

Dynamic Prototype Convolution Network for Few-Shot Semantic Segmentation

Jie Liu^{1*}, Yanqi Bao^{2*}, Guo-Sen Xie^{3,4†}, Huan Xiong⁴, Jan-Jakob Sonke⁵, Efstratios Gavves¹

¹University of Amsterdam, Netherlands ²Northeastern University, China ⁵The Netherlands Cancer Institute, Netherlands

³Nanjing University of Science and Technology, China ⁴Mohamed bin Zayed University of Artificial Intelligence, UAE

Abstract

The key challenge for few-shot semantic segmentation (FSS) is how to tailor a desirable interaction among support and query features and/or their prototypes, under the episodic training scenario. Most existing FSS methods implement such support/query interactions by solely leveraging plain operations – e.g., cosine similarity and feature concatenation – for segmenting the query objects. However, these interaction approaches usually cannot well capture the intrinsic object details in the query images that are widely encountered in FSS, e.g., if the query object to be segmented has holes and slots, inaccurate segmentation almost always happens. To this end, we propose a dynamic prototype convolution network (DPCN) to fully capture the aforementioned intrinsic details for accurate FSS. Specifically, in DPCN, a dynamic convolution module (DCM) is firstly proposed to generate dynamic kernels from support foreground, then information interaction is achieved by convolution operations over query features using these kernels. Moreover, we equip DPCN with a support activation module (SAM) and a feature filtering module (FFM) to generate pseudo mask and filter out background information for the query images, respectively. SAM and FFM together can mine enriched context information from the query features. Our DPCN is also flexible and efficient under the k -shot FSS setting. Extensive experiments on PASCAL-5ⁱ and COCO-20ⁱ show that DPCN yields superior performances under both 1-shot and 5-shot settings.

1. Introduction

Semantic segmentation has achieved tremendous success due to the advancement of deep convolutional neural networks [9, 10, 22]. Nevertheless, most leading image semantic segmentation models rely on large amounts of training images with pixel-wise annotations, which require huge human efforts. Semi- and weakly-supervised segmentation methods [15, 24, 28] are proposed to alleviate such expen-

*Equal contribution.

†Corresponding author.

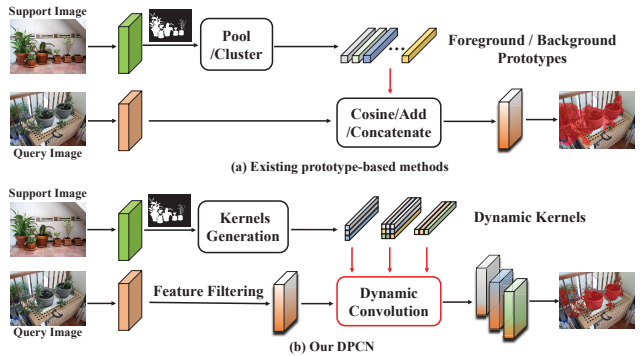


Figure 1. Comparison between (a) existing prototype-based methods and (b) proposed DPCN. (a) Existing methods usually adopt mask average pooling or clustering on support foreground to get multiple foreground/background prototypes. Then the information interaction between support and query can be plain operations, e.g., cosine similarity, element-wise sum, and channel-wise concatenation. However, this paradigm fails to segment the intrinsic details of *plants*, because this insufficient interaction cannot well address the appearance and shape variations (e.g., holes and slots in the *plants*) in FSS. (b) Our DPCN can well segment *plants* and capture the intrinsic subtle details. This benefits from dynamic convolution on query features with dynamic kernels generated from the support foreground features.

sive annotation cost. However, both semi- and weakly-supervised methods have to face a significant performance drop when only a few annotated samples for a novel object are available. In such a case, few-shot semantic segmentation (FSS) [20] is introduced to allow for dense pixel-wise prediction on novel object categories given only a few annotated samples.

Typically, most FSS methods adopt an episode-based meta-learning strategy [30], and each episode is composed of the support set and the query set, which share the same object class. The support set contains a few support images with pixel-wise annotation. FSS models are expected to learn to predict the segmentation mask for images in the query set with the guidance of the support set. Learning is based on episodes available with annotations during training. At test time, the model is expected to segment a query

image with respect to a class of interest, again provided with corresponding support set, only this time the class of interest in the query and support set is novel and not previously seen.

Currently, the most leading FSS methods are the prototype-based ones [12, 25]. As in Fig. 1(a), prototype-based paradigm typically generates multiple foreground and/or background prototypes by utilizing mask average pooling and/or clustering over support features. These prototypes are supposed to contain representative information of the target object in the support images, thus their interactions with query features by cosine similarity, element-wise summation, and feature concatenation can produce necessary predictions for the objects in the query image. However, the predictions achieved by solely relying on such limited prototypes and *plain* operations are inevitably losing some intrinsic object details in the query image. For instance, as in Fig. 1(a), since the objects *plants* have holes and slots which are intrinsic details, the segmented objects cannot well cover these details, i.e., a defective over-segmentation is achieved in this case. In addition, in the presence of large object variations (e.g., appearance and scale) in FSS, it is usually difficult to comprehensively encode the adequate patterns of the target objects by solely considering the support information as in most previous prototype-based methods.

To address the above challenges, we propose a dynamic prototype convolution network (DPCN) to fully capture the intrinsic object details for accurate FSS. DPCN belongs to prototype-based methods yet with several elegant extensions and merits. Specifically, we first propose a dynamic convolution module (DCM) to achieve more adequate interaction between support and query features, thus leading to more accurate prediction for the query objects. As in Fig. 1(b), we leverage three dynamic kernels, i.e., a square kernel and two asymmetric kernels, generated from the support foreground features. Then three convolution operations are employed in parallel onto the query features using these dynamic kernels. This interaction strategy is simple yet important to comprehensively tackle large object variations (e.g., appearance and scales) and can capture the intrinsic object details. Intuitively, the square kernel is capable of capturing the main information of an object (e.g., main body of the *plant* in Fig. 1(b)); By contrast, asymmetric kernels (i.e., kernel with size $d \times 1$ or $1 \times d$ aim to capture the subtle object details, e.g., leaves in Fig. 1(b)). As such, DPCN equipped with DCM can better handle the intrinsic object details using an extremely simple way.

Moreover, to comprehensively encode the adequate patterns of the target objects, we propose a support activation module (SAM) and a feature filtering module (FFM) to mine as much object-related context information from query image as possible. Specifically, SAM generates sup-

port activation maps and initial pseudo query mask using high-level support and query features. Then the support prototypes and pseudo query foreground features are fused to generate a refined pseudo mask for the query image in FFM. Compared with the original pseudo query mask, the refined one contains more object foreground context while filtering some noise information. Therefore, rich object-related context information from both support and query images are aggregated to the final feature, leading to better segmentation performance. Our main contributions are as follows:

- We propose a dynamic prototype convolution network (DPCN) to capture the intrinsic object details for accurate FSS. To the best of knowledge, we are the first one to do this in the FSS domain.
- We propose a novel dynamic convolution module (DCM) to achieve adequate support-query interactions. DCM can serve as a plug-and-play component to improve existing prototype learning methods.
- We propose a support activation module (SAM) and a feature filtering module (FFM) to mine complementary information of target objects from query images.

2. Related work

2.1. Semantic Segmentation

Semantic segmentation is a classical computer vision task which aims to give pixel-wise prediction for an input image. Recently, various networks [14] have been actively designed to further improve the semantic segmentation results. For capturing more contextual informations, dilated convolution [33], pyramid pooling [40], and deformable convolution [3], are proposed to enlarge the receptive field. Meanwhile, some models leverage attention mechanisms [6, 27, 34, 37] to capture long-distance dependencies for semantic segmentation, which reach state-of-the-art performances. However, these semantic segmentation approaches still fail to preserve their initial performances when insufficient training data is provided.

2.2. Few-Shot Semantic Segmentation

Few-shot semantic segmentation (FSS) learns to segment target objects in query image given few pixel-wise annotated support images. Most existing FSS methods adopt two-branch architecture which implements meta-training on the base classes and then conducts meta-testing on the disjoint novel classes. OSLSM [20] is the first two-branch FSS model. Next, PL [4] introduces prototype learning paradigm, which generates prototypes from support images to guide the segmentation of query objects. Recently, many prototype-based FSS methods emerge in the research community, such as CANet [36], SG-One [39], PANet [26], PMMs [31], PFENet [25], and ASGNet [12]. The key idea

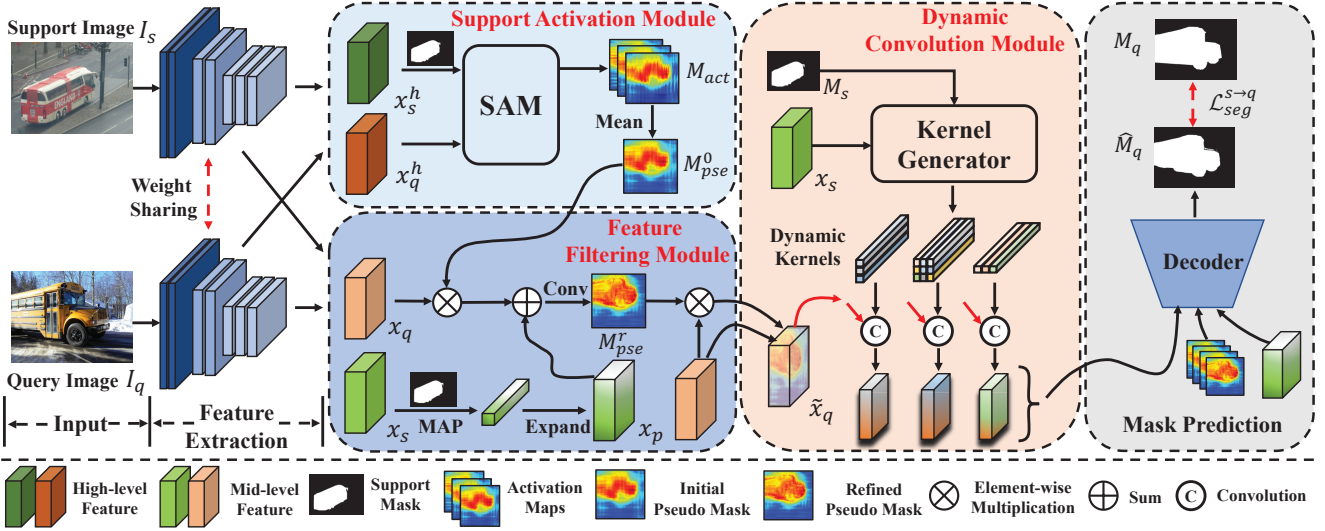


Figure 2. **Overall architecture of our proposed dynamic prototype convolution network (DPCN).** Firstly, **support activation module (SAM)** is introduced to generate activation maps and initial pseudo mask of target object in query image using high-level support and query feature. Then, **feature filtering module (FFM)** takes input mid-level support and query feature as well as initial pseudo mask to produce refined pseudo mask, which is leveraged to filter most background information in query feature. Meanwhile, **dynamic convolution module (DCM)** implements three groups of dynamic convolution over query features in parallel using kernels (multiple prototypes) generated from support foreground features, to propagate rich context information from support to query features. Finally, the updated features are concatenated and fed into a decoder for the final query segmentation mask prediction.

of these methods lies in generating or rearranging representative prototypes using different strategies, then the interaction between prototypes with query features can be formulated as a few-to-many matching problem. However, these prototype learning methods inevitably cause information loss due to limited prototypes. Therefore, graph-based methods have thrived recently as they try to preserve structural information with many-to-many matching mechanism. For instance, PGNet [36] applies attentive graph reasoning to propagate label information from support data to query data. SAGNN [30] constructs graph nodes using multi-scale features and performs k-step reasoning over nodes to capture cross-scale information. Most recently, HSNet [17] proposes to tackle the FSS task from the perspective of visual correspondence. It implements efficient 4D convolutions over multi-level feature correlation and achieves great success. Different from previous methods, we try to perform sufficient interaction between support and query features using dynamic convolution, and mine as much complementary target information from both support and query features.

2.3. Dynamic Convolution Networks

Dynamic convolution networks aim to generate diverse kernels and implement convolution over input feature with these kernels. Many previous works have explored the effectiveness of dynamic convolution in deep neural networks. DFN [11] proposes a dynamic filter network where filters are generated dynamically conditioned on input and achieves state-of-the-art performance on video and stereo

prediction task. [2] aggregates multiple parallel convolution kernels dynamically based upon their attentions, and it boosts both image classification and keypoint detection accuracy. Dynamic convolution is also used in DMNet [8] to adaptively capture multi-scale contents for predicting pixel-level semantic labels. The core of these methods is constructing multiple kernels from input features. Most recently, dynamic convolution is introduced into the few-shot object detection task by [38], which generates various kernels from the object regions in support image and then implements convolution over query feature using these kernels, leading to a more representative query feature. In this paper, we propose to generate dynamic kernels from foreground support feature to interact with query feature by convolution. Instead of only using square kernels as in [38], we also introduce asymmetric kernels to capture subtle object details. Experiments in §4.3 demonstrate well the effectiveness of our method.

3. Method

3.1. Problem Setting

We adopt the standard FSS setting, *i.e.*, following the episode-based meta-learning paradigm [23]. We start from classes \mathcal{C}_{tr} and \mathcal{C}_{ts} for the training set \mathcal{D}_{tr} and the test set \mathcal{D}_{ts} , respectively. The key difference between FSS and general semantic segmentation task is that \mathcal{C}_{tr} and \mathcal{C}_{ts} in FSS are disjoint, $\mathcal{C}_{tr} \cap \mathcal{C}_{ts} = \emptyset$. Both \mathcal{D}_{tr} and \mathcal{D}_{ts} consist of thousands of randomly sampled episodes, and each episode $(\mathcal{S}, \mathcal{Q})$ includes a support set \mathcal{S} , and a query set \mathcal{Q}

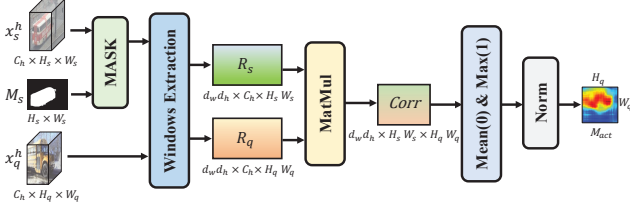


Figure 3. Illustration of the support activation module (SAM).

for a specific class c . For the k -shot setting, the support set that contains k image-mask pairs can be formulated as $\mathcal{S} = \{(I_s^i, M_s^i)\}_{i=1}^k$, where I_s^i represents i th support image and M_s^i indicates corresponding binary mask. Similarly, we define the query set as $\mathcal{Q} = \{(I_q, M_q)\}$, where I_q is the query image and its binary mask M_q is only available in the model training phase. In the meta-training stage, the FSS model takes as input \mathcal{S} and I_q from a specific class c and generates a predicted mask \hat{M}_q for the query image. Then the model can be trained with the supervision of a binary cross-entropy loss between M_q and \hat{M}_q . Finally, the model takes multiple randomly sampled episodes $(S_i^{ts}, Q_i^{ts})_{i=1}^{N_{ts}}$ from \mathcal{D}_{ts} for evaluation. Next, the 1-shot setting is adopted to illustrate our method for simplicity.

3.2. Overview

As in Fig. 2, our dynamic prototype convolution network (DPCN) consists of three key modules, i.e., support activation module (SAM), feature filtering module (FFM), and dynamic convolution module (DCM). Specifically, given the support and query images, I_s and I_q , we use a common backbone with shared weights to extract both mid-level and high-level features. We then have the SAM whose task is to generate an initial pseudo mask M_{pse}^0 for the target object in the query image. After SAM, a FFM follows, which aims to refine the pseudo mask and filter out irrelevant background information in the query feature. To incorporate relevant contextual information, we then employ the DCM, which learns to generate custom kernels from support foreground feature and employ dynamic convolution over query feature. We then feed the pseudo masks and features computed by the dynamic convolutions into a decoder to predict the final segmentation mask \hat{M}_q for the query image. Next, we describe each of the aforementioned modules in detail.

3.3. Support Activation Module

Inspired by PFENet [16, 25], recent FSS models [29, 30] usually leverage high-level features (e.g., conv5 of ResNet50) from the support and query set to generate the prior mask indicating the rough location of the target object in the query image. As this prior mask is usually obtained by element-to-element or square region-based matching between feature maps, a holistic context is not taken into account.

To counter this, with the support activation module we generate multiple activation maps of the target object in the query image using holistic region-to-region matching. Specifically, as in Fig. 3, SAM takes as input the high-level support feature $x_s^h \in \mathbb{R}^{C_h \times H_s \times W_s}$, the corresponding binary mask $M_s \in \mathbb{R}^{H_s \times W_s}$, as well as the high-level query feature $x_q^h \in \mathbb{R}^{C_h \times H_q \times W_q}$, where C_h is the channel dimension, H_s, W_s, H_q, W_q are the height and width of support and query feature, respectively.

To perform holistic matching, we first need to generate region features R_s and R_q with a fixed window operation \mathcal{W} sliding on the support and query features, respectively.

$$\begin{aligned} R_s &= \mathcal{W}(x_s^h \otimes M_s) \in \mathbb{R}^{d_h d_w \times C_h \times H_s W_s}, \\ R_q &= \mathcal{W}(x_q^h) \in \mathbb{R}^{d_h d_w \times C_h \times H_q W_q}, \end{aligned} \quad (1)$$

where \otimes stands for the Hadamard product and d_h, d_w are the window height and width. In our experiments we opt for symmetrical and asymmetrical windows, i.e., $(d_h, d_w) \in \{(5, 1), (3, 3), (1, 5)\}$ that are comprehensive and holistic regions, to account for possible object geometry variances. Having the region features, we proceed with matching by computing their cosine similarity, which generates the regional matching map $Corr \in \mathbb{R}^{d_h d_w \times H_s W_s \times H_q W_q}$. Notably, we utilize both square window (3,3) and asymmetrical windows (i.e., (5,1) and (1,5)), where square window can introduce more contextual information on regular part of target objects like the main body of object *plants*, asymmetrical windows can incorporate contextual details of slender part (e.g., leaves of *plants*).

We generate the final activation map $M_{act} \in \mathbb{R}^{H_q \times W_q}$ by taking the mean value among all regions and the maximal value among all support features followed by normalization operation. As we have three windows, we have three activation maps, $\{M_{act}^i\}_{i=1}^3$. In the end, we obtain the initial pseudo-mask $M_{pse}^0 \in \mathbb{R}^{H_q \times W_q}$, which indicates the rough location of target objects, by a mean operation.

3.4. Feature Filtering Module

As in Fig. 2, the feature filtering module is constructed on mid-level support and query features, i.e., $x_s \in \mathbb{R}^{C \times H \times W}$ and $x_q \in \mathbb{R}^{C \times H \times W}$ where C, H, W are channel, height, and width, respectively. Given x_s, x_q , and the initial pseudo mask M_{pse}^0 , the feature filtering module refines the pseudo mask, which is used to filter out irrelevant background information in the query image. We first apply masked average pooling on the features from the support set to get prototype vector $p \in \mathbb{R}^{C \times 1 \times 1}$:

$$p = \text{average}(x_s \otimes \mathcal{R}(M_s)), \quad (2)$$

where \mathcal{R} reshapes support mask M_s to be the same shape as x_s . Then, we expand the support prototype vector p to match the dimensions of the feature maps, $x_p \in \mathbb{R}^{C \times H \times W}$,

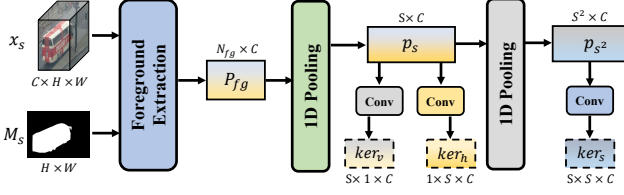


Figure 4. Illustration of kernel generator in DCM.

and combine the target object information from both the support and query features. We refine the pseudo mask with the help of a smaller network \mathcal{F} composed of a 2D convolution layer followed by a sigmoid function,

$$M_{pse}^r = \mathcal{F}((x_q \otimes \mathcal{R}(M_{pse}^0)) \oplus x_p) \in \mathbb{R}^{H \times W}, \quad (3)$$

where \oplus stands for the element-wise sum. Compared with M_{pse}^0 , M_{pse}^r gives more accurate estimation of the object location in the query image. Lastly, we obtain the final filtered query feature that discards irrelevant background by combining the feature x_q with the prior mask :

$$\tilde{x}_q = (x_q \otimes M_{pse}^r) \oplus x_q \in \mathbb{R}^{C \times H \times W}. \quad (4)$$

3.5. Dynamic Convolution Module

In the previous step we obtain a foreground feature from the query, which is minimally affected by irrelevant background. Still, the operations so far have been so that to provide a rough estimate of the location of the target object. For accurate segmentation, however, much finer pixel-level predictions are required. In the absence of significant data to train our filters on, we introduce dynamic convolutions. We illustrate DCM in Fig. 2, and Fig. 4 depicts the details of the kernel generator.

Dynamic convolutions rely on meta-learning to infer what are the optimal kernel parameters given a subset of features, agnostic to the unknown underlying class semantics. Specifically, we input the mid-level support feature x_s and the corresponding mask M_s to a kernel generator, which generates dynamic kernels, i.e., one group square kernel and two groups of asymmetrical kernels. Then, we carry out three convolution operations over the filtered query feature \tilde{x}_q using dynamic kernels. Firstly, we extract foreground vectors P_{fg} from support feature:

$$P_{fg} = \mathcal{F}_e(x_s \otimes M_s) \in \mathbb{R}^{N_{fg} \times C}, \quad (5)$$

where \mathcal{F}_e is the foreground extraction function without any learnable parameters, N_{fg} represents the number of foreground vectors. Next, two consecutive 1D pooling operations with kernel size S and S^2 are leveraged to obtain two groups of prototypes $p_s \in \mathbb{R}^{S \times C}$ and $p_{s^2} \in \mathbb{R}^{S^2 \times C}$:

$$p_s = \text{pool}_s(P_{fg}), p_{s^2} = \text{pool}_{s^2}(p_s). \quad (6)$$

As discussed above, we achieve dynamic convolution over query feature using a square kernel and two asymmetric kernels. As such, we use three parallel convolutional neural networks whose outputs are the generated kernel weights:

$$\begin{aligned} ker_v &= \mathcal{F}_{conv1}(p_s) \in \mathbb{R}^{S \times 1 \times C}, \\ ker_h &= \mathcal{F}_{conv2}(p_s) \in \mathbb{R}^{1 \times S \times C}, \\ ker_s &= \mathcal{F}_{conv3}(p_{s^2}) \in \mathbb{R}^{S \times S \times C}, \end{aligned} \quad (7)$$

where ker_v, ker_h, ker_s are the vertical, horizontal, and square kernel weights, respectively. \mathcal{F}_{conv1} , \mathcal{F}_{conv2} , and \mathcal{F}_{conv3} represent corresponding convolution sub networks, which are achieved by two consecutive 1D convolution layers. We emphasize that the above parameter generating networks *do not* share parameters. Given the vertical kernel ker_v , the query feature \tilde{x}_q can be enhanced as $\tilde{x}_q^v \in \mathbb{R}^{C \times H \times W}$:

$$\tilde{x}_q^v = \mathcal{F}_{dc}(\tilde{x}_q | ker_v), \quad (8)$$

where \mathcal{F}_{dc} denotes the dynamic convolution operation, and ker_v works as the kernel weight. Similarly, we can obtain other enhanced query features $\tilde{x}_q^h \in \mathbb{R}^{C \times H \times W}$ and $\tilde{x}_q^s \in \mathbb{R}^{C \times H \times W}$ with horizontal kernel ker_h and square kernel ker_s , respectively. With the sufficient interaction between query feature and dynamic support kernels, the object context in the generated query feature are enhanced.

Then, the enhanced query features $\tilde{x}_q^v, \tilde{x}_q^h, \tilde{x}_q^s$, support foreground feature x_p , support activation maps $\{M_{act}^i\}_{i=1}^3$, and refined pseudo mask M_{pse}^r are all reshaped to the same spatial size and concatenated to a representative feature $x_{out} \in \mathbb{R}^{(4C+4) \times H \times W}$:

$$x_{out} = \mathcal{F}_{cat}(\tilde{x}_q^v, \tilde{x}_q^h, \tilde{x}_q^s, x_p, \{M_{act}^i\}_{i=1}^3, M_{pse}^r), \quad (9)$$

where \mathcal{F}_{cat} is the concatenation operation in channel dimension. Finally, x_{out} is fed into a decoder to generate segmentation mask \hat{M}_q for query image I_q :

$$\hat{M}_q = \mathcal{F}_{cls}(\mathcal{F}_{ASPP}(\mathcal{F}_{conv}(x_{out}))), \quad (10)$$

where \mathcal{F}_{conv} , \mathcal{F}_{ASPP} , and \mathcal{F}_{cls} are three consecutive modules that constitute the decoder.

3.6. Extension to k -shot setting

So far we have focused on the one-shot setting, summarized in Fig. 2. For the k -shot setting, where more than one support images are available, most existing methods choose attention-based fusion or feature averaging. However, such a simple strategy does not make full use of the support information. By contrast, we can easily extend the dynamic convolutions to the k -shot setting and achieve substantial performance improvement. Specifically, given each support image-mask pair, we extract foreground vectors with Eq.

(5). By collecting all foreground vectors together, we get the overall support foreground vectors P_{fg} from k shots:

$$P_{fg} = (P_{fg}^1, P_{fg}^2, \dots, P_{fg}^k) \in \mathbb{R}^{N_{fg} \times C}, \quad (11)$$

where the number of foreground vectors is $N_{fg} = \sum_{i=1}^k N_{fg}^i$. By doing so, the kernel generator in DCM can generate more robust dynamic kernels, thus leading to more adequate interaction and accurate query mask estimation.

3.7. Training Loss

Our dynamic prototype convolution network is trained in an end-to-end manner with the binary cross-entropy loss (BCE). Given predicted mask \hat{M}_q and ground-truth mask M_q for query image I_q , we formulate the BCE loss between \hat{M}_q and M_q as our main loss:

$$\mathcal{L}_{seg}^{s \rightarrow q} = \frac{1}{hw} \sum_{i=1}^h \sum_{j=1}^w BCE(\hat{M}_q(i, j), M_q(i, j)). \quad (12)$$

Inspired by [26], we implement another branch to estimate the support mask using query image I_q and its corresponding predicted mask \hat{M}_q . Similar with Eq. (10), we get the predicted support mask \hat{M}_s . Then we get another loss by calculating BCE loss between \hat{M}_s and M_s :

$$\mathcal{L}_{seg}^{q \rightarrow s} = \frac{1}{h_s w_s} \sum_{i=1}^{h_s} \sum_{j=1}^{w_s} BCE(\hat{M}_s(i, j), M_s(i, j)), \quad (13)$$

where h_s and w_s are the height and width of ground-truth mask M_s for support image I_s . Note that both query and support mask prediction process share the same structure and parameters. In summary, the final loss is:

$$\mathcal{L} = \mathcal{L}_{seg}^{s \rightarrow q} + \lambda \mathcal{L}_{seg}^{q \rightarrow s}, \quad (14)$$

where λ is weight to balance the contribution of each branch and set to 1.0 in all experiments.

4. Experiments

4.1. Experimental Settings

Datasets. We follow [25] and adopt PASCAL-5ⁱ [20] and COCO-20ⁱ [18] benchmarks for evaluation. PASCAL-5ⁱ comes from PASCAL VOC2012 [5] and additional SBD [7] annotations. It contains 20 object classes split into 4 folds, which are used for 4-fold cross validation. For each fold, 5 classes are used for testing and the remaining 15 classes for training. COCO-20ⁱ is a more challenging benchmark, which is created from MSCOCO [13] and contains 80 object classes. Similarly, we split classes in COCO-20ⁱ into 4 folds with 20 classes per fold, for each fold, we utilize 20 classes for testing and the remaining 60 classes for training.

Metrics and Evaluation. Following [17, 25, 30], mean intersection over union (mIoU) and foreground-background IoU (FB-IoU) are adopted as our metrics for evaluation. While FB-IoU neglects object classes and directly averages foreground and background IoU, mIoU averages IoU values of all classes in a fold. During evaluation, 1,000 episodes are sampled from the test set for metrics calculation, and multi-scale testing is used like most existing FSS methods.

Implementation Details. We use ResNet-50 [9] and VGG-16 [22] pre-trained on ImageNet as our backbone networks. The backbone weights are fixed except for layer4, which is required to learn more robust activation maps on PASCAL-5ⁱ, meanwhile, all these weights are fixed for COCO-20ⁱ to pursue a faster training time. The model is trained with the SGD optimizer on PASCAL-5ⁱ for 200 epochs and COCO-20ⁱ for 50 epochs. The learning rate is initialized as 0.005 with batch size 8 (0.002 with batch size 32) for PASCAL-5ⁱ (COCO-20ⁱ). Data augmentation strategies in [25] are adopted in the training stage, and all images are cropped to 473 × 473 patches for two benchmarks. In addition, the window sizes in SAM are set to {5 × 1, 3 × 3, 1 × 5} for PASCAL-5ⁱ; since COCO-20ⁱ contains much larger amounts of training images with plentiful types (already having holistic contexts compared with images in PASCAL-5ⁱ), we only utilize the initial prior mask for speeding up our model training. The kernel sizes in DCM are set as 5 × 1, 5 × 5, and 1 × 5 for two datasets. We implement our model with PyTorch1.7 and conduct all the experiments with Nvidia Tesla V100 GPUs.

4.2. Comparisons with State-of-the-Arts

PASCAL-5ⁱ. We report the mIoU and FB-IoU under both 1-shot and 5-shot settings in Table A5. It can be seen that (i) DPCN achieves state-of-the-art performance under both 1-shot and 5-shot settings. Especially for the 1-shot setting, we surpass HSNet [17], which holds the previous state-of-the-art results, by 2.0% and 2.7% with VGG16 and ResNet50 as backbone networks, respectively. In addition, DPCN also presents comparable performance with HSNet under 5-shot setting while using less mid-level features. (ii) DPCN outperforms its baseline method with a large margin (e.g., mIoU 66.7% versus 61.4% with ResNet50 backbone for 1-shot setting), which is implemented with same architecture except for proposed components (i.e., SAM, FFM, and DCM). The results further demonstrate that DPCN can effectively mine complementary information from both support and query features to facilitate query image segmentation.

COCO-20ⁱ. COCO-20ⁱ is a more challenging benchmark which usually contains multiple objects and exhibits great variance. Table 2 presents the performance comparison of mIoU and FB-IoU on COCO-20ⁱ dataset. As can be seen, using VGG16 and ResNet50 as backbone, our model

Methods	Backbone	1-shot						5-shot					
		Fold-0	Fold-1	Fold-2	Fold-3	Mean	FB-IoU	Fold-0	Fold-1	Fold-2	Fold-3	Mean	FB-IoU
OSLSM (BMVC'17) [20]	VGG16	33.6	55.3	40.9	33.5	40.8	61.3	35.9	58.1	42.7	39.1	43.9	61.5
co-FCN (ICLRW'18) [19]	VGG16	36.7	50.6	44.9	32.4	41.1	60.1	37.5	50.0	44.1	33.9	41.4	60.2
AMP-2 (ICCV'19) [21]	VGG16	41.9	50.2	46.7	34.7	43.4	61.9	40.3	55.3	49.9	40.1	46.4	62.1
PFENet (TPAMI'20) [25]	VGG16	56.9	68.2	54.4	52.4	58.0	72.0	59.0	69.1	54.8	52.9	59.0	72.3
HSNet (ICCV'21) [17]	VGG16	59.6	65.7	59.6	54.0	59.7	73.4	64.9	69.0	64.1	58.6	64.1	76.6
PFENet (TPAMI'20) [25]	ResNet50	61.7	69.5	55.4	56.3	60.8	73.3	63.1	70.7	55.8	57.9	61.9	73.9
RePRI (CVPR'21) [1]	ResNet50	59.8	68.3	62.1	48.5	59.7	-	64.6	71.4	71.1	59.3	66.6	-
SAGNN (CVPR'21) [30]	ResNet50	64.7	69.6	57.0	57.3	62.1	73.2	64.9	70.0	57.0	59.3	62.8	73.3
SCL (CVPR'21) [35]	ResNet50	63.0	70.0	56.5	57.7	61.8	71.9	64.5	70.9	57.3	58.7	62.9	72.8
MLC (ICCV'21) [32]	ResNet50	59.2	71.2	65.6	52.5	62.1	-	63.5	71.6	71.2	58.1	66.1	-
MMNet (ICCV'21) [29]	ResNet50	62.7	70.2	57.3	57.0	61.8	-	62.2	71.5	57.5	62.4	63.4	-
HSNet (ICCV'21) [17]	ResNet50	64.3	70.7	60.3	60.5	64.0	76.7	70.3	73.2	67.4	67.1	69.5	80.6
Baseline	VGG16	58.4	68.0	58.0	50.9	58.8	71.2	60.7	68.8	60.2	52.2	60.4	74.3
DPCN	VGG16	58.9	69.1	63.2	55.7	61.7	73.7	63.4	70.7	68.1	59.0	65.3	77.2
Baseline	ResNet50	61.1	69.8	58.4	56.3	61.4	71.5	63.7	70.9	58.7	57.4	62.7	73.7
DPCN	ResNet50	65.7	71.6	69.1	60.6	66.7	78.0	70.0	73.2	70.9	65.5	69.9	80.7

Table 1. Comparison with state-of-the-arts on PASCAL-5ⁱ dataset under both 1-shot and 5-shot settings. mIoU of each fold, and averaged mIoU & FB-IoU of all folds are reported. Baseline results are achieved by removing three modules (i.e., SAM, FFM, and DCM) in DPCN.

Methods	Backbone	1-shot						5-shot					
		Fold-0	Fold-1	Fold-2	Fold-3	Mean	FB-IoU	Fold-0	Fold-1	Fold-2	Fold-3	Mean	FB-IoU
FWB (ICCV'19) [18]	VGG16	18.4	16.7	19.6	25.4	20.0	-	20.9	19.2	21.9	28.4	22.6	-
PFENet (TPAMI'20) [25]	VGG16	33.4	36.0	34.1	32.8	34.1	60.0	35.9	40.7	38.1	36.1	37.7	61.6
SAGNN (CVPR'21) [30]	VGG16	35.0	40.5	37.6	36.0	37.3	61.2	37.2	45.2	40.4	40.0	40.7	63.1
RePRI (CVPR'21) [1]	ResNet50	31.2	38.1	33.3	33.0	34.0	-	38.5	46.2	40.0	43.6	42.1	-
MLC (ICCV'21) [32]	ResNet50	46.8	35.3	26.2	27.1	33.9	-	54.1	41.2	34.1	33.1	40.6	-
MMNet (ICCV'21) [29]	ResNet50	34.9	41.0	37.2	37.0	37.5	-	37.0	40.3	39.3	36.0	38.2	-
HSNet (ICCV'21) [17]	ResNet50	36.3	43.1	38.7	38.7	39.2	68.2	43.3	51.3	48.2	45.0	46.9	70.7
SAGNN (CVPR'21) [30]	ResNet101	36.1	41.0	38.2	33.5	37.2	60.9	40.9	48.3	42.6	38.9	42.7	63.4
SCL (CVPR'21) [35]	ResNet101	36.4	38.6	37.5	35.4	37.0	-	38.9	40.5	41.5	38.7	39.9	-
Baseline	VGG16	32.1	36.1	35.2	32.3	33.9	60.1	35.0	40.1	37.1	36.5	37.2	61.8
DPCN	VGG16	38.5	43.7	38.2	37.7	39.5	62.5	42.7	51.6	45.7	44.6	46.2	66.1
Baseline	ResNet50	32.3	38.3	34.9	32.5	34.5	57.7	35.0	41.0	37.3	35.5	37.2	59.2
DPCN	ResNet50	42.0	47.0	43.2	39.7	43.0	63.2	46.0	54.9	50.8	47.4	49.8	67.4

Table 2. Comparison with state-of-the-arts on COCO-20ⁱ dataset under both 1-shot and 5-shot settings. mIoU of each fold, and averaged mIoU & FB-IoU of all folds are reported. Baseline results are achieved by removing three modules (i.e., SAM, FFM, and DCM) in DPCN.

SAM	FFM	DCM	1-shot mIoU					FB-IoU
			Fold-0	Fold-1	Fold-2	Fold-3	Mean	
✓	✓	✓	63.6	69.7	65.0	59.6	64.5	75.2
✓	✓	✓	67.1	71.1	63.2	60.0	65.4	76.0
✓	✓	✓	63.5	70.8	65.7	58.9	64.7	75.8
✓	✓	✓	65.7	71.6	69.1	60.6	66.7	78.0

Table 3. Ablation studies of main model components.

DPCN significantly outperforms recent methods under both 1-shot and 5-shot settings. With the ResNet50 backbone, DPCN achieves 3.8% and 2.9% of mIoU improvement over HSNet [17] (previous SOTA) under 1-shot and 5-shot settings, respectively. In addition, DPCN gains significant improvement over the baseline models. For example, DPCN with VGG16 backbone achieves 5.6% and 9.0% mIoU improvement over the baseline model, which proves the superiority of our model in such challenging scenarios.

Qualitative Results. In Fig. 5, we report some quantitative results generated from our DPCN and baseline model on the PASCAL-5ⁱ and COCO-20ⁱ benchmarks. Compared with the baseline, we can see that DPCN exhibits better per-

Kernel Size	1-shot mIoU					FB-IoU
	Fold-0	Fold-1	Fold-2	Fold-3	Mean	
3	65.2	70.4	68.5	59.4	65.9	77.5
5	65.7	71.6	69.1	60.6	66.7	78.0
7	65.5	70.7	69.3	59.0	66.1	77.5
9	65.9	70.8	68.8	59.7	66.3	77.7

Table 4. Ablation studies on kernel size.

formance in capturing object details. For instance, more tiny details are preserved in the segmentation of *plants* and *bike* (first two columns in Fig. 5). Refer to supplementary materials for more qualitative results.

4.3. Ablations

We conduct following ablation studies with ResNet-50 backbone under the 1-shot setting on PASCAL-5ⁱ dataset.

Components Analysis. DPCN contains three major components, i.e., support activation module (SAM), feature filtering module (FFM), and dynamic convolution module (DCM). Table 3 presents our validation on the effectiveness of each component. DCM, which is the most impor-

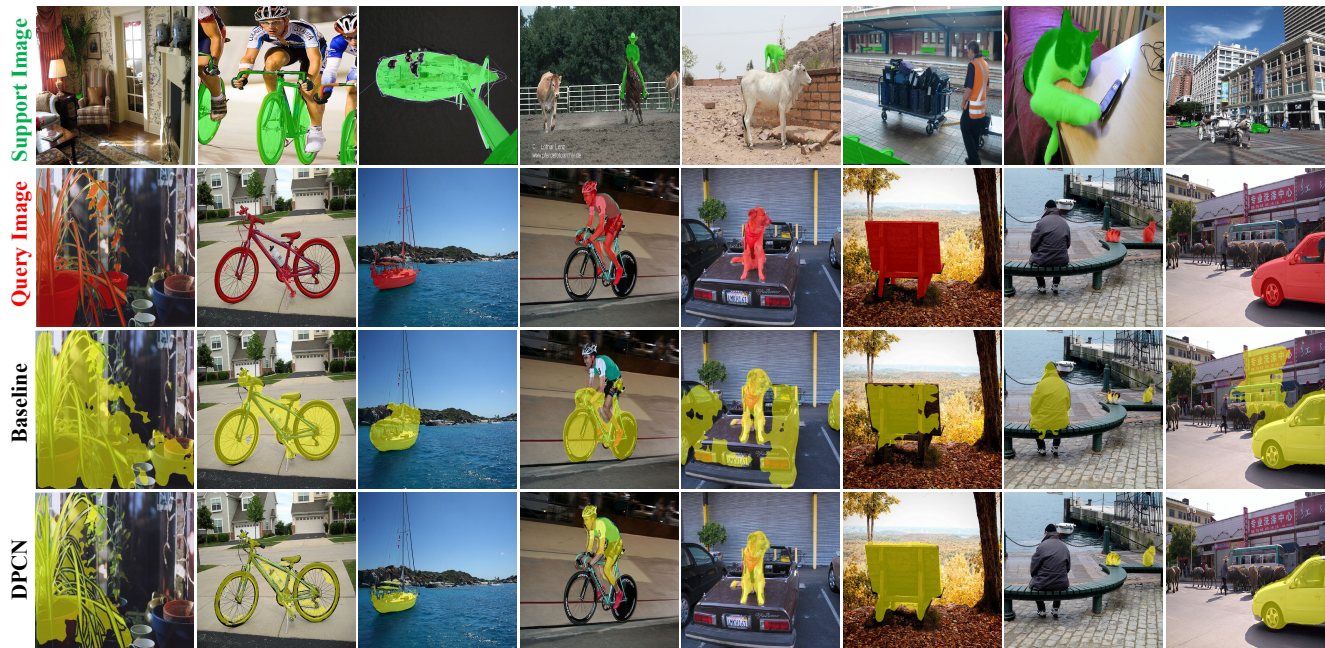


Figure 5. Qualitative results of our method DPCN and baseline model on PASCAL-5¹ and COCO-20² benchmarks. Zoom in for details.

Kernel variants	1-shot mIoU					FB-IoU
	Fold-0	Fold-1	Fold-2	Fold-3	Mean	
w/o DCM	63.6	69.7	65.0	59.6	64.5	75.2
5 × 5	64.7	71.2	65.3	58.7	65.0	76.2
1 × 5	64.9	71.0	65.2	59.7	65.2	75.4
1 × 5, 5 × 1	63.2	71.4	64.8	59.2	64.7	75.4
1 × 5, 5 × 5, 5 × 1	65.7	71.6	69.1	60.6	66.7	78.0

Table 5. Ablation studies on different kernel variants of DCM.

Methods	1-shot mIoU					FB-IoU
	Fold-0	Fold-1	Fold-2	Fold-3	Mean	
CANet	53.5	65.9	51.3	51.9	55.4	66.2
CANet+DCM	64.7	66.8	51.8	51.9	58.8	69.3
PFENet	61.7	69.5	55.4	56.3	60.8	73.3
PFENet+DCM	62.2	69.6	59.2	58.0	62.3	73.5

Table 6. Generalization ability of the proposed DCM.

tant component in our model, brings 2.2% improvement in mIoU. Meantime, SAM and FFM are also indispensable. By combing all three modules, DPCN achieves state-of-the-art performance.

Kernel Size in DCM. We take the kernel size from $\{3, 5, 7, 9\}$ to investigate the performance of our DPCN model. We can see from Table 4 that DPCN achieves the best and second-best performance when kernel sizes are 5 and 9, respectively. By the way, the performance drops slightly when the kernel sizes are 3 and 7. Therefore, the kernel size is set as 5 in all of our experiments.

Kernel Variants in DCM. Dynamic kernels are important components in DCM, so we evaluate the effectiveness of different kernel variants. As shown in Table 5, square kernel and asymmetric kernels achieve almost similar results. However, DPCN yields better performance when

choosing both square kernel and asymmetric kernels as dynamic kernels.

Generalization of DCM. DCM can be utilized as a plug-and-play module to further improve current prototype-based methods. To verify this, we apply DCM to CANet [36] and PFENet [25]. As shown in Table 6, DCM brings significant improvements on both CANet and PFENet.

5-shot Fusion Strategies. We compare our 5-shot fusion strategy (discussed in §3.6) with voting strategy [17], average [20] and OR [36] strategies on masks in Table ?? . We can see that our 5-shot fusion strategy achieves 3.2% mIoU improvement and outperforms other fusion strategies.

5. Conclusion

We propose a dynamic prototype convolution network (DPCN) with three major components (i.e., SAM, FFM, and DCM) to address the challenging FSS task. To better mine information from query image, we propose SAM and FFM to generate pseudo query mask and filter background information, respectively. Moreover, a plug-and-play module DCM is designed to implement sufficient interaction between support and query features. Extensive experiments demonstrate that DPCN achieves state-of-the-art results.

Acknowledgement

This work was partially funded by Elekta Oncology Systems AB and a RVO public-private partnership grant (PPS2102).

References

- [1] Malik Boudiaf, Hoel Kervadec, Ziko Imtiaz Masud, Pablo Piantanida, Ismail Ben Ayed, and Jose Dolz. Few-shot segmentation without meta-learning: A good transductive inference is all you need? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13979–13988, 2021. [7](#)
- [2] Yinpeng Chen, Xiyang Dai, Mengchen Liu, Dongdong Chen, Lu Yuan, and Zicheng Liu. Dynamic convolution: Attention over convolution kernels. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11030–11039, 2020. [3](#), [10](#)
- [3] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. Deformable convolutional networks. In *Proceedings of the IEEE international conference on computer vision*, pages 764–773, 2017. [2](#)
- [4] Nanqing Dong and Eric P Xing. Few-shot semantic segmentation with prototype learning. In *BMVC*, volume 3, 2018. [2](#)
- [5] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88(2):303–338, 2010. [6](#)
- [6] Jun Fu, Jing Liu, Haijie Tian, Yong Li, Yongjun Bao, Zhiwei Fang, and Hanqing Lu. Dual attention network for scene segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3146–3154, 2019. [2](#)
- [7] Bharath Hariharan, Pablo Arbeláez, Lubomir Bourdev, Subhransu Maji, and Jitendra Malik. Semantic contours from inverse detectors. In *2011 International Conference on Computer Vision*, pages 991–998. IEEE, 2011. [6](#)
- [8] Junjun He, Zhongying Deng, and Yu Qiao. Dynamic multi-scale filters for semantic segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3562–3572, 2019. [3](#)
- [9] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. [1](#), [6](#)
- [10] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4700–4708, 2017. [1](#)
- [11] Xu Jia, Bert De Brabandere, Tinne Tuytelaars, and Luc V Gool. Dynamic filter networks. *Advances in neural information processing systems*, 29:667–675, 2016. [3](#)
- [12] Gen Li, Varun Jampani, Laura Sevilla-Lara, Deqing Sun, Jonghyun Kim, and Joongkyu Kim. Adaptive prototype learning and allocation for few-shot segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8334–8343, June 2021. [2](#)
- [13] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. [6](#)
- [14] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3431–3440, 2015. [2](#)
- [15] Wenfeng Luo and Meng Yang. Semi-supervised semantic segmentation via strong-weak dual-branch network. In *European Conference on Computer Vision*, pages 784–800. Springer, 2020. [1](#)
- [16] Xiaoliu Luo, Zhuotao Tian, Taiping Zhang, Bei Yu, Yuan Yan Tang, and Jiaya Jia. Pfenet++: Boosting few-shot semantic segmentation with the noise-filtered context-aware prior mask. *arXiv preprint arXiv:2109.13788*, 2021. [4](#)
- [17] Juhong Min, Dahyun Kang, and Minsu Cho. Hypercorrelation squeeze for few-shot segmentation. *arXiv preprint arXiv:2104.01538*, 2021. [3](#), [6](#), [7](#), [8](#)
- [18] Khoi Nguyen and Sinisa Todorovic. Feature weighting and boosting for few-shot segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 622–631, 2019. [6](#), [7](#)
- [19] Kate Rakelly, Evan Shelhamer, Trevor Darrell, Alyosha Efros, and Sergey Levine. Conditional networks for few-shot semantic segmentation. 2018. [7](#)
- [20] Amirreza Shaban, Shray Bansal, Zhen Liu, Irfan Essa, and Byron Boots. One-shot learning for semantic segmentation. *arXiv preprint arXiv:1709.03410*, 2017. [1](#), [2](#), [6](#), [7](#), [8](#)
- [21] Mennatullah Siam, Boris N Oreshkin, and Martin Jagersand. Amp: Adaptive masked proxies for few-shot segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5249–5258, 2019. [7](#)
- [22] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. [1](#), [6](#), [10](#)
- [23] Jake Snell, Kevin Swersky, and Richard S Zemel. Prototypical networks for few-shot learning. *arXiv preprint arXiv:1703.05175*, 2017. [3](#)
- [24] Guolei Sun, Wenguan Wang, Jifeng Dai, and Luc Van Gool. Mining cross-image semantics for weakly supervised semantic segmentation. In *European conference on computer vision*, pages 347–365. Springer, 2020. [1](#)
- [25] Zhuotao Tian, Hengshuang Zhao, Michelle Shu, Zhicheng Yang, Ruiyu Li, and Jiaya Jia. Prior guided feature enrichment network for few-shot segmentation. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, (01):1–1, 2020. [2](#), [4](#), [6](#), [7](#), [8](#), [10](#), [11](#), [12](#)
- [26] Kaixin Wang, Jun Hao Liew, Yingtian Zou, Daquan Zhou, and Jiashi Feng. Panet: Few-shot image semantic segmentation with prototype alignment. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9197–9206, 2019. [2](#), [6](#), [10](#)
- [27] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7794–7803, 2018. [2](#)
- [28] Yude Wang, Jie Zhang, Meina Kan, Shiguang Shan, and Xilin Chen. Self-supervised equivariant attention mecha-

- nism for weakly supervised semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12275–12284, 2020. 1
- [29] Zhonghua Wu, Xiangxi Shi, Guosheng Lin, and Jianfei Cai. Learning meta-class memory for few-shot semantic segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 517–526, 2021. 4, 7
- [30] Guo-Sen Xie, Jie Liu, Huan Xiong, and Ling Shao. Scale-aware graph neural network for few-shot semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5475–5484, 2021. 1, 3, 4, 6, 7
- [31] Boyu Yang, Chang Liu, Bohao Li, Jianbin Jiao, and Qixiang Ye. Prototype mixture models for few-shot semantic segmentation. In *European Conference on Computer Vision*, pages 763–778. Springer, 2020. 2
- [32] Lihe Yang, Wei Zhuo, Lei Qi, Yinghuan Shi, and Yang Gao. Mining latent classes for few-shot segmentation. *arXiv preprint arXiv:2103.15402*, 2021. 7
- [33] Fisher Yu and Vladlen Koltun. Multi-scale context aggregation by dilated convolutions. *arXiv preprint arXiv:1511.07122*, 2015. 2
- [34] Yuhui Yuan, Xilin Chen, and Jingdong Wang. Object-contextual representations for semantic segmentation. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part VI 16*, pages 173–190. Springer, 2020. 2
- [35] Bingfeng Zhang, Jimin Xiao, and Terry Qin. Self-guided and cross-guided learning for few-shot segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8312–8321, 2021. 7
- [36] Chi Zhang, Guosheng Lin, Fayao Liu, Rui Yao, and Chunhua Shen. Canet: Class-agnostic segmentation networks with iterative refinement and attentive few-shot learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 2, 3, 8, 10, 12
- [37] Hang Zhang, Kristin Dana, Jianping Shi, Zhongyue Zhang, Xiaogang Wang, Amrbrish Tyagi, and Amit Agrawal. Context encoding for semantic segmentation. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 7151–7160, 2018. 2
- [38] Lu Zhang, Shuigeng Zhou, Jihong Guan, and Ji Zhang. Accurate few-shot object detection with support-query mutual guidance and hybrid loss. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14424–14432, 2021. 3
- [39] Xiaolin Zhang, Yunchao Wei, Yi Yang, and Thomas S Huang. Sg-one: Similarity guidance network for one-shot semantic segmentation. *IEEE transactions on cybernetics*, 50(9):3855–3865, 2020. 2
- [40] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2881–2890, 2017. 2
- [41] Qi Zhao, Binghao Liu, Shuchang Lyu, Xu Wang, and Yifan Yang. A self-distillation embedded supervised affinity attention model for few-shot segmentation. *arXiv preprint arXiv:2108.06600*, 2021. 11

Appendix A. Implementation details

We employ ResNet-50 [2] (VGG-16 [22]) pre-trained on ImageNet as our backbone networks. For ResNet-50, the dilation convolution (with stride size = 1) is introduced to ensure that the feature receptive fields of layer2, layer3, and layer4 preserve the same spatial resolution. The backbone weights are fixed except for layer4, which is required to learn more robust activation maps. The model is trained with a SGD optimizer for 200 and 50 epochs on the PASCAL-5ⁱ and the COCO-20ⁱ benchmarks, respectively. The learning rates are initialized as 0.005 and 0.002 with a poly learning rate schedule on PASCAL-5ⁱ and COCO-20ⁱ, respectively. The batch size is set as 8 on PASCAL-5ⁱ and 32 on COCO-20ⁱ. Our entire network is trained with the same learning rate during each epoch, except for layer4 of the backbone network, whose parameters starts back-propagation after training for multiple epochs to ensure a lower learning rate for fine-tuning. Data augmentation strategies in [25] are adopted in the training stage, and all images are cropped to 473 × 473 patches for two benchmarks. Besides, we leverage multi-scale testing strategy used in the most few-shot semantic segmentation methods for the model evaluation, and the original Groudtruth of the evaluated query image without any resize operations is adopted for the metric calculation. In addition, the window sizes in SAM are set to {5 × 1, 3 × 3, 1 × 5}, and the kernel sizes in DCM are set as 5 × 1, 5 × 5, and 1 × 5, respectively. We implement our model with PyTorch 1.7.0 and conduct all the experiments with Nvidia Tesla V100 GPUs and CUDA11.3.

Appendix B. Additional results and analyses

Kernel Generation Variants. The dynamic convolution module (DCM), in which we generate dynamic kernels from the support foreground and employ convolution over the query feature, is an essential component in our proposed method. Therefore, here we presents two kernel generation variants: (i) we generate both asymmetric and symmetric kernel weights in parallel (ii) we first generate asymmetric kernel weights, and the asymmetric kernel weights are further used to generate symmetric kernel weights. With a kernel size 5, we term these two variants as 5/25 (in parallel) and 5 → 25 (serial). As seen in Table A1, these two variants achieve similar performance (66.0 vs 66.7), which demonstrates the robustness of the DCM. **Experiments with Bounding Box Annotations.** Following PANet [26] and CANet [36], we evaluate our model with weakly-supervised annotation (i.e, bounding box instead of pixel-wise annotation) on the support set. As shown in Table A2, our method with the bounding box annotation achieves slightly inferior performance than that with expensive pixel-wise annotation. In addition, using bounding box annotation as supervi-

Method	1-shot mIoU					FB-IoU
	Fold0	Fold1	Fold2	Fold3	Mean	
5/25	65.9	70.6	66.9	60.5	66.0	77.0
5→25	65.7	71.6	69.1	60.6	66.7	78.0

Table A1. Ablation studies for the kernel generation variants.

Method	1-shot mIoU				
	Fold0	Fold1	Fold2	Fold3	Mean
PANet (Box)	-	-	-	-	45.1
CANet (Box)	-	-	-	-	52.0
Ours (Box)	59.8	70.5	63.2	55.5	62.3
Ours (Pixel)	65.7	71.6	69.1	60.6	66.7

Table A2. Comparison with the existing methods under the bounding box supervision under 1-shot setting.

Method	Backbone	PASCAL-5 ⁱ		COCO-20 ⁱ	
		MIoU	FB-IoU	MIoU	FB-IoU
<i>w/o</i> MS	VGG16	61.3	72.7	39.2	61.9
<i>w</i> MS	VGG16	61.7	73.7	39.5	62.5
<i>w/o</i> MS	ResNet50	65.7	77.4	41.5	62.7
<i>w</i> MS	ResNet50	66.7	78.0	43.0	63.2

Table A3. Effectiveness of multi-scale testing under the 1-shot setting on the PASCAL-5ⁱ and the COCO-20ⁱ benchmarks.

sion, our method also significantly outperforms both PANet and CANet. This experiment indicates the potential of our method in the segmentation task with weak supervision.

Multi-Scale Testing. As a post-processing method, multi-scale testing [41] is widely adopted in many semantic segmentation tasks. Many few-shot semantic segmentation methods as well as our proposed method DPCN also use this strategy to improve segmentation performance. We present our results with / without multi-scale testing (termed as *w* MS and *w/o* MS, respectively) in Table A3. The scales in our experiments are set as 473 and 641. We can see that (i) our model without the multi-scale testing also achieves state-of-the-art results using VGG16 and ResNet50 backbones, which further demonstrates the effectiveness of our proposed method (ii) multi-scale testing brings 1% mIoU improvement for our model with ResNet50 backbone, while a slight improvement when we use VGG16 as the backbone network.

Generalization Ability of DCM. The dynamic convolution module (DCM) can be used as a plug-and-play module to improve current prototype-based few-shot segmentation methods. We merge DCM into CANet and PFENet, and present the corresponding results under both 1-shot and 5-shot settings in Table A4. DCM further improve the performance of CANet and PFENet under both 1-shot and 5-shot settings, which shows the effectiveness of the DCM.

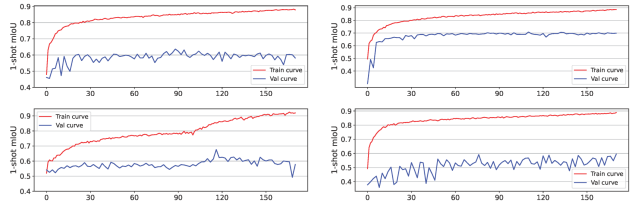


Figure A1. Training and Validation curves (x-axis: epochs, y-axis: 1-shot mIoU) on PASCAL-5ⁱ benchmark.

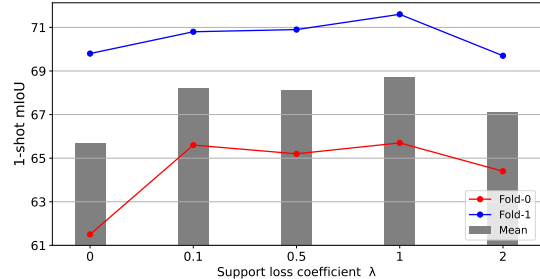


Figure A2. Ablation study for Support loss coefficient λ on fold0 and fold1 of PASCAL-5ⁱ benchmark.

Evaluation using Original Groundtruth. As in PFENet [25], we also evaluate our model with original groundtruth of the query image and the non-original one resized to the same size as training image size (473×473). We can find from Table A5 that our proposed model obtains similar performance with the original or non-original query groundtruth.

Training and validation performance. We present the performance (mIoU) changing process as the training epoch increases in Fig. A1. As can be seen, the mIoU of the training process are much better than that of the validation in each fold. Besides, the validation mIoU in fold0 and fold1 are relatively stable, while the validation mIoU in fold2 and fold3 fluctuate as the training process goes on.

Support Loss Coefficient λ . During the model training, we use the predicted query mask as a pseudo mask for predicting the support mask, which requires a support loss $\mathcal{L}_{seg}^{q \rightarrow s}$ for supervision. For the support loss coefficient λ , we take its value from $\{0, 0.1, 0.5, 1, 2\}$ to study its influence on our model. The performance of the fold0 and fold1 as well as their mean on the PASCAL-5ⁱ benchmark are used for illustration. As shown in Fig. A2, our model achieves best results when the support loss coefficient is set as 1. And λ is set as 1 in all our experiments.

Appendix C. Additional qualitative results

In this section, we present more qualitative results of our proposed DPCN and its baseline to demonstrate its

Methods	1-shot						5-shot					
	Fold-0	Fold-1	Fold-2	Fold-3	Mean	FB-IoU	Fold-0	Fold-1	Fold-2	Fold-3	Mean	FB-IoU
CANet (CVPR19) [36]	53.5	65.9	51.3	51.9	55.4	66.2	55.5	67.8	51.9	53.2	57.1	69.3
CANet (CVPR19)+DCM	64.7	66.8	51.8	51.9	58.8	69.3	65.3	67.2	52.7	52.9	59.5	70.1
PFENet (TPAMI'20) [25]	61.7	69.5	55.4	56.3	60.8	73.3	63.1	70.7	55.8	57.9	61.9	73.9
PFENet (TPAMI'20)+DCM	62.2	69.6	59.2	58.0	62.3	73.5	63.1	70.0	60.0	58.5	62.9	73.6

Table A4. Generalization ability of proposed DCM under both 1-shot and 5-shot settings.

Methods	Backbone	1-shot						5-shot					
		Fold-0	Fold-1	Fold-2	Fold-3	Mean	FB-IoU	Fold-0	Fold-1	Fold-2	Fold-3	Mean	FB-IoU
Original Groudtruth [36]	VGG16	58.9	69.1	63.2	55.7	61.7	73.7	63.4	70.7	68.1	59.0	65.3	77.2
Non-original Groudtruth	VGG16	59.3	69.5	63.3	55.8	62.0	73.8	64.0	71.2	68.4	59.0	65.7	77.2
Original Groudtruth [25]	ResNet50	65.7	71.6	69.1	60.6	66.7	78.0	70.0	73.2	70.9	65.5	69.9	80.7
Non-original Groudtruth	ResNet50	65.6	71.8	69.2	60.5	66.8	78.0	70.0	73.2	70.9	65.5	69.9	80.6

Table A5. Comparison with Original Groudtruth and Non-original Groudtruth

few-shot segmentation performance. Appearance and scale variations (more obvious in the COCO-20ⁱ benchmark) are the innate difficulty of the few-shot semantic segmentation task. Therefore, we first show some examples sampled from COCO-20ⁱ benchmark with large object appearance and scale variations in the Fig. A3 and Fig. A4, respectively. As can be seen, our model DPCN exhibits great superiority in alleviating appearance and scale variations. Besides, we also sample some examples from PASCAL-5ⁱ benchmark, and the qualitative results are presented in Fig. A5. Fur-

thermore, DPCN occasionally predicts more accurate segmentation than human-annotated ground-truth (Fig. A6), which further demonstrates the effectiveness of our method. Finally, we give some visualization of the support activation maps, initial pseudo mask as well as the refined pseudo mask in Fig. A7. We can see that the support activation maps can capture complementary object details in the query image, the initial pseudo mask gives rough pixel-wise location estimation of object, while the refined pseudo mask can estimate more accurate object location in the query image.

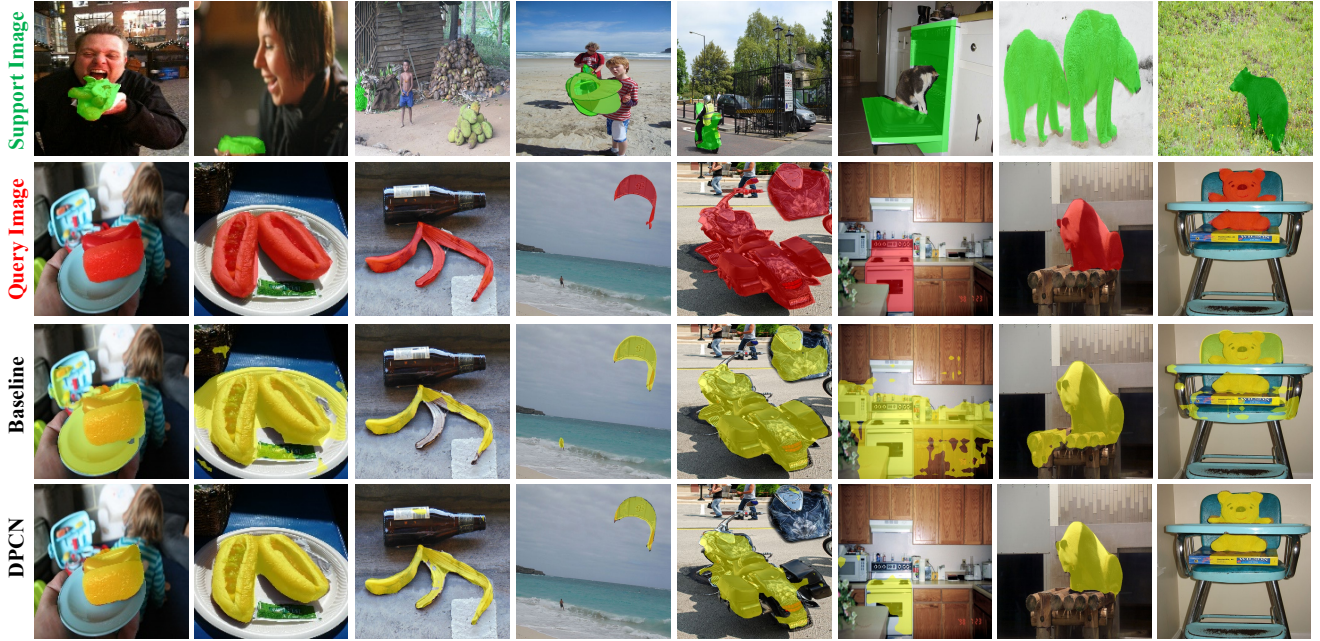


Figure A3. Qualitative results of our method DPCN and baseline model on COCO-20ⁱ benchmark with **large object appearance variations**. Zoom in for details.

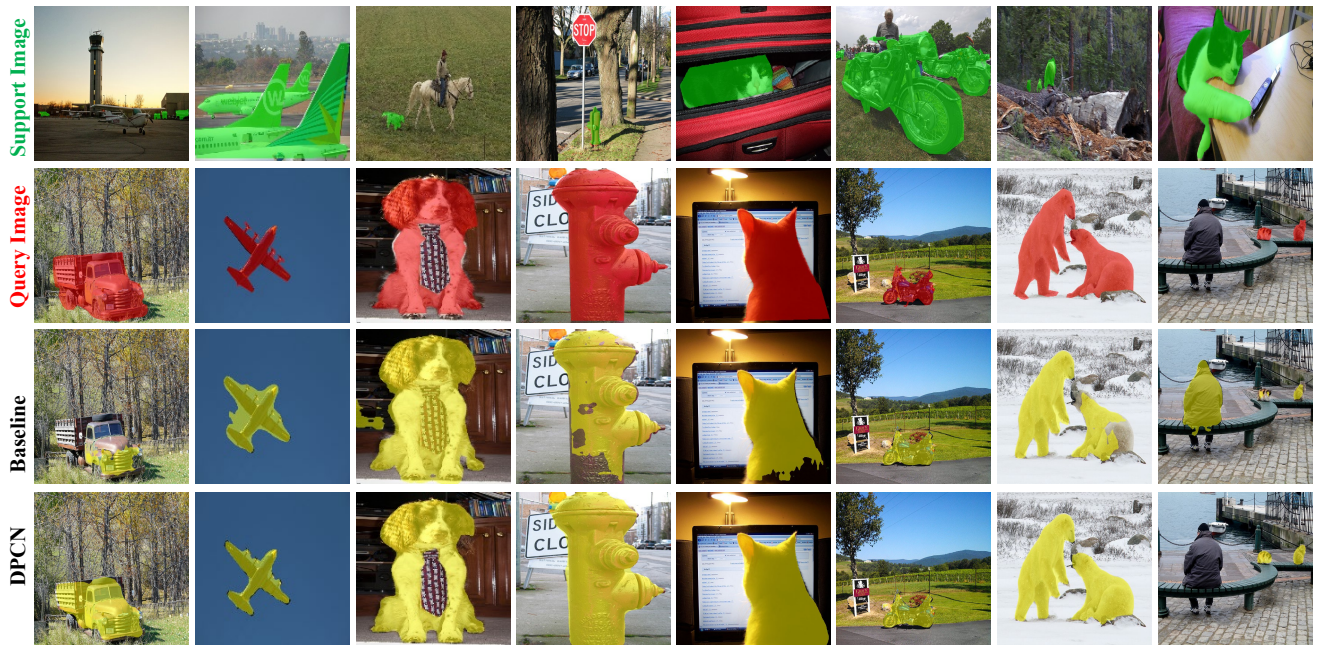


Figure A4. Qualitative results of our method DPCN and baseline model on COCO-20ⁱ benchmark with **large object scale variations**. Zoom in for details.

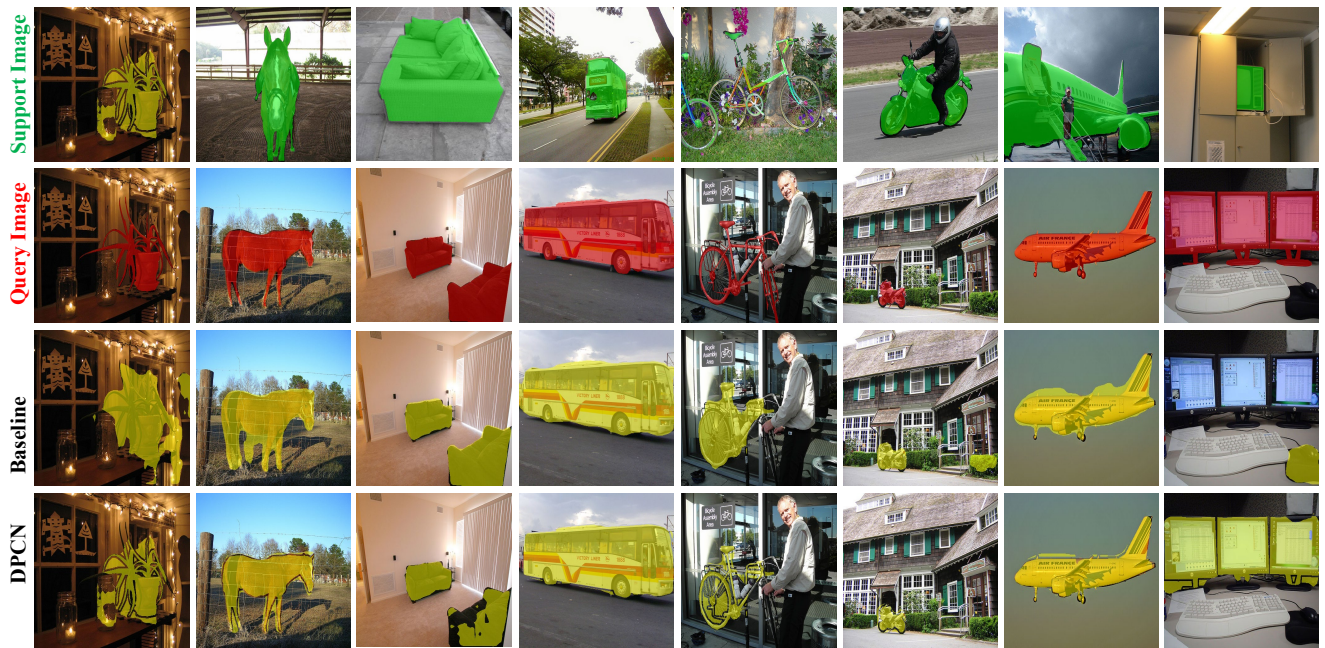


Figure A5. Qualitative results of our method DPCN and baseline model on PASCAL-5ⁱ benchmark. Zoom in for details.

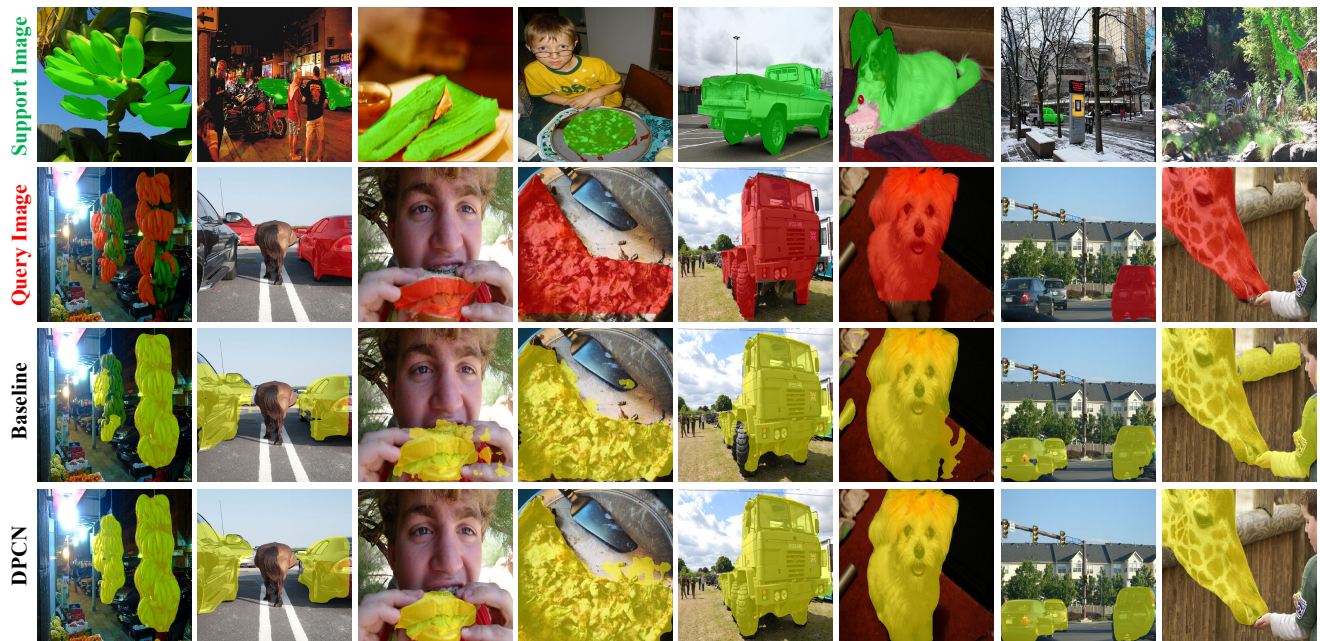


Figure A6. Our proposed DPCN occasionally predicts more accurate segmentation masks than human-annotated ground-truths. Examples are sampled from both PASCAL-5ⁱ and COCO-20ⁱ. Zoom in for details.

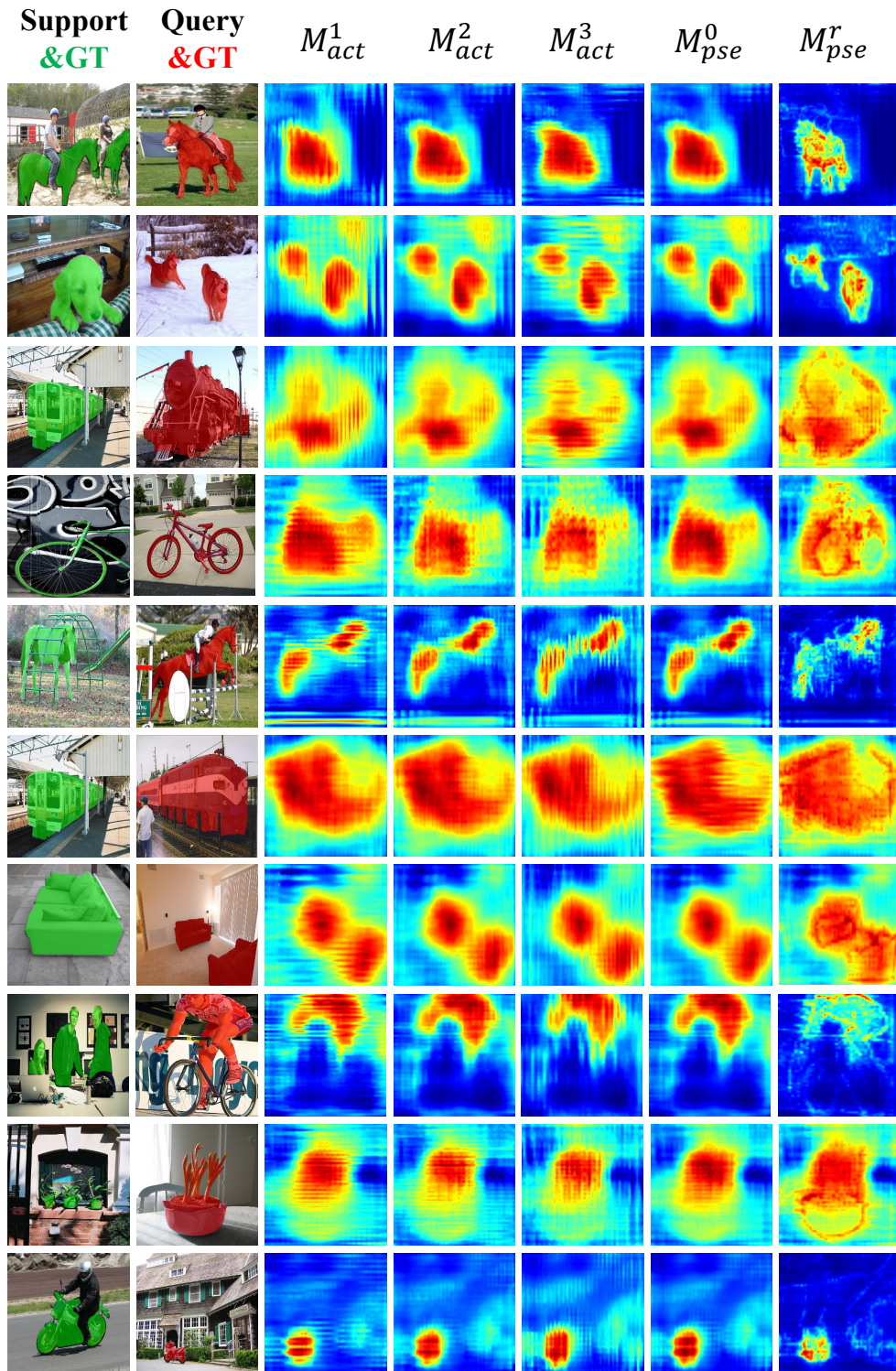


Figure A7. Visualization of the support activation maps $\{M_{act}^i\}_{i=1}^3$ and the initial pseudo mask M_{pse}^0 in the support activation module (SAM), as well as the refined pseudo mask M_{pse}^r in the feature filtering module (FFM). Zoom in for details.