

Semi-supervised few-shot learning for medical image segmentation

Abdur R Fayjie^{1,2}, Reza Azad¹, Marco Pedersoli¹, Claude Kauffman², Ismail Ben Ayed¹, and Jose Dolz¹

¹ ETS Montreal, Montreal, Canada

² CRCHUM Montreal, Montreal, Canada
abdur-razzaq.fayjie.1@ens.etsmtl.ca

Abstract. Recent years have witnessed the great progress of deep neural networks on semantic segmentation, particularly in medical imaging. Nevertheless, training high-performing models require large amounts of pixel-level ground truth masks, which can be prohibitive to obtain in the medical domain. Furthermore, training such models in a low-data regime highly increases the risk of overfitting. Recent attempts to alleviate the need for large annotated datasets have developed training strategies under the few-shot learning paradigm, which addresses this shortcoming by learning a novel class from only a few labeled examples. In this context, a segmentation model is trained on episodes, which represent different segmentation problems, each of them trained with a very small labeled dataset. In this work, we propose a novel few-shot learning framework for semantic segmentation, where unlabeled images are also made available at each episode. To handle this new learning paradigm, we propose to include surrogate tasks that can leverage very powerful supervisory signals –derived from the data itself– for semantic feature learning. We show that including unlabeled surrogate tasks in the episodic training leads to more powerful feature representations, which ultimately results in better generability to unseen tasks. We demonstrate the efficiency of our method in the task of skin lesion segmentation in two publicly available datasets. Furthermore, our approach is general and model-agnostic, which can be combined with different deep architectures.

Keywords: Few-shot learning · Semantic segmentation · CNN

Semantic segmentation is of vital importance in medical imaging, as it can assist in the treatment, diagnosis and follow-up of many diseases. Despite the automation of this task has been widely studied in the last decades, the recent advances in deep learning are driving progress on this problem. Particularly, Convolutional Neural Networks (CNN) have achieved state-of-the-art performance in a breadth of medical image segmentation problems, such as brain tissue [4,7], heart structures [3], and abdominal organs [8]. Nevertheless, a main limitation of these high-performing models is the strong need of large labeled datasets for training, which hampers their scalability to novel or rare categories.

To alleviate this issue, few-shot learning [17] has recently emerged as an efficient alternative to traditional fully supervised learning strategies. In this context, the CNN is trained to learn novel categories with only a few labeled images, which are typically referred to as *support* images. Then, the knowledge derived from the *support* images is employed to guide the segmentation of images containing the novel classes, known as *queries*. To reproduce the scenario found during testing, the network is trained on the labeled examples following the episodic training paradigm [23]. This is, at each step we sample k labeled examples from n novel categories to form an episode, which is used to train the segmentation network. By doing this, the network is trained to extract useful information for all the different episodes but prevents the specialization of a particular novel task. Recent techniques to improve the generability on this scenario integrate a mask average pooling strategy, that masks out irrelevant features based on the *support* masks [15,19,28]. Wang et al. [24] further improves generability to new classes with an additional novel prototype alignment regularization between support and query images. In other recent works [12,26], deep attention has been exploited to learn attention weights between support and query images for further label propagation. Nevertheless, these methods do not leverage unlabeled data during training, which may help in low-labeled data scenarios.

Furthermore, despite the satisfactory results achieved by this new learning paradigm on the segmentation of natural images [1,15,19,24], its use in medical images remains scarce. Few-shot segmentation on medical images was first introduced by the work in [14]. Authors proposed to leverage adversarial learning to segment brain images based on 1 or 2 labeled brain images, inspired by the success of prior semi-supervised approaches [21]. Roy et al. [18] presented a different approach, where the architecture was composed of a conditioner and a segmenter arm. To strength the information exchange between both arms, they integrated *squeeze & excite* modules [11], which facilitated the gradient flow. More recently, one-shot medical image segmentation was addressed by synthesizing realistic training examples [29], a more elegant way of performing data augmentation. Nevertheless, these methods present some limitations. First, these approaches are based on the assumption that each shot is a whole 3D image, which contains many 2D slices. And second, they integrate large architectures (e.g., encoder, decoder and discriminator in [14] or conditioner and segmenter arms in [18]), which incur in complex and potentially unstable models.

Inspired by this, our work investigates the role of unsupervised data, via surrogate tasks, in the task of segmenting medical images in a few-shot learning scenario. Particularly, we take advantage of the success of few-shot segmentation works in natural images, which is based on the episodic training paradigm. To the best of our knowledge, this is the first attempt to tackle the few-shot segmentation task in medical imaging from an episodic perspective. To further improve the performance of our model, we leverage unlabeled data as a supervisory signal of surrogate tasks. Integration of unlabeled images into auxiliary tasks, has already shown to improve the generalization capabilities of deep models in sev-

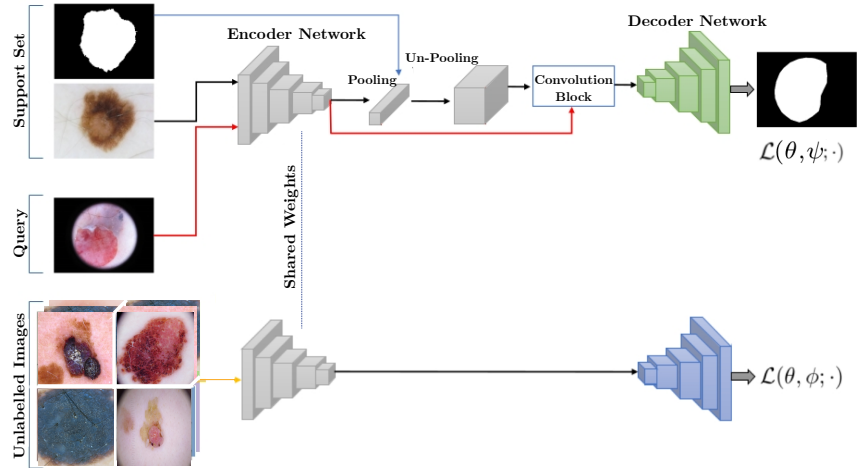


Fig. 1: Proposed few-shot semantic segmentation method for medical images using auxiliary tasks to leverage abundant unlabeled medical imaging data.

eral visual recognition tasks in other domains [6,9]. To evaluate the effectiveness of the proposed approach we resort to the task of skin cancer segmentation in two publicly available medical datasets. Results demonstrate that our approach brings an improvement of 6-7% over the baseline without incurring in extra-costs due to data annotation.

1 Methodology

1.1 Few shot semantic segmentation

In standard few-shot semantic segmentation, we typically have three datasets: a training set $D_{train} = \{(X_i^t, Y_i^t)\}_{i=1}^{N_{train}}$, a support set $D_{support} = \{(X_i^s, Y_i^s)\}_{i=1}^{N_{support}}$, and a test set $D_{test} = \{(X_i^q)\}_{i=1}^{N_{test}}$. In this setting, we denote an input image as $X_i \in \mathbb{R}^{H \times W \times Ch}$, where H , W and Ch represent the height, width and the number of channels, respectively, and $Y_i \in \{0, 1\}^{H \times W}$ its corresponding pixel-level mask. In total, each dataset contains N images, defined by N_{train} , $N_{support}$, and N_{test} , with C different classes. While classes are shared among support and test sets, they are disjoint with the training set, i.e., $\{C_{train}\} \cap \{C_{support}\} = \emptyset$.

Few-shot learning aims at training a neural network $f_{(\theta, \psi)}(\cdot)$ on the training set to have the ability to segment a novel class $c \notin C_{train}$ on the test set based on k references from $D_{support}$. Note that we employ θ and ψ to refer to the learnable parameters of the encoder and decoder, respectively. To reproduce this procedure, training on the base dataset D_{train} follows the episodic learning paradigm proposed in [23], where each episode instantiates a c -way k -shot learning task. Specifically, each episode is generated by sampling two elements. First, we build a support training set for each class c , denoted as $D_{train}^S = \{(X_s^t, Y_s^t(c))\}_{s=1}^k \subset$

D_{train} , where $Y_s^t(c)$ is the binary mask for the class c corresponding to the support image X_s^t . And second, a query set $D_{train}^Q = \{X_q^t, Y_q^t(c)\} \subset D_{train}$, where X_q^t is the query image and $Y_q^t(c)$ its corresponding binary mask for the class c . To estimate the segmentation mask of a given class c in the query image the model gets the support training set and the query image, which can be expressed as $\hat{Y}_q^t(c) = f_{(\theta, \psi)}(D_{train}^S, X_q^t)$.

Concretely, we first employ a CNN to encode the support and query images into the feature space, resulting in $f_s \in \mathbb{R}^{W' \times H' \times M}$ and $f_q \in \mathbb{R}^{W' \times H' \times M}$, respectively. The variables W' , H' and M denote the width, height and feature dimensionality on the feature space, respectively. To obtain the class prototypes, we apply mask average pooling on the feature representation over the known foreground regions of the support mask $Y_s^t(c)$. This can be formulated as:

$$p_s = \frac{1}{|\tilde{Y}_s^t(c)|} \sum_{i=1}^{W' \times H'} f_s \tilde{Y}_s^t(c) \quad (1)$$

where $\tilde{Y}_s^t(c) \in \{0, 1\}^{H' \times W'}$ denotes the down-sampled version of the support mask $Y_s^t(c)$ and $|\tilde{Y}_s^t(c)| = \sum_i \tilde{Y}_{s,i}^t(c)$ is the number of foreground locations in $\tilde{Y}_s^t(c)$. Note that p_s is a unidimensional vector with M elements, $p_s \in \mathbb{R}^{1 \times 1 \times M}$. Then, we unpool each prototype to the same spatial resolution as the query features f_q and convolve the upsampled prototypes with f_q (See Fig. 1 for the whole pipeline). The model parameters (θ, ψ) are then optimized by employing an objective function between $Y_q^t(c)$ and $\hat{Y}_q^t(c)$. In the few-shot segmentation literature, this function is typically the standard cross-entropy. Nevertheless, any other loss function could be used. Last, during testing, the model $f_{(\theta, \psi)}(\cdot)$ is evaluated on the test set D_{test} given k images from the support set $D_{support}$.

1.2 Surrogate task

A major challenge that we can encounter in few-shot learning is how to force the feature extractor to learn image features that can be readily employed on novel classes with just a handful of labeled samples. To address this issue, we propose to integrate an auxiliary loss during training, which compensates for the lack of annotated data on the novel categories. We can define formally this loss as $\mathcal{L}(\theta, \phi; \cdot)$, where ϕ denotes the parameters of the network related to the surrogate task.

Integrating unlabeled data. Several techniques have been proposed for self-supervised training, including predicting rotation [10], solving jigsaw puzzles [16] or filling removed parts of an image [27]. Among these strategies, we consider the task of image denoising, since some other techniques are not suited to our application. For example, if an image containing a car is rotated 180 degrees, it will be easy to predict its rotation, as the wheels would typically not be above the car roof. In contrast, an image of a skin tumor can be rotated in many directions, resulting in a large range of feasible rotations. To achieve this, we employ an

additional dataset, $D_{train}^{\mathcal{U}} = \{(X_i^{\mathcal{U}})\}_{i=1}^{N_{\mathcal{U}}}$, which contains $N_{\mathcal{U}}$ unlabeled images. We add random noise to these images, resulting in a larger dataset, where each image $X_i^{\mathcal{U}}$ generates multiple corrupted images $Y_i^{\mathcal{U}}$. The goal is that the encoder-decoder architecture learns a mapping between noised images and their original counterparts. To this end, we use the cross-entropy loss between the original and denoised images, which can be defined as:

$$\mathcal{L}_{sur}(\theta, \phi; Y^{\mathcal{U}}) = -\frac{1}{N_{\mathcal{U}}} \sum_i^{N_{\mathcal{U}}} \sum_j^{H \times W} X_{i,j}^{\mathcal{U}} \log \hat{Y}_{i,j}^{\mathcal{U}} \quad (2)$$

where $\hat{Y}_i^{\mathcal{U}} = f_{(\theta, \psi)}(Y_i^{\mathcal{U}})$ is the denoised image. We employed the CE instead of L_2 distance as the reconstruction loss as we empirically observed that it provided better results. This is in line with recent literature in Variational Auto-Encoders [2], which also employ CE as the reconstruction loss.

1.3 Joint objective

The final objective optimized during training is thus composed by the two terms included in the previous sections. The first term $\mathcal{L}_{few}(\theta, \psi; \cdot)$ is a function of the parameters θ and ψ of the encoder and decoder specialized on the few-shot segmentation task. The second term, $\mathcal{L}_{sur}(\theta, \phi; \cdot)$, depends on the encoder parameters θ and on the parameters ϕ of a network only dedicated to the surrogate task. Thus, the training process reduces to minimize the following function:

$$\min_{\theta, \psi, \phi} \mathcal{L}_{few}(\theta, \psi; D_{train}) + \lambda \mathcal{L}_{sur}(\theta, \phi; Y^{\mathcal{U}}) \quad (3)$$

where λ is employed to weight the importance of the surrogate task.

2 Experiments

We conduct a series of experiments to evaluate the proposed model for few-shot segmentation. We present below the datasets and experimental settings.

2.1 Dataset

We employ the FSS-1000 dataset as the base training set (D_{train}) in our experiments. To evaluate our method, we employ two publicly available medical datasets, i.e., ISIC and PH², which are described below. For the surrogate tasks, we only employ images from both FSS-1000 and ISIC datasets.

FSS-1000 Class Dataset: FSS-1000 class dataset [25] is a large-scale dataset specially designed for few-shot segmentation. It consists on 1000 classes, where each class contains 10 images with their corresponding pixel level ground truth annotations. The official training split (760 classes) is used as the base dataset for training, while the testing set (240 classes) is exploited on the surrogate tasks.

ISIC dataset: The ISIC 2018 dataset [5,22] is provided by the International 2018 Skin Imaging Collaboration Grand Challenge. The dataset contains 2594

dermoscopic images in RGB with their respective masks. An independent set of 1000 images is kept for evaluation purposes. The remaining images are employed for the auxiliary tasks, except the images used as k -shots in the support set. We resize all the images to 224×224 pixels to match the resolution of the images in the FSS-1000 dataset.

PH² dataset: The PH² dataset [13] contains a total of 200 RGB dermoscopic images of melanocytic lesions obtained at the Dermatology Service of Hospital, Pedro Hispano (Matosinhos, Portugal). All these images are used only during testing. Similarly to ISIC, we resized the images to 224×224 pixels.

2.2 Experimental Set-up

Network Details: Our encoder consists of four encoding blocks from VGG [20] pre-trained on ImageNet. We removed the max-pooling layers of the last two blocks and used atrous convolutions with dilation rate of 2 to enlarge the receptive field. Size of input images is equal to 224×224 and their encoded representations are down-scaled to 1/4 of the original input resolution. Our decoder network consists of upsampling, convolution, batch normalization and activation layers. To generate the final segmentation mask, we apply two decoding blocks.

Evaluation Protocol: In our work, we evaluate the performance of our model based on the DSC score, widely used in medical imaging to evaluate the segmentation performance. Given two segmentation masks A and B, the DSC can be defined as $DSC = \frac{2|A \cap B|}{|A| + |B|}$.

Implementation Details: The code is written in Keras with Tensorflow as backend. The tests are carried out in a server equipped with a Nvidia Titan X GPU. Both main and surrogate tasks are trained end-to-end using Adam optimization with learning rate 10^{-4} . In whole setting we trained the model for 30K iterations and evaluated on 500 episodes. The value of λ in eq. 3 is set 1.

2.3 Results

Quantitative results. Table 1 reports the segmentation results on ISIC and PH² datasets in the 1-shot scenario. First, to show the effectiveness of few-shot learning, we start by comparing this setting with standard batch-wise training, referred to as *Regular*. In this case, the model is trained on FSS-1000 dataset and directly tested on both ISIC and PH² sets. We can observe that across the two datasets, resorting to the episodic training paradigm results in an improvement of nearly 6%. Looking at the results obtained by our model, these indicate that leveraging unsupervised data via a surrogate task consistently improves the performance on both datasets. Particularly, the proposed approach outperforms the few-shot baseline by a margin between 6-7%. Compared to the upperbound, i.e., model trained and tested on the same dataset, the proposed models obtain promising results, with around 15% of difference on the PH² dataset, but only 1 target image segmented. Furthermore, we can also observe that increasing the number of additional samples on the episodes typically leads to better results.

This suggests that the model efficiently extracts semantic information from the unlabeled surrogate task to enhance the performance of the downstream task.

Model	Additional Samples	ISIC	PH ²
Regular (<i>Lower-bound</i>)	—	48.32	64.72
Few-shot	—	54.07	68.13
Few-shot + unlabeled	5	61.38	74.12
(<i>Denoising</i>) (Ours)	10	61.40	74.67
	20	60.79	74.77
Regular (<i>Upper-bound</i>)	—	86.65	89.94

Table 1: Quantitative results of the evaluated settings on ISIC and PH² datasets. Best results (different from the upperbound) highlighted in bold.

Results for the 5-shot scenario are reported in Table 2. In this setting, we only evaluate the model when 10 additional images are integrated in the episodes. Similarly to the 1-shot case, we observe an increase in performance with respect to the few-shot model trained without the surrogate task. Concretely, the proposed model obtains an improvement of 3-4%, whereas the gain was nearly to 6-7% in the case of 1-shot learning.

Model	ISIC	PH ²
Few-shot	59.63	71.15
Few-shot + Unlabeled (<i>Denoising</i>) (Ours)	62.40	75.54

Table 2: Quantitative results of 5-shot settings on ISIC and PH² datasets. Best results highlighted in bold.

These results suggest that adding a self-supervised surrogate task, i.e., denoising, improves the few-shot segmentation performance. Results also indicate that the performance improvement is more significant in the context of very few labeled samples, such as 1-shot *vs.* 5-shot.

Qualitative results. Visual segmentations on both datasets are depicted in Fig. 2. We can first observe that, by training the network in a batch-wise manner (*right column*) the segmentation results are not satisfactory, largely over-segmenting the target. If episodic training is used instead, the network shows a stronger capability to learn more general features. This is reflected in the better segmentation results provided by the few-shot baseline (*third column*). Last, leveraging unlabeled data through a surrogate task further improves the representation power of the model, resulting in richer and more generic features

(*fourth column*). Learning better features ultimately results in more reliable segmentations, as demonstrated in these visual examples. Specifically, the proposed model improves the segmentation by increasing the number of true positives (e.g., whole in the target of last row), while reducing the amount of false positives (e.g., isolated pixels on last row, and over-segmentations of other examples).

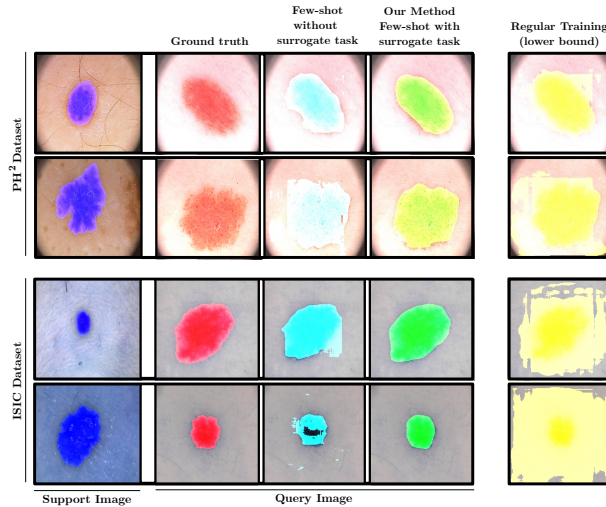


Fig. 2: Visual segmentations of the analyzed methods in the scenario of one-shot segmentation, on PH^2 and $ISIC$ datasets.

3 Conclusion

Inspired by the success of few-shot learning on visual recognition tasks we have proposed the first attempt to integrate episodic training in few-shot semantic segmentation on medical images. Furthermore, motivated by the close connection between few-shot and self-supervised learning, we have investigated the use of surrogate tasks to further improve current few-shot segmentation approaches. By solving a non-trivial proxy task that can be supervised trivially, such as denoising images with noise, the encoder network is encouraged to learn rich and generic image features which can be transferable to other ensuing tasks such as image segmentation. These enriched features lead to more powerful representations, which ultimately results in better generability to unseen tasks. Our experiments on two public skin cancer segmentation datasets revealed that exploiting unlabeled images through self-supervised tasks results in significant improvements on the few-shot segmentation performance.

References

1. Azad, R., Fayjie, A.R., Kauffman, C., Ben Ayed, I., Pedersoli, M., Dolz, J.: On the texture bias for few-shot CNN segmentation. arXiv preprint arXiv:2003.04052 (2020)
2. Basu, S., Wagstyl, K., Zandifar, A., Collins, L., Romero, A., Precup, D.: Early prediction of alzheimers disease progression using variational autoencoders. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 205–213. Springer (2019)
3. Bernard, O., Lalande, A., Zotti, C., Cervenansky, F., Yang, X., Heng, P.A., Cetin, I., Lekadir, K., Camara, O., Ballester, M.A.G., et al.: Deep learning techniques for automatic MRI cardiac multi-structures segmentation and diagnosis: is the problem solved? IEEE transactions on medical imaging **37**(11), 2514–2525 (2018)
4. Chen, H., Dou, Q., Yu, L., Qin, J., Heng, P.A.: VoxResNet: Deep voxelwise residual networks for brain segmentation from 3D MR images. NeuroImage **170**, 446–455 (2018)
5. Codella, N., Rotemberg, V., Tschandl, P., Celebi, M.E., Dusza, S., Gutman, D., Helba, B., Kalloo, A., Liopyris, K., Marchetti, M., et al.: Skin lesion analysis toward melanoma detection 2018: A challenge hosted by the international skin imaging collaboration (isic). arXiv preprint arXiv:1902.03368 (2019)
6. Doersch, C., Gupta, A., Efros, A.A.: Unsupervised visual representation learning by context prediction. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 1422–1430 (2015)
7. Dolz, J., Desrosiers, C., Ayed, I.B.: 3D fully convolutional networks for subcortical segmentation in MRI: A large-scale study. NeuroImage **170**, 456–470 (2018)
8. Fechter, T., Adebahr, S., Baltas, D., Ayed, I.B., Desrosiers, C., Dolz, J.: Esophagus segmentation in CT via 3D fully convolutional neural network and random walk. Medical physics **44**(12), 6341–6352 (2017)
9. Gidaris, S., Bursuc, A., Komodakis, N., Pérez, P., Cord, M.: Boosting few-shot visual learning with self-supervision. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 8059–8068 (2019)
10. Gidaris, S., Singh, P., Komodakis, N.: Unsupervised representation learning by predicting image rotations. In: ICLR (2018)
11. Hu, J., Shen, L., Sun, G.: Squeeze-and-excitation networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 7132–7141 (2018)
12. Hu, T., Yang, P., Zhang, C., Yu, G., Mu, Y., Snoek, C.G.: Attention-based multi-context guiding for few-shot semantic segmentation. In: AAAI. vol. 33, pp. pp. 8441–8448 (2019)
13. Mendonça, T., Ferreira, P.M., Marques, J.S., Marcal, A.R., Rozeira, J.: Ph 2-a dermoscopic image database for research and benchmarking. In: 2013 35th annual international conference of the IEEE engineering in medicine and biology society (EMBC). pp. 5437–5440. IEEE (2013)
14. Mondal, A.K., Dolz, J., Desrosiers, C.: Few-shot 3D multi-modal medical image segmentation using generative adversarial learning. arXiv preprint arXiv:1810.12241 (2018)
15. Nguyen, K., Todorovic, S.: Feature weighting and boosting for few-shot segmentation. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 622–631 (2019)
16. Noroozi, M., Favaro, P.: Unsupervised learning of visual representations by solving jigsaw puzzles. In: European Conference on Computer Vision. pp. 69–84. Springer (2016)

17. Ravi, S., Larochelle, H.: Optimization as a model for few-shot learning. In: ICLR (2017)
18. Roy, A.G., Siddiqui, S., Pölsterl, S., Navab, N., Wachinger, C.: squeeze & excite-guided few-shot segmentation of volumetric images. *Medical image analysis* **59**, 101587 (2020)
19. Siam, M., Oreshkin, B.N., Jagersand, M.: AMP: Adaptive masked proxies for few-shot segmentation. In: ICCV. pp. pp. 5249–5258 (2019)
20. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014)
21. Souly, N., Spampinato, C., Shah, M.: Semi supervised semantic segmentation using generative adversarial network. In: *Proceedings of the IEEE International Conference on Computer Vision*. pp. 5688–5696 (2017)
22. Tschandl, P., Rosendahl, C., Kittler, H.: The ham10000 dataset: A large collection of multi-source dermatoscopic images of common pigmented skin lesions; 2018. Preprint. Available from: <https://arxiv.org/ftp/arxiv/papers/1803/1803.10417.pdf>. Cited **4** (2019)
23. Vinyals, O., Blundell, C., Lillicrap, T., Wierstra, D., et al.: Matching networks for one shot learning. In: *Advances in neural information processing systems*. pp. 3630–3638 (2016)
24. Wang, K., Liew, J.H., Zou, Y., Zhou, D., Feng, J.: Panet: Few-shot image semantic segmentation with prototype alignment. In: *Proceedings of the IEEE International Conference on Computer Vision*. pp. 9197–9206 (2019)
25. Wei, T., Li, X., Chen, Y.P., Tai, Y.W., Tang, C.K.: Fss-1000: A 1000-class dataset for few-shot segmentation. *arXiv preprint arXiv:1907.12347* (2019)
26. Zhang, C., Lin, G., Liu, F., Yao, R., Shen, C.: CANet: Class-agnostic segmentation networks with iterative refinement and attentive few-shot learning. In: *CVPR*. pp. 5217–5226 (2019)
27. Zhang, R., Isola, P., Efros, A.A.: Split-brain autoencoders: Unsupervised learning by cross-channel prediction. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 1058–1067 (2017)
28. Zhang, X., Wei, Y., Yang, Y., Huang, T.: SG-one: Similarity guidance network for one-shot semantic segmentation. *arXiv preprint arXiv:1810.09091* (2018)
29. Zhao, A., Balakrishnan, G., Durand, F., Guttag, J.V., Dalca, A.V.: Data augmentation using learned transformations for one-shot medical image segmentation. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 8543–8553 (2019)