better generalization. We use SGD to optimize our model with the 0.9 momentum and 1e-3 initial learning rate, which decays by 10 times every 2,000 iterations. The model is trained for 6,000 iterations where each training batch contains 4 support-query pairs. Both images and masks are resized and cropped into (473, 473) and augmented with random horizontal flipping. The evaluation is performed on the original image.

### 4.1  Comparison with State-of-the-Arts

To validate the effectiveness of our method, we conduct extensive comparisons with SOTA methods under different backbone networks and few-shot settings.
**PASCAL-$5^i$.** We present the results of our self-support method and the improved version with one extra self-support refinement. As shown in Table 3, our method substantially outperforms MLC [82] by a large margin in the 5-shot setting, with the improvement jumping from 2.7% to 3.7% with the ResNet-50 backbone replaced by the stronger ResNet-101 network. In the 1-shot setting, our slightly inferior performance is remedied by using the stronger ResNet-101 backbone, where we surpass MLC [82] by 1.4% improvement. We can further promote the overall performance on PASCAL-$5^i$ up to 73.1% with the self-support refinement, which is a simple and straightforward extension by repeating the self-support procedure. It surpasses the previous SOTA [61] by 2.7%. Note that our method is non-parametric and thus our model uses fewest parameters while achieving the best performance.

---

[1] The "$\uparrow$" ("$\downarrow$") means that the higher (lower) is better.

**Table 3.** Quantitative comparison results on PASCAL-$5^i$ dataset. The **best** and <u>second best</u> results are highlighted with **bold** and <u>underline</u>, respectively.

| Method | Backbone | 1-shot | | | | | 5-shot | | | | | Params |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | fold0 | fold1 | fold2 | fold3 | **Mean** | fold0 | fold1 | fold2 | fold3 | **Mean** | |
| PANet [75] | Res-50 | 44.0 | 57.5 | 50.8 | 44.0 | 49.1 | 55.3 | 67.2 | 61.3 | 53.2 | 59.3 | <u>23.5 M</u> |
| PPNet [56] | | 48.6 | 60.6 | 55.7 | 46.5 | 52.8 | 58.9 | 68.3 | 66.8 | 58.0 | 63.0 | 31.5 M |
| PFENet [73] | | <u>61.7</u> | 69.5 | 55.4 | <u>56.3</u> | 60.8 | 63.1 | 70.7 | 55.8 | 57.9 | 61.9 | 34.3 M |
| CWT [59] | | 56.3 | 62.0 | 59.9 | 47.2 | 56.4 | 61.3 | 68.5 | 68.5 | 56.6 | 63.7 | - |
| HSNet [61] | | **64.3** | <u>70.7</u> | 60.3 | **60.5** | **64.0** | **70.3** | **73.2** | 67.4 | **67.1** | **69.5** | 26.1 M |
| MLC [82] | | 59.2 | **71.2** | <u>65.6</u> | 52.5 | <u>62.1</u> | 63.5 | 71.6 | 71.2 | 58.1 | 66.1 | **8.7 M** |
| SSP (Ours) | | 61.4 | 67.2 | 65.4 | 49.7 | 60.9 | <u>68.0</u> | 72.0 | <u>74.8</u> | 60.2 | 68.8 | **8.7 M** |
| SSP$_{refine}$ | | 60.5 | 67.8 | **66.4** | 51.0 | 61.4 | 67.5 | <u>72.3</u> | **75.2** | <u>62.1</u> | <u>69.3</u> | **8.7 M** |
| FWB [62] | Res-101 | 51.3 | 64.5 | 56.7 | 52.2 | 56.2 | 54.8 | 67.4 | 62.2 | 55.3 | 59.9 | <u>43.0 M</u> |
| PPNet [56] | | 52.7 | 62.8 | 57.4 | 47.7 | 55.2 | 60.3 | 70.0 | 69.4 | 60.7 | 65.1 | 50.5 M |
| PFENet [73] | | 60.5 | 69.4 | 54.4 | 55.9 | 60.1 | 62.8 | 70.4 | 54.9 | 57.6 | 61.4 | 53.4 M |
| CWT [59] | | 56.9 | 65.2 | 61.2 | 48.8 | 58.0 | 62.6 | 70.2 | 68.8 | 57.2 | 64.7 | - |
| HSNet [61] | | **67.3** | **72.3** | 62.0 | **63.1** | **66.2** | **71.8** | 74.4 | 67.0 | **68.3** | 70.4 | 45.2 M |
| MLC [82] | | 60.8 | <u>71.3</u> | 61.5 | <u>56.9</u> | 62.6 | 65.8 | 74.9 | 71.4 | 63.1 | 68.8 | **27.7 M** |
| SSP (Ours) | | <u>63.7</u> | 70.1 | <u>66.7</u> | 55.4 | 64.0 | 70.3 | <u>76.3</u> | <u>77.8</u> | 65.5 | <u>72.5</u> | **27.7 M** |
| SSP$_{refine}$ | | 63.2 | 70.4 | **68.5** | 56.3 | <u>64.6</u> | <u>70.5</u> | **76.4** | **79.0** | <u>66.4</u> | **73.1** | **27.7 M** |

**Table 4.** Quantitative comparison results on COCO-$20^i$ dataset. $^\star$ denotes the results are evaluated on the HSNet's evaluation protocol.

| Method | Backbone | 1-shot | | | | | 5-shot | | | | | Params |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | fold0 | fold1 | fold2 | fold3 | **Mean** | fold0 | fold1 | fold2 | fold3 | **Mean** | |
| PANet [75] | Res-50 | 31.5 | 22.6 | 21.5 | 16.2 | 23.0 | 45.9 | 29.2 | 30.6 | 29.6 | 33.8 | <u>23.5 M</u> |
| PPNet [56] | | 36.5 | 26.5 | 26.0 | 19.7 | 27.2 | 48.9 | 31.4 | 36.0 | 30.6 | 36.7 | 31.5 M |
| CWT [59] | | 32.2 | 36.0 | 31.6 | 31.6 | 32.9 | 40.1 | 43.8 | 39.0 | 42.4 | 41.3 | - |
| MLC [82] | | **46.8** | 35.3 | 26.2 | 27.1 | 33.9 | **54.1** | 41.2 | 34.1 | 33.1 | 40.6 | **8.7 M** |
| SSP (Ours) | | <u>46.4</u> | 35.2 | 27.3 | 25.4 | 33.6 | <u>53.8</u> | 41.5 | 36.0 | 33.7 | 41.3 | **8.7 M** |
| HSNet$^\star$ [61] | | 36.3 | **43.1** | **38.7** | **38.7** | **39.2** | 43.3 | **51.3** | **48.2** | **45.0** | **46.9** | 26.1 M |
| SSP$^\star$ (Ours) | | 35.5 | <u>39.6</u> | <u>37.9</u> | <u>36.7</u> | <u>37.4</u> | 40.6 | <u>47.0</u> | <u>45.1</u> | <u>43.9</u> | <u>44.1</u> | **8.7 M** |
| PMMs [81] | Res-101 | 29.5 | 36.8 | 28.9 | 27.0 | 30.6 | 33.8 | 42.0 | 33.0 | 33.3 | 35.5 | <u>38.6 M</u> |
| CWT [59] | | 30.3 | 36.6 | 30.5 | 32.2 | 32.4 | 38.5 | 46.7 | 39.4 | 43.2 | 42.0 | - |
| MLC [82] | | <u>50.2</u> | 37.8 | 27.1 | 30.4 | 36.4 | <u>57.0</u> | 46.2 | 37.3 | 37.2 | 44.4 | **27.7 M** |
| SSP (Ours) | | **50.4** | 39.9 | 30.6 | 30.0 | 37.7 | **57.8** | 47.0 | 40.2 | 39.9 | 46.2 | **27.7 M** |
| HSNet$^\star$ [61] | | 37.2 | <u>44.1</u> | <u>42.4</u> | **41.3** | <u>41.2</u> | 45.9 | <u>53.0</u> | **51.8** | <u>47.1</u> | <u>49.5</u> | 45.2 M |
| SSP$^\star$ (Ours) | | 39.1 | **45.1** | **42.7** | <u>41.2</u> | **42.0** | 47.4 | **54.5** | <u>50.4</u> | **49.6** | **50.2** | **27.7 M** |

**COCO-$20^i$.** This is a very challenging dataset whose images usually contain multiple objects against a complex background. As shown in Table 4, our method obtains comparable or best results with the ResNet-50 backbone. When equipped with the stronger ResNet-101 backbone, our method significantly outperforms MLC [82] with 1.3/1.8% improvements in 1/5-shot settings. To fairly compare to HSNet [61], we adopt their evaluation protocol to evaluate our method. Our method achieves SOTA when using the ResNet-101 backbone. Our method also performs best on FSS-1000 [48], shown in the supplementary material.

**Table 5.** Self-support model ablation results. "SSM" denotes the self-support module (containing the self-support foreground/background prototypes) , "SSL" denotes the self-support loss and "ASBP" denotes the adaptive self-support background prototype.

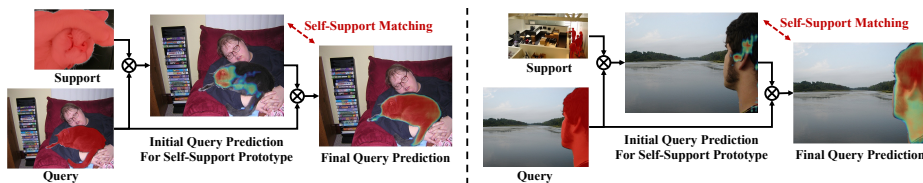| SSM | SSL | ASBP | fold0 | fold1 | fold2 | fold3 | **Mean** |
|-----|-----|------|-------|-------|-------|-------|----------|
|     |     |      | 62.2 | 70.5 | 70.7 | 55.7 | 64.8 |
| ✓   |     |      | 65.3 | 71.1 | 73.6 | 59.2 | $67.3_{\uparrow 2.5}$ |
|     | ✓   |      | 63.6 | 71.0 | 71.7 | 56.3 | $65.7_{\uparrow 0.9}$ |
| ✓   | ✓   |      | 67.0 | **72.4** | 72.9 | 59.9 | $68.1_{\uparrow 3.3}$ |
| ✓   |     | ✓    | 67.0 | 71.4 | 74.7 | 59.8 | $68.2_{\uparrow 3.4}$ |
| ✓   | ✓   | ✓    | **68.0** | 72.0 | **74.8** | **60.2** | $\mathbf{68.8_{\uparrow 4.0}}$ |



**Fig. 4.** Visualization for the working mechanism of our self-support matching. We omit the original support in self-support matching and the first row caption for clarity.

Note that our method benefits more improvement from stronger backbones and more supports because they provide better self-support prototypes, which will be validated later in Table 8.

### 4.2    Ablation Studies

As shown in Table 5, our self-support module significantly improves the performance by 2.5 %. The self-support loss further facilitates the self-support procedure and promotes the performance to 68.1%. The baseline model also benefits from the extra supervision of self-support loss. After equipped with the adaptive self-support background prototype, the self-support module can obtain extra 0.9% gain. Integrating all modules, our self-support method significantly improves the performance from 64.8% to 68.8% based on the strong baseline.

### 4.3    Self-Support Analysis

We conduct extensive experiments and analysis to understand our method.
**Self-support working mechanism.**  As shown in Figure 4, we first generate the *Initial* query predictions using the support prototype (as in Equation 2), and leverage the confident predictions to extract query features to generate self-support prototype (as in Equation 3). Then we use the self-support prototype to match with query features (as in Equation 16) and produce the final output. Note that because of the large inter-object/inter-background variation, the *Init* predictions usually only capture some small representative regions, *e.g.*, the cat/dog heads. Notwithstanding, our self-support method can handle well these hard cases by bridging the gap between query and support prototypes.
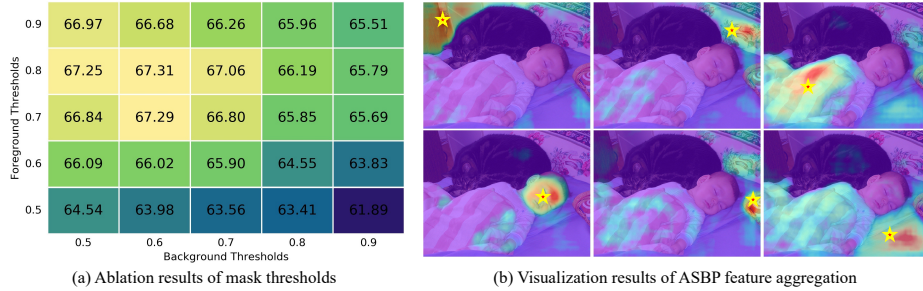
| | | | | |
|---|---|---|---|---|
| 66.97 | 66.68 | 66.26 | 65.96 | 65.51 |
| 67.25 | 67.31 | 67.06 | 66.19 | 65.79 |
| 66.84 | 67.29 | 66.80 | 65.85 | 65.69 |
| 66.09 | 66.02 | 65.90 | 64.55 | 63.83 |
| 64.54 | 63.98 | 63.56 | 63.41 | 61.89 |

(a) Ablation results of mask thresholds     (b) Visualization results of ASBP feature aggregation

**Fig. 5.** (a) Results of mask threshold variations for self-support prototypes. (b) Visualization of the feature aggregation for adaptive self-support background prototypes (ASBP) at each star-marked position. They are aggregated from the activated background regions.

**Table 6.** Ablation results of self-support module (SSM) by respectively removing foreground support prototype (FP), background support prototype (BP), self-support foreground prototype (SFP) and self-support background prototype (SBP).

| SSM | w/o FP | w/o BP | w/o SFP | w/o SBP |
|---|---|---|---|---|
| 67.3 | $66.0_{\downarrow 1.3}$ | $67.2_{\downarrow 0.1}$ | $66.5_{\downarrow 0.8}$ | $65.6_{\downarrow 1.7}$ |

**Mask threshold.** The threshold $\tau$ controls the query feature selection for self-support prototype generation (as in Equation 3). While we need to select high-confidence features for the foreground prototype, the background prototype requires more query features with a relative low threshold. This is because foreground pixels exhibit strong similarities with relatively low noise tolerance, while background is cluttered and the aggregated diverse features should tolerate more noises. Figure 5 (a) summarizes the model performance on Pascal dataset with different thresholds, where a good balance between foreground and background thresholds are respectively $\tau_{fg} \in [0.7, 0.9]$ and $\tau_{bg} \in [0.5, 0.7]$.

**Prototype ablation.** We investigate the effect of each of the prototypes by respectively removing them from the overall prototype. Table 6 summarizes the results. Both our self-support foreground and background prototypes play a critical role to account for good matching performance. The support foreground prototype is also essential for the prototype quality thanks to its foreground semantic information aggregated from multiple support images. The support background prototype can be discarded with slight impact because of the large background variation between query and support.

**Distinction from self-attention.** Readers may compare our self-support method with self-attention mechanisms. Our self-support method shares some concepts but is different from self-attention. Self-attention augments the image feature at each position by weighted aggregation of the features from all positions according to the affinity matrix. In contrast, our self-support method leverages representative query features to generate prototypes according to the query-support matching results. In Table 7 we experiment with multiple self-attention modules on the baseline. Unfortunately, all of them impose various degrees of

**Table 7.** Comparison with self-attention modules. "$\dagger$" means the improved version by removing the transformation layer.

| Baseline | NL [76] | NL$^\dagger$ [76] | GCNet [8] | Our SSM |
|----------|---------|---------|-----------|---------|
| 64.8 | $62.1_{\downarrow 2.7}$ | $64.3_{\downarrow 0.5}$ | $63.9_{\downarrow 0.9}$ | $\mathbf{67.3_{\uparrow 2.5}}$ |

**Table 8.** Comparison with other methods on performance improvement across different backbones and support shots.

| | PFENet [73] | ReRPI [7] | CWT [59] | MLC [82] | Ours |
|---|---|---|---|---|---|
| R50 $\rightarrow$ R101 | $-0.5$ | $-1.2$ | $+1.0$ | $+2.7$ | $+\mathbf{3.7}$ |
| 1shot $\rightarrow$ 5shot | $+1.3$ | $+6.2$ | $+6.7$ | $+6.2$ | $+\mathbf{8.5}$ |

harm on the matching performance, which are resulted by their self-attention augmentation which can destroy feature similarity between query and supports.
**Adaptive self-support background prototype.** This is designed to address the background clutter problem by adaptively aggregating background prototypes for each position. As shown in Figure 5 (b), the target cat is lying on a cluttered background consisting of the wardrobe, bed, quilt, baby, pillow and sheet. For each star-marked query position, the self-support background prototypes are aggregated from the corresponding semantic regions. Note that this adaptive background prototype generation is specifically designed for self-support prototypes, which cannot be directly applied to support prototype generation because it will collapse to trivial solutions by greedily aggregating similar pixels without semantic consideration.

### 4.4 Self-Support Advantages

Our self-support method has many desirable properties.
**Benefits from backbones and supports.** As shown before, our self-support method benefits more improvement from stronger backbones and more supports. Table 8 summarizes the improvement of different methods. When switching the backbone from ResNet-50 to ResNet-101, our method obtains 3.7% performance improvement, while other methods obtain at most 2.7% improvement or even performance degradation. Our method also obtains the largest improvement of 8.5% by increasing support images from 1-shot to 5-shot. The behind reason is that our self-support method benefits from the Matthew effect [60] of accumulated advantages, where better predictions induce better self-support prototypes and produce better predictions.
**High-confident predictions.** Our self-support method not only improves hard segmentation results with 0-1 labels, but also improves the soft confidence scores to produce high-confident predictions. As shown in Table 9, our self-support method significantly reduces the Mean Absolute Error (MAE) by 4.9% compared to the baseline. We further evaluate the MAE in the truth positive (TP) regions for a fair comparison, where the MAE can still be largely reduced by 5.0%. These results demonstrate that our self-support method can significantly

**Table 9.** Results of prediction quality in MAE ($\downarrow$) metric. "All/TP" means evaluating models on all/truth positive regions of the image.

|     | Baseline | SSM | SSM+SSL | SSM+ASBP | Full |
|-----|----------|-----|---------|----------|------|
| All | 17.6 | $14.6_{\downarrow3.0}$ | $14.8_{\downarrow2.8}$ | $12.9_{\downarrow4.7}$ | $\mathbf{12.7_{\downarrow\mathbf{4.9}}}$ |
| TP  | 13.2 | $9.6_{\downarrow3.6}$ | $10.1_{\downarrow3.1}$ | $\mathbf{7.8_{\downarrow\mathbf{5.4}}}$ | $8.2_{\downarrow5.0}$ |

**Table 10.** Results of using weak support annotations.

|          | Mask | Scribble | Bounding Box |
|----------|------|----------|--------------|
| Baseline | 64.8 | $63.3_{\downarrow1.5}$ | $61.7_{\downarrow3.1}$ |
| Ours     | 68.8 | $68.0_{\downarrow0.8}$ | $66.9_{\downarrow2.1}$ |

**Table 11.** Results of applying our method to other models.

| PANet [75] | PANet + Ours | PPNet [56] | PPNet + Ours |
|------------|--------------|------------|--------------|
| 55.7 | $58.3_{\uparrow2.6}$ | 62.0 | $64.2_{\uparrow2.2}$ |

improve the output quality by producing high-confident predictions, a desirable property for many real-world applications.

**Robust to weak support labels.** As shown in Table 10, when replacing the support mask with bounding box or scribble annotations for prototype generation, our self-support method still works very well with high robustness against support noises. This is because our method mainly relies on self-support prototypes and thus is less affected by from noisy support prototypes.

**Generalized to other methods.** Our self-support method is general and can be applied to other methods. As shown in Table 11, equipped with our self-support module, both the strong PANet [75] and PPNet [56] report further boost in their performance by a large improvement.

**High efficiency.** Our self-support method is very efficient, which is a non-parametric method with few extra computation and $\sim$28 FPS running speed on a Tesla V100 GPU (with the ResNet-50 backbone in the 1-shot setting).

## 5   Conclusion

In this paper, we address the critical intra-class appearance discrepancy problem inherent in few-shot segmentation, by leveraging the query feature to generate self-support prototypes and perform self-support matching with query features. This strategy effectively narrows down the gap between support prototypes and query features. Further, we propose an adaptive self-support background prototype and a self-support loss to facilitate the self-support procedure. Our self-support network has various desirable properties, and achieves SOTA on multiple benchmarks. We have thoroughly investigated the self-support procedure with extensive experiments and analysis to substantiate its effectiveness and deepen our understanding on its working mechanism.

## 6    More Implementation Details

Our baseline model is adopted from MLC [82] with a metric learning framework consisting of only an encoder.

Our improved model with self-support refinement is to repeat the self-support procedure based on the predicted mask $\mathcal{M}_2$ produced by our self-support network. Specifically, the refined self-support foreground prototype $\mathcal{P}_{q,f}^r$ generation can be formulated as:

$$\mathcal{P}_{q,f}^r = MAP(\widetilde{\mathcal{M}}_{2,q,f}, \mathcal{F}_q), \tag{12}$$

where $\mathcal{F}_q$ is the query feature. Similarly, we can generate the refined self-support background prototype $P_{q,b}^{\star,r}$ by

$$\mathcal{P}_{q,b}^{\star,r} = ASBP(\widetilde{\mathcal{M}}_{2,q,b}, \mathcal{F}_q), \tag{13}$$

where $ASBP$ is the adaptive self-support background prototype generation module. The $\widetilde{\mathcal{M}}_{2,q,f} = \mathbb{1}(\mathcal{M}_{2,f} > \tau_{fg})$ and $\widetilde{\mathcal{M}}_{2,q,b} = \mathbb{1}(\mathcal{M}_{2,b} > \tau_{bg})$ are defined according to the estimated query mask $\mathcal{M}_2 = \{\mathcal{M}_{2,f}, \mathcal{M}_{2,b}\}$ by Equation 8 in the main paper. We use the same $\{\tau_{fg} = 0.7, \tau_{bg} = 0.6\}$ settings as in the main paper.

Finally, we weight-combine the original support prototype $\mathcal{P}_s = \{\mathcal{P}_{s,f}, \mathcal{P}_{s,b}\}$, self-support prototype $\mathcal{P}_q = \{\mathcal{P}_{q,f}, \mathcal{P}_q^\star\}$ and refined self-support prototype $\mathcal{P}_q^r = \{\mathcal{P}_{q,f}^r, \mathcal{P}_{q,b}^{\star,r}\}$:

$$\mathcal{P}_s^{\star,r} = \alpha_1\mathcal{P}_s + \alpha_2\mathcal{P}_q + \alpha_3\mathcal{P}_q^r, \tag{14}$$

where $\alpha_1$, $\alpha_2$ and $\alpha_3$ are the tuning weights which are set as $\alpha_1 = 0.5$, $\alpha_2 = 0.2$ and $\alpha_3 = 0.3$ in our experiments. Then we compute the cosine distance between the augmented support prototype with self-support refinement $\mathcal{P}_s^{\star,r}$ and query feature $\mathcal{F}_q$ to generate the matching prediction output $\mathcal{M}_3$:

$$\mathcal{M}_3 = \text{softmax}(\text{cosine}(\mathcal{P}_s^{\star,r}, \mathcal{F}_q)). \tag{15}$$

The final output $M_{final}$ is the weighted combination of $\mathcal{M}_2$ and $\mathcal{M}_3$ for good performance:

$$\mathcal{M}_{final} = \beta_1 M_2 + \beta_2 M_3, \tag{16}$$

where $\beta_1$ and $\beta_2$ are the tuning weights and we set $\beta_1 = 0.3$ and $\beta_2 = 0.7$ in our experiments.

## 7    More Quantitative Results

In the main paper, we repeat the evaluation procedure of all our experiments by 5 times with different random seeds to obtain stable results.

$^\dagger$ Corresponding author.

**Table 12.** Quantitative comparison results on FSS-1000 dataset with the mIoU metric of positive labels in a binary segmentation map.

| Method | Publication | 1-shot | 5-shot |
|--------|-------------|--------|--------|
| OSLSM [69] | BMVC'17 | 70.3 | 73.0 |
| GNet [66] | Arxiv'18 | 71.9 | 74.3 |
| FSS [48] | CVPR'20 | 73.5 | 80.1 |
| DoG-LSTM [3] | WACV'21 | 80.8 | 83.4 |
| DAN [74] | ECCV'20 | 85.2 | 88.1 |
| SSP (Ours) | - | 86.9 | 88.2 |
| $SSP_{refine}$ | - | 87.3 | 88.6 |

**Table 13.** Performance improvement on Pascal VOC dataset of our self-support method on baseline models with different backbones and support shots.

| Backbone | Shot | Baseline | Ours | $\Delta$ |
|----------|------|----------|------|----------|
| ResNet-50 | 1-shot | 57.8 | 60.9 | +3.1 |
|  | 5-shot | 64.8 | 68.8 | +4.0 |
| ResNet-101 | 1-shot | 60.1 | 64.0 | +3.9 |
|  | 5-shot | 67.8 | 72.5 | +4.7 |

To further validate the effectiveness of our self-support method, we evaluate our model on FSS-1000 [48], which is a recently proposed large-scale few-shot segmentation dataset containing 1000 classes. The dataset is split into train/val/test sets with 520, 240, 240 classes respectively. We follow the common practice [48] to evaluate performance on FSS-1000 using the intersection-over-union (IoU) of positive labels in a binary segmentation map. The evaluation procedure is conducted on 2400 randomly sampled support-query pairs. As shown in Table 12, our self-support matching model outperforms other methods. And the self-support refinement step can further promote our performance to 87.3/88.6 mIoU in 1/5-shot settings.

As shown in Table 13, we present the performance improvement on Pascal VOC dataset [19] of our self-support method on the baseline models. Our method can consistently improve the performance by a large margin with different backbone models and support shots. We can also observe more performance gains on the stronger backbone model and more support shots, which is consistent with the advantage conclusion in the main paper.

## 8    More Qualitative Results

We present more qualitative results in 1-shot setting with ResNet-50 backbone for better visualization. As shown in Figure 6, Figure 7, Figure 8, and Figure 9,

objects in support and query images have large appearance discrepancy even belonging to the same class. Thus the initial predictions generated by the traditional matching network can cover only a small region of the target object. On the other hand, equipped with our self-support method, the model can produce satisfactory results with substantial qualitative improvement. Note that the initial masks are obtained by setting foreground and background thresholds on the original mask prediction $\mathcal{M}_1$.
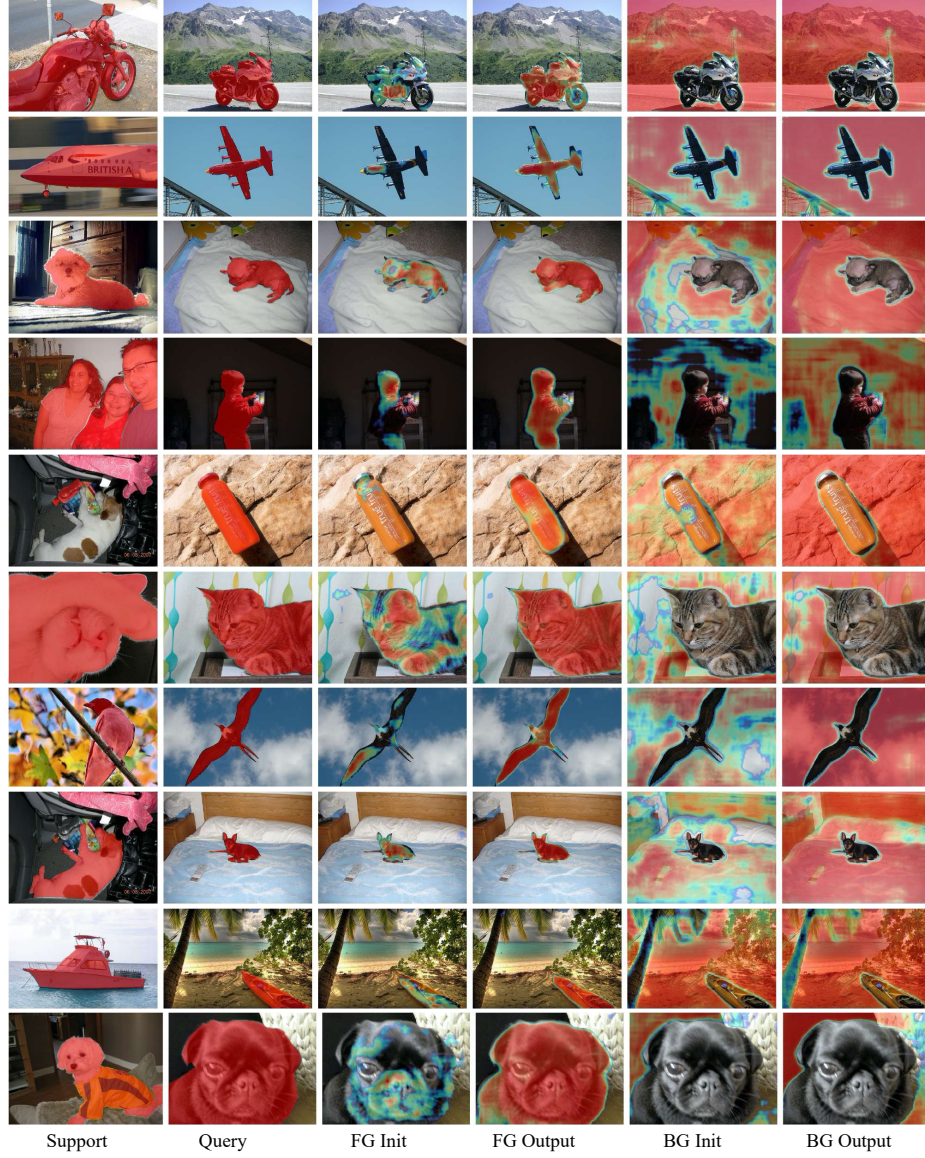
**Fig. 6.** The 1-shot visualization results of our model containing the *Init* and final outputs.

Support        Query        FG Init        FG Output        BG Init        BG Output

**Fig. 7.** The 1-shot visualization results of our model containing the *Init* and final outputs.

**Fig. 8.** The 1-shot visualization results of our model containing the *Init* and final outputs.
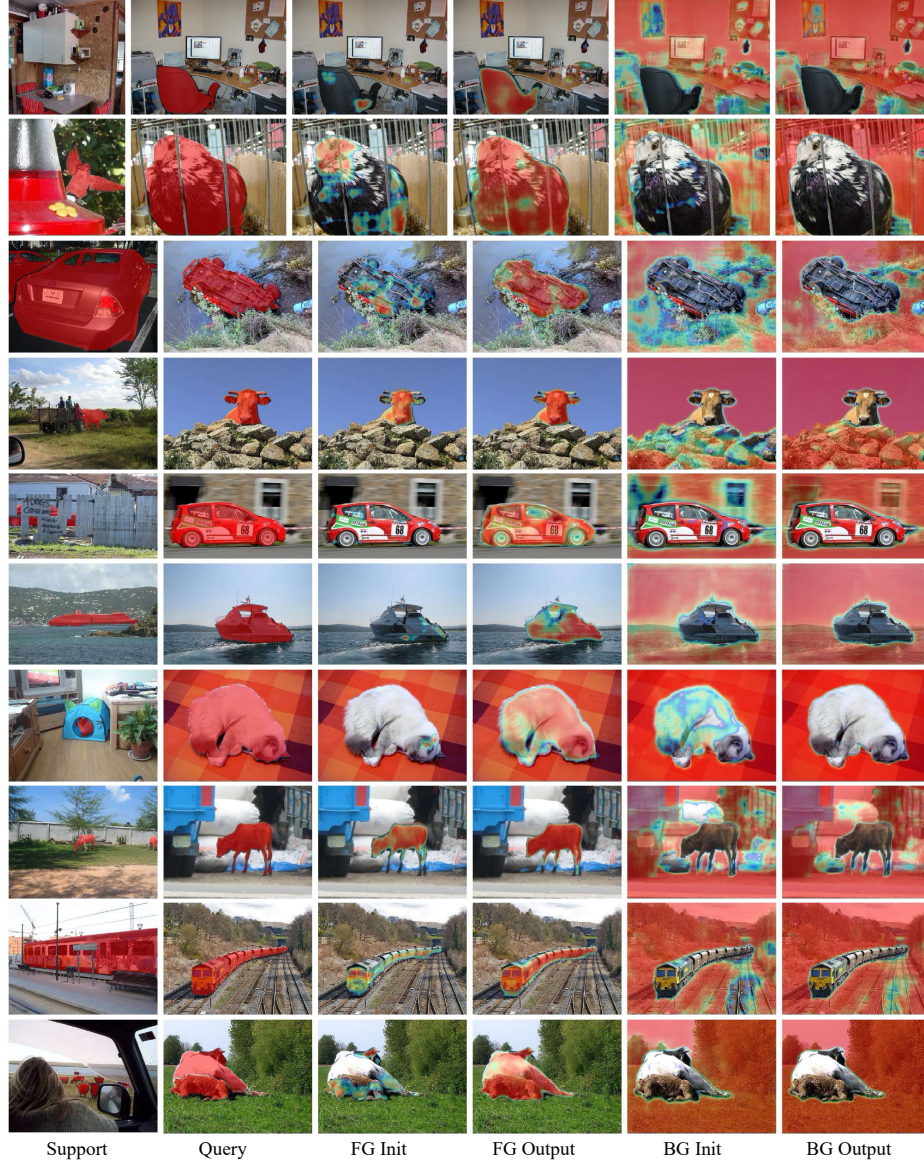
| Support | Query | FG Init | FG Output | BG Init | BG Output |

**Fig. 9.** The 1-shot visualization results of our model containing the *Init* and final outputs.

# References

1. Allen, K., Shelhamer, E., Shin, H., Tenenbaum, J.: Infinite mixture prototypes for few-shot learning. In: ICML (2019) 4
2. Antoniou, A., Edwards, H., Storkey, A.: How to train your maml. In: ICLR (2019) 4
3. Azad, R., Fayjie, A.R., Kauffmann, C., Ben Ayed, I., Pedersoli, M., Dolz, J.: On the texture bias for few-shot cnn segmentation. In: WACV (2021) 16
4. Badrinarayanan, V., Kendall, A., Cipolla, R.: Segnet: A deep convolutional encoder-decoder architecture for image segmentation. IEEE TPAMI (2017) 4
5. Benenson, R., Popov, S., Ferrari, V.: Large-scale interactive object segmentation with human annotators. In: CVPR (2019) 1
6. Bertinetto, L., Henriques, J.F., Torr, P.H., Vedaldi, A.: Meta-learning with differentiable closed-form solvers. In: ICLR (2019) 4
7. Boudiaf, M., Kervadec, H., Masud, Z.I., Piantanida, P., Ben Ayed, I., Dolz, J.: Few-shot segmentation without meta-learning: A good transductive inference is all you need? In: CVPR (2021) 13
8. Cao, Y., Xu, J., Lin, S., Wei, F., Hu, H.: Gcnet: Non-local networks meet squeeze-excitation networks and beyond. In: CVPRW (2019) 13
9. Chen, L.C., Papandreou, G., Kokkinos, I., Murphy, K., Yuille, A.L.: Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. IEEE TPAMI (2017) 4
10. Chen, L.C., Yang, Y., Wang, J., Xu, W., Yuille, A.L.: Attention to scale: Scale-aware semantic image segmentation. In: CVPR (2016) 4
11. Chen, L.C., Zhu, Y., Papandreou, G., Schroff, F., Adam, H.: Encoder-decoder with atrous separable convolution for semantic image segmentation. In: ECCV (2018) 4
12. Chen, W.Y., Liu, Y.C., Kira, Z., Wang, Y.C.F., Huang, J.B.: A closer look at few-shot classification. In: ICLR (2019) 4
13. Cheng, B., Chen, L.C., Wei, Y., Zhu, Y., Huang, Z., Xiong, J., Huang, T.S., Hwu, W.M., Shi, H.: Spgnet: Semantic prediction guidance for scene parsing. In: ICCV (2019) 4
14. Dai, J., Qi, H., Xiong, Y., Li, Y., Zhang, G., Hu, H., Wei, Y.: Deformable convolutional networks. In: ICCV (2017) 4
15. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: CVPR (2009) 1, 5, 9
16. Dhillon, G.S., Chaudhari, P., Ravichandran, A., Soatto, S.: A baseline for few-shot image classification. In: ICLR (2019) 4
17. Doersch, C., Gupta, A., Zisserman, A.: Crosstransformers: spatially-aware few-shot transfer. In: NeurIPS (2020) 4
18. Dong, N., Xing, E.P.: Few-shot semantic segmentation with prototype learning. In: BMVC (2018) 4
19. Everingham, M., Van Gool, L., Williams, C.K., Winn, J., Zisserman, A.: The pascal visual object classes (voc) challenge. IJCV (2010) 5, 9, 16
20. Fan, Q., Fan, D.P., Fu, H., Tang, C.K., Shao, L., Tai, Y.W.: Group collaborative learning for co-salient object detection. In: CVPR (2021) 7
21. Fan, Q., Ke, L., Pei, W., Tang, C.K., Tai, Y.W.: Commonality-parsing network across shape and appearance for partially supervised instance segmentation. In: ECCV (2020) 7

22. Fan, Q., Tang, C.K., Tai, Y.W.: Few-shot video object detection. arXiv preprint arXiv:2104.14805 (2021) 4
23. Fan, Q., Zhuo, W., Tang, C.K., Tai, Y.W.: Few-shot object detection with attention-rpn and multi-relation detector. In: CVPR (2020) 4
24. Finn, C., Abbeel, P., Levine, S.: Model-agnostic meta-learning for fast adaptation of deep networks. In: ICML (2017) 4
25. Fu, J., Liu, J., Tian, H., Li, Y., Bao, Y., Fang, Z., Lu, H.: Dual attention network for scene segmentation. In: CVPR (2019) 4
26. Fu, J., Liu, J., Wang, Y., Li, Y., Bao, Y., Tang, J., Lu, H.: Adaptive context network for scene parsing. In: ICCV (2019) 4
27. Gairola, S., Hemani, M., Chopra, A., Krishnamurthy, B.: Simpropnet: Improved similarity propagation for few-shot image segmentation. In: IJCAI (2020) 2, 4
28. Gidaris, S., Komodakis, N.: Dynamic few-shot visual learning without forgetting. In: CVPR (2018) 4
29. Goodfellow, I., Bengio, Y., Courville, A.: Deep Learning. MIT Press (2016) 1
30. Gordon, J., Bronskill, J., Bauer, M., Nowozin, S., Turner, R.: Meta-learning probabilistic inference for prediction. In: ICLR (2019) 4
31. Grant, E., Finn, C., Levine, S., Darrell, T., Griffiths, T.: Recasting gradient-based meta-learning as hierarchical bayes. In: ICLR (2018) 4
32. He, H., Zhang, J., Thuraisingham, B., Tao, D.: Progressive one-shot human parsing. In: AAAI (2021) 2, 4
33. He, J., Deng, Z., Zhou, L., Wang, Y., Qiao, Y.: Adaptive pyramid context network for semantic segmentation. In: CVPR (2019) 4
34. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: CVPR (2016) 1, 5, 9
35. Hou, R., Chang, H., Ma, B., Shan, S., Chen, X.: Cross attention network for few-shot classification. In: NeurIPS (2019) 4
36. Huang, Z., Wang, X., Huang, L., Huang, C., Wei, Y., Liu, W.: Ccnet: Criss-cross attention for semantic segmentation. In: ICCV (2019) 4
37. Kang, B., Liu, Z., Wang, X., Yu, F., Feng, J., Darrell, T.: Few-shot object detection via feature reweighting. In: ICCV (2019) 4
38. Kim, S., Chikontwe, P., Park, S.H.: Uncertainty-aware semi-supervised few shot segmentation. In: IJCAI (2021) 2
39. Kirillov, A., Girshick, R., He, K., Dollár, P.: Panoptic feature pyramid networks. In: CVPR (2019) 4
40. Koch, G., Zemel, R., Salakhutdinov, R.: Siamese neural networks for one-shot image recognition. In: ICMLW (2015) 4
41. Koffka, K.: Principles of Gestalt psychology. Routledge (1935) 2, 5
42. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In: NeurIPS (2012) 1
43. Lee, K., Maji, S., Ravichandran, A., Soatto, S.: Meta-learning with differentiable convex optimization. In: CVPR (2019) 4
44. Lee, Y., Choi, S.: Gradient-based meta-learning with learned layerwise metric and subspace. In: ICML (2018) 4
45. Li, G., Jampani, V., Sevilla-Lara, L., Sun, D., Kim, J., Kim, J.: Adaptive prototype learning and allocation for few-shot segmentation. In: CVPR (2021) 2, 3, 4
46. Li, H., Eigen, D., Dodge, S., Zeiler, M., Wang, X.: Finding task-relevant features for few-shot learning by category traversal. In: CVPR (2019) 4
47. Li, W., Wang, L., Xu, J., Huo, J., Gao, Y., Luo, J.: Revisiting local descriptor based image-to-class measure for few-shot learning. In: CVPR (2019) 4

48. Li, X., Wei, T., Chen, Y.P., Tai, Y.W., Tang, C.K.: Fss-1000: A 1000-class dataset for few-shot segmentation. In: CVPR (2020) 4, 10, 16
49. Lin, G., Milan, A., Shen, C., Reid, I.: Refinenet: Multi-path refinement networks for high-resolution semantic segmentation. In: CVPR (2017) 4
50. Lin, G., Shen, C., Van Den Hengel, A., Reid, I.: Efficient piecewise training of deep structured models for semantic segmentation. In: CVPR (2016) 4
51. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: ECCV (2014) 9
52. Liu, B., Ding, Y., Jiao, J., Ji, X., Ye, Q.: Anti-aliasing semantic reconstruction for few-shot semantic segmentation. In: CVPR (2021) 4
53. Liu, C., Fu, Y., Xu, C., Yang, S., Li, J., Wang, C., Zhang, L.: Learning a few-shot embedding model with contrastive learning. In: AAAI (2021) 4
54. Liu, L., Cao, J., Liu, M., Guo, Y., Chen, Q., Tan, M.: Dynamic extension nets for few-shot semantic segmentation. In: ACM MM (2020) 2, 4
55. Liu, W., Zhang, C., Lin, G., Liu, F.: Crnet: Cross-reference networks for few-shot segmentation. In: CVPR (2020) 3, 4
56. Liu, Y., Zhang, X., Zhang, S., He, X.: Part-aware prototype network for few-shot semantic segmentation. In: ECCV (2020) 2, 10, 14
57. Liu, Y., Zhang, X., Zhang, S., He, X.: Part-aware prototype network for few-shot semantic segmentation. In: ECCV (2020) 3, 4
58. Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. In: CVPR (2015) 3
59. Lu, Z., He, S., Zhu, X., Zhang, L., Song, Y.Z., Xiang, T.: Simpler is better: Few-shot semantic segmentation with classifier weight transformer. In: ICCV (2021) 4, 10, 13
60. Merton, R.K.: The matthew effect in science: The reward and communication systems of science are considered. Science (1968) 13
61. Min, J., Kang, D., Cho, M.: Hypercorrelation squeeze for few-shot segmentation. In: ICCV (2021) 4, 9, 10
62. Nguyen, K., Todorovic, S.: Feature weighting and boosting for few-shot segmentation. In: ICCV (2019) 2, 4, 10
63. Noh, H., Hong, S., Han, B.: Learning deconvolution network for semantic segmentation. In: ICCV (2015) 4
64. Ouyang, C., Biffi, C., Chen, C., Kart, T., Qiu, H., Rueckert, D.: Self-supervision with superpixels: Training few-shot medical image segmentation without annotation. In: ECCV (2020) 2, 4
65. Qi, H., Brown, M., Lowe, D.G.: Low-shot learning with imprinted weights. In: CVPR (2018) 4
66. Rakelly, K., Shelhamer, E., Darrell, T., Efros, A.A., Levine, S.: Few-shot segmentation propagation with guided networks. arXiv preprint arXiv:1806.07373 (2018) 16
67. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: MICCAI (2015) 4
68. Rusu, A.A., Rao, D., Sygnowski, J., Vinyals, O., Pascanu, R., Osindero, S., Hadsell, R.: Meta-learning with latent embedding optimization. In: ICLR (2019) 4
69. Shaban, A., Bansal, S., Liu, Z., Essa, I., Boots, B.: One-shot learning for semantic segmentation. In: BMVC (2017) 4, 16
70. Siam, M., Doraiswamy, N., Oreshkin, B.N., Yao, H., Jagersand, M.: Weakly supervised few-shot object segmentation using co-attention with visual and semantic embeddings. In: IJCAI (2020) 2, 4

71. Siam, M., Oreshkin, B.N., Jagersand, M.: Amp: Adaptive masked proxies for few-shot segmentation. In: ICCV (2019) 2, 3, 4
72. Tian, P., Wu, Z., Qi, L., Wang, L., Shi, Y., Gao, Y.: Differentiable meta-learning model for few-shot semantic segmentation. In: AAAI (2020) 4
73. Tian, Z., Zhao, H., Shu, M., Yang, Z., Li, R., Jia, J.: Prior guided feature enrichment network for few-shot segmentation. IEEE TPAMI (2020) 3, 4, 9, 10, 13
74. Wang, H., Zhang, X., Hu, Y., Yang, Y., Cao, X., Zhen, X.: Few-shot semantic segmentation with democratic attention networks. In: ECCV (2020) 2, 4, 16
75. Wang, K., Liew, J.H., Zou, Y., Zhou, D., Feng, J.: Panet: Few-shot image semantic segmentation with prototype alignment. In: ICCV (2019) 3, 4, 10, 14
76. Wang, X., Girshick, R., Gupta, A., He, K.: Non-local neural networks. In: CVPR (2018) 13
77. Wu, Z., Shi, X., Lin, G., Cai, J.: Learning meta-class memory for few-shot semantic segmentation. In: ICCV (2021) 4
78. Xie, G.S., Liu, J., Xiong, H., Shao, L.: Scale-aware graph neural network for few-shot semantic segmentation. In: CVPR (2021) 4
79. Xie, G.S., Xiong, H., Liu, J., Yao, Y., Shao, L.: Few-shot semantic segmentation with cyclic memory network. In: ICCV (2021) 4
80. Yan, X., Chen, Z., Xu, A., Wang, X., Liang, X., Lin, L.: Meta r-cnn : Towards general solver for instance-level low-shot learning. In: ICCV (2019) 4
81. Yang, B., Liu, C., Li, B., Jiao, J., Ye, Q.: Prototype mixture models for few-shot semantic segmentation. In: ECCV (2020) 2, 3, 4, 10
82. Yang, L., Zhuo, W., Qi, L., Shi, Y., Gao, Y.: Mining latent classes for few-shot segmentation. In: ICCV (2021) 2, 3, 4, 9, 10, 13, 15
83. Yang, X., Wang, B., Chen, K., Zhou, X., Yi, S., Ouyang, W., Zhou, L.: Brinet: Towards bridging the intra-class and inter-class gaps in one-shot segmentation. In: BMVC (2020) 2, 4
84. Yu, F., Koltun, V., Funkhouser, T.: Dilated residual networks. In: CVPR (2017) 4
85. Yuan, Y., Chen, X., Wang, J.: Object-contextual representations for semantic segmentation. In: ECCV (2020) 4
86. Zhang, B., Xiao, J., Qin, T.: Self-guided and cross-guided learning for few-shot segmentation. In: CVPR (2021) 4
87. Zhang, C., Lin, G., Liu, F., Guo, J., Wu, Q., Yao, R.: Pyramid graph networks with connection attentions for region-based one-shot semantic segmentation. In: ICCV (2019) 2, 4
88. Zhang, C., Lin, G., Liu, F., Yao, R., Shen, C.: Canet: Class-agnostic segmentation networks with iterative refinement and attentive few-shot learning. In: CVPR (2019) 2, 4
89. Zhang, F., Chen, Y., Li, Z., Hong, Z., Liu, J., Ma, F., Han, J., Ding, E.: Acfnet: Attentional class feature network for semantic segmentation. In: ICCV (2019) 4
90. Zhang, H., Zhang, L., Qi, X., Li, H., Torr, P.H., Koniusz, P.: Few-shot action recognition with permutation-invariant attention. In: ECCV (2020) 4
91. Zhao, H., Shi, J., Qi, X., Wang, X., Jia, J.: Pyramid scene parsing network. In: CVPR (2017) 4
92. Zhao, H., Zhang, Y., Liu, S., Shi, J., Loy, C.C., Lin, D., Jia, J.: Psanet: Point-wise spatial attention network for scene parsing. In: ECCV (2018) 4
93. Zhou, B., Zhao, H., Puig, X., Fidler, S., Barriuso, A., Torralba, A.: Scene parsing through ade20k dataset. In: CVPR (2017) 1
94. Zhu, K., Zhai, W., Zha, Z.J., Cao, Y.: Self-supervised tuning for few-shot segmentation. In: IJCAI (2020) 4

95. Zhu, Z., Xu, M., Bai, S., Huang, T., Bai, X.: Asymmetric non-local neural networks for semantic segmentation. In: ICCV (2019) 4
96. Zhuge, Y., Shen, C.: Deep reasoning network for few-shot semantic segmentation. In: ACM MM (2021) 2, 4