

FSS-1000: A 1000-Class Dataset for Few-Shot Segmentation

Xiang Li¹ Tianhan Wei¹ Yau Pun Chen¹ Yu-Wing Tai² Chi-Keung Tang¹

¹HKUST

²Tencent

{xli@de, twei@ab, ypchen}@connect.ust.hk, yuwingtai@tencent.com, cktang@cs.ust.hk

Abstract

Over the past few years, we have witnessed the success of deep learning in image recognition thanks to the availability of large-scale human-annotated datasets such as PASCAL VOC, ImageNet, and COCO. Although these datasets have covered a wide range of object categories, there are still a significant number of objects that are not included. Can we perform the same task without a lot of human annotations? In this paper, we are interested in few-shot object segmentation where the number of annotated training examples are limited to 5 only. To evaluate and validate the performance of our approach, we have built a few-shot segmentation dataset, FSS-1000, which consists of 1000 object classes with pixelwise annotation of ground-truth segmentation. Unique in FSS-1000, our dataset contains significant number of objects that have never been seen or annotated in previous datasets, such as tiny daily objects, merchandise, cartoon characters, logos, etc.

We build our baseline model using standard backbone networks such as VGG-16, ResNet-101, and Inception. To our surprise, we found that training our model from scratch using FSS-1000 achieves comparable and even better results than training with weights pre-trained by ImageNet which is more than 100 times larger than FSS-1000. Both our approach and dataset are simple, effective, and easily extensible to learn segmentation of new object classes given very few annotated training examples. Dataset is available at <https://github.com/HKUSTCV/FSS-1000>

1. Introduction

Although unprecedented in the number of object categories when first released, contemporary image datasets for training deep neural networks such as PASCAL VOC [5] (19,740 images, 20 classes), ILSVRC [29] (1,281,167 images, 1,000 classes), and COCO [21] (204,721 images, 80 classes) are actually quite limited for visual recognition tasks in the real world: a rough estimate of the number of different objects on the Earth falls in the range of 500,000

to 700,000, following the total number of nouns in the English language. While the exact total number of visual object categories is smaller than these numbers, these large-scale datasets contribute less than 1% in total. Extending a new object category to existing datasets is a major undertaking because a lot of human annotation effort is required: in ImageNet, the mean number of images in a given class is 650. More importantly, observe that the number of images within each object category in ImageNet for instance can vary significantly, ranging from 1 to 3,047. This inevitably introduces undesirable biases which may have a detrimental effect on important tasks solely relying on pre-trained weights obtained using a dataset that is biased in both the choice of object classes (small number) and images within a given class (uneven distribution). Biases in existing datasets have also been recently reported [9, 20].

Thus, Few-Shot Learning has emerged as an attractive alternative for important computer vision tasks, especially when the given new dataset is very small and dissimilar so relying on the aforementioned pre-trained weights may not work well. Particularly relevant is image segmentation which requires extremely labor-intensive, pixelwise labeling for supervised learning. In few-shot segmentation, given an input consisting of a small support image set with labels (5 in this paper) and a query image set without labels, the learned model should properly segment the query images, even the pertinent objects belong to an object class unseen before.

There is *no* large-scale object dataset for few-shot segmentation. Previous research on few-shot segmentation relies on a manual split of the PASCAL VOC dataset to train and evaluate a new model [31, 26], but only 20 and 80 classes in the PASCAL VOC and COCO datasets respectively contain pixelwise segmentation information. Thus, building a large-scale object segmentation dataset is necessary to extensively and objectively evaluate the performance of our and future few-shot models.

FSS-1000 is the first large-scale dataset for few-shot segmentation with built-in object category hierarchy which emphasizes the number of object *classes* rather than the number of images. FSS-1000 is highly scalable: 10 new images with ground-truth segmentation are all it takes for new object class extension.

This research is supported in part by Tencent and the Research Grant Council of the Hong Kong SAR under grant no. 1620818.

Dataset	Images	Classes	Classification	Detection	Segmentation	Mean	Stddev
SUN [37]	131,067	3,819				39.22	717.68
ImageNet	3,200,000	5,247				650.02	526.03
Open Image	9,052,839	7,186				1409.62	14429.29
PASCAL VOC 2012	19,740	20				215.90	164.07
MS COCO	204,721	80				4492.13	7487.38
FSS-1000	10,000	1,000				10	0

Table 1. Large-scale datasets comparison. Mean and standard deviation are based on the expected number of images in each class.

Figure 1. Normalized image distribution. To make these datasets comparable, we normalize each dataset respectively in the total number of images (y-axis) and in the total number of object super-categories (x-axis) such that the area under each curve is 1 to make them comparable. All existing datasets are biased toward a number of object categories except FSS-1000 (red).

Our baseline network architecture is constructed by appending a decoder module to the relation network [33], which is a simple and elegant deep model effective and originally designed for few-shot image classification only. Reshaping the relation network into a fully-convolutional U-Net architecture [24], our extensive experimental results show that this baseline model trained from scratch on FSS-1000, which is less than 1% of the size of contemporary large-scale datasets, outperforms the model fine-tuned from weights pre-trained on ImageNet/COCO dataset. In addition, without any fine-tuning / re-training, our trained baseline network can be applied to any unseen classes directly with decent performance. With its excellent segmentation performance as well as extensibility, FSS-1000 and our baseline model are expected to make a lasting contribution to few-shot image segmentation. Please also refer to the supplemental materials for our extensive experimental results.

2. Related Work

We first review the relationship and difference between FSS-1000 and modern datasets aiming to solve image segmentation and few-shot classification. Then we review contemporary research on few-shot learning and semantic segmentation and discuss how we relate the few-shot segmentation to previous research.

Large-Scale Datasets When deep learning had started to become a dominating tool for computer vision, the importance of building large-scale datasets was emphasized for training deep networks. The PASCAL VOC [5] was the first to provide a challenging image dataset for object class recognition and semantic segmentation. The latest version VOC2012 contains 20 object classes and 9,993 images with segmentation annotations. Despite the absence of segmen-

tation labels, the Imagenet [4] is built upon the backbone of WordNet and provides image-level labels for 5,247 classes for training, out of which a subset of 1,000 categories are split out to form the ILSVRC [29] dataset. This challenge has made a significant impact on the rapid progress in visual recognition task and computer vision in recent years. The latest Open Image dataset [17] contains 7,186 trainable distinct object classes for classification and 600 classes for detection, making it the largest existing dataset with object classes and location annotations. Following the PASCAL VOC and ImageNet, the COCO segmentation dataset [21] includes more than 200,000 images with instance-wise semantic segmentation labels. There are 80 object classes and over 1.5 million object instances in COCO dataset.

In this paper, we instead focus on broadening the number of object classes in a segmentation dataset rather than increasing dataset size. Our FSS-1000 consists of 1,000 object classes, wherein each class we label 10 images with binary segmentation annotation. So in total, our dataset contains 10,000 images with pixelwise segmentation labels. We are particularly interested in segmentation due to its obvious benefits: segmentation captures the essential feature of an object without background; instance level segmentation can be ready from segmentation. The structure of our dataset is similar to widely-used datasets for few-shot visual recognition. For example, the Omniglot dataset [18] consists of 1,623 different handwritten characters of 50 different alphabets, which is equivalent to 1,623 object classes with 50 images in each class. The MiniImageNet, first proposed in [35], consists of 60,000 images with 100 classes each having 600 examples. But none of these few-shot learning datasets incorporate dense pixelwise segmentation labels, which is essential in training a deep network model for semantic segmentation.

Few-Shot Learning Recent research in few-shot classification can be classified into 1) learn a good initial condition for the network to be fine-tuned on extremely small training set, as proposed in [8, 27]; 2) rely on memory properties of RNN, introduced in [23, 30]; 3) learn a metric between few-shot samples and queries, as in [2, 10, 18, 16, 33]. We choose to extend the relation network [33] for few-shot segmentation because it is a simple, general and working framework. By concatenating the CNN feature maps between support images and query images, the relation module can consider the hidden relationship between these two sets of images guided by the loss function. In the original relation network, it uses the MSE loss to compare the fi-

Figure 2. Example images and their corresponding segmentation in FSS-1000. For the 12 super-categories here, 5 examples are shown, where the ground-truth segmentation map is overlaid in red in the corresponding image.

nal probability vector to the ground truth. In this paper, we simply modify the loss to calculate pixelwise differences between the segmentation ground truth and heatmap. In OSLSM [31], the authors proposed a two-branch network to solve few-shot segmentation. The network is quite complex, and their training set was limited to the PASCAL VOC dataset with only 20 object classes. Consequently, their feature extractor may suffer severe bias making it hard to be generalized to other objects. The guided network [26] can also suffer the same limitation on their dataset choice. Though point annotation can be used to guide the training of few-shot segmentation, the sparse annotation can seriously hamper accuracy.

Semantic Image Segmentation Previous research exploiting CNN to make dense prediction often relied on patchwise training [3, 6, 25] and pre- and post-processing of superpixels [6, 11]. In [22] the authors first proposed a simple and elegant fully convolutional network (FCN) to solve semantic segmentation. Notably, this is the first work which was trained end-to-end on a fully convolutional network for dense pixel prediction, which showed that the last layer feature maps from a good backbone network such as VGG-16 contain sufficient foreground features which can be decoded by the upsampling network to produce segmentation results. Intuitively, that is also the guiding principle behind our modification on relation network architecture. Though modern network architectures [12, 14, 19] achieve high accuracy in the COCO challenge by adding complex network modules and branches, these models cannot be adapted easily to segment new classes with few training examples.

3. FSS-1000

Recent few-shot datasets [18, 35] support few-shot classification but there is no large-scale few-shot segmentation dataset. In this section, we first introduce the details of data collection and annotation, then discuss the properties of FSS-1000. Table 1 and Figure 1 compare FSS-1000 with existing popular datasets. FSS-1000 targets at solving general objects few-shot segmentation problem. So datasets only focusing on sub-domain object categories in the world (e.g. handwritten characters, human faces and road scenes) are not included in the comparison.

3.1. Data Collection

Object Classes We first referred to the classes in ILSVRC [29] in our choice of object categories for FSS-1000. Consequently, FSS-1000 has 584 classes out of its 1,000 classes overlap with the classes in the ILSVRC dataset. We find ILSVRC dataset heavily biases toward animals, both in terms of the distribution of categories and number of images. Therefore, we fill in the other 486 by new classes unseen in any existing datasets. Specifically, we include more daily objects so that network models trained on FSS-1000 can learn from diverse artificial and man-made objects/features in addition to natural and organic objects/features where the latter was emphasized by existing large-scale datasets. Our diverse 1,000 object classes are further arranged in a hierarchy to be detailed in section 3.2.

Raw Images To avoid bias, the raw images were retrieved by querying object keywords on three different Internet search engines, namely, Google, Bing and Yahoo. We downloaded the first 100 results returned (or less if less than 100 images were returned) from a given search engine. No special criteria or assumption was used to select the candidates, however, due to the bias of Internet search engines, a large number of the images returned contain a single object photographed with sharp focus. In the final step, we intentionally included some images with a relatively small object, multiple objects or other objects in the background to balance the easy and hard examples of the dataset.

Images with aspect ratio larger than 2 or smaller than 0.5 were excluded. Since all images and their segmentation maps were to be resized to 224×224 , bad aspect ratio would destroy important geometric properties after the resize operation. For the same reason, images with height or width less than 224 pixels were discarded because they would trigger upsampling which would affect the image quality after resizing.

Pixelwise Segmentation Annotation We used Photoshop's "quick selection" tool which allows users to loosely select an object automatically, and refined or corrected the selected area to produce the desired segmentation. Figure 2 shows example images overlaid with their corresponding segmentation maps in FSS-1000.

3.2. Properties

This section summarizes the three desirable properties of FSS-1000:

Figure 3. Hierarchy of FSS-1000. Arrow represents “is a subclass of” relationship.

Figure 4. Example of instance annotation in the FSS-1000 dataset.

Scalability To extend FSS-1000 to include a new class, all it takes are 10 images with pixelwise binary segmentation labels for the new class. This is significantly easier than other datasets such as PASCAL VOC and COCO. First, the mean number of images in a given class is much larger than 10 in these datasets. Second, in these large-scale datasets the object classes need to be first pre-defined. Thus we believe binary annotation is a better annotation strategy in few-shot learning datasets, since it allows easy expansion of new object classes without concerning old object classes that have already been annotated.

Hierarchy Figure 3 shows examples of one sub-category for each given super-category in the dataset to illustrate the hierarchical structure of FSS-1000. The object classes are arranged hierarchically following a 3-level structure, while not every bottom-level subclass has a middle-level superclass. The top of the object hierarchy consists of 12 super-categories while the bottom contains the 1,000 classes as the leaf nodes. Note that this is strictly not a tree structure because a given class may belong to more than one superclass (e.g., an apple is both “fruit” and “food”).

Instance FSS-1000 dataset supports instance-level segmentation with instance segmentation labels in 758 out of the 1,000 classes in the dataset, which are significantly more classes than PASCAL VOC and MS COCO. One major difference between our dataset and PASCAL VOC / MS COCO instance level segmentation is that our dataset only annotates one type of objects in one image, despite there may be other object categories appearing in the background. We annotate at most 10 instances in a single image, which follows the same instance annotation principle adopted by COCO. Figure 4 shows examples of instance annotations in the dataset.

4. Methodology

4.1. Problem Formulation

In few-shot learning, the train-test split is on *object categories*, thus, all testing categories are unseen during training. In both training and testing, the input is divided

Figure 5. Our baseline network architecture using VGG-16 as backbone. The relation module is adapted from [33] where a decoder module is appended to produce the segmentation map. Both support and query features are concatenated to the decoder module via skip connection. More details of this standard architecture are available in supplemental materials.

into two sets, namely, the support set and the query set. The support set consists of samples with annotation, while the query set contains samples without annotation. In few-shot classification, the support set usually includes C classes and K training examples. This setting is defined as C -way- K -shot classification [7, 33]. In few-shot segmentation, we adopt this notation but extend the query output to be per-pixel classification of the query image, rather than a single class label. Specifically, in few-shot segmentation, the input-output pair is given by (X, Y) , where

$$L = \{l_{(i,j)}; i \in \{1, 2, \dots, C\}\}$$

$$X = \{(I_s, L_s, I_q); s \in \{1, 2, \dots, K\}\}$$

$$Y = \{y_{(i,j)}; i \in \{1, 2, \dots, C\}\}$$

$l_{(i,j)}$ is the ground-truth class label and $y_{(i,j)}$ represents the predicted class label for pixel (i, j) in a given image. I_s is the 3-channel RGB support image. For each support input X with image and label pair (I_s, L_s) , the model predicts a pixelwise classification map over query image I_q . Following the annotation strategy of FSS-1000, we set $C = 2$ and only focus on few-shot binary segmentation problem in this paper. However, a general C -way- K -shot segmentation could be solved by a union of C binary segmentation tasks.

4.2. Network Architecture

Pipeline Our network consists of three sub-modules: an encoder module E , a relation module R and a decoder module D . For a given input X to the network, the encoder E encodes the support and query images respectively into feature maps $E(I_s)$ and $E(I_q)$. For K -shot forwarding, we perform element-wise averaging over the depth channels of support feature maps, so that the encoder

module always produces support feature maps of the same depth regardless of the size of the support set.

The support and query feature maps are then combined in the relation module R . We choose channel-wise concatenation as the combination operation, while other choices such as parameter regression and nearest neighbors are possible and discussed in [26]. The relation module generates coarse segmentation results in low-resolution based on the concatenated feature maps. Finally, the coarse result is fed into the decoder module to restore the prediction map to the same resolution of the input. Figure 5 shows the entire workflow. In summary, the output is defined by

$$Y = D \left(R \left(\bigoplus_{s=1}^K E(I_s), E(I_q) \right) \right).$$

Loss function We use the cross entropy loss between the query prediction output and the ground-truth annotation to train our model. Specifically, under our binary few-shot segmentation setting, binary cross entropy (BCE) loss is adopted to optimize the parameters in the network:

$$\argmin_{\theta} \sum_{i,j} -L_{(i,j)} \log y_{(i,j)} + (1 - L_{(i,j)}) \log(1 - y_{(i,j)})$$

Mean square error (MSE) is also a widely used objective function for semantic segmentation task. Different from BCE loss, MSE models the problem as regression to the target output. Our experiments show that BCE and MSE loss achieve similar performance under our network setting.

4.3. Network Module Details

One can design his/her own or choose any popular feature extraction backbone such as VGG-16 [32], ResNet [13] and Inception [34] as the encoder module inside the network. The support and query features compose the combined feature map whose depth is twice the channel number of the last-layer output of the encoder. The relation module utilizes two 1×1 convolutional layers on the combined feature map to embed the relationship between the support features and query features. The decoder module is designed according to the number of downscale operations in the encoder module, which applies equivalent upsample blocks to restore the resolution back to the original input. In each upsample block stands a nearest neighbor upsampling layer and a convolutional layer. Skip connection is adopted between encoder and decoder feature maps, following the scheme proposed by U-Net [24]. We find it helpful to produce fine details in segmentation when information in the encoder feature maps are fused to the decoder module by channel-wise concatenation. ReLU activation is applied throughout the deep network except for the last layer's activation where Sigmoid is used in order to scale the output to a suitable range to calculate cross-entropy loss. More detail parameters of our architecture are provided in the supplemental materials.

Method	MeanIoU
VGG-16-BCEloss	80.12%
VGG-16-MSEloss	79.66%
ResNet-101-BCEloss	79.43%
ResNet-101-MSEloss	79.12%
InceptionV3-BCEloss	79.02%
InceptionV3-MSEloss	79.22%

Table 2. Different network settings to explore the best settings for our network architecture.

Method	MeanIoU
OSLSM-1shot [31]	70.29%
OSLSM-5shot	73.02%
Guided Network-1shot [26]	71.94%
Guided Network-5shot	74.27%
Ours-1shot	73.47%
Ours-5shot	80.12%

Table 3. Different few-shot segmentation networks trained and tested on FSS-1000.

Method	PASCAL-5 ⁰	PASCAL-5 ¹	PASCAL-5 ²	PASCAL-5 ³	Mean
OSLSM [31]	34.2%	57.9%	43.2%	37.8%	43.3%
GN [26]	33.1%	58.9%	44.3%	39.9%	44.1%
Ours	37.4%	60.9%	46.6%	42.2%	46.8%
PANet [36]	51.8%	64.6%	59.8%	46.5%	55.7%
CANet [39]	55.5%	67.8%	51.9%	53.2%	57.1%
Ours*	50.6%	70.3%	58.4%	55.1%	58.6%

Table 4. Comparison of different models on PASCAL-5¹. GN is Guided Network and Ours* is our model trained on FSS-1000. All models are using 5-shot setting.

5. Experiments

We conduct experiments to evaluate the practicability of FSS-1000 and the performance of our method under few-shot learning settings. We evaluate models with the same network architecture but trained on different datasets to show that FSS-1000 is effective for few-shot segmentation task. Different support sets and their influence on query results will be discussed. Finally we illustrate that models trained on FSS-1000 are capable to generalize the few-shot segmentation knowledge to new unseen classes. The metric we use is the intersection-over-union (IoU) of positive labels in a binary segmentation map. IoU is a standard metric and widely adopted in evaluating image segmentation methods. All the networks are implemented in PyTorch. We use Adam solver [15] to optimize the parameters. The learning rate is initially set to 10^{-3} (10^{-4} for fine-tuning) and halved for every 50,000 episodes. We train all the networks for 500,000 episodes.

Network setting To explore the best settings for our network, we train different models using a combination of different backbones and loss functions on FSS-1000. Table 2 tabulates the respective performance on VGG-16, ResNet-101 and InceptionNet as backbone, and BCE and MSE as loss function. Based on the result, we choose VGG-16 as feature extractor and use BCE loss in our model throughout the experimental section.

5.1. Benchmarks

5.1.1 FSS-1000

We train OSLSM and Guided Network on FSS-1000 to provide benchmarks and justify our dataset. Table 3 shows that

No.	ImageNet	FSS	fsCOCO	FSS (test set)	fsCOCO (test set)
I				71.34%	42.11%
II				79.30%	47.99%
III				80.12%	48.31%
IV				82.66%	50.56%

Table 5. Comparison of models trained and tested on different datasets. Each model (row) shows the training stages, e.g., model I uses the pre-trained weights from ImageNet then fine-tuned on fsCOCO’s training classes, and finally tested on the novel test classes in both FSS and fsCOCO. All learning rates are initially set to 10^{-4} except the model trained without using ImageNet pre-trained weights, which is set to 10^{-3} .

Figure 6. MeanIoU of superclasses in FSS-1000 tested with models trained on fsPASCAL, fsCOCO and FSS-1000. Bars at the bottom indicate the percentage of the number of categories overlapping with FSS-1000 in the corresponding dataset.

our adapted relation network achieves the best results on FSS-1000. Moreover, ours is the only model whose 5-shot training boosts the accuracy by over 10% compared to the 1-shot case. We believe that embedding multiple support images at the input end of the network and encouraging the feature extractor to consider correlation between multiple support images and the query image is the appropriate way to design k-shot ($k > 1$) segmentation network, rather than simply combining 1-shot prediction [31] or merging high-level features of multiple supports [26].

5.1.2 PASCAL-5ⁱ

To compare with previous few-shot methods, we train and test our network on PASCAL-5ⁱ [31]. Table 4 shows that our simple baseline model (Ours) marginally outperforms OSLSM and Guided Network. More importantly, our model trained only on FSS-1000 without fine-tuning on PASCAL-5ⁱ (Ours*) achieves much better results compared to models trained on PASCAL-5ⁱ (Ours), exceeding the state-of-the-art performance of the very recent [39, 36].

5.2. Effect of Pre-training

We compare our network model trained on different datasets to demonstrate the effectiveness of FSS-1000 in few-shot segmentation. Since there are no publicly available few-shot image segmentation datasets, we convert PASCAL VOC 2012 and COCO datasets by setting the desired foreground class label as positive and all others as negative, followed by the identical clean-up stage described in section 3.1 to the binarized labels. Two new datasets

Figure 7. Image results of our baseline model respectively trained on fsPASCAL, fsCOCO and FSS-1000. Support labels and predicted segmentation are overlaid in red in corresponding support images and query images. Ground truth labels for query images are in green. The classes in the first two rows are present in fsPASCAL and fsCOCO whereas the rest are unique in FSS-1000.

Figure 8. MeanIoU of superclasses in FSS-1000 tested with k-shot models ($k = 1, 3, 5, 7$).

are thus produced: fsPASCAL and fsCOCO. There are respectively 4,318 image and label pairs in 20 object classes in fsPASCAL, which consists of 15 training classes and 5 test classes, and 48,015 image and label pairs in 80 object classes in fsCOCO, containing 60 training classes and 20 test classes. The generation of these datasets are in line with the settings in [39].

For FSS-1000, we build the validation/test set by randomly sampling 20 distinct sub-categories from the 12 super-categories; the other images and labels are used in training. The train/validation/test split used in the experiments consists of 5,200/2,400/2,400 image and label pairs. Each test set of fsPASCAL, fsCOCO and FSS are designed to be disjoint with all the training sets in terms of classes for fair comparison.

Table 5 tabulates the performance of different models. For each model (row), the marks in sequence indicate the dataset(s) used in pre-training stages with the last mark indicating the dataset used in fine-tuning. Model III has only one indicating that it is exclusively trained on the dataset.

Using the pre-trained weights from ImageNet, Model II trained on FSS-1000 outperforms the fsCOCO-trained model I on both test sets by a large margin of 8% and 5.8%, which is due to the FSS training set containing the COCO

Figure 9. Effect of different support sets. The leftmost support of each row is used to generate 1-shot results. For each class, we show the result of a good support set followed by a bad support set in the next row.

	Human (PS)	Human (GrabCut)	CPU	GPU
Time	180m32s	53m22s	9m13s	16.9s
95%+ IOU	100%	71.4%	58.4%	58.4%
90%+ IOU	100%	80.4%	70.4%	70.4%
80%+ IOU	100%	91.0%	87.4%	87.4%
70%+ IOU	100%	95.8%	90.2%	90.2%

Table 6. 500 test images are randomly sampled from FSS-1000 to compare time and accuracy performance of labeling segmentation data between humans and few-shot model.

training classes, but with more variety. Notably, *without* using any pre-trained weights Model III achieves slightly better results compared to Model II, which substantiate our claim that bias in feature extractor does exist in models pre-trained and/or trained on a dataset unevenly distributed in object categories and images within each class.

Interestingly, Model IV pre-trained on FSS-1000 and fine-tuned on fsCOCO achieves the best result on both test sets, outperforming Model III exclusively trained on FSS and the model I pre-trained on ILSRVC fine-tuned on fsCOCO. We believe the former is due to the addition of more data, and the latter is due to the difference in requirement of feature maps ideal for classification and segmentation task. Intuitively, semantic segmentation requires more accurate low-level features to produce fine details in segmentation map, while classification focuses on high-level features for image understanding. Therefore, we argue that pre-training with FSS-1000 serves as a good alternative for ImageNet pre-training in few-shot semantic segmentation.

Overall, models trained on fsCOCO produce quite good results in test classes that are similar to COCO training classes. For these classes, sometimes their segmentation results are better in local details compared to the results produced by models trained on FSS-1000 due to more variations in the training set. However, it failed in classes significantly different from the 60 COCO training classes. The somewhat limited variation in object categories in existing datasets makes it hard for models trained on them to generalize to more unseen classes under the few-shot setting.

On the other hand, models trained on FSS-1000 classes can handle these cases. Quantitative results and qualitative results are shown in Figure 6 and Figure 7 respectively. Results on fsPASCAL and further comparisons are provided in supplementary material.

5.3. Effect of Support Set

We train four different models, using 1, 3, 5 and 7 support images respectively, to study how different number of support images influence the accuracy of few-shot segmentation. Two important observations can be summarized from Figure 8.

First, more support images generally boost the segmentation accuracy because more variations of color, pose, and scale of the object are included. However, the performance increase becomes negligible when more than 5 support images are given. Due to this bottleneck effect, we set up most of the experiments under the 5-shot setting.

Second, the accuracy boost is different among different classes. For easy cases (e.g. rigid objects), the improvement is not obvious because a single support image is enough for the deep network to capture and distinguish strong features of the object. For hard cases (e.g. deformable objects), more support images are essential for the network to learn the complex shapes to make correct segmentation.

Figure 9 demonstrates the effect of support set, which shows that scale and pose of the object to be segmented are the most important characteristics to guide few-shot semantic segmentation on FSS-1000. Since FSS-1000 does not explicitly consider scale variations (future work), a tiny or oversized object in the support set is not a good reference for segmentation. Significant differences in scales can mislead the network to capture wrong feature contents in the query. Besides, significantly different poses in support and query sets can result in bad segmentation results, due to the intrinsic fragility to rotation in CNN features.

5.4. Auto-Labeling on Novel and Unseen Classes

Traditionally a large number of human-annotated images are required to train a deep network for segmenting a new class. Table 6 tabulates the tradeoff in time and accuracy for annotating 500 test images in FSS-1000 by humans (using Photoshop and GrabCut [28] algorithm) and our few-shot segmentation.

With its good accuracy and time tradeoff, despite the current limitations in scale invariance aforementioned, FSS-1000 allows us to automatically segment a novel object category by just providing a few support examples *without* re-training or fine-tuning a given model. We pick a number of very novel classes unseen by FSS-1000, and label 5 images of each class serving as the support set. Figure 10 shows the test results which demonstrates that our model trained on FSS-1000 is capable of generalizing to these unseen classes. More extensive results on novel classes are included in supplementary materials.

Figure 10. Test results for unseen classes. From top to bottom: *android robot*; the *river* from UC Merced Land Use Dataset [38]; a large *cell* image cropped into patches; herds of *sheep*, *penguin* from Oxford penguin counting dataset [1]; flock of wild *goose*; different images of fields of *sunflower* depict various scales in the presence of occlusion and perspective distortion.

For example, *android robot* is an unreal object unseen in FSS-1000. In cartography from satellite images which often come in overlapping image tiles, cartographers need to label only 5 images or tiles and our system can automatically segment the rest, such as recognizing *river* in our example where saliency detection does not work in general. The *cell* example shows the good potential of FSS-1000 in instance segmentation which significantly contributes to cell counting in medical image analysis where, for instance, a patient’s health directly correlates to his or her red blood cell count. With the advance of whole slide images (WSI) in which the width and height often exceed 100,000 pixels (and thus many cells to count), using our few-shot segmentation trained on FSS-1000, pathologists only need to label 5 image relevant regions and then the rest of the WSI will be automatically labeled. Although manual corrections for missed or wrong cells may still be necessary given the current accuracy, comparing with exhaustive labeling which requires hours or even days to complete, the potential contribution of FSS-1000 is substantial. Similarly, the related

Figure 11. Iterative few-shot segmentation. Left and right show respectively the support sets and results before and after including corrected failure cases in the support set. Complete testing set of *Eiffel Tower* is available in the supplemental material.

animal examples of *sheep*, *penguin* and wild *goose* show FSS-1000’s potential for large-scale instance segmentation. Finally, our baseline backbone network is not very robust to scale variance, occlusion and background noises (future work). In *sunflower*, the segmentation results for instances too big or too small (especially for images with depth of field where faraway sunflowers are out of focus) become incomplete or even totally omitted. Despite that, FSS-1000 still reports limited success.

5.5. Iterative Few-Shot Segmentation

Our few-shot segmentation successively benefits from support sets improved easily by including failure cases after correction in each pass. Consider the Eiffel Tower unseen by FSS-1000 in Figure 11 where we manually label 200 images for quantitative evaluation (IoU). The first support set (left) did not have sufficient view and scale variations and did not see clearly the bottom part of the tower which resulted in its incomplete segmentation in some test cases. After mining a few of such hard cases, correcting and including them in the second support set (right), the previous hard cases could now be correctly segmented. We believe that few-shot segmentation performed in stages can offer an immediate performance boost.

6. Conclusion

Few-shot learning/segmentation is an emerging attractive alternative where only a few training examples are required. However, there is no existing large-scale dataset for few-shot segmentation. In this paper, we address the limitation of existing large-scale datasets in their biases and lack of scalability, and build the first few-shot segmentation dataset FSS-1000 emphasizing class diversity rather than dataset size. We adapt the relation network architecture to few-shot segmentation. This baseline few-shot segmentation model, trained exclusively on FSS-1000 without using pre-trained weights, achieves higher accuracy than previous methods including on test sets unseen by FSS-1000. We further demonstrated the efficacy and potential of FSS-1000 in large-scale segmentation on totally unseen classes without re-training or fine-tuning, and showed its promise on few-shot instance segmentation and iterative few-shot recognition tasks.

References

- [1] Carlos Arteta, Victor Lempitsky, and Andrew Zisserman. Counting in the wild. In *ECCV*, 2016. 8
- [2] Luca Bertinetto, João F. Henriques, Jack Valmadre, Philip H. S. Torr, and Andrea Vedaldi. Learning feed-forward one-shot learners. In *NIPS*, 2016. 2
- [3] Dan C. Ciresan, Alessandro Giusti, Luca Maria Gambardella, and Jürgen Schmidhuber. Deep neural networks segment neuronal membranes in electron microscopy images. In *NIPS*, 2012. 3
- [4] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR*, 2009. 2
- [5] Mark Everingham, Luc Van Gool, Christopher K. I. Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. In *IJCV*, 2010. 1, 2
- [6] Clément Farabet, Camille Couprie, Laurent Najman, and Yann LeCun. Learning hierarchical features for scene labeling. In *TPAMI*, 2013. 3
- [7] Li Fei-Fei, Rob Fergus, and Pietro Perona. One-shot learning of object categories. In *TPAMI*, 2006. 4
- [8] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *ICML*, 2017. 2
- [9] Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A. Wichmann, and Wieland Brendel. Imagenet-trained CNNs are biased towards texture; increasing shape bias improves accuracy and robustness. In *ICLR*, 2019. 1
- [10] Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett. Prototypical networks for few-shot learning. In *NIPS*, 2017. 2
- [11] Bharath Hariharan, Pablo Andrés Arbeláez, Ross B. Girshick, and Jitendra Malik. Simultaneous detection and segmentation. In *ECCV*, 2014. 3
- [12] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross B. Girshick. Mask r-cnn. In *ICCV*, 2017. 3
- [13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 5
- [14] Dai Jifeng, He Kaiming, and Sun Jian. Instance-aware semantic segmentation via multi-task network cascades. In *CVPR*, 2016. 3
- [15] Diederik P. Kingma and Jimmy Lei Ba. Adam: A Method for Stochastic Optimization. In *ICLR*, 2015. 5
- [16] Gregory R. Koch. Siamese neural networks for one-shot image recognition. In *ICML Workshop*, 2015. 2
- [17] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Mallocci, Tom Duerig, and Vittorio Ferrari. The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale. In *arXiv:1811.00982*, 2018. 2
- [18] Brenden M. Lake, Ruslan Salakhutdinov, and Joshua B. Tenenbaum. Human-level concept learning through probabilistic program induction. In *Science*, 2015. 2, 3
- [19] Yi Li, Haozhi Qi, Jifeng Dai, Xiangyang Ji, and Yichen Wei. Fully convolutional instance-aware semantic segmentation. In *CVPR*, 2017. 3
- [20] Tsung-Yi Lin, Priya Goyal, Ross B. Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *ICCV*, 2017. 1
- [21] Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays, Pietro Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollár. Microsoft COCO: Common Objects in Context. In *ECCV*, 2014. 1, 2
- [22] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *TPAMI*, 2017. 3
- [23] Tsendsuren Munkhdalai and Hong Yu. Meta networks. In *ICML*, 2017. 2
- [24] Philipp Fischer Olaf Ronneberger and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, 2015. 2, 5
- [25] Pedro H. O. Pinheiro and Ronan Collobert. Recurrent convolutional neural networks for scene labeling. In *ICML*, 2014. 3
- [26] Kate Rakelly, Evan Shelhamer, Trevor Darrell, Alexei A. Efros, and Sergey Levine. Few-shot segmentation propagation with guided networks. In *ICLR Workshop*, 2018. 1, 3, 5, 6
- [27] Sachin Ravi and Hugo Larochelle. Optimization as a model for few-shot learning. In *ICLR*, 2017. 2
- [28] Carsten Rother, Vladimir Kolmogorov, and Andrew Blake. "grabcut": interactive foreground extraction using iterated graph cuts. In *SIGGRAPH*, 2004. 7
- [29] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. In *IJCV*, 2015. 1, 2, 3
- [30] Adam Santoro, Sergey Bartunov, Matthew Botvinick, Daan Wierstra, and Timothy P. Lillicrap. Meta-learning with memory-augmented neural networks. In *ICML*, 2016. 2
- [31] Amirreza Shaban, Shray Bansal, Zhen Liu, Irfan Essa, and Byron Boots. One-shot learning for semantic segmentation. In *BMVC*, 2017. 1, 3, 5, 6
- [32] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015. 5
- [33] Flood Sung, Yongxin Yang, Li Zhang, Tao Xiang, Philip HS Torr, and Timothy M Hospedales. Learning to compare: Relation network for few-shot learning. In *CVPR*, 2018. 2, 4
- [34] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott E. Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *CVPR*, 2015. 5
- [35] Oriol Vinyals, Charles Blundell, Timothy P. Lillicrap, Koray Kavukcuoglu, and Daan Wierstra. Matching networks for one shot learning. In *NIPS*, 2016. 2, 3
- [36] Kaixin Wang, Jun Hao Liew, Yingtian Zou, Daquan Zhou, and Jiashi Feng. Panet: Few-shot image semantic segmentation with prototype alignment. 5, 6
- [37] Jianxiong Xiao, Krista A. Ehinger, James Hays, Antonio Torralba, and Aude Oliva. Sun database: Exploring a large collection of scene categories. In *IJCV*, 2014. 2

- [38] Yi Yang and Shawn Newsam. Bag-of-visual-words and spatial extensions for land-use classification. In *ACM GIS*, 2010. 8
- [39] 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 *CVPR*, 2019. 5, 6