

A New Few-shot Segmentation Network Based on Class Representation

Yuwei Yang¹, Fanman Meng^{2*}, Hongliang Li³, King N.Ngan⁴, Qingbo Wu⁵

School of Information and Communication Engineering

University of Electronic Science and Technology of China

Chengdu, China

¹ywyang@std.uestc.edu.cn, {²fmmeng, ³hlli, ⁵qbwu}@uestc.edu.cn, ⁴knngan@ee.cuhk.edu.hk

Abstract—This paper studies few-shot segmentation, which is a task of predicting foreground mask of unseen classes by a few of annotations only, aided by a set of rich annotations already existed. The existing methods mainly focus the task on “*how to transfer segmentation cues from support images (labeled images) to query images (unlabeled images)*”, and try to learn efficient and general transfer module that can be easily extended to unseen classes. However, it is proved to be a challenging task to learn the transfer module that is general to various classes. This paper solves few-shot segmentation in a new perspective of “*how to represent unseen classes by existing classes*”, and formulates few-shot segmentation as the representation process that represents unseen classes (in terms of forming the foreground prior) by existing classes precisely. Based on such idea, we propose a new class representation based few-shot segmentation framework, which firstly generates class activation map of unseen class based on the knowledge of existing classes, and then uses the map as foreground probability map to extract the foregrounds from query image. A new two-branch based few-shot segmentation network is proposed. Moreover, a new CAM generation module that extracts the CAM of unseen classes rather than the classical training classes is raised. We validate the effectiveness of our method on Pascal VOC 2012 dataset, the value FB-IoU of one-shot and five-shot arrives at 69.2% and 70.1% respectively, which outperforms the state-of-the-art method.

Index Terms—Few-shot Segmentation, Class Activation Map, Classification

I. INTRODUCTION

Deep learning has obviously improved the performance of many computer vision tasks such as classification [1], object detection [2] and segmentation [3]. However, its drawback is the serious dependence on manual annotations that are very time-consuming to be generated, especially for dense prediction tasks such as image segmentation. To this end, weakly-supervised [4] and semi-supervised manner [5], [6], [7] attract researchers’ attention.

Few-shot segmentation is a new semi-supervised segmentation task that predicts the foreground mask of unseen classes based on few of annotations only, aided by the annotations of classes that are already existed. The key step of such task is to learn the general knowledge from known classes that can be easily extended to unseen classes. The existing methods focus on solving the problem of *how to transfer segmentation cues from support images (labeled images) to query images (unlabeled images)*, and try to learn a general transformation module that has the capacity of transferring the segmentation

cues from support image to query image for various classes, so that the transferred cues can be used directly to guide the segmentation of query image for unseen classes. Based on such strategy, the existing few-shot segmentation framework is build as two-branch based segmentation network, where the two branches such as support branch and query branch are used to generate features for support image and query image respectively, and transformation module is added between the two branches to transfer the segmentation cues between the two branches. Based on such framework, the existing methods try to design new transformation module that is more general and efficient, and several types of transformation modules have been proposed [8], [9]. It is proved that the segmentation results can be enhanced by improving the transformation module. However, learning general transformation module is also proved to be a challenging task [8].

Different from the existing strategy, we solve few-shot segmentation in a new perspective of “*how to represent unseen classes by existing classes*”. The idea is based on the assumption that every class can be formed by a basic attribute set A . By learning the basic element set A from known classes, the representation of each unseen classes can also be obtained. Therefore, the prior of unseen classes can be established, and is further used to achieve the segmentation of unseen classes. In other words, unseen classes is firstly represented by existing classes with a representation module. Then, the representation is used to segment regions of unseen classes more efficiently.

Motivated by this, this paper proposes a new few-shot segmentation network based on the strategy of representation. A new representation module in terms of class activation map is proposed. The idea is to represent the images of unseen classes based on its activation regions by the classification model of known classes. A two-branch based few-shot segmentation network is proposed. Different from the classical two branches such as support and query branch in the existing few-shot segmentation framework, the first branch is the prior generation branch that generates object prior of query image in terms of class activation map by the CAM generation module. The second branch is the segmentation branch that segments the foreground of query image based on the prior. A new CAM generation module based on the task of highlighting unseen classes rather than training classes is proposed, which firstly learns the CAM extraction module

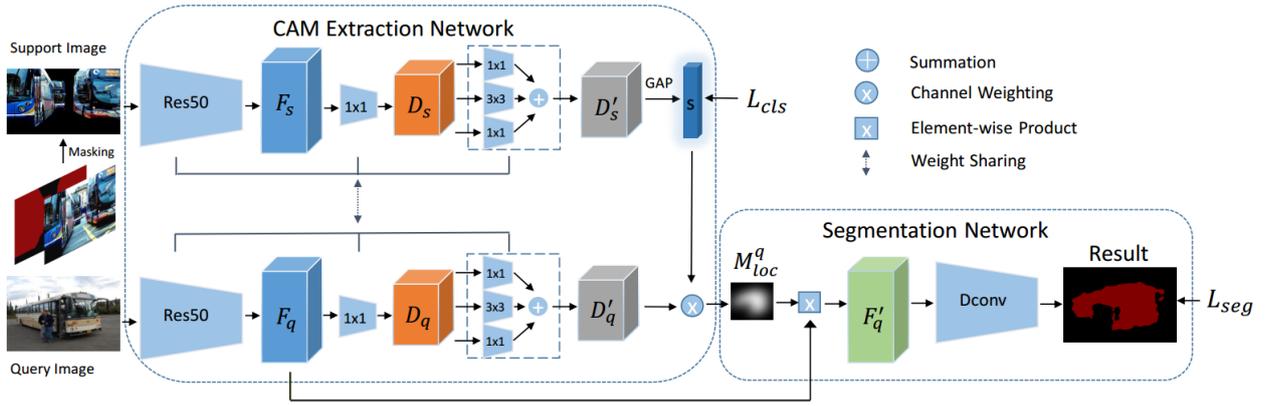


Fig. 1. The pipeline of the proposed method. The support image and manual annotation are sent to the classification sub-network to get the classification score S firstly. The query image is then sent to the classification sub-network to get class activation maps D_q for each existing class. Afterwards, a refined block is adopted to get maps D'_q . The class activation map M_{loc}^q of query image is obtained by averaging the maps D'_q weighted by S , which is then serves as the foreground probability map to update the original query features F_q to F'_q . Finally, the deconvolution module is used to output the segmentation result.

from support images of unseen classes, and then applies the CAM extraction module on query image to extract the prior map. We verify the proposed method on Pascal VOC 2012 dataset, the value FB-IoU of one-shot and five-shot arrives at 69.2% and 70.1% respectively, which outperforms several recent comparison methods.

II. PROPOSED METHOD

A. Problem Definition

Let $P = \{(I_s^i, Y_s^i)\}_{i=1}^k$ be support images and manual annotations for unseen classes, where k is the number of support images, Y_s^i is the binary annotation mask for image I_s^i . The goal of few-shot segmentation is to build a model $f(I_q, P)$ that outputs the binary mask \hat{M}_q for query images I_q based on P , aided by a set of existing training dataset $P' = \{(I_j', Y_j')\}_{j=1}^{n_s}$.

B. Overview

The proposed network is shown in Fig. 1, which consists of two sub-networks such as the CAM generation sub-network and the segmentation sub-network. Given a support image and a query image, the first sub-network is used to generate class activation map of query image, based on the classification model of known classes. A CAM generation module is proposed. The second one outputs the segmentation mask of query image. We next detail the two sub-networks.

C. CAM Generation Module for Unseen Class

Our goal is to represent unseen classes based on the existing classes in terms of class activation map, i.e., generating class activation map for unseen classes. Note that the classical CAM generation methods can not be used directly, as unseen classes is not considered in the classification model. Therefore, a new CAM generation module is proposed. Different from the classical CAM generation methods that use the gradient of back propagation to form the CAM, we form and learn the CAM extraction directly. We firstly learn the weight vector

$S = \{s_1, s_2, \dots, s_n\}$ for unseen classes, where s_i is the weight (similarity) of the i th class to unseen classes. Then, the probability map of query image is obtained by averaging the CAMs of query image highlighted by different known classes, weighted by the vector S . The detailed structure can be found in Fig. 1, where the proposed module consists of two steps such as learning the weight vector by support images, and the CAM extraction for query image based on S .

1) *Learning Weights by Support Images*: We intend to sufficiently use the manual annotations and support images to generate accurate CAM. A new CAM extraction module derived from the classical classification network for extracting CAM of unseen classes is proposed. The structure is shown in Fig. 1. The support image is used to obtain the weight S . Specifically, we firstly set zero to the background pixels in order to consider the foregrounds of the manual mask only. Then, Res50 is used to extract the convolution feature of support image. Based on the last deep convolution feature of Res50, a 1×1 convolution operation is applied to reduce the channel dimension of the convolution feature to the number of classes n , where i th channel means the class activation map of query image for i th class. We set the obtained features as D_s . Afterwards, a multi-scale feature extraction block is implemented to obtain final class activation map denoted as D'_s that has the same size to D_s . The refined block consists of feature extraction step and feature combination step. One 3×3 convolution operation and two 1×1 convolution operation are used to extract multi-scale features, and the multi-scale features are combined to obtain the refined class activation map. Finally, a global average pooling is applied on the multiple-scale class activation map to get the weight vector S .

Note that the proposed CAM extraction module is based on the classification network pre-trained on known classes. Here, we use the loss function Eq.1 to supervise the learning of the

classification network, i.e.,

$$L_{cls} = \log(1 + \exp(-\hat{S} \cdot S)) \quad (1)$$

where \hat{S} is the class-level labels, and S is the classification score.

It is seen that the weight S is very important to the CAM generation for unseen class. Here, although the extraction of weight is learned automatically, it can be simply considered as the similarities between unseen class and training classes. Therefore, the CAM of unseen class can be obtained by the sum of the CAMs of training classes weighted by their similarities.

2) *Extracting CAM for Query Images:* Given a query image of unseen class, we forward it to the classification network (Res50) to obtain the convolution feature F_q . A 1×1 convolution layer is then applied on F_q to obtain D_q with channel number n . After that, we implement the multi-scale feature extraction block to obtain feature $D'_q = \{D_1, D_2, \dots, D_n\}$ with the same size to D_q . Then, each channel feature of D'_q is weighted with the weight vector S , and the weighted channel features are summed to obtain the foreground prior M_{loc}^q of query image. Such process can be represented by

$$M_{loc}^q = \sum_{i=1}^n D_i \cdot s_i \quad (2)$$

The foreground prior M_{loc}^q is then forwarded to the segmentation network for the foreground prediction. It is seen that the object regions of unseen class is obtained based on the relationships between unseen class and known classes.

Some class activation maps of unseen class can be found in Fig. 2. It is seen that the object regions of unseen class are highlighted successfully, which demonstrates the effectiveness of the proposed CAM based representation module.

D. The Network of Few-shot Segmentation

After obtaining the class activation map M_{loc}^q , we next normalize it into $[0, 1]$ by

$$\bar{M}_{loc}^q = \frac{M_{loc}^q - \min(M_{loc}^q)}{\max(M_{loc}^q) - \min(M_{loc}^q)} \quad (3)$$

Then, the normalized class activation map serves as an attention module to weight the feature F_q of query image extracted from the backbone network. The filtered feature is then forwarded to a simple deconvolution block to obtain the final segmentation result \hat{M}_q . Such process can be represented as

$$\hat{M}(q) = Dconv(F_q * \underbrace{Cat(\bar{M}_{loc}^q, \dots, \bar{M}_{loc}^q)}_{\text{total of } n_F \text{ items}}) \quad (4)$$

where n_F is the number of channels of F_q . $Cat(\cdot)$ is the concatenation operation. The loss function L_{seg} is set to the cross-entropy loss.

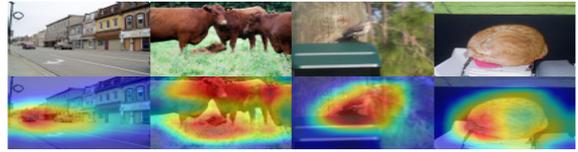


Fig. 2. The class activation maps of images from unseen classes by the proposed method. It is seen that the proposed method can generate class activation maps of unseen classes successfully.

TABLE I
THE DETAILS FOR CLASSIFYING THE 20 CLASSES INTO FOUR SUB-DATASETS. THERE ARE 4 SUB-DATASETS, AND $PASCAL - 5^i$ REPRESENTS THE i TH SUBSET, WHERE $i = \{0, 1, 2, 3\}$.

sub-dataset	corresponding classes
$PASCAL - 5^0$	aeroplane, bicycle, bird, boat, bottle
$PASCAL - 5^1$	bus, car, cat, chair, cow
$PASCAL - 5^2$	diningtable, dog, horse, motorbike, person
$PASCAL - 5^3$	potted plant, sheep, sofa, train, tv/monitor

E. Training and Inference

In training stage, because the proposed network is an end-to-end network, we train the network based on known classes directly. The details of the training setting can be found in Section III-A. It is worth noting that the class activation map is implicitly represented by feature D and D' here. Therefore, the CAM extraction manner for unseen class is learned automatically and directly without back propagation.

In the reference stage, the segmentation result is obtained directly by the network without fine-tuning.

III. EXPERIMENTS

We implement the proposed network on Pytorch. Adam optimizer is used to update parameters. One Nvidia Titan XP GPU is used. We set learning rate to $1e-4$ which decays 0.7 times per 10 epochs. Our backbone is set to Res50 pre-trained on ImageNet, and the top three layers is frozen during training. The size of input image is 320×320 .

A. Implementation Details

We implement experiment on Pascal VOC 2012 [10] dataset and its augmentation dataset SBD [11]. Similar to [8], we split the 20 classes into 4 sub-datasets, each of which contains 5 classes. Details can be found in Table I. For the four sub-datasets, one is selected as the unseen dataset for evaluation, the other three are used as known datasets for training. The image pairs for training are randomly selected from the training dataset. For fair comparison with the existing methods, we use the same seed for random sampling, and select the same 1000 image pairs in the testing stage. In the training stage, we use two of the three sub-datasets (ten classes) to train our classification network, and the rest one sub-dataset (five classes) as unseen classes to train the proposed network.

The FB-IoU [9] that calculates mean intersection over union of both foreground and background is used for objective evaluation.

TABLE II

THE FB-IOU VALUES OF ONE-SHOT AND FIVE-SHOT SEGMENTATION BY THE PROPOSED METHOD AND THE COMPARISON METHODS ON PASCAL VOC 2012 DATASET. THE BEST RESULTS ARE IN BOLD.

Methods	FG-BG[9]	OSLSM[8]	co-FCN[9]	PL[12]	SG-One[13]	A-MCG[14]	CA-Net[15]	Ours
One-shot	55.1	61.3	60.1	61.2	63.1	61.2	66.2	69.2
Five-shot	55.6	61.5	60.2	62.3	65.9	62.2	69.6	70.1

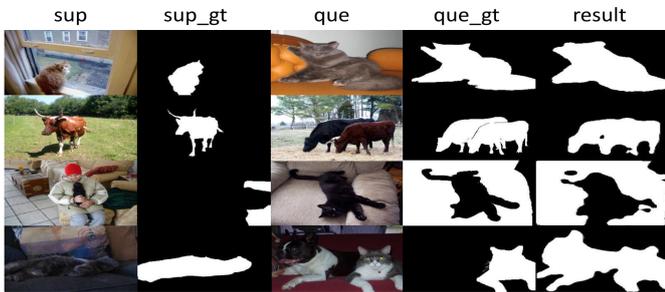


Fig. 3. The subjective results of the proposed method. From left to right: support images, ground-truth of support images, query images, ground-truth of query images and segmentation results, respectively.

B. Subjective Results

The subjective results of the proposed method are shown in Fig. 3. The support images, the ground-truth of support images, the query images, the ground-truth of query images and the segmentation results are displayed from left column to right column, respectively. The first three rows show successful results. It is seen that the proposed method segments objects from these images successfully. Meanwhile, the last row displays some case of failures, where the region of “Dog” is wrongly segmented as “Cat”. This is caused by the fact that “Cat” and “Dog” are very similar so that it intends to segment both of the object regions as foreground.

C. Objective Results and the Comparisons with Benchmarks

We next display the objective results in terms of FB-IoU value. In addition, we compare the proposed method with several recent few-shot segmentation methods. The results are displayed in Table II, where One-shot and Five-shot segmentation are considered. It is seen that the FB-IoU of One-shot segmentation on the four evaluation sub-dataset is 69.2%, which is better than the comparison methods. In addition, the value of the FB-IoU of Five-shot is 70.1%, which also outperforms the comparison methods. This demonstrates the effectiveness of the proposed method.

IV. CONCLUSION

In this paper, a new few-shot segmentation strategy based on class representation is proposed. A novel few-shot segmentation network is established. The proposed segmentation network consists of two branches. One is CAM generation network that obtains the class activation map of query image based on the classification model pre-trained on known image and support image of unseen class. The other is segmentation network that segments foreground from query image based

on the class activation map. A new CAM generation module for unseen class is proposed. The proposed method is verified on Pascal VOC dataset. Experimental results demonstrate the effectiveness of our proposed method with larger FB-IoU values.

ACKNOWLEDGMENT

This work was supported in part by the National Natural Science Foundation of China under Grant 61871087, Grant 61502084, Grant 61831005, and Grant 61601102, and supported in part by Sichuan Science and Technology Program under Grant 2018JY0141.

REFERENCES

- [1] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [2] S. Ren, K. He, R. Girshick, and S. Jian, “Faster r-cnn: Towards real-time object detection with region proposal networks,” in *International Conference on Neural Information Processing Systems*, 2015.
- [3] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam, “Encoder-decoder with atrous separable convolution for semantic image segmentation,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 801–818.
- [4] J. Ahn, S. Cho, and S. Kwak, “Weakly supervised learning of instance segmentation with inter-pixel relations,” 2019.
- [5] B. Lake, R. Salakhutdinov, J. Gross, and J. Tenenbaum, “One shot learning of simple visual concepts,” in *Proceedings of the Annual Meeting of the Cognitive Science Society*, vol. 33, no. 33, 2011.
- [6] J. Snell, K. Swersky, and R. Zemel, “Prototypical networks for few-shot learning,” in *Advances in Neural Information Processing Systems*, 2017, pp. 4077–4087.
- [7] A. Alfassy, L. Karlinsky, A. Aides, J. Shtok, S. Harary, R. Feris, R. Giryes, and A. M. Bronstein, “Laso: Label-set operations networks for multi-label few-shot learning,” 2019.
- [8] Z. L. I. E. B. Boots, A. Shaban, and S. Bansal, “One-shot learning for semantic segmentation,” *BMVC*, 2017.
- [9] T. D. A. E. S. Levine, K. Rakelly, and E. Shelhamer, “Conditional networks for few-shot semantic segmentation,” in *ICLR workshop*, 2018.
- [10] M. Everingham, S. M. A. Eslami, L. V. Gool, C. K. I. Williams, J. Winn, and A. Zisserman, “The pascal visual object classes challenge: A retrospective,” *International Journal of Computer Vision*, vol. 111, no. 1, pp. 98–136, 2015.
- [11] B. Hariharan, P. Arbeláez, L. Bourdev, S. Maji, and J. Malik, “Semantic contours from inverse detectors,” in *2011 International Conference on Computer Vision*. IEEE, 2011, pp. 991–998.
- [12] N. Dong and E. Xing, “Few-shot semantic segmentation with prototype learning,” in *BMVC*, vol. 1, 2018, p. 6.
- [13] X. Zhang, Y. Wei, Y. Yang, and T. Huang, “Sg-one: Similarity guidance network for one-shot semantic segmentation,” *CoRR*, vol. abs/1810.09091, 2018. [Online]. Available: <http://arxiv.org/abs/1810.09091>
- [14] T. Hu, P. Yang, C. Zhang, G. Yu, Y. Mu, and C. G. M. Snoek, “Attention-based multi-context guiding for few-shot semantic segmentation,” in *AAAI*, 2019.
- [15] C. Zhang, G. Lin, F. Liu, R. Yao, and C. Shen, “CANet: Class-agnostic segmentation networks with iterative refinement and attentive few-shot learning,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR’19)*, 2019.